

Abstract Interpretation Over Bitvectors

by

Tushar Sharma

A dissertation submitted in partial fulfillment of
the requirements for the degree of

Doctor of Philosophy

(Computer Sciences)

at the

UNIVERSITY OF WISCONSIN–MADISON

2017

Date of final oral examination: 08/11/2017

The dissertation is approved by the following members of the Final Oral Committee:

Thomas Reps, Professor, Computer Sciences

Somesh Jha, Professor, Computer Sciences

Benjamin Liblit, Associate Professor, Computer Sciences

Aws Albarghouthi, Assistant Professor, Computer Sciences

Vadim Shapiro, Professor, Mechanical Engineering and Computer Sciences

© Copyright by Tushar Sharma 2017
All Rights Reserved

Dedicated to mummy, papa, didi, and bhaiya.

ACKNOWLEDGMENTS

My graduate student journey was a memorable one, and several people made it a meaningful experience. Without the help and support of mentors, co-workers, friends, and family, this dissertation would not have been possible.

First and foremost, I am indebted to my advisor, Prof. Thomas Reps, who has motivated and advised me throughout my journey in the doctoral program. Without his immense knowledge and expertise in program analysis and abstract interpretation, none of my research projects would have been successful. His passion and enthusiasm for research is unparalleled. Tom has given me so much freedom to pursue my research interests, and helped me in converting my half-baked ideas into useful publications. Over the years, he cultivated me into an independent researcher by teaching me the fine art of asking interesting research questions, reading interesting related work, and building the research ideas into concrete projects worthy of publication. Most importantly, he has treated me more like a friend, than as a subordinate, and for that, I will be forever indebted. I hope you enjoyed working with me, as much I did working with you.

I thank my dissertation committee for taking the time to read my dissertation. I thank Prof. Somesh Jha, Prof. Ben Liblit, and Prof. Shan Lu for reading my prelim document, attending my prelim talk, and providing useful and constructive suggestions that guided me through my years as a dissertator. I thank my undergrad mentor, Prof. Sundar Balasubramaniam, who got me excited about static analysis and type systems, and motivated me to pursue graduate school.

Over the last few years, I had the fortune of working and collaborating with the best colleagues. I want to thank Junghee Lim, Matt Elder, Evan Driscoll, and Aditya Thakur, for mentoring me at the earlier years of my graduate school, and answering a whole variety of questions from the

complex subtleties of static analysis to basic coding questions. I want to thank Vijay Chidambaram, Venkatesh Srinivasan, and Tycho Andersen for collaborating with me on several projects and making my graduate school experience much more fun. I have had so many wonderful conversations with Venkatesh Srinivasan, Peter Ohmann, and Jason Breck. It has been a pleasure to have you guys around, on whom I can vent my pessimism and nihilism. I am appreciative for David Bingham, Stephen Lee, Venkatesh Srinivasan, Samuel Drews, Aditya Thakur, Evan Driscoll, Drew Davidson and Jason Breck, who patiently listened to my boring presentations and provided me constructive feedback.

The employees at GrammaTech have helped me tremendously over the years to answer my queries, fix bugs, provide useful technical support in the tools and infrastructure support for machine-code analysis. In particular, Junghee Lim, Evan Driscoll, Suan Yong, Brian Alliet, Tom Johnson, Alexey Loginov and David Melski have promptly answered my questions and provided helpful comments to improve my tools.

I thank my friends for making my life in Madison an exciting and happy one. Suyash Singh was a friend before we started the graduate journey and over our common graduate years at Madison, our friendship has grown stronger (and quirky). Thank you for the laughs and the support! Hemanth Pikkili, Gurdaman Khaira, Kushal Sinha, and Suyash Singh have been amazing roommates over the years, and I have never felt bored in Madison because of you guys. Thank you Suyash, Rishabh, Fulya, and Kushal for being the 'chemical brothers/sisters'; Raja, Rohit, and Sankar for being the 'doods'; Saily, Etienne, Cherry, Bang, Dipto, Saurabh, Keive, Anand, and Amy for the 'spicy balls'; Selah, Obi, Manali, Jason, and Janina for being the voice of reason; and Thanu, Venkatesh, Goldy and Dipto for listening to my graduate school rants.

My mom and dad have unconditionally loved and supported me throughout my life. They have always supported me in my graduate

school journey, even though we have been thousands of miles away. My sister, Sugandha, always provides me with moral support. Her innocence, charisma, and kindness is unparalleled. My brother has been my role model, my mentor, my ally, and my competition. He has always believed in me, and challenged me to dream bigger, higher, and better. Words would never be enough to describe the gratitude I have for my family. I want to thank Mrs. and Mr. Kulkarni and Shefali for being my family in US. Thanks to them, thanksgiving and christmas holidays have been full of life.

Finally, I want to thank the love of my life, Sali Kulkarni, for being there for me, for understanding me, and for believing in me, even when I didn't believe in myself.

This dissertation is supported, in part, by a gift from Rajiv and Ritu Batra; by DARPA under cooperative agreement HR0011-12-2-0012; by NSF under grant CCF-{0810053, 0904371}; by ONR under grants N00014-09-1-0510, 11-C-0447; by AFRL under contract FA9550-09-1-0279 and FA8650-10-C-7088; by ARL under grant W911NF-09-1-0413, by AFRL under DARPA MUSE award FA8750-14-2- 0270, DARPA STAC award FA8750-15-C-0082; and by the UW-Madison Office of the Vice Chancellor for Research and Graduate Education with funding from the Wisconsin Alumni Research Foundation. Any opinions, findings, and conclusions or recommendations expressed in this dissertation are those of the author, and do not necessarily reflect the views of the sponsoring agencies. The author's advisor Thomas Reps has an ownership interest in GrammaTech, Inc., which has licensed elements of the technology reported in this dissertation.

CONTENTS

Contents	v
List of Tables	vi
List of Figures	vii
Abstract	viii
1 Introduction	1
1.1 <i>Motivation.</i>	4
1.2 <i>Thesis Contributions</i>	6
1.3 <i>Thesis Organization</i>	14
2 Background	15
2.1 <i>Concrete Semantics</i>	15
2.2 <i>Abstract Interpretation</i>	17
2.3 <i>Creation of Abstract Transformers.</i>	21
2.4 <i>Fixed-point computation</i>	31
2.5 <i>Weighted Pushdown Systems</i>	32
3 Abstract Domains of Affine Relations	35
3.1 <i>Abstract Domains for affine-relation analysis</i>	37
3.2 <i>Relating AG and KS Elements</i>	46
3.3 <i>Relating KS and MOS</i>	48
3.4 <i>Using KS for Interprocedural Analysis</i>	55
3.5 <i>Experiments</i>	71
3.6 <i>Related Work</i>	89
3.7 <i>Chapter Notes</i>	94
4 A New Abstraction Framework for Affine Transformers	96

4.1	<i>Preliminaries</i>	99
4.2	<i>Overview</i>	101
4.3	<i>Affine-Transformer-Abstraction Framework</i>	106
4.4	<i>Discussion and Related Work</i>	118
4.5	<i>Chapter Notes</i>	121
5	An Abstract Domain for Bit-vector Inequalities	122
5.1	<i>Overview</i>	124
5.2	<i>Terminology</i>	126
5.3	<i>Base Domains</i>	127
5.4	<i>The View-Product Combinator</i>	128
5.5	<i>Synthesizing Abstract Operations for Reduced-Product Domains</i>	132
5.6	<i>Experimental Evaluation</i>	136
5.7	<i>Related Work</i>	137
5.8	<i>Chapter Notes</i>	139
6	Sound Bit-Precise Numerical Domains Framework for Inequalities	140
6.1	<i>Terminology</i>	142
6.2	<i>Overview</i>	145
6.3	<i>The BVsFD Abstract-Domain Framework</i>	151
6.4	<i>Experimental Evaluation</i>	157
6.5	<i>Chapter Notes</i>	162
7	Conclusion and Future Work	163
7.1	<i>Bit-vector-precise Equality Domains</i>	165
7.2	<i>Bit-vector-precise Inequality Domains</i>	166
A	Domain Conversions	168
A.1	<i>Soundness of MOS to AG transformation</i>	168

<i>A.2 Soundness of KS Without Pre-State Guards to MOS transformation</i>	169
<i>A.3 Soundness of KS Without Pre-State Guards to MOS transformation</i>	169
B Howell Properties	172
C Correctness of KS Join	175
D Soundness of the Abstract-Domain Operations for Affine-Transformers-Abstraction Framework	178
E Soundness of Abstract Composition for Affine-Transformers-Abstraction Framework	180
<i>E.1 Non-Relational Base Domain</i>	180
<i>E.2 Weakly-Convex Base Domain</i>	182
F Soundness of The Merge Operation	188
References	191

LIST OF TABLES

2.1	Abstract-domain operations.	20
2.2	Truth table for multiplication and addition of two parity values.	25
2.3	Snapshots in the fixed-point analysis for Ex. 2.3.	32
2.4	Semiring operators in terms of abstract-domain operations.	34
4.1	Abstract-domain operations.	100
4.2	Example demonstrating two ways of relating MOS and AG.	101
4.3	Base abstract-domain operations.	107
4.4	Abstract-domain operations for the $\text{ATA}[\mathcal{B}]$ -domain.	108
4.5	Foundation-domain operations.	113
5.1	Machine-code analysis using \mathcal{BVJ} . Columns 6–9 show the times (in seconds) for the $\mathcal{E}_{\mathbb{Z}_{2^w}}$ -based analysis, and for the \mathcal{BVJ} -based analysis; and the degree of improvement in precision measured as the number of control points at which \mathcal{BVJ} -based analysis gave more precise invariants compared to $\mathcal{E}_{\mathbb{Z}_{2^w}}$ -based analysis, and the number of procedures for which \mathcal{BVJ} -based analysis gave more precise summaries compared to $\mathcal{E}_{\mathbb{Z}_{2^w}}$ -based analysis.	136
6.1	Snapshots in the fixed-point analysis for Ex. 6.1 using the $\text{BV}SFD_2(\mathcal{OCT})$ domain. B_{v_1, v_2, \dots, v_k} are the bounding constraints for the variables v_1, v_2, \dots, v_k	150
6.2	Information about the loop benchmarks containing true assertions, a subset of the SVCOMP benchmarks.	159

LIST OF FIGURES

1.1	Each + represents a solution of the indicated inequality in 4-bit unsigned bit-vector arithmetic.	5
2.1	Concrete semantics for $\mathcal{L}(\text{SLANG})$	18
2.2	Reinterpretation semantics for $\mathcal{L}(\text{SLANG})$ for domain \mathcal{A}_{PAR}	24
3.1	(a) The King-Søndergaard algorithm for symbolic abstraction ($\hat{\alpha}_{\text{KS}}^{\uparrow}(\varphi)$). (b) The Thakur-Elder-Reps bilateral algorithm for symbolic abstraction, instantiated for the KS domain: $\hat{\alpha}_{\text{TER}[\text{KS}]}^{\uparrow}(\varphi)$. In both algorithms, <i>lower</i> is maintained in Howell form throughout.	66
3.2	Some of the characteristics of the corpus of 19,066 (non-privileged, non-floating point, non-mmx) instructions.	71
3.3	Program information. All nine utilities are from Microsoft Windows version 5.1.2600.0, except setup, which is from version 5.1.2600.5512. The columns show the number of instructions (Instrs); the number of procedures (Procs); the number of basic blocks (BBs); the number of branch instructions (Branches); and the number of Δ_0 , Δ_1 , and Δ_2 rules in the WPDS encoding (WPDS Rules).	73
3.4	A fragment of the TSL specification of the concrete semantics of the Intel IA32 instruction set.	76
3.5	Comparison of the performance of MOS-reinterpretation and KS-reinterpretation for x86 instructions.	79
3.6	Comparison of the precision of MOS-reinterpretation and KS-reinterpretation for x86 instructions.	79

3.7	Performance of WPDS-based interprocedural analysis. The times, in seconds, for WPDS construction, performing interprocedural dataflow analysis (i.e., running <code>post*</code> and performing path-summary) and finding one-vocabulary affine relations at branch instructions, using MOS-reinterpretation, KS-reinterpretation, $\hat{\alpha}_{KS}$, and $\hat{\alpha}_{KS}^+$ to generate weights. The columns labeled “t/o” report the number of WPDS rules for which weight generation timed out during symbolic abstraction.	83
3.8	Comparison of the precision of the WPDS weights computed using MOS-reinterpretation, KS-reinterpretation, and $\hat{\alpha}_{KS}$. (E.g., <code>KS-reinterp < MOS-reinterp</code> reports the number of rules for which the KS-reinterp weight was more precise than the MOS-reinterp weight.)	85
3.9	Comparison of the precision of the one-vocabulary affine relations identified to hold at branch points via interprocedural analysis, using weights created using MOS-reinterpretation, KS-reinterpretation, and $\hat{\alpha}_{KS}$. (E.g., <code>KS-reinterp < MOS-reinterp</code> reports the number of branch points at which the KS-reinterp results were more precise than the MOS-reinterp results.)	86
3.10	Comparison of the precision of the two-vocabulary affine relations identified to hold at procedure-exit points via interprocedural analysis, using weights created using MOS-reinterpretation, KS-reinterpretation, and $\hat{\alpha}_{KS}$. (E.g., <code>KS-reinterp < MOS-reinterp</code> reports the number of procedure-exit points at which the KS-reinterp results were more precise than the MOS-reinterp results.)	87
3.11	Simplified version of an example that caused KS results to be less precise than MOS results, due to compose not distributing over join in the KS domain.	88

4.1	Abstract transformers and snapshots in the fixpoint analysis with the MOS domain for Ex. 4.4.	103
4.2	Abstract transformers and fixpoint analysis with the $\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]$ domain for Ex. 4.4.	105
5.1	Example snippet of Intel x86 machine code.	125
6.1	Wrap-around on variable x , treated as an unsigned char. . .	143
6.2	Lazy abstract transformers with the $BVSFD_2(\mathcal{POLY})$ domain for Ex. 6.1. ID refers to the identity transformation.	149
6.3	Abstract-domain interface for \mathcal{A}	152
6.4	Reinterpretation semantics for $\mathcal{L}(\text{ELANG})$	155
6.5	Precision and performance numbers for SV-COMP loop benchmarks.	159
6.6	Precision and performance numbers for SV-COMP array benchmarks with \mathcal{POLY} as the base domain.	161

ABSTRACT

Most critical applications, such as airplane and rocket controllers, need correctness guarantees. Usually these correctness guarantees can be described as safety properties in the form of assertions. Verifying an assertion amounts to showing that the assertion holds true for all possible runs of an application. Abstract interpretation is a method to automatically verify a program by soundly abstracting the concrete executions of the program to elements in an abstract domain, and checking the correctness guarantees using the abstraction. However, traditional abstract domains treat the machine integers as mathematical integers. As a result, the conclusions drawn from such an abstract interpretation are, in general, unsound. In other words, the assertions shown to be true by traditional abstract interpretation approaches might actually be false because the underlying point space does not faithfully model bit-vector arithmetic.

This dissertation advances the field of abstract interpretation by providing sound abstraction techniques and abstract-domain frameworks that faithfully model bit-vector semantics. We focus on numerical abstract domains for bitvectors, which can provide equality and inequality invariants.

The first part of the dissertation focuses on abstract domains capable of expressing bit-vector-sound *equality* invariants. The performance and precision of two equality domains is compared, and sound inter-conversion methods are provided. Furthermore, we generalize one of the equality domains to develop a new abstract-domain framework that is capable of expressing a certain class of disjunctions over bit-vector-sound equality constraints. This framework can be instantiated with any relational or non-relational base abstract domain over bitvectors.

The second part of the dissertation focuses on abstract domains capable of expressing bit-vector-sound *inequality* invariants. We develop an

abstract domain that is capable of expressing a certain class of bit-vector-sound inequalities over memory variables and registers. Furthermore, we develop an abstract-domain framework that takes an abstract domain that is sound with respect to mathematical integers, and creates an abstract domain whose operations and abstract transformers are sound with respect to machine integers.

1 INTRODUCTION

Nowadays, humans rely on computer systems in almost every aspect of life. Reliance on critical systems such as automated banking systems, flight-control systems, automated cruise control and braking systems, and traffic-light control systems has increased as a result. The presence of bugs in these critical systems have lead to embarrassing issues like the Excel 2007 bug [98] and overflow in youtube view limit [6], unnecessary human intervention [31], losses in millions of dollars [38, 36] or worse, loss of human life [28, 32].

Most critical applications, such as airplane and rocket controllers, need correctness guarantees. Usually these correctness guarantees can be described as safety properties in the form of assertions. Verifying an assertion amounts to showing that the assertion's condition holds true for all possible runs of an application. Proving an assertion is, in general, an undecidable problem. Nevertheless, there exist static-analysis techniques that are able to verify automatically some kinds of program assertions. One such technique is abstract interpretation [20], which soundly abstracts the concrete executions of the program to elements in an abstract domain, and checks the correctness guarantees using the abstraction.

Abstract interpretation has seen a lot of progress in last 40 years since it was introduced by Cousot and Cousot in 1977. Many numeric abstract domains have been introduced to verify the correctness of assertions by the programmer, as well as divide-by-zero, array-bounds checks, buffer-overflow behavior, etc. Arguably, abstract interpretation has made the most progress in numeric program analysis, that is, discovering numerical properties about the program. The material in this dissertation is centered on discovering numerical properties on the machine integers through abstract interpretation.

The building block for the abstraction of the program is an *abstract*

domain. An abstract domain defines the level and detail of the abstraction. An abstract domain can be seen as a logic fragment [86] that is capable of expressing only a certain class of constraints. For instance, an interval value $v \in [0, 5]$, where v is a variable, expresses the constraint: $0 \leq v \leq 5$. The richer the abstract domain, the more intricate the constraints it can express. Richer abstract domains can prove more assertions possibly at the cost of worse performance.

Mathematically, abstract domains are (usually) partial orders, where the order is given by set containment: the smaller the set denoted by an abstract-domain element, the more precise it is. For instance, the interval $i_1 : v \in [0, 5]$ is more precise than the interval $i_2 : v \in [-10, 10]$.

Often \mathcal{A} is really a family of abstract domains, in which case $\mathcal{A}[\mathcal{U}]$ denotes the specific instance of \mathcal{A} that is defined over vocabulary \mathcal{U} . The two important steps in abstract interpretation (AI) are:

1. Abstraction: The abstraction of the program is constructed using the abstract domain and abstract semantics.
2. Fixpoint analysis: Iteration until a fixpoint is reached is performed on the abstraction of the program to identify invariants of the program.

In the typical setup of AI, the set of states that can arise at each point in the program is safely represented by the abstract-domain element found for that program point in the fixed point. This setup can be used to prove assertions. Let us denote the set of variables in a program by V . In this case, the abstract-domain elements abstract the set of values at a program point, thus they are abstracting the values for V , and therefore $\mathcal{U} = V$. For the purpose of this thesis, abstract-domain elements are abstract transformers, that is, they are abstractions of concrete transformers describing the control-flow graph edges in the program. When abstract-domain elements are abstract transformers, the results provide function summaries or loop summaries [19, 92]. In principle, summaries can be computed offline for large libraries of code so that client static analyses can use them

to obtain verification results more efficiently. Abstract transformers are abstracting the set of values at two program points: one is represented by the pre-transformation variables V , and the other is represented by the post-transformation variables V' , and therefore $U = V \cup V'$. A static analyzer needs a way to construct abstract transformers for the concrete operations in the programs. Semantic reinterpretation (see §2.3.1) and symbolic abstraction (see §2.3.2) provide automatic ways to construct abstract transformers.

Past literature has introduced many non-relational and relational abstract domains. Examples of non-relational domains include constant propagation [51] ($v_i = \alpha_i$), signs [19] ($\pm v_i \leq 0$), intervals [19] ($\alpha_i \leq v_i \leq \beta_i$), congruences [39] ($v_i \equiv \alpha_i \pmod{\beta_i}$), and interval congruences [65] ($v_i \in [\alpha_i, \beta_i] \pmod{\gamma_i}$). Here, v_i refers to a variable in the program, and the symbols α_i , β_i and γ_i refer to constants.

The classical example of a relational domain is the domain of linear equalities introduced by Karr [49] ($\sum_i \alpha_{ij} v_i = \beta_j$). Granger [40] introduced the linear-congruence domain ($\sum_i \alpha_{ij} v_i = \beta_j \pmod{\gamma_j}$). One of the most widely used relational abstract domains is the polyhedral domain [23], which is capable of expressing relational affine inequalities ($\sum_i \alpha_{ij} v_i \leq \beta_j$). While the polyhedral domain is useful, it is also very slow, and not scalable to large systems. With that in mind, previous research [68, 69, 97, 58, 89, 73] has also provided weaker variants of the polyhedra domain that are capable of expressing some; but not all; affine inequalities. For instance, the octagon abstract domain [68] can express only relational inequalities involving at most two variables where the coefficients on the variables are only allowed to be plus or minus one ($\pm v_i \pm v_j \leq \alpha_{ij}$).

1.1 Motivation.

These abstract domains suffer from one huge limitation: they treat the program variables as mathematical integers or rationals. However, the native machine-integer data-types used in programs (e.g., `int`, `unsigned int`, `long`, etc.) perform bit-vector arithmetic, and arithmetic operations wrap around on overflow. Thus, the underlying point space used in the aforementioned abstract domains does not faithfully model bit-vector arithmetic, and consequently the conclusions drawn from an analysis based on these domains are, in general, unsound, unless special steps are taken [96][15].

Example 1.1. *The following C-program fragment incorrectly computes the average of two `int`-valued variables [9]:*

```
int low, high, mid;
assume(0 <= low <= high);
mid = (low + high)/2;
assert(0 <= low <= mid <= high);
```

A static analysis based on polyhedra or octagons would draw the wrong conclusion that the assertion always holds. In particular, assuming 32-bit `ints`, when the sum of `low` and `high` is greater than $2^{31}-1$, the sum overflows, and the resulting value of `mid` is smaller than `low`. Consequently, there exist runs in which the assertion fails. These runs are overlooked when the polyhedral domain is used for static analysis because the domain fails to take into account the bit-vector semantics of program variables. \square

The problem that we wish to solve is not one of merely *detecting* overflow—e.g., to restrain an analyzer from having to explore what happens after an overflow occurs. On the contrary, our goal is to be able to track soundly the effects of arithmetic operations, including wrap-around effects of operations that overflow. This ability is useful, for instance, when analyzing code generated by production code generators, such as dSPACE TargetLink [25], which use the “compute-through-overflow” tech-

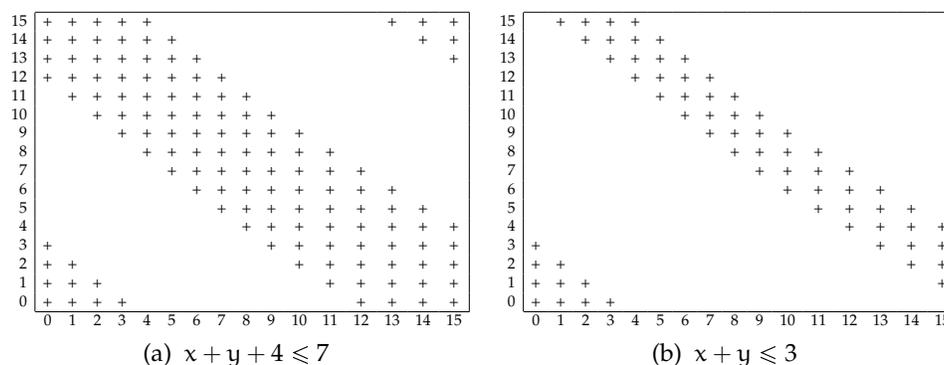


Figure 1.1: Each + represents a solution of the indicated inequality in 4-bit unsigned bit-vector arithmetic.

nique [35]. Furthermore, clever idioms for bit-twiddling operations, such as the ones explained in [102], sometimes rely on overflow [24].

1.1.1 Challenges in dealing with bit-vectors.

Some of the ideas used in designing an inequality domain for reals do not carry over to ones designed for bit-vectors. First, in bit-vector arithmetic, additive constants cannot be canceled on both sides of an inequality, as illustrated in the following example.

Example 1.2. Let x and y be 4-bit unsigned integers. Fig. 1.1 depict the solutions in bit-vector arithmetic of the inequalities $x + y + 4 \leq 7$ and $x + y \leq 3$. Although $x + y + 4 \leq 7$ and $x + y \leq 3$ are syntactically quite similar, their solution spaces are quite different. In particular, because of wrap-around of values computed on the left-hand sides using bit-vector arithmetic, one cannot just subtract 4 from both sides to convert the inequality $x + y + 4 \leq 7$ into $x + y \leq 3$. \square

Second, in bit-vector arithmetic, positive constant factors cannot be canceled on both sides of an inequality; for example, if x and y are 4-bit bit-vectors, then $(4, 4)$ is in the solution set of $2x + 2y \leq 4$, but not of $x + y \leq 2$.

These properties do not carry over because bitvectors are not a field, but rather a ring \mathbb{Z}_m of integers modulo $m = 2^w$ for $w > 1$. Moreover, the ring has zero divisors, and consequently even numbers in bit-vector arithmetic do not have a multiplicative inverse. Because of these properties, results from linear algebra over fields do not apply to the abstract domains over \mathbb{Z}_m .

While some simple domains do exist that are capable of representing certain kinds of inequalities over bit-vectors (e.g., intervals with a congruence constraint, sometimes called “strided-intervals” [83, 91, 5, 72]), such domains are non-relational; that is, they are not capable of expressing relations among several variables. On the other hand, there exist relational bit-precise domains such as the domains introduced by Müller-Olm/Seidl (denoted by MOS) [78] and King/Søndergaard (denoted by KS) [53], but they cannot express inequalities. Moreover, it is not clear how MOS and KS relate to each other.

1.2 Thesis Contributions

In this thesis, we provide sound abstract domains, and abstract-transformer construction techniques, that allow one to prove assertions and provide function summaries for programs with machine integers. The communication with existing analyzers is done strictly through a generic abstract-domain interface. Thus, our work should be compatible with most existing analysis engines. As a result, our abstraction techniques can be applied to source-code analysis as well as machine-code analysis.

1.2.1 Quick Summary

Broadly, our contributions to the numeric analysis of programs with machine integers fall into two categories: i) Bit-Vector Precise Equality Ab-

stract Domains, ii) Bit-Vector Precise Inequality Abstract Domains. The thesis makes the following contributions:

- **Bit-Vector Precise Equality Abstract Domains**
 - **Abstract Domains of Affine Relations [29, 30]:** An affine relation is a *linear-equality* constraint over a given set of variables that hold machine integers. In this work, we compare the MOS and KS abstract domains, along with several variants. These domains are capable of expressing affine relations over bitvectors. We show that MOS and KS are, in general, *incomparable* and give sound interconversion methods for KS and MOS. We introduce a third domain for representing affine relations, called AG, which stands for *affine generators*. Furthermore, we present an experimental study comparing the precision and performance of analyses with the KS and MOS domains. (See §1.2.2 for more details.)
 - **Abstraction Framework for Affine Transformers [93]:** In this work, we define the Affine-Transformers Abstraction Framework, which represents a new family of numerical abstract domains. This framework is parameterized by a base numerical abstract domain, and allows one to represent a set of affine transformers (or, alternatively, certain disjunctions of transition formulas). Specifically, this framework is a generalization of the MOS domain. The choice of the base abstract domain allows the client to have some control over the performance/precision trade-off. (See §1.2.3 for more details.)
- **Bit-Vector Precise Inequality Abstract Domains**
 - **An Abstract Domain for Bit-vector Inequalities [95]:** This work describes the design and implementation of a new abstract domain, called the *Bit-Vector Inequality* domain, which is capable of capturing certain inequalities over bit-vector-valued

variables (which represent a program’s registers and/or its memory variables). This domain tracks properties of the values of selected registers and portions of memory via *views*, and provides automatic heuristics to gather equality and inequality views from the program. Furthermore, experiments are provided to show the usefulness of the Bit-Vector Inequality domain. (See §1.2.4 for more details.)

- **Sound Bit-Precise Numerical Domains Framework for Inequalities [94]:** This work introduces a class of abstract domains, parameterized on a base domain, that is sound with respect to bitvectors whenever the base domain is sound with respect to mathematical integers. The base domain can be any numerical abstract domain. We also describe how to create abstract transformers for this framework that incorporate lazy wrap-around to achieve more precision, without sacrificing soundness with respect to bitvectors. We use a finite number of disjunctions of base-domain elements to help retain precision. Furthermore, we present experiments to empirically demonstrate the usefulness of the framework. (See §1.2.5 for more details.)

1.2.2 Abstract Domains of Affine Relations

As mentioned before, the abstract domains MOS and KS can express relational equalities over bitvectors. An element in the King/Søndergaard domain (KS) is an affine-closed set of linear-equality constraints over bitvectors ($\sum_i \alpha_{ij} v_i + \sum_i \alpha'_{ij} v'_i = \beta_j$, where $\alpha_{ij}, \alpha'_{ij}, \beta_j \in \mathbb{Z}_m$). An element in the Müller-Olm/Seidl domain (MOS) is an affine-closed set of affine transformers. An affine transformer is a relation on states, defined by $\vec{v}' = \vec{v} \cdot C + \vec{d}$, where \vec{v}' and \vec{v} are row vectors that represent the post-transformation state and the pre-transformation state, respectively.

C is the linear component of the transformation and \vec{d} is a constant vector. For example, $[x' \ y'] = [x \ y] \begin{bmatrix} 1 & 0 \\ 2 & 0 \end{bmatrix} + [10 \ 0]$ denotes the affine transformation $(x' = x + 2y + 10 \wedge y' = 0)$ over variables $\{x, y\}$. We denote such an affine transformation by $C : \vec{d}$. An element of these domains represents a set of points that satisfies an affine relation over variables that hold machine integers.

Chapter 3 considers MOS and KS abstract domains, along with several variants, and studies how they relate to each other. We show that MOS and KS are, in general, *incomparable*. In particular, we show that KS can express transformers with affine guards, which MOS cannot express. On the other hand, MOS can express some non-affine-closed sets, which are not expressible by KS. We give sound interconversion methods for KS and MOS. That is, we give an algorithm to convert a KS element to an over-approximating MOS element, as well as an algorithm to convert an MOS element to an over-approximating KS element.

We introduce a third domain for representing affine relations, called AG, which stands for *affine generators*. Whereas an element in the KS domain consists of a set of *constraints* on the values of variables, AG represents a collection of allowed values of variables via a set of *generators*. We show that AG is the generator counterpart of KS: a KS element can be converted to an AG element, and vice versa, with no loss of precision.

Furthermore, Chapter 3 presents an experimental study with the Intel IA32 (x86) instruction set in which the symbolic abstraction (§2.3.2) method and two reinterpretation methods (§2.3.1)—KS-reinterpretation and MOS-reinterpretation—are compared in terms of their performance and precision. The precision comparison is done by comparing the affine invariants obtained at branch points, as well as the affine procedure summaries obtained for procedures. For KS-reinterpretation and MOS-reinterpretation, we also compare the abstract transformers generated for individual x86 instructions.

1.2.3 A New Abstraction Framework for Affine Transformers

Abstractions of affine transformers can be used to obtain affine-relation invariants at each program point in the program [75]. An affine relation is a linear-equality constraint between numeric-valued variables of the form $\sum_{i=1}^n a_i v_i + b = 0$. For a given set of variables $\{v_i\}$, affine-relation analysis (ARA) identifies affine relations that are invariants of a program. The results of ARA can be used to determine a more precise abstract value for a variable via *semantic reduction* [21], or detect the relationship between program variables and loop-counter variables. While the MOS domain is useful for finding affine-relation invariants in a program, the join operation used at confluence points can lose precision in many cases, leading to imprecise function summaries. Furthermore, the analysis does not scale well as the number of variables in the vocabulary increases. In other words, it has one baked-in performance-versus-precision aspect.

Chapter 4 provide analysis techniques to abstract the behavior of the program as a set of affine transformations over bit-vectors. Remember that an affine transformation has the following form: $\vec{v}' = \vec{v} \cdot C + \vec{d}$). In this work, we generalize the ideas used in the MOS domain—in particular, to have an abstraction of *sets of affine transformers*—but to provide a way for a client of the abstract domain to have some control over the performance/precision trade-off. Toward this end, we define a new family of numerical abstract domains, denoted by $\text{ATA}[\mathcal{B}]$. (ATA stands for Affine-Transformers Abstraction.) Following our observation, $\text{ATA}[\mathcal{B}]$ is parameterized by a base numerical abstract domain \mathcal{B} , and allows one to represent a set of affine transformers (or, alternatively, certain disjunctions of transition formulas).

The overall contribution of our work is the framework $\text{ATA}[\mathcal{B}]$, for which we present

- methods to perform basic abstract-domain operations, such as equality and join.
- a method to perform abstract composition, which is needed to perform abstract interpretation.
- a faster method to perform abstract composition when the base domain is non-relational.

1.2.4 An Abstract Domain for Bit-vector Inequalities

Two of the biggest challenges in machine-code analysis are: (1) identifying inequality invariants while handling overflow in arithmetic operations over bit-vector data-types, and (2) identifying invariants that capture properties of values in memory.

When analyzing machine code, memory is usually modeled as a flat array. When analyzing Intel x86 machine code, for instance, memory is modeled as a map from 32-bit bit-vectors to 8-bit bit-vectors. Consequently, an analysis has to deal with complications arising from the little-endian addressing mode and aliasing, as we now illustrate.

Example 1.3. Consider the following machine-code snippet:

```
mov  eax, [ebp]
mov  [ebp+2], ebx
```

The first instruction loads the four bytes pointed to by register `ebp` into the 32-bit register `eax`. Suppose that the value in register `ebp` is A . After the first instruction, the bytes of `eax` contain, in least-significant to most-significant order, the value at memory location A , the value at location $A + 1$, the value at location $A + 2$, and the value at location $A + 3$. The second instruction stores the value in register `ebx` into the memory pointed to by `ebp+2`. Due to this instruction, the values at memory locations $A + 2$ through $A + 5$ are overwritten, after which the value in register `eax` no longer equals (the little-endian interpretation of) the bytes in memory pointed to by `ebp`. □

These challenges lead to the following problem statement:

Design an abstract domain of relational bit-vector affine-inequalities over memory-values and registers.

Chapter 5 expands the set of techniques available for abstract interpretation and model checking of machine code. We describe the design and implementation of a new abstract domain, called the *Bit-Vector Inequality* (\mathcal{BVJ}) domain, that is capable of capturing certain inequalities over bit-vector-valued variables. We also consider some variants of the \mathcal{BVJ} domain.

The key insight used to design \mathcal{BVJ} domain (and its variants) involves a new domain combinator (denoted by \mathcal{V}), called the *view-product combinator*. \mathcal{V} constructs a reduced product of two domains [21], using shared *view-variables* to communicate information between the domains.

The *Bit-Vector Memory-Equality Domain* $\mathcal{BVM}\mathcal{E}$, a domain of bit-vector affine-equalities over variables and memory-values, is created by applying the view-product combinator \mathcal{V} to the bit-vector memory domain (§5.3) and the bit-vector equality domain. The *Bit-Vector Inequality Domain* \mathcal{BVJ} , a domain of bit-vector affine-inequalities over variables, is created by applying \mathcal{V} to the bit-vector equality domain and a bit-vector interval domain. The *Bit-Vector Memory-Inequality Domain* $\mathcal{BVM}\mathcal{I}$, a domain of relational bit-vector affine-inequalities over variables and memory, is then created by applying \mathcal{V} to the $\mathcal{BVM}\mathcal{E}$ domain and the bit-vector interval domain. The latter construction illustrates that \mathcal{V} composes: the $\mathcal{BVM}\mathcal{I}$ domain is created via two applications of \mathcal{V} .

This work makes the following contributions:

- The bit-vector memory domain, a non-relational memory domain capable of expressing invariants involving memory accesses.
- The view-product combinator \mathcal{V} , a general procedure to construct more expressive domains.
- Three domains for machine-code analysis constructed using \mathcal{V} :

- The bit-vector memory-equality domain $\mathcal{BVM}\mathcal{E}$, which captures equality relations among bit-vector variables and memory.
- The bit-vector inequality domain $\mathcal{BV}\mathcal{I}$, which captures inequality relations among bit-vector variables.
- The bit-vector memory-inequality domain $\mathcal{BVM}\mathcal{I}$, which captures inequality relations among bit-vector variables and memory.
- A procedure for synthesizing best abstract operations for reduced products of domains that meet certain requirements.
- Experimental results that illustrate the effectiveness of the $\mathcal{BV}\mathcal{I}$ domain applied to machine-code analysis. On average (geometric mean), our $\mathcal{BV}\mathcal{I}$ -based analysis is about 3.5 times slower than an affine-equality-based analysis, while finding improved (more-precise) invariants at 29% of the branch points.

1.2.5 Sound Bit-Precise Numerical Domains Framework for Inequalities

Chapter 6 describes a second approach to designing and implementing a bit-precise relational domain capable of expressing inequality invariants. This work presents the design and implementation of a new framework for abstract domains, called the *Bit-Vector-Sound Finite-Disjunctive (BVSEFD)* domains, which are capable of capturing useful program invariants such as inequalities over bit-vector-valued variables.

We introduce a class of abstract domains, called $BVS(\mathcal{A})$, that is sound with respect to bitvectors whenever \mathcal{A} is sound with respect to mathematical integers. The \mathcal{A} domain can be any numerical abstract domain. For example, it can be the polyhedral domain, which can represent useful program invariants as inequalities. We also describe how to create abstract transformers for $BVS(\mathcal{A})$ that are sound with respect to bitvectors. For $v \subseteq V$ and $av \in \mathcal{A}$, we denote the result by $WRAP_v(av)$; the operation

performs wraparound on av for variables in v . We give an algorithm for $WRAP_v(av)$ that works for any relational abstract domain. We use a finite number of disjunctions of \mathcal{A} elements—captured in the domain $\mathcal{FD}_k(\mathcal{A})$ —to help retain precision. Our contributions are the following:

- We propose a framework for abstract domains, called $BVSFD_d(\mathcal{A})$, to express bit-precise relational invariants by performing wrap-around over abstract domain \mathcal{A} and using disjunctions to retain precision. This abstract domain is parameterized by a positive value d , which provides the maximum number of disjunctions that the abstract domain can make use of.
- We provide a generic technique via reinterpretation to create the abstract transformer for the path through a basic block to a given successor, such that the transformer incorporates lazy wrap-around.
- We present experiments to show how the performance and precision of $BVSFD_d$ analysis changes with the tunable parameter d .

1.3 Thesis Organization

The remainder of the thesis is organized as follows. Chapter 2 provides background on concrete semantics, abstract interpretation and automatic generation of abstract transformers via semantic reinterpretation and symbolic abstraction. Chapter 3 presents the theoretical and experimental comparison of MOS and KS domains. Chapter 4 presents the Affine-Transformers Abstraction framework, which generalizes the MOS domain. Chapter 5 describe the design and implementation of the Bit-Vector Inequality (BVI) domain, that is capable of capturing certain inequalities over bit-vector-valued variables (which represent a program’s registers and/or its memory variables). Chapter 6 describes the design and implementation of the Bit-Vector-Sound Finite-Disjunctive framework. Chapter 7 concludes and describes possible future work.

2 BACKGROUND

This chapter presents a gentle introduction to concrete semantics (§2.1), abstract interpretation (§2.2, §2.3, and §2.4) and weighted pushdown systems used to perform abstract interpretation (§2.5). The notion of the best abstract transformer is also introduced (§2.3.2).

2.1 Concrete Semantics

A semantics is a mathematical characterization of program behavior. The concrete semantics is the most precise semantics describing the actual execution of the program. Let us consider a very simple concrete language to illustrate concrete semantics. An SLang program is a sequence of basic blocks with execution starting from the first block. Each basic block consists of a sequence of statements and a list of control-flow instructions.

The statements are restricted to an assignment of a linear expression. A control-flow instruction consists of either a jump statement or a conditional statement. The guard in the conditional statement can be true, false, negation of a condition, conjunction/disjunction of two conditions, an equality/inequality on two bitvector expressions, or a modulus check operation. For the sake of simplicity, we assume that each variable in the program is a 32-bit unsigned integer.

```

⟨SLang⟩ :: (Block)*
⟨Block⟩ :: †: (⟨Stmt⟩;)* ⟨Next⟩
⟨Next⟩ :: jump †;
        | if  $\phi_{\text{Cond}}$  then jump †1 else jump †2
⟨Expr⟩ :: n | n * v + ⟨Expr⟩
⟨Stmt⟩ :: v = ⟨Expr⟩
⟨Cond⟩ :: True | False | Not(⟨Cond⟩) | ⟨BCond⟩ | ⟨ECond⟩
⟨ECond⟩ :: ⟨Expr⟩ ⟨ExprOp⟩ ⟨Expr⟩ | ⟨Expr⟩ % ⟨Expr⟩ == ⟨Expr⟩
⟨BCond⟩ :: ⟨Cond⟩ ⟨CondOp⟩ ⟨Cond⟩
⟨CondOp⟩ :: And | Or
⟨ExprOp⟩ :: == | < | ≤ | %

```

A concrete state σ for an SLang program is a tuple of concrete values, $\sigma : \prod_{v \in V} BV^{s(v)}$, where $s(v)$ is the size of variable v in bits and BV^b is a bitvector with b bits. We use k to denote the size of the vocabulary V . We use the row vector $\vec{v} = (v_1, v_2, \dots, v_k)$ to denote the tuple of variables in the vocabulary V (which makes explicit the implicit order in " $\sigma : \prod_{v \in V} BV^{s(v)}$ "). A concrete state σ is described by the tuple of bitvectors $(bv_1, bv_2, \dots, bv_k)$, where the value of variable v_i is bv_i . The concrete states described by a tuple of length k are one-vocabulary concrete states. Let \mathbb{C} be the set of all one-vocabulary concrete states.

The concrete transformer τ_E for a basic-block edge $E = B \rightarrow B'$ is described by $\llbracket E \rrbracket_{\text{Block}}$. A concrete transformer is a *two-vocabulary* transition relation, denoted by $R[\vec{v}; \vec{v}']$, where $\vec{v}' = (v'_1, v'_2, \dots, v'_k)$ is a row vector of length k . \vec{v} and \vec{v}' represent the pre-transformation state and the post-transformation state, respectively. τ_E is set of tuples of bitvectors of the form $(bv_1, bv_2, \dots, bv_k, bv'_1, bv'_2, \dots, bv'_k)$. Thus, a concrete transformer is a set of two-vocabulary concrete states. A (two-vocabulary) concrete state is sometimes called an *assignment* to the variables of the pre-state and the post-state vocabulary. Let \mathbb{C} be the set of all concrete transformers. The concrete domain $\mathcal{C} = \mathcal{P}[\mathbb{C}]$ is the powerset of the set of concrete transformers.

Fig. 2.1 provides the concrete semantics for the SLang program. The concrete semantics for an edge $B \rightarrow B'$ from the basic block B to basic block B' is denoted by $\llbracket B \rightarrow B' \rrbracket$, where $\llbracket \cdot \rrbracket : \mathbb{C}$ represents the concrete evaluator. Rule 2.1 in Fig. 2.1 specifies how the concrete evaluation for basic-block pairs feeds into concrete evaluation on a sequence of statements. The concrete evaluator takes a one-vocabulary concrete state σ as input. Note that l' is the label corresponding to the basic block B' . Rule 2.2 states that the concrete evaluation of a sequence of instructions can be broken down into a concrete evaluation of a smaller sequence of instructions, by recursively performing statement-level concrete evaluation

$\llbracket \cdot \rrbracket_{\text{stmt}}$ on the first instruction in the sequence. Rules 2.3, 2.4, 2.5, and 2.6 handle control-flow statements. Rule 2.3 delegates the responsibility of executing the last instruction in the statement sequence to $\llbracket \cdot \rrbracket_{\text{next}}$. Rule 2.4 deals with unconditional-jump instructions. Rules 2.5 and 2.6 deal with conditional-jump instructions; both delegate to the assume evaluation function $\llbracket \cdot \rrbracket_{\text{assume}}$. The assume is given the guard condition ϕ_{Cond} for the true case and $\text{Not}(\phi_{\text{Cond}})$ for the false case. Rule 2.7 deals with assumes. Assume operations filters out the concrete states that do not match the assume condition. If the evaluation of ϕ on the concrete state σ is true, then the state σ is returned; otherwise, a special empty concrete state $\{\}$ is returned. The result $\llbracket E \rrbracket_{\text{Block}} \sigma = \{\}$ specifies that the concrete state σ cannot concretely execute along the edge E .

Rule 2.8 handles assignment statements. Assignment $\llbracket v = \text{exp} \rrbracket_{\text{stmt}}$ returns a concrete state that is the same as σ except the variable v is updated with the evaluated value of exp . Rules 2.9-2.12 specify the concrete evaluation of linear expressions. Similarly, 2.13-2.17 specify the concrete evaluation of conditional expressions.

2.2 Abstract Interpretation

Abstract Interpretation is a process of discovering properties about the program by “running” a safe approximation of the program [20]. This safe approximation of the program is called an *abstraction* of the program. The building block for the abstraction of the program is an *abstract domain*, denoted by \mathcal{A} . An abstract domain defines the level and detail of the abstraction. The program properties inferred by means of abstract interpretation are a safe approximation of actual program properties, and hence they are invariants for the program.

Basic Block:

$$\llbracket B \rightarrow B' \rrbracket_{\text{Block}} \sigma = \llbracket [s_1; \dots; s_n; \text{next}] \rrbracket_{\text{Seq}}(\sigma, l') \quad (2.1)$$

$$\llbracket [s_1; \dots; s_n; \text{next}] \rrbracket_{\text{Seq}}(\sigma, l') = \llbracket [s_2; \dots; s_n; \text{next}] \rrbracket_{\text{Seq}}(\llbracket [s_1] \rrbracket_{\text{Stmt}} \sigma, l') \quad (2.2)$$

Control Flow:

$$\llbracket [\text{next}] \rrbracket_{\text{Seq}}(\sigma, l') = \llbracket \text{next} \rrbracket_{\text{Next}}(\sigma, l') \quad (2.3)$$

$$\llbracket \text{jump } l' \rrbracket_{\text{Next}}(\sigma, l') = \sigma \quad (2.4)$$

$$\begin{aligned} \llbracket \text{if } \phi_{\text{Cond}} \text{ then jump } l_1 \text{ else jump } l_2 \rrbracket_{\text{Next}}(\sigma, l_1) = \\ \llbracket \text{assume}(\phi_{\text{Cond}}) \rrbracket_{\text{Assume}} \sigma \end{aligned} \quad (2.5)$$

$$\begin{aligned} \llbracket \text{if } \phi_{\text{Cond}} \text{ then jump } l' \text{ else jump } l_2 \rrbracket_{\text{Next}}(\sigma, l_2) = \\ \llbracket \text{assume}(\text{Not}(\phi_{\text{Cond}})) \rrbracket_{\text{Assume}} \sigma \end{aligned} \quad (2.6)$$

$$\llbracket \Phi \rrbracket_{\text{Assign}} \sigma = \begin{cases} \sigma & \text{if } \llbracket \Phi \rrbracket_{\text{Cond}} \sigma = \text{True} \\ \{\} & \text{otherwise} \end{cases} \quad (2.7)$$

Assignments:

$$\llbracket v = \text{exp} \rrbracket_{\text{Stmt}} \sigma = \sigma[v \leftarrow \llbracket \text{exp} \rrbracket_{\text{Expr}} \sigma] \quad (2.8)$$

Expressions:

$$\llbracket n \rrbracket_{\text{Expr}} \sigma = n \quad (2.9)$$

$$\llbracket v \rrbracket_{\text{Expr}} \sigma = \sigma[v] \quad (2.10)$$

$$\llbracket \text{exp}_1 * \text{exp}_2 \rrbracket_{\text{Expr}} \sigma = \llbracket \text{exp}_1 \rrbracket_{\text{Expr}} \sigma * \llbracket \text{exp}_2 \rrbracket_{\text{Expr}} \sigma \quad (2.11)$$

$$\llbracket \text{exp}_1 + \text{exp}_2 \rrbracket_{\text{Expr}} \sigma = \llbracket \text{exp}_1 \rrbracket_{\text{Expr}} \sigma + \llbracket \text{exp}_2 \rrbracket_{\text{Expr}} \sigma \quad (2.12)$$

Conditions:

$$\llbracket \text{True/False} \rrbracket_{\text{Cond}} \sigma = \text{True/False} \quad (2.13)$$

$$\llbracket \text{exp}_1 \text{ Op } \text{exp}_2 \rrbracket_{\text{Cond}} \sigma = \llbracket \text{exp}_1 \rrbracket_{\text{Expr}} \sigma \text{ Op } \llbracket \text{exp}_2 \rrbracket_{\text{Expr}} \sigma \quad (2.14)$$

$$\llbracket \text{Not}(\text{exp}) \rrbracket_{\text{Cond}} \sigma = \text{Not}(\llbracket \text{exp} \rrbracket_{\text{Expr}} \sigma) \quad (2.15)$$

$$\llbracket \text{exp}_1 \% \text{exp}_2 == \text{exp}_3 \rrbracket_{\text{Cond}} \sigma = \llbracket \text{exp}_1 \rrbracket_{\text{Expr}} \sigma \% \llbracket \text{exp}_2 \rrbracket_{\text{Expr}} \sigma == \llbracket \text{exp}_3 \rrbracket_{\text{Expr}} \sigma \quad (2.16)$$

$$\llbracket \text{cond}_1 \text{ Op } \text{cond}_2 \rrbracket_{\text{Cond}} \sigma = \llbracket \text{cond}_1 \rrbracket_{\text{Cond}} \sigma \text{ Op } \llbracket \text{cond}_2 \rrbracket_{\text{Cond}} \sigma \quad (2.17)$$

Figure 2.1: Concrete semantics for $\mathcal{L}(\text{SLANG})$.

2.2.1 Abstract Domain

As mentioned in §2.1, the *concrete domain*, denoted by \mathcal{C} , is the powerset of the set of concrete transformers. Often, but not always, the concrete

domain \mathcal{C} and the abstract domain \mathcal{A} form a Galois connection $\mathcal{G} = \mathcal{C} \xrightleftharpoons[\alpha]{\gamma} \mathcal{A}$ [20]. γ takes an abstract-domain element A , and gives back the concrete-domain element $C \in \mathcal{C}$, that A represents. α takes a concrete domain value $C \in \mathcal{C}$, and abstracts C to the least abstract-domain element A such that the set of concrete elements represented by A is a superset of C . For a given abstract domain \mathcal{A} , $\mathcal{A}[V]$ denotes the specific instance of \mathcal{A} that is defined over vocabulary V .

Parity Abstract Domain

As an illustrative example of an abstract domain, consider the *parity* abstract domain that only tracks whether each variable in the program is even or odd. A parity map is a mapping from variables to their parity, $\text{pm} : V \rightarrow \{e, o\}$, where e represents the set of even integers and o represents the set of odd integers. Thus, $\gamma(e) = \{0, 2, 4, \dots, 2^{32} - 2\}$ and $\gamma(o) = \{1, 3, 5, \dots, 2^{32} - 1\}$ for 32-bit unsigned integers. The concretization $\gamma(\text{pm}) = \{(bv_1, bv_2, \dots, bv_k) : \bigwedge_{0 \leq i \leq k} bv_i \in \gamma(\text{pm}(v_i))\}$. Let us denote the set of all parity maps by \mathcal{PM} . We define the the parity abstract domain, denoted by \mathcal{A}_{Par} , as the powerset of \mathcal{PM} . Therefore, $\mathcal{A}_{\text{Par}} = \mathcal{P}[\mathcal{PM}]$. The concretization $\gamma(a)$ for an element $a \in \mathcal{A}_{\text{Par}}$ is defined as the union over the concretizations of the parity maps in a .

$$\gamma(a) = \bigcup_{\text{pm} \in a} \gamma(\text{pm})$$

Example 2.1. Consider $V = \{v, v'\}$ and $a \in \mathcal{A}_{\text{Par}}[V]$, such that $a = \{(e, e), (o, o)\}$. The abstract element a represent all concrete states such that $\text{parity}(v) = \text{parity}(v')$. Consider the abstraction for the concrete state set $\text{css} = \{(1, 3), (11, 13), (211, 215)\}$. $\alpha(\text{css}) = \{(o, o)\}$. Notice that some precision is lost in the abstraction process. For instance, the concrete state $(1, 1) \in \gamma(\alpha(\text{css}))$, but $(1, 1) \notin \text{css}$.

Table 2.1: Abstract-domain operations.

Result Type	Operation	Description
\mathcal{A}	\perp	<i>bottom element</i>
bool	$(\mathbf{a}_1 == \mathbf{a}_2)$	<i>equality</i>
\mathcal{A}	$(\mathbf{a}_1 \sqcup \mathbf{a}_2)$	<i>join</i>
\mathcal{A}	$(\mathbf{a}_1 \nabla \mathbf{a}_2)$	<i>widen</i>
\mathcal{A}	Id	<i>identity element</i>
\mathcal{A}	$(\mathbf{a}_1 \circ \mathbf{a}_2)$	<i>composition</i>

2.2.2 Abstract-Domain Operations

A static analyzer needs a way to construct abstract transformers for the concrete operations in the programs. *Semantic reinterpretation* and *symbolic abstraction* provide automatic ways to construct abstract transformers. We provide a more detailed discussion of these automatic techniques in §2.3.

For an analysis that provides function summaries or loop summaries, the fixpoint analysis is performed using equality, join (\sqcup), and abstract-composition (\circ) operations on abstract transformers. We provide a more detailed description of the fixpoint analysis in §2.4.

For the analyses discussed in this dissertation, the program is abstracted to a control-flow graph, where each edge in the graph is labeled with an abstract transformer. As described in §2.1, an abstract transformer is a *two-vocabulary* transition relation $R[\vec{v}; \vec{v}']$, where \vec{v} and \vec{v}' are row vectors of length k that represent the pre-transformation state and the post-transformation state, respectively.

Tab. 2.1 lists the abstract-domain operations needed to generate the program abstraction and perform fixpoint analysis on it. The bottom element is the abstract-domain element at the bottom of the abstract domain. The concretization of the bottom element is the empty set: $\gamma(\perp) = \emptyset$. Equality is the standard lattice equality operation. Join is the least upper bound operation on the abstract-domain elements, that is, $\mathbf{a}_1 \sqcup \mathbf{a}_2 \sqsupseteq \mathbf{a}_1, \mathbf{a}_2$ and $\forall \mathbf{a}_3 : \mathbf{a}_3 \sqsupseteq \mathbf{a}_1, \mathbf{a}_2 \Rightarrow \mathbf{a}_3 \sqsupseteq \mathbf{a}_1 \sqcup \mathbf{a}_2$. The widen operation is needed

for domains with infinite ascending chains to ensure termination [23]. For abstract domains with finite ascending chains, the widen operation can be safely replaced with the join operation. Id is the identity element, which represents the identity transformation ($\bigwedge_{i=1}^k v'_i = v_i$). Finally, the abstract-composition operation $\alpha_1 \circ \alpha_2$ returns a sound overapproximation of the composition of abstract transformer α_1 with abstract transformer α_2 .

For the parity abstract domain, join is set union and meet is set intersection. The empty set $\{\}$ is the bottom element for \mathcal{A}_{Par} . Consider $\vec{v} = (v, v')$. The abstract-domain element $\{(e, e), (o, o)\}$ is the identity transformer, specifying that the parity of the variable v is preserved.

Example 2.2. Consider $\alpha_1 = \{(e, e), (o, o)\}$ and $\alpha_2 = \{(e, o), (o, e)\}$. The abstract transformer α_1 states that the parity of the variable v is preserved. Similarly, the abstract transformer α_2 states that the parity of variable v_2 is flipped. The join of α_1 and α_2 , denoted by $\alpha_1 \sqcup \alpha_2$ is $\{(e, e), (o, o), (e, o), (o, e)\}$. $\alpha_1 \sqcup \alpha_2$ represents the set of all concrete states because all the permutations of parity transformations are possible. The meet of α_1 and α_2 , denoted by $\alpha_1 \sqcap \alpha_2$ is $\{\}$.

Consider the abstract composition of $\alpha_1 \circ \alpha_2$. Intuitively, the abstract transformer $\alpha_1 \circ \alpha_2$ takes the parity of the variable v and,

1. applies the abstract transformation α_1 to v 's parity (which preserves the parity) to get a temporary parity t_p ,
2. applies the abstract transformation α_2 to the parity t_p (which flips the parity).

The overall effect is that the parity of the variable v is flipped. Thus, $\alpha_1 \circ \alpha_2 = \{(e, o), (o, e)\} = \alpha_2$.

2.3 Creation of Abstract Transformers.

To perform program analysis, the program-state transitions that are associated with the edges of a control-flow graph also need to be abstracted. We will use the map $\llbracket \cdot \rrbracket^\# : E \rightarrow \mathcal{A}$ to denote the map that specifies abstract trans-

formers for each edge in the CFG. We say that $\llbracket e \rrbracket^\#$ is a sound approximation of $\llbracket e \rrbracket$ if the following condition holds for every edge $e \in E$:

$$\gamma(\llbracket e \rrbracket^\#) \supseteq \gamma(\llbracket e \rrbracket)$$

Example 2.3. *The following example is used to illustrate abstract interpretation with the parity abstract domain:*

```

L0:   v=v+1
L1:   while (*) {
L2:     v=v+2
      }
L3:   if (v%4==2) {
L4:     v=v+1
L5:   }
END:

```

For instance, the abstract transformer for the concrete operation $v = v + 2$ starting from node L2 and ending at node L1, denoted by $\tau_{L2 \rightarrow L1}^\#$, is defined as $\{(e, e), (o, o)\}$. Thus, the abstract transformer $\tau_{L2 \rightarrow L1}^\#$ is the identity function on parity maps because $v = v + 2$ does not change the parity of variable v .

More interesting is the abstract transformer for $v = v + 1$ which flips the parity of v . $\tau_{L0 \rightarrow L1}^\#$ is defined as $\{(e, o), (o, e)\}$. \square

2.3.1 Semantic Reinterpretation

In this section, we describe how the abstract transformers are generated using semantic reinterpretation [47, 79, 81, 64, 60, 30]. Semantic reinterpretation is an automatic and efficient method to create abstract transformers. Semantic reinterpretation makes use of abstract components of each of the concrete operations and concrete data types used to specify the semantics of a programming language. These abstract components are collectively referred to as the semantic core [60] (sometimes called a semantic algebra [90]). Specifically, the reinterpretation consists of a domain of abstract transformers $\mathcal{A}[V, V']$, a domain of abstract integers $\mathcal{A}^{\text{INT}}[t, V]$, and operations to lookup a variable's value in the post-state

of an abstract transformer and to create an updated version of a given abstract transformer. Here t denotes a temporary variable not in V or V' . Given blocks $B : [l : s_1; \dots; s_n; \text{next}]$ and $B' : [l' : s'_1; \dots; s'_n; \text{next}']$ in an SLang program (see §2.1), where B' is a successor of B , reinterpretation of B can provide an abstract transformer for the transformation that starts from the first instruction in B and ends in the first instruction in B' , denoted by $B \rightarrow B'$.

Fig. 2.2 provide the reinterpretation semantics for the \mathcal{A}_{Par} domain. Rule 2.18 specifies how abstract-transformer evaluation for basic-block pairs feeds into abstract-transformer evaluation on a sequence of statements. The evaluation on a sequence of statements starts with the identity abstract transformer, denoted by Id . Rule 2.19 states that the abstract transformer for a sequence of instruction can be broken down into an abstract transformer for a smaller sequence of instruction, by recursively performing statement-level abstract interpretation $\llbracket \cdot \rrbracket_{\text{Stmt}}^\#$ on the first instruction in the sequence. In this rule and subsequent $\llbracket \cdot \rrbracket_{\text{Next}}^\#$ and $\llbracket \cdot \rrbracket_{\text{Stmt}}^\#$ rules, “ a ” denotes the intermediate abstract transformer value. It starts as Id at the beginning of the instruction sequence, and gets updated or accessed by assignment and control-flow statements in the sequence.

Rules 2.20, 2.21, 2.22 and 2.23 handle control-flow statements. Rule 2.20 delegates the responsibility of executing the last instruction in the statement sequence to $\llbracket \cdot \rrbracket_{\text{Next}}^\#$. Rule 2.21 deals with unconditional-jump instructions. Rules 2.22 and 2.23 handle conditional branching by delegating the responsibility to $\llbracket \text{assume}(\phi_{\text{Cond}}) \rrbracket_{\text{Assume}}^\# a$ and $\llbracket \text{assume}(\text{Not}(\phi_{\text{Cond}})) \rrbracket_{\text{Assume}}^\# a$ respectively. Rule 2.24 handles assume by delegating the abstract execution of the conditional ϕ to $\llbracket \cdot \rrbracket_{\text{Cond}}^\#$ (see rules 2.32-2.35) and then performing a meet operation.

Rule 2.25 handles the assignment statement by merely performing a post-state-vocabulary update on the current abstract transformer “ a .”

Rules 2.26-2.29 handle reinterpretation of expressions. The abstract

Basic Block:

$$\llbracket B \rightarrow B' \rrbracket_{\text{Block}}^\sharp = \llbracket [s_1; \dots; s_n; \text{nxt}] \rrbracket_{\text{Seq}}^\sharp(\text{Id}, l') \quad (2.18)$$

$$\llbracket [s_1; \dots; s_n; \text{nxt}] \rrbracket_{\text{Seq}}^\sharp(a, l') = \llbracket [s_2; \dots; s_n; \text{nxt}] \rrbracket_{\text{Seq}}^\sharp(\llbracket [s_1] \rrbracket_{\text{Stmnt}}^\sharp(a), l') \quad (2.19)$$

Control Flow:

$$\llbracket [\text{nxt}] \rrbracket_{\text{Seq}}^\sharp(l') = \llbracket \text{nxt} \rrbracket_{\text{Next}}^\sharp(l') \quad (2.20)$$

$$\llbracket \text{jump } l' \rrbracket_{\text{Next}}^\sharp(a, l') = a \quad (2.21)$$

$$\begin{aligned} \llbracket \text{if } \phi_{\text{Cond}} \text{ then jump } l_1 \text{ else jump } l_2 \rrbracket_{\text{Next}}^\sharp(a, l_1) = \\ \llbracket \text{assume}(\phi_{\text{Cond}}) \rrbracket_{\text{Assume}}^\sharp a \end{aligned} \quad (2.22)$$

$$\begin{aligned} \llbracket \text{if } \phi_{\text{Cond}} \text{ then jump } l' \text{ else jump } l_2 \rrbracket_{\text{Next}}^\sharp(a, l_2) = \\ \llbracket \text{assume}(\text{Not}(\phi_{\text{Cond}})) \rrbracket_{\text{Assume}}^\sharp a \end{aligned} \quad (2.23)$$

$$\llbracket \phi \rrbracket_{\text{Assume}}^\sharp a = a \sqcap \llbracket \phi \rrbracket_{\text{Cond}}^\sharp a \quad (2.24)$$

Assignments:

$$\llbracket v = \text{exp} \rrbracket_{\text{Stmnt}}^\sharp a = \text{update}(a, v', \llbracket \text{exp} \rrbracket_{\text{Expr}}^\sharp a) \quad (2.25)$$

Expressions:

$$\llbracket n \rrbracket_{\text{Expr}}^\sharp a = \{(\text{parity}(n), p_1, p_2, \dots, p_k) : p_i \in \{e, o\}\} \quad (2.26)$$

$$\llbracket v \rrbracket_{\text{Expr}}^\sharp a = \text{lookup}(a, v') \quad (2.27)$$

$$\llbracket \text{exp}_1 * \text{exp}_2 \rrbracket_{\text{Expr}}^\sharp a = \llbracket \text{exp}_1 \rrbracket_{\text{Expr}}^\sharp a *^\sharp \llbracket \text{exp}_2 \rrbracket_{\text{Expr}}^\sharp a \quad (2.28)$$

$$\llbracket \text{exp}_1 + \text{exp}_2 \rrbracket_{\text{Expr}}^\sharp a = \llbracket \text{exp}_1 \rrbracket_{\text{Expr}}^\sharp a +^\sharp \llbracket \text{exp}_2 \rrbracket_{\text{Expr}}^\sharp a \quad (2.29)$$

Abstract Integers:

$$\begin{aligned} i *^\sharp i' = \{(\text{parity}(p_t) *^\sharp \text{parity}(p'_t), p_1, p_2, \dots, p_k) : \\ (\text{parity}(p_t), p_1, p_2, \dots, p_k) \in i, (\text{parity}(p'_t), p_1, p_2, \dots, p_k) \in i'\} \end{aligned} \quad (2.30)$$

$$\begin{aligned} i +^\sharp i' = \{(\text{parity}(p_t) +^\sharp \text{parity}(p'_t), p_1, p_2, \dots, p_k) : \\ (\text{parity}(p_t), p_1, p_2, \dots, p_k) \in i, (\text{parity}(p'_t), p_1, p_2, \dots, p_k) \in i'\} \end{aligned} \quad (2.31)$$

Conditions:

$$\llbracket \text{True/False} \rrbracket_{\text{Cond}}^\sharp a = \top/\perp \quad (2.32)$$

$$\begin{aligned} \llbracket \text{exp \% 2} == 0 \rrbracket_{\text{Cond}}^\sharp a = \{(p_1, \dots, p_k, p_1, \dots, p_k) : \\ (e, p_1, \dots, p_k) \in \llbracket \text{exp} \rrbracket_{\text{Expr}}^\sharp a\} \end{aligned} \quad (2.33)$$

$$\begin{aligned} \llbracket \text{exp \% 2} == 1 \rrbracket_{\text{Cond}}^\sharp a = \{(p_1, \dots, p_k, p_1, \dots, p_k) : \\ (o, p_1, \dots, p_k) \in \llbracket \text{exp} \rrbracket_{\text{Expr}}^\sharp a\} \end{aligned} \quad (2.34)$$

$$\llbracket \phi_{\text{Other}} \rrbracket_{\text{Cond}}^\sharp a = \top \quad (2.35)$$

Variable lookup and update:

$$\text{lookup}(a, v'_i) = \{(p'_i, p_1, \dots, p_k) : (p_1, \dots, p_k, p'_1, \dots, p'_k) \in a\} \quad (2.36)$$

$$\begin{aligned} \text{update}(a, v'_i, ai) = \{(p_1, p_2, \dots, p_k, p'_1, p'_2, \dots, p'_{i-1}, p_t, p_{i+1}, \dots, p'_k) : \\ (p_1, p_2, \dots, p_k, p'_1, p'_2, \dots, p'_k) \in a, (p_t, p_1, p_2, \dots, p_k) \in ai\} \end{aligned} \quad (2.37)$$

Figure 2.2: Reinterpretation semantics for $\mathcal{L}(\text{SLANG})$ for domain \mathcal{A}_{PAR} .

p_1	p_2	$*^\#$	$+^\#$
e	e	e	e
e	o	e	o
o	e	e	o
o	o	o	e

Table 2.2: Truth table for multiplication and addition of two parity values.

interpretation of an expression gives back an abstract integer. An abstract integer is an abstract-domain element over (t, \vec{v}) , where t is the temporary variable not in \vec{v} referring to the value of abstract integer. For example, the abstract integer $i = \{(e, e), (o, o)\}$, with $\vec{v} = \{v_1\}$, is an abstract integer whose parity is the same as $\{v_1\}$. Rule 2.26 creates an abstract integer from a constant integer n by setting the *parity* of the t variable to $parity(n)$ in all possible pre-vocabulary-state parity combinations. Rule 2.27 creates an abstract integer from a variable by performing a variable “lookup” operation in α . Rules 2.28 and 2.29 delegate computation to the corresponding abstract-integer operations.

Rules 2.30 and 2.31 are abstract-multiplication and abstract-addition operations respectively. The operators $*^\#$ and $+^\#$ for $\{e, o\}$ are defined in Tab. 2.2.

Rules 2.32-2.35 provide abstract interpretation of conditionals. The result of these operations is an abstract transformer that overapproximates the set of values satisfied by the conditional for the given input abstract transformer α . Rule 2.32 handles the true and false conditionals. Rules 2.33 and 2.34 handles the conditionals which check for the parity of an expression “ exp ”. Rule 2.33 filters the tuples in the abstract integers for “ exp ”, where the parity for “ exp ” is even. Similarly, rule 2.34 only considers tuples where the parity for “ exp ” is odd. Rule 2.35 safely handles all the other expression to return \top . For brevity, we do not explicitly give rules for the Boolean operators AND, OR and NOT, but they can be implemented

as set intersection, set union, and set complementation, respectively.

Rules 2.36 and 2.37 are variable lookup and update operations. Lookup takes an abstract transformer α and a variable v_i , and returns the abstract integer α_i such that the relationship of t with \vec{v} in α_i is the same as the relationship of v'_i with \vec{v} in α . The variable-update operation works in the opposite direction. Update takes an abstract transformer α , a variable v'_i , and α_i , and returns α' such that the relationship of v'_i with \vec{v} in α' is the same as the relationship of t with \vec{v} in α_i , and all the other relationships that do not involve v' remain the same.

Example 2.4. *For instance, consider the abstract transformer for the concrete operation $v = v + 1$, denoted by $\tau_{L0 \rightarrow L1}^\#$. Remember that the desired answer is $\{(e, o), (o, e)\}$. Semantic reinterpretation can arrive at this answer by performing $\llbracket [v = v + 1; \text{jmp } L1] \rrbracket_{\text{Seq}}^\#(\text{Id}, L1)$, where $\text{Id} = \{(e, e), (o, o)\}$. First, it will perform $\llbracket v = v + 1 \rrbracket_{\text{Stmnt}}^\#$, which in turn will obtain $\llbracket v \rrbracket_{\text{Expr}}^\# = \{(e, e), (o, o)\}$ (through rule 2.36) and $\llbracket 1 \rrbracket_{\text{Expr}}^\# = \{(o, e), (o, o)\}$ (through rule 2.26). Using these values, $\llbracket v + 1 \rrbracket_{\text{Expr}}^\# = \llbracket v \rrbracket_{\text{Expr}}^\# +^\# \llbracket 1 \rrbracket_{\text{Expr}}^\# = \{(e +^\# o, e), (o +^\# o, o)\} = \{(o, e), (e, o)\}$. Then, $\text{update}(\text{Id}, v', \{(o, e), (e, o)\})$ is performed to obtain $\llbracket v = v + 1 \rrbracket_{\text{Stmnt}}^\# = \{(e, o), (o, e)\}$. Finally, this value is left unchanged by the unconditional jump to L1 to obtain $\tau_{L0 \rightarrow L1}^\# = \{(e, o), (o, e)\}$.*

2.3.2 Symbolic Abstraction

Symbolic abstraction (denoted by $\hat{\alpha}$) [84] provides an automatic way of obtaining the best abstract transformer, given the concrete semantics of the transformer as a formula (for abstract domains and logics that meet certain properties [84]).

Best Abstract Transformer For every concrete operation (concrete transformer) in the actual program, the corresponding over-approximating abstract operation (abstract transformer) needs to be provided to perform

abstract interpretation (§2.2). Cousot and Cousot [21] give the notion of a *best* abstract transformer, but do not provide an algorithm to either (i) apply the best transformer, or (ii) create a representation of the best transformer. Their notion is as follows: If τ is a concrete transformer, let τ^s be τ lifted to operate on a set pointwise (i.e., $\tau^s(S) = \{\tau(s) | s \in S\}$). Then the best abstract transformer τ^\sharp equals $\alpha \circ \tau^s \circ \gamma$.

An abstract transformer can be viewed as an element in an abstract domain in which each element represents an input-output relation. The advantage of adopting this viewpoint is that it allows symbolic abstraction — and in particular *algorithms* for symbolic abstraction — to be applied to the problem of creating best abstract transformers. That is, if φ_τ is a formula in \mathcal{L} that captures the semantics of concrete transformer τ , \mathcal{A} is the abstract domain of abstract transformers, and $\hat{\alpha}_{\mathcal{A}}$ is a symbolic-abstraction algorithm for mapping an \mathcal{L} formula to the smallest over-approximating value in \mathcal{A} , then $\hat{\alpha}_{\mathcal{A}}(\varphi_\tau)$ yields the best abstract transformer for τ . $\hat{\alpha}_{\mathcal{A}}(\varphi_\tau)$ is the symbolic version of the α function, referred as the symbolic-abstraction operation. Similarly, the symbolic-concretization operation, denoted by $\hat{\gamma}_{\mathcal{A}}(a)$, is the symbolic version of the γ function; it takes an abstract value $a \in \mathcal{A}$ and returns a formula $\varphi \in \mathcal{L}$ describing the abstract value symbolically as a logical formula. Thus, a and $\hat{\gamma}_{\mathcal{A}}(a)$ represent the same set of concrete states (i.e., $\gamma_{\mathcal{A}}(a) = \llbracket \hat{\gamma}_{\mathcal{A}}(a) \rrbracket$). Usually, symbolic concretization is straightforward to implement for abstract domains. For instance,

$$\hat{\gamma}_{\mathcal{A}_{\text{par}}}(a) = \bigvee_{(p_1, p_2, \dots, p_k, p'_1, p'_2, \dots, p'_k) \in a} \left(\bigwedge_{i=0}^k (v_i \% 2 = (p_i = e) ? 0 : 1) \wedge (v'_i \% 2 = (p'_i = e) ? 0 : 1) \right)$$

I now discuss two approaches to symbolic abstraction, namely, *symbolic abstraction from below* [84] and the *bilateral approach to symbolic abstraction* [100].

Symbolic Abstraction From Below

Symbolic Abstraction From Below ($\hat{\alpha}^{\text{be}}$) [84] starts from \perp (the abstract value that denotes the empty set), and climbs up the abstract-domain lattice using models obtained from a succession of calls to a Satisfiability Modulo Theory (SMT) solver (Alg. 1). A satisfiability modulo theories (SMT) solver [74, 26] is a decision procedure for logical formulas with respect to combinations of background theories expressed in classical logic. For instance, the logic can be Quantifier Free Bit-Vector logic (QFBV). The SMT solver takes a logical formula and returns SAT (or UNSAT) if the formula is satisfiable (or unsatisfiable). The beta (β) function, also referred as the representation function, obtains an abstract-domain element given a concrete state or model (M). The algorithm refines the search space after each SMT call to ensure progress by using $\neg\hat{\gamma}(\text{ans})$ as a blocking clause. Once the SMT call returns unsatisfiable, ans contains an over-approximation of $\llbracket\varphi\rrbracket$.

This algorithm is not guaranteed to terminate when the abstract-domain lattice has infinite ascending chains. Moreover, because ans is initialized to \perp and always has a value that is \sqsubseteq the desired final answer, the algorithm is unable to give back a non-trivial abstract value, i.e., a value other than \top , in case the SMT solver times out.

Algorithm 1 Symbolic Abstraction From Below $\hat{\alpha}^{\text{be}}(\varphi)$

```

1:  $\text{ans} \leftarrow \perp$ 
2: while  $r = \text{SMTCall}(\varphi \wedge \neg\hat{\gamma}(\text{ans}))$  is Sat do
3:    $M \leftarrow \text{GetModel}(r)$ 
4:    $\text{ans} \leftarrow \text{ans} \sqcup \beta(M)$ 
5: return  $\text{ans}$ 

```

Example 2.5. Consider the abstraction of the concrete transformer $\tau_{L3 \rightarrow L4}$. The symbolic formula for this transformer is $\varphi_{\tau_{L3 \rightarrow L4}} = (v' = v) \wedge (\text{mod}(v', 4) = 2)$, where mod is the modulus operator. First, $\text{ans} = \perp$ and SMTCall returns true.

Suppose the model M is $(v = 2, v' = 2)$. $\beta(v = 2, v' = 2) = (e, e)$, because that is the best value that overapproximates model M . Thus, $\text{ans} = \{(e, e)\}$ and $\hat{\gamma}(\text{ans}) = (\text{mod}(v', 2) = 0 \wedge \text{mod}(v, 2) = 0)$ and consequently, the next SMT query is given as $(v' = v) \wedge (\text{mod}(v', 4) = 2) \wedge \neg(\text{mod}(v', 2) = 0 \wedge \text{mod}(v, 2) = 0)$. This time the query is unsatisfiable and $\text{ans} = \{(e, e)\}$ is returned.

Note that the value obtained via semantic reinterpretation is $\llbracket L3 \rightarrow L4 \rrbracket_{\text{Block}}^{\#} = \top$. This answer is returned because the reinterpretation semantics is myopic and greedy, and is incapable of obtaining the best abstract transformer in the general case.

Bilateral Approach to Symbolic Abstraction

The bilateral algorithm for symbolic abstraction ($\hat{\alpha}^{\text{bi}}$) [100] starts from both below (\perp) and above (\top). (See Alg. 2.) At any point in the computation, it maintains two values, *lower* (an under-approximation of the best abstract transformer) and *upper* (an over-approximation of the best abstract transformer). When *lower* is equal to *upper*, the algorithm stops and returns *upper*. For *lower* and *upper* such that $\text{lower} \sqsubset \text{upper}$, an abstract consequence is any value p , such that $p \sqsupseteq \text{lower}$ and $p \sqsubseteq \text{upper}$. These restrictions on p guarantee that the algorithm makes progress on each step. For conjunctive domains, p can be chosen as any conjunct in *lower* not implied by *upper*.

Like $\hat{\alpha}^{\text{be}}$, this algorithm is not guaranteed to terminate when the abstract-domain lattice has infinite ascending chains. The key advantage of this algorithm over $\hat{\alpha}^{\text{be}}$ is that $\hat{\alpha}^{\text{bi}}$ can return *upper* as a safe over-approximation in case of a timeout. This safe over-approximation of the symbolic abstraction ($\hat{\alpha}$) is denoted by $\tilde{\alpha}$. Thus, for any given φ , $\hat{\alpha}(\varphi) \sqsubseteq \tilde{\alpha}(\varphi)$.

Example 2.6. Consider the abstraction of the concrete transformer $\tau_{L3 \rightarrow L4}$ again. First, $\text{lower} = \perp$, $\text{upper} = \top$ and SMTCall returns true. Suppose that the model

Algorithm 2 Bilateral Symbolic Abstraction $\hat{\alpha}^{\text{bi}}(\varphi)$

```

1:  $lower \leftarrow \perp$ 
2:  $upper \leftarrow \top$ 
3: while  $lower \neq upper$  do
4:    $p \leftarrow \text{AbstractConsequence}(lower, upper)$ 
5:    $r \leftarrow \text{SMTCall}(\varphi \wedge \neg \hat{\gamma}(p))$ 
6:   if  $r$  is Sat then
7:      $M \leftarrow \text{GetModel}(r)$ 
8:      $lower \leftarrow lower \sqcup \beta(M)$ 
9:   else
10:     $upper \leftarrow upper \sqcap p$ 
11: return  $upper$ 

```

M is $(v = 2, v' = 2)$. Thus, $\beta(v = 2, v' = 2) = (e, e)$, $lower = \{(e, e)\}$ and $\hat{\gamma}(lower) = (\text{mod}(v', 2) = 0 \wedge \text{mod}(v, 2) = 0)$. The bilateral algorithm calls the `AbstractConsequence` operation. Note that the abstract consequence p can be created automatically by adding 0 or more tuples to $lower$ such that p is a subset of $upper$. Suppose that `AbstractConsequence` adds the tuple (e, o) , so that $p = \{(e, o), (e, e)\}$. Consequently, the next SMT query is given as $(v' = v) \wedge (\text{mod}(v', 4) = 2) \wedge \neg(\text{mod}(v, 2) = 0)$. This time the query is unsatisfiable and $upper$ is set to $\top \sqcap p = \{(e, o), (e, e)\}$. Because $upper$ is still not equal to $lower$, `AbstractConsequence` p is called again. This time, there is only one option for p ; namely $\{(e, e)\}$. The next SMT query is given as $(v' = v) \wedge (\text{mod}(v', 4) = 2) \wedge \neg(\text{mod}(v', 2) = 0 \wedge \text{mod}(v, 2) = 0)$. The query is unsatisfiable and $upper$ is set to $\{(e, o), (e, e)\} \sqcap p = \{(e, e)\}$. $upper$ is now equal to $lower$, and $upper$ is returned.

Note that the advantage of this approach is that a sound non-trivial value of $upper = \{(e, o), (e, e)\}$ could be returned if the third SMT call timed out.

2.4 Fixed-point computation

To obtain program summaries, an iterative fixed-point computation needs to be performed. Tab. 2.3 provides the fixpoint analysis for Ex. 2.3. Each row in the table shows the value of path summaries $PS_n(L0 \rightarrow L^*)$ starting from the beginning of the program $L0$ to different program points in the program, where n is the iteration number of the path summary. The intermediate path summary at the first iteration (i), denoted by $PS_{(i)}^\sharp(L0 \rightarrow L1)$, is obtained from $\tau_{L0 \rightarrow L1}^\sharp$. The intermediate path summary at the first iteration (ii), denoted by $PS_{(ii)}^\sharp(L0 \rightarrow L1)$ is calculated as the join of the path summary at the previous iteration $PS_{(i)}^\sharp(L0 \rightarrow L1)$ with the abstract composition of $\tau_{L2 \rightarrow L1}^\sharp$ and $PS_{(i)}^\sharp(L0 \rightarrow L1)$. Quiescence is discovered during the second iteration, because $PS_{(ii)}^\sharp(L0 \rightarrow L1) \sqsubseteq PS_{(i)}^\sharp(L0 \rightarrow L1)$. Therefore, $PS_{(ii)}^\sharp(L0 \rightarrow L1)$ is the path summary from $L0$ to $L1$. Using this path summary for program point $L1$, path summary at $L3$, denoted by $PS_{(i)}^\sharp(L0 \rightarrow L3)$ can be calculated by abstract composition of $\tau_{L1 \rightarrow L3}^\sharp$ and $PS_{(ii)}^\sharp(L0 \rightarrow L1)$. In a similar manner, $PS_{(i)}^\sharp(L0 \rightarrow L4)$ is calculated as $\tau_{L3 \rightarrow L4}^\sharp \circ PS_{(i)}^\sharp(L0 \rightarrow L3)$. Because node $L5$ is the confluence point of edges $L3 \rightarrow L5$ and $L4 \rightarrow L5$, the fixpoint iteration needs to perform a join operation of the path summaries obtained from abstract composition along these two edges. Finally, $PS_{(i)}^\sharp(L0 \rightarrow END) = \{(e, o), (o, e), (o, o)\}$ is obtained by composing $PS_{(i)}^\sharp(L0 \rightarrow L5)$ with the abstract semantics $\tau_{L5 \rightarrow END}$ of the edge from $L5$ to END , which has the identity transformation Id .

The abstract interpretation with the parity abstract domain \mathcal{A}_{par} is able to show that if the input value of v is even, then the output value of v will always be odd; if the input value is odd, then nothing is known about the parity of the output value of v .

Table 2.3: Snapshots in the fixed-point analysis for Ex. 2.3.

Iteration	Equation	Value
$PS_{(i)}^\sharp(L0 \rightarrow L1)$	$\tau_{L0 \rightarrow L1}^\sharp$	$\{(e, o), (o, e)\}$
$PS_{(ii)}^\sharp(L0 \rightarrow L1)$	$PS_{(i)}^\sharp(L0 \rightarrow L1) \sqcup$ $(\tau_{L1 \rightarrow L2}^\sharp \circ PS_{(i)}^\sharp(L0 \rightarrow L1))$	$\{(e, o), (o, e)\} \sqcup \{(e, o), (o, e)\}$ $= \{(e, o), (o, e)\}$
$PS_{(i)}^\sharp(L0 \rightarrow L3)$	$\tau_{L1 \rightarrow L3}^\sharp \circ PS_{(ii)}^\sharp(L0 \rightarrow L1)$	$\{(e, e), (o, o)\} \circ \{(e, o), (o, e)\}$ $= \{(e, o), (o, e)\}$
$PS_{(i)}^\sharp(L0 \rightarrow L4)$	$\tau_{L3 \rightarrow L4}^\sharp \circ PS_{(i)}^\sharp(L0 \rightarrow L3)$	$\{(e, e)\} \circ \{(e, o), (o, e)\}$ $= \{(o, e)\}$
$PS_{(i)}^\sharp(L0 \rightarrow L5)$	$(\tau_{L4 \rightarrow L5}^\sharp \circ PS_{(i)}^\sharp(L0 \rightarrow L4)) \sqcup$ $(\tau_{L3 \rightarrow L5}^\sharp \circ PS_{(i)}^\sharp(L0 \rightarrow L3))$	$\{((e, o), (o, e)) \circ \{(o, e)\}\} \sqcup$ $\{(e, e), (o, o)\} \circ \{(e, o), (o, e)\}$ $= \{(o, o)\} \sqcup \{(e, o), (o, e)\}$ $= \{(e, o), (o, e), (o, o)\}$
$PS_{(i)}^\sharp(L0 \rightarrow END)$	$(\tau_{L5 \rightarrow END}^\sharp \circ PS_{(i)}^\sharp(L0 \rightarrow L5))$	$Id \circ \{(e, o), (o, e), (o, o)\}$ $= \{(e, o), (o, e), (o, o)\}$

2.5 Weighted Pushdown Systems

Weighted PushDown Systems (WPDSs) are a modern formalism for solving flow-sensitive, context-sensitive, interprocedural dataflow-analysis problems [10, 85].

A weighted pushdown system is a pushdown system whose rules are given values from some domain of weights. The weight domains of interest are bounded idempotent semirings.

2.5.1 Bounded Idempotent Semiring

Definition 2.7. A *bounded idempotent semiring* is a quintuple $(D, \oplus, \otimes, 0, 1)$, where D is a set, 0 and 1 are elements of D , and \oplus (the join operation) and \otimes (the extend operation) are binary operators on D such that

1. (D, \oplus) is a commutative monoid with 0 as its neutral element, and where \oplus is idempotent (i.e., for all $a \in D$, $a \oplus a = a$).
2. (D, \otimes) is a monoid with the neutral element 1 .

3. \otimes distributes over \oplus , i.e., for all $a, b, c \in D$ we have

$$\begin{aligned} a \otimes (b \oplus c) &= (a \otimes b) \oplus (a \otimes c) \text{ and} \\ (a \oplus b) \otimes c &= (a \otimes c) \oplus (b \otimes c). \end{aligned}$$

4. 0 is an annihilator with respect to \otimes , i.e., for all $a \in D$, $a \otimes 0 = 0 = 0 \otimes a$.

5. In the partial order \sqsubseteq defined by: $\forall a, b \in D, a \sqsubseteq b$ iff $a \oplus b = b$, there are no infinite ascending chains.

Defn. 2.7(1) and Defn. 2.7(5) mean that (D, \oplus) is a **join semilattice with no infinite ascending chains**.

Note that the distributivity property (Property 3 in Defn. 2.7) holds only if the extend operation distributes over the join operation. The analysis is still sound if the semiring is not distributive. If the semiring is distributive, then the WPDS-based dataflow analysis gives the join-over-all-paths solution.

2.5.2 Application of WPDS to Interprocedural Dataflow Analysis

For our interprocedural-analysis experiments, we use the WALi system [50] for WPDSs to perform abstract interpretation (§2.2). To construct the WPDS from a given control flow graph (CFG), each edge in the CFG (which represents a concrete transformer) is converted to a rule in WPDS. The semiring weight on an edge in the CFG is calculated by abstracting the concrete transformer for the edge to an abstract transformer. Hence, the WPDS weights are abstract transformers, represented as semiring weights. The asymptotic cost of weight generation is linear in the size of the program: to generate the weights, each basic block in the program is visited once, and a weight is generated by the relevant method.

Table 2.4: Semiring operators in terms of abstract-domain operations.

Semiring Operator	Abstract-domain Operation	Description
0	\perp	<i>bottom element</i>
1	Id	<i>identity element</i>
$(\alpha_1 \oplus \alpha_2)$	$(\alpha_1 \sqcup \alpha_2)$	<i>join</i>
$(\alpha_1 \oplus \alpha_2)$	$(\alpha_1 \nabla \alpha_2)$	<i>widen at loop headers</i>
$(\alpha_1 \otimes \alpha_2)$	$(\alpha_2 \circ \alpha_1)$	<i>composition</i>

Tab. 2.4 provides the relationship between the semiring operators and the corresponding abstract-domain operations. 0 is the bottom element. 1 is the identity operation. The join operation $(\alpha_1 \oplus \alpha_2)$ is implemented as the join of the abstract transformers, except at the loop headers where it is implemented as widening operation to ensure termination. The extend operation $(\alpha_1 \otimes \alpha_2)$ is implemented as abstract-composition operation $(\alpha_2 \circ \alpha_1)$. Notice that the order of the arguments in abstract composition has changed.

Semiring EWPDS merge functions [56] are used to preserve caller-save and callee-save registers across call sites. Running a WPDS-based analysis to find the least fixed-point value for a given set of program points involves calling two operations, “post*” and “path summary”, as detailed in [85]. The post* queries use the FWPDS algorithm [55].

3 ABSTRACT DOMAINS OF AFFINE RELATIONS

This chapter considers some known abstract domains for affine-relation analysis, along with several variants, and studies how they relate to each other. An affine relation is a *linear-equality* constraint over a given set of variables that hold machine integers. An abstract-domain element represents the set of states that satisfies a conjunction of affine relations. ARA finds, for each point in the program, a domain element that over-approximates the set of states that can arise at that point. ARA generalizes such analyses as linear-constant propagation [87] and induction-variable analysis.

We show that the abstract domains for ARA of Müller-Olm/Seidl (MOS) and King/Søndergaard (KS) are, in general, incomparable. However, we give sound interconversion methods. That is, we give an algorithm to convert a KS element v_{KS} to an over-approximating MOS element v_{MOS} —i.e., $\gamma(v_{KS}) \subseteq \gamma(v_{MOS})$ —as well as an algorithm to convert an MOS element w_{MOS} to an over-approximating KS element w_{KS} —i.e., $\gamma(w_{MOS}) \subseteq \gamma(w_{KS})$.

We describe our own variant of the KS domain, which is inspired by—but different from, and arguably easier to use than—the version of KS developed by King and Søndergaard. Our version is presented in §3.4.

For MOS, it was not previously known how to perform $\tilde{\alpha}_{MOS}(\varphi)$ in a non-trivial fashion (i.e., other than defining $\tilde{\alpha}_{MOS}$ to be $\lambda f. \top$). In contrast, [53, Fig. 2] gave an algorithm for $\hat{\alpha}_{KS}$, which led us to examine more closely how MOS and KS are related.

Contributions This chapters contributions include the following:

- We introduce a third domain for representing affine relations, called AG, which stands for *affine generators* (§3.1.3). Whereas an element in the KS domain consists of a set of *constraints* on the values of variables, AG represents a collection of allowed values of variables

via a set of *generators*. We show that AG is the generator counterpart of KS: a KS element can be converted to an AG element, and vice versa, with no loss of precision (§3.2).

- We show that MOS and KS/AG are, in general, *incomparable* (§3.3.1). In particular, we show that KS and AG can express transformers with affine guards, which MOS cannot express.
- We give sound interconversion methods between MOS and KS/AG (§3.3.2–§3.3.4):
 - We show that an AG element v_{AG} can be converted to an over-approximating MOS element v_{MOS} —i.e., $\gamma(v_{AG}) \subseteq \gamma(v_{MOS})$.
 - We show that an MOS element w_{MOS} can be converted to an over-approximating AG element w_{AG} —i.e., $\gamma(w_{MOS}) \subseteq \gamma(w_{AG})$.
- Consequently, by means of the conversion path $\varphi \rightarrow \text{KS} \rightarrow \text{AG} \rightarrow \text{MOS}$, we obtain a method for performing $\tilde{\alpha}_{MOS}(\varphi)$ (§3.3.5).

Experiments §3.5 presents an experimental study with the Intel IA32 (x86) instruction set in which the $\hat{\alpha}_{KS}$ method and two greedy, operator-by-operator reinterpretation methods—KS-reinterpretation (§2.3.1) and MOS-reinterpretation [62, §4.1.2]—are compared in terms of their performance and precision. The precision comparison is done by comparing the affine invariants obtained at branch points, as well as the affine procedure summaries obtained for procedures. For KS-reinterpretation and MOS-reinterpretation, we also compare the abstract transformers generated for individual x86 instructions. The experiments were designed to answer the following questions:

- Which method of obtaining abstract transformers is fastest: $\hat{\alpha}_{KS}$, KS-reinterpretation, or MOS-reinterpretation?
- Does MOS-reinterpretation or KS-reinterpretation yield more precise abstract transformers for machine instructions?
- For what percentage of branch points and procedures does $\hat{\alpha}_{KS}$ pro-

duce more precise answers than KS-reinterpretation?

Organization The chapter is organized as follows: §3.1 summarizes relevant features of the various ARA domains. §3.2 presents the AG domain, and shows how an AG element can be converted to a KS element, and vice versa. §3.3 presents our results on the incomparability of the MOS and KS domains, but gives sound methods to convert a KS element into an over-approximating MOS element, and vice versa. §3.4 explains how to use the KS domain for interprocedural analysis. §3.5 presents experimental results. §3.6 discusses related work. Proofs can be found in the appendices.

3.1 Abstract Domains for affine-relation analysis

The two existing ARA domains—one defined by [76, 78] (MOS) and one defined by [52, 53] (KS) are introduced in this section. Both MOS and KS are based on an extension of linear algebra to modules over a ring [44, 43, 1, 99, 76, 78]. §3.1.1 introduces terminology used in describing the ARA domains. §3.1.2 introduces Howell form, which can be used as a normal form in the MOS and KS domains. §3.1.3 defines the affine-generator (AG) domain, which is a generator representation of the KS domain. §3.1.4 and §3.1.5 define the KS domain and MOS domain, respectively. Finally, we discuss the domain heights of AG, KS, and MOS domains in §3.1.6.

3.1.1 Terminology

All numeric values are integers in \mathbb{Z}_{2^w} for some bit width w . That is, values are w -bit machine integers with the standard operations for machine addition and multiplication. Addition and multiplication in \mathbb{Z}_{2^w} form a

ring, not a field, so some facets of standard linear algebra do not apply, and thus we must be cautious about carrying over intuition from standard linear algebra. In particular, each odd element in \mathbb{Z}_{2^w} has a multiplicative inverse (which may be found in time $O(\log w)$ [102, Fig. 10-5]), but no even element has a multiplicative inverse. The *rank* of a value $x \in \mathbb{Z}_{2^w}$ is the maximum integer $p \leq w$ such that 2^p divides evenly into x [76, 78]. For example, $\text{rank}(1) = 0$, $\text{rank}(24) = 3$, and $\text{rank}(0) = w$.

Throughout the paper, k is the size of the *vocabulary*, the variable-set under analysis. A *two-vocabulary* relation is a transition relation between values of variables in the *pre-state* vocabulary and values of variables in the *post-state* vocabulary.

Matrix addition and multiplication are defined as usual, forming a matrix ring. We denote the transpose of a matrix M by M^t . A *one-vocabulary matrix* is a matrix with $k + 1$ columns. A *two-vocabulary matrix* is a matrix with $2k + 1$ columns. In each case, the “+1” is for technical reasons (which vary according to what kind of matrix we are dealing with). I denotes the (square) identity matrix (whose size can be inferred from context). The rows of a matrix M are numbered from 1 to $\text{rows}(M)$; the columns of M are numbered starting from 1.

Because the MOS domain inherently involves pre-state-vocabulary to post-state-vocabulary transformers (see §3.1.5), our definitions of the AG and KS domains (§3.1.3 and §3.1.4, respectively) are also two-vocabulary domains. Technically, AG and KS can have an arbitrary number of vocabularies, including just a single vocabulary. To be able to give simpler examples, some of the AG and KS examples use one-vocabulary domain elements.

States in the various abstract domains are represented by row vectors of length $k + 1$. The *row space* of a matrix M is defined by $\text{row } M \stackrel{\text{def}}{=} \{\vec{v} \mid \exists w: wM = \vec{v}\}$. When we speak of the “null space” of a matrix, we actually mean the set of row vectors whose transposes are in the traditional

null space of the matrix. Thus, we define $\text{null}^t M \stackrel{\text{def}}{=} \{ \vec{v} \mid M\vec{v}^t = 0 \}$.

3.1.2 Matrices in Howell Form

One way to appreciate how linear algebra in rings differs from linear algebra in fields is to see how certain issues are finessed when converting a matrix to *Howell form* [44]. The Howell form of a matrix is an extension of reduced row-echelon form [67] suitable for matrices over \mathbb{Z}_n . Because Howell form is *canonical* for matrices over principal ideal rings [44, 99], it provides a way to test whether two abstract-domain elements are equal—i.e., whether they represent the same set of concrete values. Such an equality test is needed during program analysis to determine whether a fixed point has been reached.

Definition 3.1. *The leftmost nonzero value in a row vector is its **leading value**. The leading value's index is the **leading index**. A matrix M is in **row-echelon form** if*

- All all-zero rows are at the bottom.
- Each row's leading index is greater than that of the row above it.

If M is in row-echelon form, let $[M]_i$ denote the matrix that consists of all rows of M whose leading index is i or greater.

A matrix M is in **Howell form** if

1. M is in row-echelon form and has no all-zero rows,
2. the leading value of every row is a power of two,
3. each leading value is the largest value in its column, and
4. for every row r of M , for any $p \in \mathbb{Z}$, if i is the leading index of $2^p r$, then $2^p r \in \text{row}([M]_i)$.

□

In Defn. 3.1, item (4) may be confusing, and thus warrants an example.

Example 3.2. Suppose that $w = 4$, so that we are working in \mathbb{Z}_{16} . Consider the following two matrices and their Howellizations:

$$M_1 \stackrel{\text{def}}{=} \begin{bmatrix} 4 & 2 & 4 \\ 0 & 8 & 0 \end{bmatrix} \quad \text{HOWELLIZE}(M_1) = \begin{bmatrix} 4 & 2 & 4 \\ 0 & 8 & 0 \end{bmatrix}$$

$$M_2 \stackrel{\text{def}}{=} \begin{bmatrix} 4 & 2 & 4 \\ 0 & 4 & 0 \end{bmatrix} \quad \text{HOWELLIZE}(M_2) = \begin{bmatrix} 4 & 2 & 4 \\ 0 & 4 & 0 \end{bmatrix}$$

First, notice that M_1 does not satisfy item (4). M_1 has only one row, $[4 \ 2 \ 4]$, and consider what happens when this row is multiplied by powers of 2:

$$2^1 \cdot [4 \ 2 \ 4] = [8 \ 4 \ 8]$$

$$2^2 \cdot [4 \ 2 \ 4] = [0 \ 8 \ 0]$$

$$2^3 \cdot [4 \ 2 \ 4] = [0 \ 0 \ 0]$$

In particular, the leading index of $2^2 \cdot [4 \ 2 \ 4] = [0 \ 8 \ 0]$ is 2; however, because $\text{row}([M]_2) = \emptyset$, $[0 \ 8 \ 0] \notin \text{row}([M]_2)$. Consequently, $[0 \ 8 \ 0]$ must be included in $\text{HOWELLIZE}(M_1)$. We say that a row like $[0 \ 8 \ 0]$ is a logical consequence of $[4 \ 2 \ 4]$ that is added to satisfy item (4) of Defn. 3.1.

In contrast, matrix M_2 satisfies item (4) (and, in fact, M_2 is already in Howell form). For matrix M_2 to fail to satisfy item (4), there would have to be some row r and power p for which (a) the leading index i of $2^p r$ is strictly greater than the leading index of r , (b) $2^p r \neq 0$, and (c) $2^p r \notin \text{row}([M]_i)$. In this example, the only interesting quantity of the form $2^p r$ is $2^2 \cdot [4 \ 2 \ 4] = [0 \ 8 \ 0]$. The leading index of $[0 \ 8 \ 0]$ is 2, but $[0 \ 8 \ 0] = 2 \cdot [0 \ 4 \ 0]$, and so $[0 \ 8 \ 0] \in \text{row}([M]_2)$. Consequently, M_2 satisfies item (4). \square

The Howell form of a matrix is unique among all matrices with the same row space (or null space) [44]. As mentioned earlier, this property of Howell form provides a way to test two MOS elements, two KS elements, or two AG elements for equality.

Algorithm 3 HOWELLIZE: Put the matrix G in Howell form.

```

1: procedure HOWELLIZE( $G$ )
2:   Let  $j = 0$   $\triangleright$   $j$  is the number of already-Howellized rows
3:   for all  $i$  from 1 to  $2k + 1$  do
4:     Let  $R = \{\text{all rows of } G \text{ with leading index } i\}$ 
5:     if  $R \neq \emptyset$  then
6:       Pick an  $r \in R$  that minimizes  $\text{rank } r_i$ 
7:       Pick the odd  $u$  and rank  $p$  so that  $u2^p = r_i$ 
8:        $r \leftarrow u^{-1}r$   $\triangleright$  Adjust  $r$ , leaving  $r_i = 2^p$ 
9:       for all  $s$  in  $R \setminus \{r\}$  do
10:        Pick the odd  $v$  and rank  $t$  so that  $v2^t = s_i$ 
11:         $s \leftarrow s - (v2^{t-p})r$   $\triangleright$  Zero out  $s_i$ 
12:        if row  $s$  contains only zeros then
13:          Remove  $s$  from  $G$ 
14:        In  $G$ , swap  $r$  with  $G_{j+1}$   $\triangleright$  Place  $r$  for row-echelon form
15:        for all  $h$  from 1 to  $j$  do  $\triangleright$  Set values above  $r_i$  to be  $0 \leq \cdot < r_i$ 
16:           $d \leftarrow G_{h,i} \gg p$   $\triangleright$  Pick  $d$  so that  $0 \leq G_{h,i} - dr_i < r_i$ 
17:           $G_h \leftarrow G_h - dr$   $\triangleright$  Adjust row  $G_h$ , leaving  $0 \leq G_{h,i} < r_i$ 
18:        if  $r_i \neq 1$  then  $\triangleright$  Add logical consequences of  $r$  to  $G$ 
19:          Add  $2^{w-p}r$  as last row of  $G$   $\triangleright$  New row has leading
    index  $> i$ 
20:      $j \leftarrow j + 1$ 

```

The notion of a *saturated set of generators* used by [78] is closely related to Howell form, but is defined for an unordered set of matrices rather than row-vectors arranged in a matrix, and has no analogue of item (3). The algorithms of Müller-Olm and Seidl do not compute multiplicative inverses (see §3.6.2), so a saturated set has no analogue of item (2). Consequently, a saturated set is not canonical among generators of the same space.

Our technique for putting a matrix in Howell form is the procedure HOWELLIZE (Alg. 3). Much of HOWELLIZE is similar to a standard Gaussian-elimination algorithm, and it has the same overall cubic-time complexity as Gaussian elimination. In particular, HOWELLIZE minus lines 15–19 puts G in row-echelon form (item (1) of Defn. 3.1) with the leading value of every row

a power of two. (Line 8 enforces item (2) of Defn. 3.1.) `HOWELLIZE` differs from standard Gaussian elimination in how the pivot is picked (line 6) and in how the pivot is used to zero out other elements in its column (lines 7–13). Lines 15–17 of `HOWELLIZE` enforce item (3) of Defn. 3.1, and lines 18–19 enforce item (4). Lines 12–13 remove all-zero rows, which is needed for Howell form to be canonical.

Alg. 3 is simple and easy to implement. For analyses over large vocabularies, one should replace Alg. 3, which has cubic-time complexity with, say, the algorithm of [99], which has the same asymptotic complexity as matrix multiplication.

3.1.3 The Affine Generator Domain

An element in the Affine Generator domain (AG) is a two-vocabulary matrix whose rows are the affine generators of a two-vocabulary relation.

An AG element is an r -by- $(2k + 1)$ matrix G , with $0 < r \leq 2k + 1$. The concretization of an AG element is

$$\gamma_{AG}(G) \stackrel{\text{def}}{=} \{(\vec{v}, \vec{v}') \mid \vec{v}, \vec{v}' \in \mathbb{Z}_{2^w}^k \wedge [1 \mid \vec{v} \ \vec{v}'] \in \text{row } G\}.$$

The AG domain captures all two-vocabulary affine spaces, and treats them as relations between pre-states and post-states.

The bottom element of the AG domain is the empty matrix, and the AG

element that represents the identity relation is the matrix
$$\left[\begin{array}{c|cc} & \vec{v} & \vec{v}' \\ \hline 1 & 0 & 0 \\ 0 & I & I \end{array} \right].$$

The AG element

$$\left[\begin{array}{c|ccccc} & v_1 & v_2 & v'_1 & v'_2 \\ \hline 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{array} \right] \quad (3.1)$$

represents the transition relation in which $v'_1 = v_1$, v_2 can have any value, and v'_2 can have any even value.

To compute the join of two AG elements, stack the two matrices vertically and Howellize the result.

3.1.4 The King/Søndergaard Domain

An element in the King/Søndergaard domain (KS) is a two-vocabulary matrix whose rows represent constraints on a two-vocabulary relation. A KS element is an r -by- $(2k + 1)$ matrix M , with $0 \leq r \leq 2k + 1$. The concretization of a KS element M is

$$\gamma_{\text{KS}}(M) \stackrel{\text{def}}{=} \{(\vec{v}, \vec{v}') \mid \vec{v}, \vec{v}' \in \mathbb{Z}_{2^w}^k \wedge [\vec{v} \ \vec{v}' \mid 1] \in \text{null}^t G\}.$$

Like the AG domain, the KS domain captures all two-vocabulary affine spaces, and treats them as relations between pre-states and post-states.

It is easy to read out affine equalities from a KS element M (regardless of whether M is in Howell form): if

$$\left[\begin{array}{cccccc|c} v_1 & \dots & v_k & v'_1 & \dots & v'_k & 1 \\ a_1 & \dots & a_k & a'_1 & \dots & a'_k & b \end{array} \right]$$

is a row of M , then $\sum_i a_i v_i + \sum_i a'_i v'_i = -b$ is a constraint on $\gamma_{\text{KS}}(M)$. The conjunction of these constraints describes $\gamma_{\text{KS}}(M)$ exactly.

The bottom element of the KS domain is the matrix $\left[\begin{array}{cc|c} \vec{v} & \vec{v}' & 1 \\ 0 & 0 & 1 \end{array} \right]$, and the KS element that represents the identity relation is the matrix

$$\left[\begin{array}{cc|c} \vec{v} & \vec{v}' & 1 \\ I & -I & 0 \end{array} \right].$$

Suppose that $w = 4$, so that we are working in \mathbb{Z}_{16} . The

KS element

$$\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ \hline 1 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 8 & 0 \end{array} \quad (3.2)$$

represents the transition relation in which $v'_1 = v_1$, v_2 can have any value, and v'_2 can have any even value. Thus, Eqns. (3.1) and (3.2) represent the same transition relation in AG and KS, respectively.

A Howell-form KS element can easily be checked for emptiness: it is empty if and only if it contains a row whose leading entry is in its last column. In that sense, an implementation of the KS domain in which all elements are kept in Howell form has redundant representations of bottom (whose concretization is \emptyset). However, such KS elements can always be detected during HOWELLIZE and replaced by the canonical representation

$$\begin{array}{cc|c} \vec{v} & \vec{v}' & 1 \\ \hline 0 & 0 & 1 \end{array}.$$

of bottom, namely,

The original King and Søndergaard paper 2008 gives polynomial-time algorithms for join and projection; projection can be used to implement composition (see §3.4.3).

3.1.5 The Müller-Olm/Seidl Domain

An element in the Müller-Olm/Seidl domain (MOS) is an affine-closed set of affine transformers, as detailed in [78]. An MOS element is represented by a set of $(k + 1)$ -by- $(k + 1)$ matrices. Each matrix T is a one-vocabulary transformer of the form $T = \left[\begin{array}{c|c} 1 & b \\ \hline 0 & M \end{array} \right]$, which represents the state transformation $\vec{v}' := \vec{v} \cdot M + b$, or, equivalently, $[1|\vec{v}'] := [1|\vec{v}] T$.

An MOS element B consists of a set of $(k + 1)$ -by- $(k + 1)$ matrices, and represents the affine span of the set, denoted by $\langle B \rangle$ and defined as follows:

$$\langle B \rangle \stackrel{\text{def}}{=} \left\{ T \mid \exists w \in \mathbb{Z}_{2^w}^{|\mathcal{B}|} : T = \sum_{B \in \mathcal{B}} w_B B \wedge T_{1,1} = 1 \right\}.$$

The meaning of B is the union of the graphs of the affine transformers in $\langle B \rangle$

$$\Upsilon_{\text{MOS}}(B) \stackrel{\text{def}}{=} \{(\vec{v}, \vec{v}') \mid \vec{v}, \vec{v}' \in \mathbb{Z}_{2^w}^k \wedge \exists T \in \langle B \rangle: [1 \mid \vec{v}] T = [1 \mid \vec{v}']\}.$$

The bottom element of the MOS domain is \emptyset , and the MOS element that represents the identity relation is the singleton set $\{I\}$. If $w = 4$, the MOS element $B = \left\{ \left[\begin{array}{c|cc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right], \left[\begin{array}{c|cc} 0 & 0 & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{array} \right] \right\}$ represents the affine span

$$\langle B \rangle = \left\{ \left[\begin{array}{c|cc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right], \left[\begin{array}{c|cc} 1 & 0 & 2 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right], \left[\begin{array}{c|cc} 1 & 0 & 4 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right], \dots, \left[\begin{array}{c|cc} 1 & 0 & 14 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right] \right\},$$

which corresponds to the transition relation in which $v'_1 = v_1$, v_2 can have any value, and v'_2 can have any even value—i.e., B represents the same transition relation as Eqns. (3.1) and (3.2).

The operations join and compose can be performed in polynomial time. If B and C are MOS elements, $B \sqcup C = \text{HOWELLIZE}(B \cup C)$ and $B;C = \text{HOWELLIZE}\{bc \mid b \in B \wedge c \in C\}$. In this setting, HOWELLIZE of a set of $(k+1)$ -by- $(k+1)$ matrices $\{M_1, \dots, M_n\}$ means “Apply Alg. 3 to a larger, n -by- $(k+1)^2$ matrix, each of whose rows is the linearization (e.g., in row-major order) of one of the M_i .”

3.1.6 Domain Heights

In all three domains, an element can be represented via an appropriate matrix in Howell form (where in the case of the MOS domain, we mean a matrix in the extended sense discussed in §3.1.5). For a fixed bit width and a fixed number of columns, there are only a constant number of Howell-form matrices. Consequently, the KS, AG, and MOS domains are all finite domains, and hence of finite height.

Domain elements need not necessarily be maintained in Howell form; instead, they could be Howellized on demand when it is necessary to check containment (see §3.4.6). Our implementation maintains domain elements in Howell form using essentially the “list of lists” sparse-matrix representation: each matrix is represented via a C++ vector of rows; each row is a vector of (column-index, nonzero-value) pairs.

3.2 Relating AG and KS Elements

AG and KS are equivalent domains. One can convert an AG element to an equivalent KS element with no loss of precision, and vice versa. In essence, these are a single abstract domain with two representations: constraint form (KS) and generator form (AG).

We use an operation similar to singular value decomposition, called diagonal decomposition:

Definition 3.3. *The **diagonal decomposition** of a square matrix M is a triple of matrices, L , D , R , such that $M = LDR$; L and R are invertible matrices; and D is a diagonal matrix in which all entries are either 0 or a power of 2. \square*

[78, Lemma 2.9] give a decomposition algorithm that nearly performs diagonal decomposition, except that the entries in their D might not be powers of 2. We can easily adapt that algorithm. Suppose that their method yields LDR (where L and R are invertible). Pick u and r so that $u_i 2^{r_i} = D_{i,i}$ with each u_i odd, and define U as the diagonal matrix where $U_{i,i} = u_i$. (If $D_{i,i} = 0$, then $u_i = 1$.) It is easy to show that U is invertible. Let $L' = LU$ and $D' = U^{-1}D$. Consequently, $L'D'R = LDR = M$, and $L'D'R$ is a diagonal decomposition.

From diagonal decomposition we derive the dualization operation, denoted by \cdot^\perp , such that the rows of M^\perp generate the null space of M , and vice versa.

Definition 3.4. The **dualization** of M , denoted by M^\perp , is defined as follows:

- $P_{AD}(M)$ is the $(2k+1)$ -by- $(2k+1)$ matrix $\begin{bmatrix} M \\ \vec{0} \end{bmatrix}$,
- L, D, R is the diagonal decomposition of $P_{AD}(M)$,
- T is the diagonal matrix with $T_{i,i} \stackrel{\text{def}}{=} 2^{w-\text{rank}(D_{i,i})}$, and
- $M^\perp \stackrel{\text{def}}{=} (L^{-1})^t T (R^{-1})^t$

□

This definition of dualization has the following useful property:

Theorem 3.5. For any matrix M , $\text{null}^t M = \text{row } M^\perp$ and $\text{row } M = \text{null}^t M^\perp$.

Proof. For any matrix M , it is a common lemma that $(M^{-1})^t = (M^t)^{-1}$. Thus, the notation M^{-t} denotes $(M^{-1})^t$.

Lemma 3.6. Let D and T be square, diagonal matrices, where $D_{ii} = 2^{p_i}$ and $T_{ii} = 2^{w-p_i}$ for all i . Then, $\text{null}^t T = \text{row } D$ and $\text{null}^t D = \text{row } T$.

Proof. Let z be any row vector. To see that $\text{null}^t T = \text{row } D$:

$$\begin{aligned} z \in \text{null}^t T &\Leftrightarrow Tz^t = 0 \Leftrightarrow \forall i: z_i 2^{w-p_i} = 0 \\ &\Leftrightarrow \forall i: 2^{p_i} | z_i \Leftrightarrow \exists v: \forall i: v_i 2^{p_i} = z_i \\ &\Leftrightarrow \exists v: vD = z \Leftrightarrow z \in \text{row } D. \end{aligned}$$

One can show that $\text{null}^t D = \text{row } T$ by essentially the same reasoning. □

Again, let L, D , and R be the diagonal decomposition of M (see Defn. 3.3, and construct T from D as in Lem. 3.6. Recall that L is invertible. To see that $\text{row } M = \text{null}^t M^\perp$,

$$\begin{aligned} \text{row } M = \text{row } LDR = \text{row } DR, \text{ so } x \in \text{row } DR &\Leftrightarrow xR^{-1} \in \text{row } D \\ &\Leftrightarrow xR^{-1} \in \text{null}^t T \Leftrightarrow TR^{-t}x^t = 0 \Leftrightarrow x \in \text{null}^t TR^{-t}. \end{aligned}$$

We know that L^{-t} is also invertible, so

$$\text{null}^t TR^{-t} = \text{null}^t L^{-t}TR^{-t} = \text{null}^t M^\perp.$$

Thus, $\text{row } M = \text{null}^t M^\perp$. One can show that $\text{null}^t M = \text{row } M^\perp$ by essentially the same reasoning. \square

We can therefore use dualization to convert between equivalent KS and AG elements. For a given (padded, square) AG matrix $G = [c|Y Y']$, we seek a KS matrix Z of the form $[X X'|b]$ such that $\gamma_{\text{KS}}(Z) = \gamma_{\text{AG}}(G)$. We construct Z by letting $[b|X X'] = G^\perp$ and permuting those columns to $Z \stackrel{\text{def}}{=} [X X'|b]$. This works by Thm. 3.5, and because

$$\begin{aligned} \gamma_{\text{AG}}(G) &= \{(\mathbf{x}, \mathbf{x}') \mid [1|\mathbf{x} \ \mathbf{x}'] \in \text{row } G\} \\ &= \{(\mathbf{x}, \mathbf{x}') \mid [1|\mathbf{x} \ \mathbf{x}'] \in \text{null}^t G^\perp\} \\ &= \{(\mathbf{x}, \mathbf{x}') \mid [\mathbf{x} \ \mathbf{x}'|1] \in \text{null}^t Z\} = \gamma_{\text{KS}}(Z). \end{aligned}$$

Furthermore, to convert from any KS matrix to an equivalent AG matrix, we reverse the process. Reversal is possible because dualization is an involution: for any matrix M , $(M^\perp)^\perp = M$.

3.3 Relating KS and MOS

3.3.1 MOS and KS are Incomparable

The MOS and KS domains are incomparable: some relations are expressible in each domain that are not expressible in the other. Intuitively, the central difference is that MOS is a domain of sets of *functions*, while KS is a domain of *relations*.

KS can capture restrictions on both the pre-state and post-state vocabularies, while MOS can capture restrictions only on its post-state vocab-

ulary. For example, when $k = 1$, the KS element for “assume $v = 2$ ” (in

un-Howellized form) is $\left[\begin{array}{cc|c} v & v' & 1 \\ 1 & 0 & -2 \\ & 1 & 0 \end{array} \right]$, i.e., “ $v = 2 \wedge v = v'$ ”.

In contrast, an MOS element cannot encode an assume statement. The smallest MOS element that over-approximates “assume $v = 2$ ” is the identity transformer $\left\{ \left[\begin{array}{c|c} 1 & 0 \\ 0 & 1 \end{array} \right] \right\}$. In general, an MOS element cannot encode a non-trivial condition on the pre-state. If an MOS element contains a single transition, it encodes that transition for every possible pre-state. Therefore, KS can encode relations that MOS cannot encode.

On the other hand, an MOS element can encode two-vocabulary relations that are not affine. One example is the matrix basis $\mathcal{B} = \left\{ \left[\begin{array}{c|cc} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{array} \right], \left[\begin{array}{c|cc} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 1 \end{array} \right] \right\}$. The set that \mathcal{B} encodes is

$$\begin{aligned} \gamma_{\text{MOS}}(\mathcal{B}) &= \left\{ \left[\begin{array}{cccc} v_1 & v_2 & v'_1 & v'_2 \end{array} \right] \left| \begin{array}{l} \exists w_0, w_1: \left[\begin{array}{c|cc} 1 & v_1 & v_2 \\ 0 & w_0 & w_0 \\ 0 & w_1 & w_1 \end{array} \right] = \left[\begin{array}{c|cc} 1 & 0 & 0 \\ 0 & w_0 & w_0 \\ 0 & w_1 & w_1 \end{array} \right] \\ \wedge w_0 + w_1 = 1 \end{array} \right. \right\} \\ &= \left\{ \left[\begin{array}{cccc} v_1 & v_2 & v'_1 & v'_2 \end{array} \right] \left| \exists w_0: v'_1 = v'_2 = w_0 v_1 + (1 - w_0) v_2 \right. \right\} \\ &= \left\{ \left[\begin{array}{cccc} v_1 & v_2 & v'_1 & v'_2 \end{array} \right] \left| \exists w_0: v'_1 = v'_2 = v_1 + (1 - w_0)(v_2 - v_1) \right. \right\} \\ &= \left\{ \left[\begin{array}{cccc} v_1 & v_2 & v'_1 & v'_2 \end{array} \right] \left| \exists p: v'_1 = v'_2 = v_1 + p(v_2 - v_1) \right. \right\} \end{aligned} \quad (3.3)$$

Affine spaces are closed under affine combinations of their elements. Thus, $\gamma_{\text{MOS}}(\mathcal{B})$ is not an affine space because some affine combinations of its elements are not in $\gamma_{\text{MOS}}(\mathcal{B})$. For instance, let $a = \left[\begin{array}{cccc} 1 & -1 & 1 & 1 \end{array} \right]$, $b = \left[\begin{array}{cccc} 2 & -2 & 6 & 6 \end{array} \right]$, and $c = \left[\begin{array}{cccc} 0 & 0 & -4 & -4 \end{array} \right]$. By Eqn. (3.3), we have $a \in \gamma_{\text{MOS}}(\mathcal{B})$ when $p = 0$ in Eqn. (3.3), $b \in \gamma_{\text{MOS}}(\mathcal{B})$ when $p = -1$, and $c \notin \gamma_{\text{MOS}}(\mathcal{B})$ (the equation “ $-4 = 0 + p(0 - 0)$ ” has no solution for p). Moreover, $2a - b = c$, so c is an affine combination of a and b . Thus, $\gamma_{\text{MOS}}(\mathcal{B})$ is not closed under affine combinations of its elements, and so $\gamma_{\text{MOS}}(\mathcal{B})$ is not an affine space. Because every KS element encodes a two-

vocabulary affine space, MOS can represent $\gamma_{\text{MOS}}(\mathcal{B})$ but KS cannot.

3.3.2 Converting MOS Elements to KS

Soundly converting an MOS element \mathcal{B} to a KS element is equivalent to stating two-vocabulary affine constraints satisfied by \mathcal{B} .

To convert an MOS element \mathcal{B} to a KS element, we

1. rewrite \mathcal{B} so that every matrix it contains has a 1 in its top-left corner,
2. build a two-vocabulary AG matrix from each one-vocabulary matrix in \mathcal{B} ,
3. join the resulting AG matrices, and
4. convert the joined AG matrix to a KS element.

For Step (1), we rewrite \mathcal{B} so that

$$\mathcal{B} = \left\{ \left[\begin{array}{c|c} 1 & \mathbf{c}_i \\ \hline 0 & \mathbf{N}_i \end{array} \right] \right\}, \text{ where } \mathbf{c}_i \in \mathbb{Z}_{2^w}^{1 \times k} \text{ and } \mathbf{N}_i \in \mathbb{Z}_{2^w}^{k \times k}.$$

If our original MOS element \mathcal{B}_0 fails to satisfy this property, we can construct an equivalent \mathcal{B} that does. Let $\mathcal{C} = \text{HOWELLIZE}(\mathcal{B}_0)$; pick the unique $\mathbf{B} \in \mathcal{C}$ such that $\mathbf{B}_{1,1} = 1$, and let $\mathcal{B} = \{\mathbf{B}\} \cup \{\mathbf{B} + \mathbf{C} \mid \mathbf{C} \in (\mathcal{C} \setminus \{\mathbf{B}\})\}$. \mathcal{B} now satisfies the property, and $\langle \mathcal{B} \rangle = \langle \mathcal{B}_0 \rangle$.

In Step (2), we construct the matrices

$$\mathbf{G}_i = \left[\begin{array}{c|cc} 1 & 0 & \mathbf{c}_i \\ \hline 0 & \mathbf{I} & \mathbf{N}_i \end{array} \right].$$

Note that, for each matrix $\mathbf{B}_i \in \mathcal{B}$, $\gamma_{\text{MOS}}(\{\mathbf{B}_i\}) = \gamma_{\text{AG}}(\mathbf{G}_i)$. In Step (3), we join the \mathbf{G}_i matrices in the AG domain to yield one matrix \mathbf{G} . Thm. 3.7 states the soundness of this transformation from MOS to AG, i.e., $\gamma_{\text{MOS}}(\mathcal{B}) \subseteq \gamma_{\text{AG}}(\mathbf{G})$. Finally, \mathbf{G} is converted in Step (4) to an equivalent KS element by the method given in §3.2.

Thm. 3.7 states that the transformation from MOS to AG is sound.

Theorem 3.7. *Suppose that \mathcal{B} is an MOS element such that, for every $B \in \mathcal{B}$, $B = \left[\begin{array}{c|c} 1 & c_B \\ \hline 0 & M_B \end{array} \right]$ for some $c_B \in \mathbb{Z}_{2^w}^{1 \times k}$ and $M_B \in \mathbb{Z}_{2^w}^{k \times k}$. Define $G_B = \left[\begin{array}{c|c} 1 & 0 & c_B \\ \hline 0 & I & M_B \end{array} \right]$ and $G = \bigsqcup_{AG} \{G_B \mid B \in \mathcal{B}\}$. Then, $\gamma_{MOS}(\mathcal{B}) \subseteq \gamma_{AG}(G)$.*

Proof. See App. A. □

Because we can easily read affine relations from KS elements (§3.1.4), this conversion method also gives an easy way to create a quantifier-free formula that over-approximates the meaning of an MOS element. In particular, the formula read out of the KS element obtained from MOS-to-KS conversion captures affine relations implied by the MOS element.

3.3.3 Converting KS Without Pre-State Guards to MOS

If a KS element is total with respect to pre-state inputs, then we can convert it to an equivalent MOS element. First, convert the KS element to an AG element G . When G expresses no restrictions on its pre-state, it has the form

$$G = \left[\begin{array}{c|cc} 1 & 0 & b \\ \hline 0 & I & M \\ 0 & 0 & R \end{array} \right], \quad (3.4)$$

where $b \in \mathbb{Z}_{2^w}^{1 \times k}$; $I, M \in \mathbb{Z}_{2^w}^{k \times k}$; and $R \in \mathbb{Z}_{2^w}^{k \times r}$ with $0 \leq r \leq k$.

Definition 3.8. *An AG matrix of the form*

$$\left[\begin{array}{c|cc} 1 & 0 & b \\ \hline 0 & I & M \end{array} \right],$$

*such as the G_i matrices discussed in §3.3.2, is said to be in **explicit form**. An AG matrix in this form represents the transition relation $x' = x \cdot M + b$. □*

Explicit form is desirable because we can immediately convert the AG matrix of Defn. 3.8 into the MOS element

$$\left\{ \left[\begin{array}{c|c} 1 & \mathbf{b} \\ \hline 0 & \mathbf{M} \end{array} \right] \right\}.$$

The matrix G in Eqn. (3.4) is not in explicit form because of the rows $[0|0 \ R]$; however, G is quite close to being in explicit form, and we can read off a *set* of matrices to create an appropriate MOS element. We produce this set of matrices via the SHATTER operation, where

$$\text{SHATTER}(G) \stackrel{\text{def}}{=} \left\{ \left[\begin{array}{c|c} 1 & \mathbf{b} \\ \hline 0 & \mathbf{M} \end{array} \right] \right\} \cup \left\{ \left[\begin{array}{c|c} 0 & \mathbf{R}_{j,*} \\ \hline 0 & 0 \end{array} \right] \mid 1 \leq j \leq r \right\},$$

where $\mathbf{R}_{j,*}$ is row j of R .

As shown in Thm. 3.9, $\gamma_{AG}(G) = \gamma_{MOS}(\text{SHATTER}(G))$. Intuitively, this property holds because the coefficients of the $[0|0 \ \mathbf{R}_{j,*}]$ rows in an affine combination of the rows of G correspond to coefficients of the $\left\{ \left[\begin{array}{c|c} 0 & \mathbf{R}_{j,*} \\ \hline 0 & 0 \end{array} \right] \right\}$ matrices in an affine combination of the matrices in $\text{SHATTER}(G)$.

Theorem 3.9. *When $G = \left[\begin{array}{c|c} 1 & 0 \ \mathbf{b} \\ \hline 0 & \mathbf{I} \ \mathbf{M} \\ 0 & 0 \ \mathbf{R} \end{array} \right]$, then $\gamma_{AG}(G) = \gamma_{MOS}(\text{SHATTER}(G))$.*

Proof. See App. A. □

3.3.4 Converting KS With Pre-State Guards to MOS

If a KS element is not total with respect to pre-state inputs, then there is no MOS element with the same concretization. However, we can find sound over-approximations within MOS for such KS elements.

We convert the KS element into an AG matrix G as in §3.3.3 and put G in Howell form. There are two ways that G can enforce guards on the

Algorithm 4 MAKEEXPLICIT: Transform an AG matrix G in Howell form to near-explicit form.

Require: G is an AG matrix in Howell form

```

1: procedure MAKEEXPLICIT( $G$ )
2:   for all  $i$  from 2 to  $k + 1$  do   ▷ Consider each col. of the pre-state
   VOC.
3:     if there is a row  $r$  of  $G$  with leading index  $i$  then
4:       if  $\text{rank } r_i > 0$  then
5:         for all  $j$  from 1 to  $2k + 1$  do ▷ Build  $s$  from  $r$ , with  $s_i = 1$ 
6:            $s_j \leftarrow r_j \gg \text{rank } r_i$ 
7:           Append  $s$  to  $G$ 
8:            $G \leftarrow \text{Howellize}(G)$ 
9:   for all  $i$  from 2 to  $k + 1$  do
10:    if there is no row  $r$  of  $G$  with leading index  $i$  then
11:      Insert, as the  $i^{\text{th}}$  row of  $G$ , a new row of all zeroes

```

pre-state vocabulary: it might contain one or more rows whose leading value is even, or it might skip some leading indexes in row-echelon form.

While we cannot put G in explicit form, we can run MAKEEXPLICIT to coarsen G so that it is close enough to the form that arose in §3.3.3. Adding extra rows to an AG element can only enlarge its concretization. Thus, to handle a leading value 2^p , $p > 0$ in the pre-state vocabulary, MAKEEXPLICIT introduces an extra, over-approximating row constructed by copying the row with leading value 2^p and right-shifting each value in the copied row by p bits (lines 4–8). After the loop on lines 2–8 finishes, every leading value in a row that generates pre-state-vocabulary values is 1. MAKEEXPLICIT then introduces all-zero rows so that each leading element from the pre-state vocabulary lies on the diagonal (lines 9–11).

Example 3.10. Suppose that $k = 3$, $w = 4$, and $G = \left[\begin{array}{c|cccccc} 1 & 0 & 2 & 0 & 0 & 0 & 0 \\ & 4 & 0 & 12 & 2 & 4 & 0 \\ & & & & & 4 & 0 & 8 \end{array} \right]$. After line 11 of MAKEEXPLICIT, all pre-state vocabulary leading values of G have been made ones, and the resulting G' has $\text{row } G' \supseteq \text{row } G$. In our case, $G' =$

$$\left[\begin{array}{c|cccccc} 1 & 0 & 2 & 0 & 0 & 0 & 0 \\ \hline & 1 & 0 & 3 & 0 & 1 & 0 \\ & & & 2 & 0 & 0 & \\ & & & & & & 8 \end{array} \right].$$
 To handle “skipped” indexes, lines 9–11 insert all-zero rows into G' so that each leading element from the pre-state vocabulary lies on

the diagonal. The resulting matrix is

$$\left[\begin{array}{c|cccccc} 1 & 0 & 2 & 0 & 0 & 0 & 0 \\ \hline & 1 & 0 & 3 & 0 & 1 & 0 \\ & & 0 & 0 & 0 & 0 & 0 \\ & & & 0 & 0 & 0 & 0 \\ & & & & 2 & 0 & 0 \\ & & & & & & 8 \end{array} \right]. \quad \square$$

Theorem 3.11. For $G \in AG$, $\gamma_{AG}(G) \subseteq \gamma_{MOS}(\text{SHATTER}(\text{MAKEEXPLICIT}(G)))$.

Proof. See App. A. □

Thus, we can use the KS-to-AG conversion method of §3.2, MAKEEXPLICIT, and SHATTER to obtain an over-approximation of a KS element in MOS.

Example 3.12. The final MOS element for Ex. 3.10 is

$$\left\{ \left[\begin{array}{c|cccc} 1 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right], \left[\begin{array}{c|cccc} 0 & 2 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right], \left[\begin{array}{c|cccc} 0 & 0 & 0 & 8 \\ \hline 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right] \right\}.$$

□

3.3.5 Symbolic Abstraction for the MOS Domain

As mentioned in the introduction, it was not previously known how to perform symbolic abstraction for MOS. Using $\hat{\alpha}_{KS}$ (see §2.3.2 and §3.4.8) in conjunction with the algorithms from §3.2 and §3.3.4, we can soundly define $\tilde{\alpha}_{MOS}(\varphi)$ as

let $G = \text{CONVERTKSTOAG}(\hat{\alpha}_{KS}(\varphi))$ **in** $\text{SHATTER}(\text{MAKEEXPLICIT}(G))$.

3.4 Using KS for Interprocedural Analysis

This section presents a two-vocabulary version of the KS abstract domain, focusing on the operations that are useful in a program analyzer. Unlike previous work by [52, 53], it is not necessary to perform bit-blasting to use the version of KS presented here. §3.4.1–§3.4.7 describe the suite of operations needed to use the KS domain in an interprocedural-analysis algorithm in the style of [92] or [54], or to use the KS domain as a weight domain in a weighted pushdown system (WPDS) [10, 85, 56, 50].

§3.4.8 discusses symbolic abstraction for the KS domain, which provides one way to create two-vocabulary KS elements that represent abstract transformers needed by a program analyzer. (§2.3.1 provides an alternative method for creating KS abstract transformers, based on the operator-reinterpretation paradigm.)

§3.4.9 shows how to compute the number of tuples that satisfy a KS element.

3.4.1 Meet

As discussed in §3.1.4, the meaning of a KS matrix X can be expressed as a formula by forming a conjunction that consists of one equality for each row of X . We can obtain a KS element that precisely represents the conjunction of any number of such formulas by stacking the rows that represent the equalities, and putting the resulting matrix in Howell form. Consequently, we can compute the meet $X \sqcap Y$ of any two KS elements X and Y by putting the block matrix $\begin{bmatrix} X \\ Y \end{bmatrix}$ into Howell form. The resulting matrix exactly represents the intersection of the meanings of X and Y :

$$\gamma_{\text{KS}}(X \sqcap Y) = \gamma_{\text{KS}}(X) \cap \gamma_{\text{KS}}(Y).$$

3.4.2 Project and Havoc

[52, §3] describe a way to project a KS element X onto a suffix x_i, \dots, x_k of its vocabulary: (i) put X in row-echelon form to create X' ; (ii) create X'' by removing from X' every row a in which any of a_1, \dots, a_{i-1} is nonzero (i.e., $X'' = [X']_i$); and (iii) remove columns $1, \dots, i-1$. (Note that the resulting matrix has only a portion of the original vocabulary; we have projected away $\{x_1, \dots, x_{i-1}\}$.) However, although their method works for Boolean-valued KS elements (i.e., KS elements over \mathbb{Z}_2^k), when the leading values of X are not all 1, as can occur in KS elements over \mathbb{Z}_2^k for $w > 1$, step (ii) is not guaranteed to produce the most-precise projection of X onto x_i, \dots, x_k , although the KS element obtained is always sound.

Example 3.13. Suppose that $X = \begin{bmatrix} v_1 & v_2 & 1 \\ 4 & 2 & 6 \end{bmatrix}$, with $w = 4$, and the goal is to project away the first column (for v_1). When the King/Søndergaard projection algorithm is applied to X , we obtain the empty matrix, which represents no constraints on v_2 —i.e., $v_2 \in \{0, 1, \dots, 15\}$. However, closer inspection reveals that v_2 cannot be even; if v_2 were even, then both of the terms $4v_1$ and $2v_2$ would both be divisible by 4, and hence both values would have at least two zeros as their least-significant bits. Such a pair of values could not sum to a value congruent to 6 because the binary representation of 6 ends with $\dots 10$. \square

Instead, we put X in Howell form before removing rows. By Thm. 3.14, step (ii) above returns the exact projection of the original KS element onto the smaller vocabulary.

Theorem 3.14. Suppose that M has c columns. If matrix M is in Howell form, $\vec{v} \in \text{null}^t M$ if and only if $\forall i: \forall y_1, \dots, y_{i-1}: \begin{bmatrix} y_1 & \dots & y_{i-1} & v_i & \dots & v_c \end{bmatrix} \in \text{null}^t([M]_i)$.

Proof. See App. B. \square

Example 3.15. The Howell form of X from Ex. 3.13 is $\left[\begin{array}{cc|c} & v_1 & v_2 & 1 \\ 4 & 2 & 6 \\ 0 & 8 & 8 \end{array} \right]$, and thus

we obtain the following answer for the projection of X onto v_2 : $\left[\begin{array}{c|c} & v_2 & 1 \\ 8 & 8 \end{array} \right]$, which represents $v_2 \in \{1, 3, \dots, 15\}$.

This example illustrates that while the answer produced in Ex. 3.13 by the King/Søndergaard projection algorithm is a sound over-approximation, it is not as precise as the most-precise answer that can be represented in the KS domain.

□

Given KS element M , it is also possible to project away a set of variables V that does not constitute a prefix of the vocabulary: create M' by permuting the columns of M so that the columns for the variables in V come first—the order chosen for the V columns themselves is unimportant—and then project away V from M' as described earlier.

The *havoc* operation removes all constraints on a set of variables V . To havoc V from KS element M , project away V and then (i) add back an all-0 column for each variable in V , and (ii) permute columns to restore the original variable order. Because of the all-0 columns, the resulting KS element has no constraints on the values of the variables in V .

Example 3.16. Suppose that we wish to havoc v_2 from the KS value

$\left[\begin{array}{cc|c} & v_1 & v_2 & 1 \\ 2 & 4 & 6 \end{array} \right]$. We permute columns and Howellize to create $\left[\begin{array}{cc|c} & v_2 & v_1 & 1 \\ 4 & 2 & 6 \\ 0 & 8 & 8 \end{array} \right]$,

project onto the vocabulary suffix v_1 , obtaining $\left[\begin{array}{c|c} & v_1 & 1 \\ 8 & 8 \end{array} \right]$, add back an all-0 col-

umn for v_2 , $\left[\begin{array}{cc|c} & v_2 & v_1 & 1 \\ 0 & 8 & 8 \end{array} \right]$, and permute columns back to the original order to

obtain $\left[\begin{array}{cc|c} & v_1 & v_2 & 1 \\ 8 & 0 & & 8 \end{array} \right]. \square$

3.4.3 Compose

[53, §5.2] present a technique to compose two-vocabulary affine relations. For completeness, that algorithm follows. Suppose that we have KS elements $Y = \left[\begin{array}{cc|c} Y_{\text{pre}} & Y_{\text{post}} & \mathbf{y} \end{array} \right]$ and $Z = \left[\begin{array}{cc|c} Z_{\text{pre}} & Z_{\text{post}} & \mathbf{z} \end{array} \right]$, where Y_{pre} , Y_{post} , Z_{pre} , and Z_{post} are k -column matrices, and \mathbf{y} and \mathbf{z} are column vectors. We want to compute the relational composition " $Y; Z$ "; i.e., find some X such that $(\vec{v}, \vec{v}'') \in \gamma_{\text{KS}}(X)$ if and only if $\exists \vec{v}': (\vec{v}, \vec{v}') \in \gamma_{\text{KS}}(Y) \wedge (\vec{v}', \vec{v}'') \in \gamma_{\text{KS}}(Z)$.

Because the KS domain has a projection operation, we can create $Y; Z$ by first constructing the three-vocabulary matrix W ,

$$W = \left[\begin{array}{ccc|c} Y_{\text{post}} & Y_{\text{pre}} & 0 & \mathbf{y} \\ Z_{\text{pre}} & 0 & Z_{\text{post}} & \mathbf{z} \end{array} \right],$$

and then projecting away the first vocabulary of W . Any element $(\vec{v}', \vec{v}, \vec{v}'') \in \gamma_{\text{KS}}(W)$ has $(\vec{v}, \vec{v}') \in \gamma_{\text{KS}}(Y)$ and $(\vec{v}', \vec{v}'') \in \gamma_{\text{KS}}(Z)$; consequently, the projection yields a matrix X such that $\gamma_{\text{KS}}(X) = \gamma_{\text{KS}}(Y); \gamma_{\text{KS}}(Z)$, as required.

Alternatively, we can think of abstract composition as happening in three steps: (i) adding 0-columns and reordering vocabularies in Y and Z ; (ii) computing the meet W of the resulting matrices; and (iii) projecting onto the initial and final vocabulary of W . Thus, because reordering, meet, and projection are all exact operations, abstract composition is also an exact operation:

$$\gamma_{\text{KS}}(Y; Z) = \gamma_{\text{KS}}(Y); \gamma_{\text{KS}}(Z).$$

Note that the steps of the abstract-composition algorithm mimic a standard

way to express the composition of concrete relations, i.e.,

$$Y;Z = \exists U: Y[\text{pre}, \text{post}] \wedge Z[\text{pre}, \text{post}] \wedge U = Y.\text{post} \wedge U = Z.\text{pre}.$$

3.4.4 Join

To join two KS elements Y and Z , we first construct the matrix $\begin{bmatrix} -Y & Y \\ Z & 0 \end{bmatrix}$ and then project onto the last $2k + 1$ columns.

[52, §3] give a method to compute the join of two KS elements by building a $(6k + 3)$ -column matrix and projecting onto its last $2k + 1$ variables. We improve their approach slightly, building a $(4k + 2)$ -column matrix and then projecting onto its last $2k + 1$ variables.

If Y and Z are considered as representing *linear* spaces, rather than affine spaces, this approach works because $\begin{bmatrix} -Y & Y \\ Z & 0 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = 0$ is true just if $(Y(v - u) = 0) \wedge (Zu = 0)$. Because $(v - u) \in \text{null } Y$, and $u \in \text{null } Z$, we know that v is the sum of values in $\text{null } Y$ and $\text{null } Z$, and so v is in their linear closure. Thm. 3.17 demonstrates the correctness of the same algorithm in affine spaces; that proof is driven by roughly the same intuition.

Theorem 3.17. *If Y and Z are both $N + 1$ -column KS matrices, and $\gamma_{\text{KS}}(Y)$ and $\gamma_{\text{KS}}(Z)$ are both non-empty sets, then $Y \sqcup Z$ is the projection of $\begin{bmatrix} -Y & Y \\ Z & 0 \end{bmatrix}$ onto its right-most $N + 1$ columns.*

Proof. See App. C. □

Join is not exact in the same sense that meet, project, and compose are above: affine spaces are not closed under union. However, this algorithm does return the least upper bound of Y and Z in the space of KS elements.

Neither meet nor compose distribute over join, as illustrated in the following examples:

Meet over join ¹

$$\begin{aligned}
“(v_1 = v'_1)” &= \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & -1 & 0 & 0 \end{array} \right] = \top \sqcap \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & -1 & 0 & 0 \end{array} \right] \\
&= \left(\left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & 0 & 0 & 0 \end{array} \right] \sqcup \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 0 & 1 & 0 & 0 \end{array} \right] \right) \sqcap \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & -1 & 0 & 0 \end{array} \right] \\
&\neq \left(\left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & 0 & 0 & 0 \end{array} \right] \sqcap \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & -1 & 0 & 0 \end{array} \right] \right) \\
&\quad \sqcup \left(\left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 0 & 1 & 0 & 0 \end{array} \right] \sqcap \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & -1 & 0 & 0 \end{array} \right] \right) \\
&= \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & 0 & 0 & 0 \end{array} \right] \sqcup \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & 0 & 0 & 0 \end{array} \right] = \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & 0 & 0 & 0 \end{array} \right] = “(v_1 = 0) \wedge (v'_1 = 0)”
\end{aligned}$$

Compose over join

$$\begin{aligned}
\text{(i) } “(v_1 = v_2)” &= \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & -1 & 0 & 0 & 0 \end{array} \right] = \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & -1 & 0 & 0 \\ 0 & 1 & 0 & -1 & 0 \end{array} \right]; \\
&\quad \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & -1 & 0 & 0 & 0 \end{array} \right] \\
&= \left(\left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & -1 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{array} \right] \sqcup \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & -1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{array} \right] \right);
\end{aligned}$$

¹In this example, we use the fact that $\top = \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & 0 & 0 & 0 \end{array} \right] \sqcup \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 0 & 1 & 0 & 0 \end{array} \right]$. Although technically we are not working with a vector space over a field, the intuition is that the

KS element $\left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 1 & 0 & 0 & 0 \end{array} \right]$ represents the “line” $v_1 = 0$, the KS element $\left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & \\ 0 & 1 & 0 & 0 \end{array} \right]$ represents the “line” $v'_1 = 0$, and their affine closure is the whole “plane” (i.e., \top).

$$\begin{aligned}
& \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & -1 & 0 & 0 & 0 \end{array} \right] \\
\neq & \left(\begin{array}{l} \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & -1 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{array} \right]; \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & -1 & 0 & 0 & 0 \end{array} \right] \\ \sqcup \\ \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & -1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{array} \right]; \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & -1 & 0 & 0 & 0 \end{array} \right] \end{array} \right) \\
= & \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{array} \right] \sqcup \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{array} \right] \\
= & \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{array} \right] = "(v_1 = 0) \wedge (v_2 = 0)"
\end{aligned}$$

$$\begin{aligned}
\text{(ii) } "(v'_1 = v'_2)" &= \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 0 & 0 & 1 & -1 & 0 \end{array} \right] \\
&= \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 0 & 0 & 1 & -1 & 0 \end{array} \right]; \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & -1 & 0 & 0 \\ 0 & 1 & 0 & -1 & 0 \end{array} \right] \\
&= \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 0 & 0 & 1 & -1 & 0 \end{array} \right]; \\
&\left(\begin{array}{l} \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{array} \right] \sqcup \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ -1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{array} \right] \end{array} \right)
\end{aligned}$$

$$\begin{aligned}
& \neq \left(\begin{array}{c} \left[\begin{array}{ccccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 & 0 \\ 0 & 0 & 1 & -1 & & 0 \end{array} \right]; \left(\begin{array}{c} \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{array} \right] \\ \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ -1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{array} \right] \end{array} \right) \\
\sqcup \\
& \left[\begin{array}{ccccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 & 0 \\ 0 & 0 & 1 & -1 & & 0 \end{array} \right]; \left(\begin{array}{c} \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{array} \right] \\ \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{array} \right] \end{array} \right) \\
& = \left[\begin{array}{ccccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 & 0 \\ 0 & 0 & 1 & 0 & & 0 \\ 0 & 0 & 0 & 1 & & 0 \end{array} \right] \sqcup \left[\begin{array}{ccccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 & 0 \\ 0 & 0 & 1 & 0 & & 0 \\ 0 & 0 & 0 & 1 & & 0 \end{array} \right] \\
& = \left[\begin{array}{ccccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 & 0 \\ 0 & 0 & 1 & 0 & & 0 \\ 0 & 0 & 0 & 1 & & 0 \end{array} \right] = "(v'_1 = 0) \wedge (v'_2 = 0)"
\end{aligned}$$

Similarly, join does not distribute over either meet or compose, as illustrated in the following examples:

Join over meet

$$\begin{aligned}
"(v_1 = v'_1)" &= \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 1 & -1 & & 0 \end{array} \right] = \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 1 & 0 & & 0 \\ 0 & 1 & & 0 \end{array} \right] \sqcup \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 1 & -1 & & 0 \end{array} \right] \\
&= \left(\left(\left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 1 & 0 & & 0 \end{array} \right] \sqcap \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 0 & 1 & & 0 \end{array} \right] \right) \sqcup \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 1 & -1 & & 0 \end{array} \right] \right) \\
&\neq \left(\left(\left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 1 & 0 & & 0 \end{array} \right] \sqcup \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 1 & -1 & & 0 \end{array} \right] \right) \sqcap \right. \\
&\quad \left. \left(\left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 0 & 1 & & 0 \end{array} \right] \sqcup \left[\begin{array}{ccc|c} v_1 & v'_1 & 1 & 0 \\ 1 & -1 & & 0 \end{array} \right] \right) \right) \\
&= \top \sqcap \top = "true"
\end{aligned}$$

Join over compose

$$\begin{aligned}
“(v'_1 = v_1 + 1)” &= \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] = \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \sqcup \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \\
&= \left(\text{Id}; \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \right) \sqcup \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \\
&= \left(\left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 0 \end{array} \right]; \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \right) \sqcup \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \\
&\neq \left(\left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 0 \end{array} \right] \sqcup \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \right); \\
&\quad \left(\left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \sqcup \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] \right) \\
&= \top; \left[\begin{array}{cc|c} v_1 & v'_1 & 1 \\ 1 & -1 & 1 \end{array} \right] = \top = \text{“true”}
\end{aligned}$$

3.4.5 Assuming Conditions

By “assuming” a condition φ on a KS element X , we mean to compute a minimal KS element Y such that

$$\gamma_{\text{KS}}(Y) \supseteq \gamma_{\text{KS}}(X) \cap \{v \mid \varphi(v)\}.$$

This operation is needed to compute the transformer for an assume edge in a program graph (i.e., the true-branch or false-branch of an if-then-else statement). It can also be used to create transformers for assignments; for instance, the transformer for the assignment $x \leftarrow 3u + 2v$ can be created by starting with the KS element for the identity relation on vocabulary

$\vec{V}, \left[\begin{array}{cc|c} \vec{V} & \vec{V}' & 1 \\ \hline I & -I & 0 \end{array} \right]$, having $x' \in \vec{V}'$, and assuming the equality $x' = 3u + 2v$.

Assuming a w -bit affine constraint is straightforward: rewrite the constraint to isolate 0 on one side; form a matrix row from resulting constraint's coefficients; append the row to the KS element X ; and Howellize. In other words, when φ is an affine constraint, we create a one-row KS element that represents φ exactly, and take the meet with X .

It is also possible to perform an assume with respect to an affine congruence of the form " $\text{lhs} = \text{rhs} \pmod{2^h}$ ", with $h < w$. We rewrite the congruence as an equivalent congruence modulo 2^w , by multiplying the modulus 2^h and all of the coefficients by 2^{w-h} , to obtain the w -bit affine constraint " $2^{w-h}\text{lhs} = 2^{w-h}\text{rhs}$ ". We then proceed as before.

3.4.6 Containment

Two KS elements X and Y are equal if their concretizations are equal: $\gamma(X) = \gamma(Y)$. However, when each KS element is in Howell form, equality checking is trivial because Howell form is unique among all matrices with the same row space (or null space) [44]. Consequently, containment can be checked using meet and equality: $X \sqsubseteq Y$ iff $X = X \sqcap Y$.

3.4.7 Merge Functions

[54] extended the Sharir and Pnueli algorithm 1981 for interprocedural dataflow analysis to handle local variables. At a site where procedure P calls procedure Q , the local variables of P are modeled as if the current incarnations of P 's locals are stored in locations that are inaccessible to Q and to procedures transitively called by Q . Because the contents of P 's locals cannot be affected by the call to Q , a *merge function* is used to combine them with the element returned by Q to create the state in P after

the call to Q has finished. Other work using merge functions includes [75] and [56].

To simplify the discussion, assume that all scopes have the same number of locals L. Each merge function is of the form

$$\text{MERGE}(\text{CallSiteVal}, \text{CalleeExitVal}) \stackrel{\text{def}}{=} \\ \text{CallSiteVal} ; \text{REPLACELOCALS}(\text{CalleeExitVal}).$$

Suppose that the i^{th} vocabulary consists of sub-vocabularies \vec{g}_i and \vec{l}_i . The operation $\text{REPLACELOCALS}(\text{CalleeExitVal})$ is defined as follows:

1. Project away vocabulary \vec{l}_2 from CalleeExitVal.
2. Insert L all-0 columns for vocabulary \vec{l}_2 . The KS element now has no constraints on the variables in \vec{l}_2 .
3. Append L rows, $\left[\begin{array}{c|c|c|c|c} \vec{g}_1 & \vec{l}_1 & \vec{g}_2 & \vec{l}_2 & 1 \\ \hline 0 & I & 0 & -I & 0 \end{array} \right]$, so that in $\text{REPLACELOCALS}(\text{CalleeExitVal})$ each variable in vocabulary \vec{l}_2 is constrained to have the value of the corresponding variable in vocabulary \vec{l}_1 .

3.4.8 Symbolic Abstraction ($\hat{\alpha}(\varphi)$)

[53, Fig. 2] gave an algorithm for the problem of *symbolic abstraction* (§2.3.2) with respect to the KS domain: given a quantifier-free bit-vector (QFBV) formula φ , the algorithm returns the best element in KS that over-approximates $\llbracket \varphi \rrbracket$. The algorithm given by King and Søndergaard, which we will denote by $\hat{\alpha}_{\text{KS}}^{\uparrow}(\varphi)$, needs the minor correction of using Howell form instead of row-echelon form for the projections that take place in its join operations, as discussed in §3.4.4.

Pseudo-code for $\hat{\alpha}_{\text{KS}}^{\uparrow}(\varphi)$ is shown in Fig. 3.1(a). The matrix *lower* is maintained in Howell form throughout. In line 6, $\hat{\alpha}_{\text{KS}}^{\uparrow}(\varphi)$ uses the *symbolic concretization* of p , denoted by $\hat{\gamma}(p)$. For most abstract domains, including

<p>Require: φ: a QFBV formula</p> <p>Ensure: $\hat{\alpha}(\varphi)$ for the KS domain</p> <pre> 1: <i>lower</i> $\leftarrow \perp$ 2: <i>i</i> $\leftarrow 1$ 3: 4: while <i>i</i> \leq rows(<i>lower</i>) do 5: <i>p</i> \leftarrow <i>lower</i>[rows(<i>lower</i>) - <i>i</i> + 1] {<i>p</i> \sqsupseteq <i>lower</i>} 6: <i>S</i> \leftarrow Model($\varphi \wedge \neg \hat{\gamma}(p)$) 7: if <i>S</i> is TimeOut then return \top 8: else if <i>S</i> is None then {$\varphi \Rightarrow \hat{\gamma}(p)$} 9: <i>i</i> $\leftarrow i + 1$ 10: 11: else {<i>S</i> $\models \hat{\gamma}(p)$} 12: <i>lower</i> \leftarrow <i>lower</i> \sqcup $\beta(S)$ 13: <i>ans</i> \leftarrow <i>lower</i> 14: return <i>ans</i> </pre> <p style="text-align: center;">(a)</p>	<p>Require: φ: a QFBV formula</p> <p>Ensure: $\hat{\alpha}(\varphi)$ for the KS domain</p> <pre> 1: <i>lower</i> $\leftarrow \perp$ 2: <i>i</i> $\leftarrow 1$ 3: <i>upper</i> $\leftarrow \top$ 4: while <i>i</i> \leq rows(<i>lower</i>) do 5: <i>p</i> \leftarrow <i>lower</i>[rows(<i>lower</i>) - <i>i</i> + 1] {<i>p</i> \sqsupseteq <i>lower</i>, <i>p</i> $\not\sqsupseteq$ <i>upper</i>} 6: <i>S</i> \leftarrow Model($\varphi \wedge \neg \hat{\gamma}(p)$) 7: if <i>S</i> is TimeOut then return <i>upper</i> 8: else if <i>S</i> is None then {$\varphi \Rightarrow \hat{\gamma}(p)$} 9: <i>i</i> $\leftarrow i + 1$ 10: <i>upper</i> \leftarrow <i>upper</i> \sqcap <i>p</i> 11: else {<i>S</i> $\models \hat{\gamma}(p)$} 12: <i>lower</i> \leftarrow <i>lower</i> \sqcup $\beta(S)$ 13: <i>ans</i> \leftarrow <i>lower</i> 14: return <i>ans</i> </pre> <p style="text-align: center;">(b)</p>
--	--

Figure 3.1: (a) The King-Søndergaard algorithm for symbolic abstraction ($\hat{\alpha}_{\text{KS}}^{\uparrow}(\varphi)$). (b) The Thakur-Elder-Reps bilateral algorithm for symbolic abstraction, instantiated for the KS domain: $\hat{\alpha}_{\text{TER[KS]}}^{\uparrow}(\varphi)$. In both algorithms, *lower* is maintained in Howell form throughout.

KS, it is easy to write a $\hat{\gamma}$ function. As mentioned in §3.1.4, affine equalities can be read out from a KS element *M* (regardless of whether *M* is in Howell form) as follows:

$$\text{If } \left[\begin{array}{cccc|c} v_1 & \dots & v_k & v'_1 & \dots & v'_k & 1 \\ a_1 & \dots & a_k & a'_1 & \dots & a'_k & b \end{array} \right] \text{ is a row of } M, \text{ then}$$

$$\sum_i a_i v_i + \sum_i a'_i v'_i = -b \text{ is a constraint on } \gamma_{\text{KS}}(M).$$

The conjunction of these constraints describes $\gamma_{\text{KS}}(M)$ exactly. Conse-

quently, $\hat{\gamma}(M)$ can be defined as follows:

$$\hat{\gamma}(M) \stackrel{\text{def}}{=} \bigwedge_{\substack{[a_1 \cdots a_k \ a'_1 \cdots a'_k | b] \\ \text{is a row of } M}} \sum_i a_i v_i + \sum_i a'_i v'_i = -b$$

The algorithm $\hat{\alpha}_{\text{KS}}^\uparrow(\varphi)$ is a successive-approximation algorithm: it computes a sequence of successively larger approximations to $\llbracket \varphi \rrbracket$. It maintains an under-approximation of the final answer in the variable “*lower*”, which is initialized to \perp on line 1. On each iteration, the algorithm selects p , a single row (constraint) of *lower* (line 5), and calls a decision procedure to determine whether there is a model that satisfies the formula “ $\varphi \wedge \neg \hat{\gamma}(p)$ ” (line 6). When $\varphi \wedge \neg \hat{\gamma}(p)$ is unsatisfiable, φ implies $\hat{\gamma}(p)$. In this case, p cannot be used to figure out how to make *lower* larger, so variable i is incremented (line 9), which means that on the next iteration of the loop, the algorithm selects the row immediately above p (line 5).

On the other hand, if the decision procedure returns a model S , the under-approximation *lower* is updated to make it larger via the join performed on the right-hand side of the assignment in line 12

$$\text{lower} \leftarrow \text{lower} \sqcup \beta(S). \quad (3.5)$$

Because KS elements represent two-vocabulary relations, S is an assignment of concrete values to both the pre-state and post-state variables:

$$S = [\dots, v_i \mapsto v_i, \dots, v'_i \mapsto v'_i, \dots],$$

or, equivalently,

$$S = [\vec{X} \mapsto \vec{v}, \vec{X}' \mapsto \vec{v}']. \quad (3.6)$$

The notation $\beta(S)$ in line 12 denotes the abstraction of the singleton state-set $\{S\}$ to a KS element. $\{S\}$ can always be represented exactly in the KS

domain as follows (where the superscript t denotes the operation of vector transpose):

$$\beta(S) \stackrel{\text{def}}{=} \left[\begin{array}{cc|c} \vec{X} & \vec{X}' & 1 \\ \hline I & 0 & (-\vec{v})^t \\ 0 & I & (-\vec{v}')^t \end{array} \right] \quad (3.7)$$

An Improvement to $\hat{\alpha}_{KS}^\uparrow$

Algorithm $\hat{\alpha}_{KS}^\uparrow$ from Fig. 3.1(a) is related to, but distinct from, an earlier $\hat{\alpha}$ algorithm, due to [84] (RSY), which applies not just to the KS domain, but to all abstract domains that meet a certain interface. (In other words, $\hat{\alpha}_{RSY}$ is the cornerstone of a *framework* for symbolic abstraction.) The two algorithms resemble one another in that they both find $\hat{\alpha}(\varphi)$ via successive approximation from below. However, there is a key difference in the nature of the satisfiability queries that are passed to the decision procedure by the two algorithms. Compared to $\hat{\alpha}_{RSY}$, $\hat{\alpha}_{KS}$ issues comparatively inexpensive satisfiability queries in which only a single affine equality is negated²—i.e., line 6 of Fig. 3.1(a) calls $\text{Model}(\varphi \wedge \neg\hat{y}(p))$, where p is a single constraint from *lower*.

This difference—together with the observation that in practice $\hat{\alpha}_{KS}$ was about ten times faster than $\hat{\alpha}_{RSY}$ when the latter was instantiated for the KS domain—led Thakur, Elder, and Reps 2012 (TER) to investigate the fundamental principles underlying $\hat{\alpha}_{RSY}$ and $\hat{\alpha}_{KS}$. They developed a new framework, $\hat{\alpha}_{TER}^\uparrow$, that transfers $\hat{\alpha}_{KS}$'s advantages from the KS domain to other abstract domains [100].

In addition to generating less expensive satisfiability queries, the second benefit of $\hat{\alpha}_{TER}^\uparrow$ is that $\hat{\alpha}_{TER}^\uparrow$ generally returns a more precise answer than $\hat{\alpha}_{RSY}$ and $\hat{\alpha}_{KS}$ when a timeout occurs. Because $\hat{\alpha}_{RSY}$ and $\hat{\alpha}_{KS}$ maintain only under-approximations of the desired answer, if the successive-

²See [100, §3] for a more extensive explanation of the differences between $\hat{\alpha}_{KS}$ and $\hat{\alpha}_{RSY}$.

approximation process takes too much time and needs to be stopped, they must return \top to be sound. In contrast, $\hat{\alpha}_{\text{TER}}^\dagger$ is *bilateral*, and can generally return a nontrivial (non- \top) element in case of a timeout. That is, $\hat{\alpha}_{\text{TER}}^\dagger$ maintains both an under-approximation and a (nontrivial) over-approximation of the desired answer, and hence is resilient to timeouts: $\hat{\alpha}_{\text{TER}}^\dagger$ returns the over-approximation if it is stopped at any point.

Fig. 3.1(b) shows the $\hat{\alpha}_{\text{TER}}^\dagger$ algorithm instantiated for the KS domain, which we call $\hat{\alpha}_{\text{TER}[\text{KS}]}^\dagger$. The differences between Fig. 3.1(a) and (b) are highlighted in gray. $\hat{\alpha}_{\text{TER}[\text{KS}]}^\dagger$ initializes *upper* to \top on line 3. At any stage in the algorithm *upper* $\sqsupseteq \hat{\alpha}(\varphi)$. It is sound to update *upper* on line 10 by performing a meet with the row *p* that was selected in line 5. Because *p*'s leading index, $LI(p)$, is less than the leading index of every row in *upper*, *p* constrains the value of variable $v_{LI(p)}$, whereas *upper* places no constraints on variable $v_{LI(p)}$. Therefore, $p \not\sqsupseteq \text{upper}$, which guarantees progress because $p \sqcap \text{upper} \sqsubset \text{upper}$. Termination is guaranteed.

In case of a decision-procedure timeout (line 7), $\hat{\alpha}_{\text{TER}[\text{KS}]}^\dagger$ returns *upper* as the answer (line 7). If the algorithm finishes without a timeout, then $\hat{\alpha}_{\text{TER}[\text{KS}]}^\dagger$ computes $\hat{\alpha}(\varphi)$; on the other hand, if a timeout occurs, the element returned is generally an over-approximation of $\hat{\alpha}(\varphi)$ —i.e., $\hat{\alpha}_{\text{TER}[\text{KS}]}^\dagger$ computes $\tilde{\alpha}(\varphi)$.

In the KS instantiation of $\hat{\alpha}_{\text{TER}}^\dagger$, *upper* can actually be represented implicitly by $\text{lower}[(\text{rows}(\text{lower}) - i + 2) \dots \text{rows}(\text{lower})]$. Consequently, the assignment $\text{upper} \leftarrow \text{lower}[(\text{rows}(\text{lower}) - i + 2) \dots \text{rows}(\text{lower})]$ need only be performed if line 7 is reached, and neither of the assignments on lines 3 and 10 need to be performed explicitly.

3.4.9 Number of Satisfying Solutions

The *size* of a KS element X with k variables over \mathbb{Z}_2^w is the number of k -tuples that satisfy X . The size computation is inexpensive; the size of X

depends on the leading values in X , and the number of rows in X . (X is assumed to be in Howell form.)

- If X is bottom, then $\text{SIZE}(X) = 0$.
- Otherwise, we can derive how to compute $\text{SIZE}(X)$ by imagining that we are building up a partial assignment for the variables, from right to left. (In what follows, for simplicity we assume that we have a one-vocabulary KS element and “right to left” means from higher-indexed variables to lower-indexed variables.) In this case, each variable v_i is constrained by the current partial assignment to the variables $\{v_j \mid i < j\}$, and by the row with leading index i :
 - If the leading value of that row is 1, then for every partial assignment to the variables $\{v_j \mid i < j\}$, there is exactly one consistent value for v_i , namely, whatever value for v_i satisfies the equation for the row when the values in the partial assignment are used for the higher-indexed variables.
 - If the leading value of a row is 2^m for some value m , then for every partial assignment to the variables $\{v_j \mid i < j\}$, there is exactly one consistent value for $2^m v_i$, namely, whatever value y for $2^m v_i$ satisfies the equation for the row when the values in the partial assignment are used for the higher-indexed variables.

However, there are 2^m different ways to choose v_i to obtain the needed value y . That is, if \bar{v} is a value such that $2^m \bar{v} = y$, then so are all 2^m values in the set

$$\{(\bar{v} + 2^{w-m}p) \pmod{2^w} \mid 0 \leq p \leq 2^m - 1\}.$$
- Finally, if there is no row with leading index i , then v_i is fully unconstrained, and can take on any of the 2^w available values.

Altogether, the product of these counts is the number of satisfying solutions of KS element X . In particular, let u be the number of indices that are not the leading index of any row of X . Then $\text{SIZE}(X)$

Instruction Characteristics		
Kind	# instruction instances	# different opcodes
ordinary	12,734	164
lock prefix	2,048	147
rep prefix	2,143	158
repne prefix	2,141	154
full corpus	19,066	179

Figure 3.2: Some of the characteristics of the corpus of 19,066 (non-privileged, non-floating point, non-mmx) instructions.

is the product of the leading values in X , times $(2^w)^u$.

Example 3.18. Consider again the KS element from Eqn. (3.2)

$$X_0 = \left[\begin{array}{cccc|c} v_1 & v_2 & v'_1 & v'_2 & 1 \\ 1 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 8 & 0 \end{array} \right],$$

where $w = 4$, so that we are working in \mathbb{Z}_{16} . Then $\text{SIZE}(X_0)$ equals $1 \times 8 \times (2^4)^2 = 2,048$.

3.5 Experiments

In this section, we present the results of experiments to evaluate the costs and benefits—in terms of time and precision—of the methods described in earlier sections. The experiments were designed to shed light on the following questions:

1. Which method of obtaining abstract transformers is fastest: $\hat{\alpha}_{\text{KS}}$ (§2.3.2 and §3.4.8), KS-reinterpretation (§2.3.1), or MOS-reinterpretation (§2.3.1, [62, §4.1.2])?

2. Does MOS-reinterpretation or KS-reinterpretation yield more precise abstract transformers for machine instructions?
3. For what percentage of program points does $\hat{\alpha}_{\text{KS}}$ produce more precise answers than KS-reinterpretation and MOS-reinterpretation? This question actually has two versions, depending on whether we are interested in
 - one-vocabulary affine relations that hold at branch points
 - two-vocabulary procedure summaries obtained at procedure-exit points.

As shown in §3.3.1, the MOS and KS domains are incomparable. To compare the final results obtained using the two domains, we converted each MOS element to a KS element, using the algorithm from §3.3.2, and then checked for equality, containment (§3.4.6), or incomparability. It might be argued that this approach biases the results in favor of KS. However, if we have run an MOS-based analysis and are interested in using affine relations in a client application, we must extract an affine relation from each computed MOS element. In §3.3.1, we showed that, in general, an MOS element \mathcal{B} does not represent an affine relation; thus, a client application needs to obtain an affine relation that over-approximates $\gamma_{\text{MOS}}(\mathcal{B})$. Consequently, the comparison method that we used *is* sensible, because it compares the precision of the affine relations that would be seen by a client application.

3.5.1 Experimental Setup

To address these questions, we performed two experiments. Both experiments were run on a single core of a single-processor 16-core 2.27 GHz Xeon computer running 64-bit Windows 7 Enterprise (Service Pack 1), configured so that a user process has 4 GB of memory.

Program Name	Measures of Size						
	Instrs	Procs	BBs	Branches	WPDS Rules		
					Δ_0	Δ_1	Δ_2
write	232	10	134	26	10	151	5
finger	532	18	298	48	18	353	20
subst	1093	16	609	74	16	728	13
chkdsk	1468	18	787	119	18	887	32
convert	1927	38	1013	161	38	1266	22
route	1982	40	931	243	40	1368	63
comp	2377	35	1261	224	35	1528	30
logoff	2470	46	1145	306	46	1648	72
setup	4751	67	1862	589	67	2847	121

Figure 3.3: Program information. All nine utilities are from Microsoft Windows version 5.1.2600.0, except setup, which is from version 5.1.2600.5512. The columns show the number of instructions (Instrs); the number of procedures (Procs); the number of basic blocks (BBs); the number of branch instructions (Branches); and the number of Δ_0 , Δ_1 , and Δ_2 rules in the WPDS encoding (WPDS Rules).

Per-Instruction Experiment On a corpus of 19,066 instances of x86 instructions, we measured (i) the time taken to create MOS and KS transformers via the operator-by-operator reinterpretation method supported by TSL [61, 62], and (ii) the relative precision of the abstract transformers obtained by the two methods.

This corpus was created using the ISAL instruction-decoder generator [62, §2.1] in a mode in which the input specification of the concrete syntax of the x86 instruction set was used to create a randomized instruction *generator*—instead of the standard mode in which ISAL creates an instruction *recognizer*. By this means, we are assured that the corpus has substantial coverage of the syntactic features of the x86 instruction set (including opcodes, addressing modes, and prefixes, such as “lock”, “rep”, and “repne”); see Fig. 3.2.

Interprocedural-Analysis Experiment We performed flow-sensitive, context-sensitive, interprocedural affine-relation analysis on the executables of nine Windows utilities, using four different sets of abstract transformers:

1. MOS transformers for basic blocks, created by performing operator-by-operator MOS-reinterpretation.
2. KS transformers for basic blocks, created by performing operator-by-operator KS-reinterpretation.
3. KS transformers for basic blocks, created by symbolic abstraction of quantifier-free bit-vector (QFBV) formulas that capture the precise bit-level semantics of register-access/update operations in the different basic blocks. We denote this symbolic-abstraction method by $\hat{\alpha}_{KS}$.
4. KS transformers for basic blocks, created by symbolic abstraction of quantifier-free bit-vector (QFBV) formulas that, in addition to register-access/update operations, also capture the precise bit-level semantics of all *memory-access/update* and *flag-access/update* operations. We denote this symbolic-abstraction method by $\hat{\alpha}_{KS}^+$.

For these programs, the generated abstract transformers were used as “weights” in a weighted pushdown system (WPDS).

Fig. 3.3 lists several size parameters of the executables (number of instructions, procedures, basic blocks, branches, and number of WPDS rules). WPDS rules can be divided into three categories, called Δ_0 , Δ_1 , and Δ_2 rules [10, 85]. The number of Δ_1 rules corresponds roughly to the total number of edges in a program’s intraprocedural control-flow graphs; the number of Δ_2 rules corresponds to the number of call sites in the program; the number of Δ_0 rules corresponds to the number of procedure-exit sites.

$\hat{\alpha}_{KS}^+$ has the potential to create more precise KS weights than $\hat{\alpha}_{KS}$ because $\hat{\alpha}_{KS}^+$ can account for transformations of register values that involve a sequence of memory-access/update and/or flag-access/update opera-

tions within a basic block \mathcal{B} . For example, suppose that \mathcal{B} contains a store to memory of register `eax`'s value, and a subsequent load from memory of that value into `ebx`. Because $\hat{\alpha}_{\text{KS}}^+$ uses a formula that captures the two memory operations, it can find a weight that captures the transformation $\text{ebx}' = \text{eax}$. A second type of example involving a store to memory followed by a load from memory within a basic block involves a sequence of the form “push *constant*; ... pop `edi`”, and thus represents the transformation $\text{edi}' = \text{constant}$. Such sequences occur in several of the programs listed in Fig. 3.3.

Fig. 3.4 shows a TSL specification for the `MOV` and `ADD` instructions of the Intel IA32 instruction set. As illustrated in line 13 of Fig. 3.4, the top-level function that is reinterpreted in TSL is `interpInstr`, which is of type

$$\text{interpInstr} : \text{instruction} \times \text{state} \rightarrow \text{state}.$$

To use semantic reinterpretation to implement `CREATEABSTRANS` for the KS domain, `interpInstr` is reinterpreted as a $\text{KS}[V; V']$ transformer; that is, $\text{interpInstr}_{\text{KS}}$ has the type

$$\text{interpInstr}_{\text{KS}} : \text{instruction} \times \text{KS}[V; V'] \rightarrow \text{KS}[V; V'].$$

Let Id denote the $\text{KS}[V; V']$ identity relation, $\left[\begin{array}{c|c} \vec{V} & \vec{V}' \\ \hline \text{I} & -\text{I} \\ \hline & 0 \end{array} \right]$. To reinterpret an individual instruction ι , one invokes $\text{interpInstr}_{\text{KS}}(\iota, \text{Id})$.

For a basic block $\mathcal{B} = [\iota_1, \dots, \iota_m]$, there are two approaches to performing $\text{KS}[V; V']$ reinterpretation:

- *Composed reinterpretation:*

$$w_{\text{KS}}^{\mathcal{B}} \leftarrow \text{interpInstr}_{\text{KS}}(\iota_1, \text{Id}) ; \text{interpInstr}_{\text{KS}}(\iota_2, \text{Id}) ; \dots ; \text{interpInstr}_{\text{KS}}(\iota_m, \text{Id}).$$

```

[1] // Abstract-syntax declarations
[2] reg: EAX() | EBX() | . . . ;
[3] flag: ZF() | SF() | . . . ;
[4] operand: Indirect(reg reg INT8 INT32) | DirectReg(reg) | Immediate(INT32) | ...;
[5] instruction: MOV(operand operand) | ADD(operand operand) | . . . ;
[6] state: State(MAP[INT32,INT8] // memory-map
[7]             MAP[reg,INT32] // register-map
[8]             MAP[flag,BOOL]); // flag-map
[9] // Interpretation functions
[10] INT32 interpOp(state S, operand op) { . . . };
[11] state updateFlag(state S, INT32 v1, INT32 v2, INT32 v3) { . . . };
[12] state updateState(state S, operand op, INT32 val) { . . . };
[13] state interpInstr(instruction I, state S) {
[14]   with(I) (
[15]     MOV(dstOp, srcOp):
[16]       let srcVal = interpOp(S, srcOp);
[17]       in ( updateState(S, dstOp, srcVal) ),
[18]     ADD(dstOp, srcOp):
[19]       let dstVal = interpOp(S, dstOp);
[20]       srcVal = interpOp(S, srcOp);
[21]       result = dstVal + srcVal;
[22]       S2 = updateFlag(S, dstVal, srcVal, result);
[23]       in ( updateState(S2, dstOp, result) ),
[24]     . . .
[25]   );
[26] };

```

Figure 3.4: A fragment of the TSL specification of the concrete semantics of the Intel IA32 instruction set.

- *Chained reinterpretation:*

$$w_{KS}^B \leftarrow \text{interpInstr}_{KS}(t_m, \dots \text{interpInstr}_{KS}(t_2, \text{interpInstr}_{KS}(t_1, \text{Id})) \dots).$$

Our experiments use chained reinterpretation for two reasons:

1. There are cases in which chained reinterpretation creates a more precise $\text{KS}[V; V']$ element. For instance, consider the following code fragment, which zeros the two low-order bytes of register `eax` and does a bitwise-or of `eax` into `ebx` (`ax` denotes the two low-order bytes of register `eax`):

```

t1 : xor  ax, ax
t2 : or   ebx, eax

```

The semantics of this code fragment can be expressed as follows:

$$\text{ebx}' = (\text{ebx} \mid (\text{eax} \ \& \ \text{xFFFF0000})) \wedge \text{eax}' = (\text{eax} \ \& \ \text{xFFFF0000}),$$

where “&” and “|” denote bitwise-and and bitwise-or, respectively. $\text{interpInstr}_{\text{KS}}(t_1, \text{Id})$ creates the KS element

$\left[\begin{array}{cccc|c} \text{eax} & \text{ebx} & \text{eax}' & \text{ebx}' & 1 \\ 0 & 1 & 0 & -1 & 0 \\ 0 & 0 & 2^{16} & 0 & 0 \end{array} \right]$, which captures $(\text{ebx}' = \text{ebx}) \wedge (2^{16}\text{eax}' = 0)$. The two approaches to reinterpretation produce the following answers:

- *Composed reinterpretation:*

$$\begin{aligned}
& \text{interpInstr}_{\text{KS}}(t_1, \text{Id}) ; \text{interpInstr}_{\text{KS}}(t_2, \text{Id}) \\
&= \left[\begin{array}{cccc|c} \text{eax} & \text{ebx} & \text{eax}' & \text{ebx}' & 1 \\ 0 & 1 & 0 & -1 & 0 \\ 0 & 0 & 2^{16} & 0 & 0 \end{array} \right] ; \left[\begin{array}{cccc|c} \text{eax} & \text{ebx} & \text{eax}' & \text{ebx}' & 1 \\ 1 & 0 & -1 & 0 & 0 \end{array} \right] \\
&= \left[\begin{array}{cccc|c} \text{eax} & \text{ebx} & \text{eax}' & \text{ebx}' & 1 \\ 0 & 0 & 2^{16} & 0 & 0 \end{array} \right] \\
&= (2^{16}\text{eax}' = 0).
\end{aligned}$$

- *Chained reinterpretation:*

$$\begin{aligned}
& \text{interpInstr}_{\text{KS}}(\iota_2, \text{interpInstr}_{\text{KS}}(\iota_1, \text{Id})) \\
&= \text{interpInstr}_{\text{KS}}(\iota_2, \left[\begin{array}{cccc|c} \text{eax} & \text{ebx} & \text{eax}' & \text{ebx}' & 1 \\ 0 & 1 & 0 & -1 & 0 \\ 0 & 0 & 2^{16} & 0 & 0 \end{array} \right]) \\
&= \left[\begin{array}{cccc|c} \text{eax} & \text{ebx} & \text{eax}' & \text{ebx}' & 1 \\ 0 & 2^{16} & 0 & -2^{16} & 0 \\ 0 & 0 & 2^{16} & 0 & 0 \end{array} \right] \\
&= (2^{16}\text{ebx}' = 2^{16}\text{ebx}) \wedge (2^{16}\text{eax}' = 0).
\end{aligned}$$

In particular, the reinterpretation of “ ι_2 : or ebx, eax” takes place in a context in which the two low-order bytes of eax are partially constant ($2^{16}\text{eax} = 0$). Because of this additional piece of information, the reinterpretation technique recovers the additional conjunct “ $2^{16}\text{ebx}' = 2^{16}\text{ebx}$ ”.

2. In the case of $\hat{\alpha}_{\text{KS}}$, a formula φ_{B} is created that captures the concrete semantics of B (via symbolic execution), and then the KS weight for B is obtained by performing $w_{\text{KS}}^{\text{B}} \leftarrow \hat{\alpha}_{\text{KS}}(\varphi_{\text{B}})$. Letting QFBV denote the type of quantifier-free bit-vector formulas, the QFBV reinterpretation of interpInstr has the type

$$\text{interpInstr}_{\text{QFBV}} : \text{instruction} \times \text{QFBV} \rightarrow \text{QFBV}.$$

Symbolic execution is performed by chained reinterpretation:

$$\varphi_{\text{B}} \leftarrow \text{interpInstr}_{\text{QFBV}}(\iota_m, \dots \text{interpInstr}_{\text{QFBV}}(\iota_2, \text{interpInstr}_{\text{QFBV}}(\iota_1, \text{Id})) \dots).$$

The $\hat{\alpha}_{\text{KS}}^+$ weight for B is created similarly, except that we also arrange for φ_{B} to encode all memory-access/update and flag-access/update operations.

Total instructions	MOS-reinterpretation time (seconds)	KS-reinterpretation time (seconds)
19,066	23.3	348.2

Figure 3.5: Comparison of the performance of MOS-reinterpretation and KS-reinterpretation for x86 instructions.

Identical precision	MOS-reinterpretation more precise	KS-reinterpretation more precise	Total
18,158	0	908	19,066

Figure 3.6: Comparison of the precision of MOS-reinterpretation and KS-reinterpretation for x86 instructions.

For our experiments, we wanted to control for any precision improvements that might be due solely to the use of chained reinterpretation; thus, we use chained reinterpretation for all of the weight-generation methods.³

Due to the high cost in §3.5.3 of constructing WPDS weights via $\hat{\alpha}_{KS}$ and $\hat{\alpha}_{KS}^+$, we ran all WPDS analyses without the code for libraries. Values are returned from x86 procedure calls in register `eax`, and thus in our experiments library functions were modeled approximately (albeit unsoundly, in general) by “`havoc(eax)`”.

To implement $\hat{\alpha}_{KS}$ and $\hat{\alpha}_{KS}^+$, we used the Yices solver [26], version 1.0.19, with the timeout for each invocation set to three seconds.

We compared the precision of the one-vocabulary affine relations at branch points, as well as two-vocabulary affine relations at procedure exits, which can be used as procedure summaries.

3.5.2 Reinterpretation of Individual Instructions

Figs. 3.5 and 3.6 summarize the results of the per-instruction experiment. They answer questions (1) and (2) posed at the beginning of §3.5.

- KS-reinterpretation created an abstract transformer that was more precise than the one created by MOS-reinterpretation for about 4.76% of the instructions. MOS-reinterpretation never created an abstract transformer that was more precise than the one created by KS-reinterpretation.
- However, MOS-reinterpretation is much faster: to generate abstract transformers for the entire corpus of instructions, MOS-reinterpretation is about 14.9 times faster than KS-reinterpretation.

Example 3.19. *One instruction for which the abstract transformer created by KS-reinterpretation is more precise than the transformer created by MOS-reinterpretation is $\iota \stackrel{\text{def}}{=} \text{“add bh, al”}$. This instruction adds the value of al , the low-order byte of register eax , to the value of bh , the second-to-lowest byte of register ebx , and stores the result in bh . The semantics of this instruction can be expressed as a QFBV formula as follows:*

$$\varphi_{\iota} \stackrel{\text{def}}{=} \text{ebx}' = \left((\text{ebx} \& 0\text{xFFF00FF}) \mid ((\text{ebx} + 256 * (\text{eax} \& 0\text{xFF})) \& 0\text{xFF00}) \right) \wedge (\text{eax}' = \text{eax}). \quad (3.8)$$

Eqn. (3.8) shows that the semantics of the instruction involves non-linear bit-masking operations.

The abstract transformer created via MOS-reinterpretation corresponds to $\text{havoc}(\text{ebx}')$; all other registers are unchanged. That is, if we only had the three registers eax , ebx , and ecx , the abstract transformer created via MOS-

³MOS-reinterpretation of a basic block is performed by chained reinterpretation, using $\text{interpInstr}_{\text{MOS}} : \text{instruction} \times \text{MOS} \rightarrow \text{MOS}$.

reinterpretation would be

$$\left\{ \left[\begin{array}{c|ccc} 1 & 0 & 0 & 0 \\ \hline 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right], \left[\begin{array}{c|ccc} 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right] \right\},$$

which captures the affine transformation “ $(\mathbf{eax}' = \mathbf{eax}) \wedge (\mathbf{ecx}' = \mathbf{ecx})$ ”. In contrast, the transformer created via KS-reinterpretation is

$$\left[\begin{array}{cccccc|c} \mathbf{eax} & \mathbf{ebx} & \mathbf{ecx} & \mathbf{eax}' & \mathbf{ebx}' & \mathbf{ecx}' & 1 \\ \hline 1 & 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & 2^{24} & 0 & 0 & -2^{24} & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & -1 & 0 \end{array} \right],$$

which corresponds to “ $(\mathbf{eax}' = \mathbf{eax}) \wedge (2^{24}\mathbf{ebx}' = 2^{24}\mathbf{ebx}) \wedge (\mathbf{ecx}' = \mathbf{ecx})$ ”. Both transformers are over-approximations of the instruction’s semantics, but the extra conjunct $(2^{24}\mathbf{ebx}' = 2^{24}\mathbf{ebx})$ in the KS element captures the fact that the low-order byte of \mathbf{ebx} is not changed by executing “`add bh, al`”.

In contrast, $\hat{\alpha}_{\text{KS}}(\varphi_{\downarrow})$, the most-precise over-approximation of φ_{\downarrow} that can be expressed as a KS element is (the Howellization of)

$$\left[\begin{array}{cccccc|c} \mathbf{eax} & \mathbf{ebx} & \mathbf{ecx} & \mathbf{eax}' & \mathbf{ebx}' & \mathbf{ecx}' & 1 \\ \hline 1 & 0 & 0 & -1 & 0 & 0 & 0 \\ 2^{24} & 2^{16} & 0 & 0 & -2^{16} & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & -1 & 0 \end{array} \right],$$

which corresponds to “ $(\mathbf{eax}' = \mathbf{eax}) \wedge (2^{16}\mathbf{ebx}' = 2^{16}\mathbf{ebx} + 2^{24}\mathbf{eax}) \wedge (\mathbf{ecx}' = \mathbf{ecx})$ ”. Multiplying by a power of 2 serves to shift bits to the left; because it is performed in arithmetic mod 2^{32} , bits shifted off the left end are unconstrained. Thus, the second conjunct captures the relationship between the low-order two bytes of \mathbf{ebx}' , the low-order two bytes of \mathbf{ebx} , and the low-order byte of \mathbf{eax} . \square

3.5.3 Interprocedural Analysis

Fig. 3.7 shows the times for WPDS construction (including constructing the weights that serve as abstract transformers) and performing interprocedural dataflow analysis by performing post* and path summary. Columns 11 and 15 of Fig. 3.7 show the number of $\hat{\alpha}$ calls for which weight generation timed out during $\hat{\alpha}_{\text{KS}}$ and $\hat{\alpha}_{\text{KS}}^+$, respectively. During WPDS construction, if Yices times out during $\hat{\alpha}_{\text{KS}}$ or $\hat{\alpha}_{\text{KS}}^+$, the implementation uses a weight that is less precise than the best transformer, but it always uses a weight that is at least as precise as the weight obtained using KS-reinterpretation.⁴ The number of WPDS rules is given in Fig. 3.3; a timeout occurred for about 1.0% of the $\hat{\alpha}_{\text{KS}}$ calls (computed as a geometric mean⁵), and for about 1.65%

⁴This footnote explains more precisely how weights were constructed in the $\hat{\alpha}_{\text{KS}}$ runs. We used the following “chained” method for generating weights:

1. KS-reinterpretation is the method of §2.3.1.
2. “Stålmarck” is the generalized-Stålmarck algorithm of [101], *starting with the element obtained via KS-reinterpretation*. The generalized-Stålmarck algorithm successively over-approximates the best transformer from above. By starting the algorithm with the element obtained via KS-reinterpretation, the generalized-Stålmarck algorithm does not have to work its way down from \top ; it merely continues to work its way down from the over-approximation already obtained via KS-reinterpretation.
3. $\hat{\alpha}_{\text{KS}}$ is actually $\hat{\alpha}_{\text{TER}[\text{KS}]}$, from Fig. 3.1(b), which maintains both an under-approximation and a (nontrivial) over-approximation of the desired answer, and hence is resilient to timeouts—i.e., it returns the over-approximation if a timeout occurs. In the chained method for generating weights, $\hat{\alpha}_{\text{TER}[\text{KS}]}$ is *started with the element obtained via the Stålmarck method as an over-approximation* as a way to accelerate its performance.

The generalized-Stålmarck algorithm is a faster algorithm than $\hat{\alpha}_{\text{TER}[\text{KS}]}$, but is not guaranteed to find the best abstract transformer [101]. $\hat{\alpha}_{\text{TER}[\text{KS}]}$ is guaranteed to obtain the best abstract transformer, except for cases in which an SMT solver timeout is reported. The use of KS-reinterpretation accelerates Stålmarck, and the use of Stålmarck accelerates $\hat{\alpha}_{\text{TER}[\text{KS}]}$. Moreover, $\hat{\alpha}_{\text{TER}[\text{KS}]} \sqsubseteq$ KS-reinterpretation is always guaranteed to hold for the weights that are computed.

⁵We use “computed as a geometric mean” as a shorthand for “computed by converting the data to ratios; finding the geometric mean of the ratios; and converting the result back to a percentage”. For instance, suppose that you have improvements of 3%, 17%, 29% (i.e., .03, .17, and .29). The geometric mean of the values .03, .17, and .29 is .113.

Prog. name	Performance of Interprocedural Analysis													
	MOS-reinterp				KS-reinterp				$\hat{\alpha}_{ks}^+$					
	WPDS	post*	query at branch points	WPDS	post*	query at branch points	WPDS	post*	query at branch points	WPDS	post*	query at branch points	t/o	
write	0.088	0.073	0.091	0.58	0.063	0.257	59.167	0.067	2.681	2(1.17%)	97.732	0.066	6.047	2(1.17%)
finger	0.25	0.889	0.275	1.405	0.192	0.366	80.294	0.208	3.61	3(0.78%)	272.939	0.178	6.217	7(1.82%)
subst	0.38	0.895	0.257	2.045	0.412	0.535	164.033	0.427	5.266	2(0.29%)	261.865	0.433	8.775	4(0.57%)
chkdsk	0.472	0.187	0.314	2.577	0.261	0.831	359.949	0.267	61.907	8(0.81%)	349.519	0.275	26.828	11(1.11%)
convert	0.672	1.643	0.7	3.481	0.647	1.091	206.25	0.629	22.409	17(1.43%)	317.142	0.627	30.959	19(1.60%)
route	1.144	3.581	1.413	6.393	0.82	2.073	396.732	0.959	87.864	4(0.32%)	704.816	0.902	75.348	28(2.27%)
comp	0.8	2.916	1.112	4.55	0.738	1.606	599.272	0.835	41.013	6(0.39%)	912.319	0.806	54.349	9(0.59%)
logoff	1.111	4.858	1.936	7.843	1.01	2.908	449.851	1.061	57.276	24(1.58%)	915.493	1.033	141.981	32(2.10%)
setup	2.054	3.259	1.907	13.598	0.521	5.731	1191.37	0.591	277.8	58(2.28%)	1852.79	0.599	464.781	93(3.65%)

Figure 3.7: Performance of WPDS-based interprocedural analysis. The times, in seconds, for WPDS construction, performing interprocedural dataflow analysis (i.e., running post* and performing path-summary) and finding one-vocabulary affine relations at branch instructions, using MOS-reinterpretation, KS-reinterpretation, $\hat{\alpha}_{ks}$, and $\hat{\alpha}_{ks}^+$ to generate weights. The columns labeled “t/o” report the number of WPDS rules for which weight generation timed out during symbolic abstraction.

of the $\hat{\alpha}_{KS}^+$ calls.

The experiment showed that the cost of constructing weights via $\hat{\alpha}_{KS}$ is high, which was to be expected because $\hat{\alpha}_{KS}$ repeatedly calls an SMT solver. Creating KS weights via $\hat{\alpha}_{KS}$ is about 81.5 times slower than creating them via KS-reinterpretation (computed as the geometric mean of the construction-time ratios).

Moreover, creating KS weights via KS-reinterpretation is itself 5.9 times slower than creating MOS weights using MOS-reinterpretation. The latter number is different from the 14.9-fold slowdown reported in §3.5.2 for two reasons: (i) §3.5.2 reported the cost of creating KS and MOS abstract transformers for individual instructions, whereas in Fig. 3.7 the transformers are for basic blocks, and (ii) the WPDS construction times in Fig. 3.7 include the cost of creating merge functions for use at procedure-exit sites, which was about the same for KS-reinterpretation and MOS-reinterpretation.

A comparison of the $\hat{\alpha}_{KS}$ columns of Fig. 3.7 against the $\hat{\alpha}_{KS}^+$ columns reveals that

- Creating KS weights via $\hat{\alpha}_{KS}^+$ is about 1.7 times slower than creating weights via $\hat{\alpha}_{KS}$ (computed as the geometric mean of the construction-time ratios). The slowdown occurs because the formula created for use by $\hat{\alpha}_{KS}^+$ is more complicated than the one created for use by $\hat{\alpha}_{KS}$: the former contains additional conjuncts that capture the effects of memory-access/update and flag-access/update operations.
- The times for performing “post*” and “path summary” are almost the same for both methods, because these phases do not involve any calls to the respective $\hat{\alpha}$ procedures.
- Answering queries at branch points was 1.4 times slower for $\hat{\alpha}_{KS}^+$ compared to $\hat{\alpha}_{KS}$. The reason for the slowdown is that this phase must

Instead, we express the original improvements as ratios and take the geometric mean of 1.03, 1.17, and 1.29, obtaining 1.158. We subtract 1, convert to a percentage, and report “15.8% improvement (computed as a geometric mean)”.

The advantage of this approach is that it handles datasets that include one or more instances of 0% improvement, as well as negative percentage improvements.

Prog. name	Δ_1 rules Rules	WPDS Weights			
		MOS-reinterp < KS-reinterp	KS-reinterp < MOS-reinterp	$\hat{\alpha}_{KS} < \text{KS-reinterp}$	$\hat{\alpha}_{KS}^+ < \hat{\alpha}_{KS}$
write	151	0(0.00%)	0(0.00%)	11(7.28%)	0(0.00%)
finger	353	0(0.00%)	0(0.00%)	29(8.22%)	1(0.28%)
subst	728	0(0.00%)	0(0.00%)	59(8.10%)	0(0.00%)
chkdsk	887	0(0.00%)	0(0.00%)	86(9.70%)	1(0.11%)
convert	1266	0(0.00%)	2(0.16%)	131(10.35%)	0(0.00%)
route	1368	0(0.00%)	3(0.22%)	142(10.38%)	0(0.00%)
comp	1528	0(0.00%)	1(0.07%)	163(10.67%)	0(0.00%)
logoff	1648	0(0.00%)	4(0.24%)	191(11.59%)	1(0.06%)
setup	2847	0(0.00%)	20(0.70%)	432(15.17%)	8(0.28%)

Figure 3.8: Comparison of the precision of the WPDS weights computed using MOS-reinterpretation, KS-reinterpretation, and $\hat{\alpha}_{KS}$. (E.g., KS-reinterp < MOS-reinterp reports the number of rules for which the KS-reinterp weight was more precise than the MOS-reinterp weight.)

call the respective $\hat{\alpha}$ procedures once for each branch point: “post*” and “path summary” return the weight that holds at the *beginning* of a basic block $B = [\iota_1, \dots, \iota_m]$. To obtain the one-vocabulary affine relation that hold just before branch point ι_m at the end of B , we need to perform an additional $\hat{\alpha}$ computation for $[\iota_1, \dots, \iota_{m-1}]$ (i.e., for B , but *without* the branch instruction at the end of B).

Figs. 3.8, 3.9, and 3.10 present three studies that compare the precision obtained via MOS-reinterpretation, KS-reinterpretation, $\hat{\alpha}_{KS}$, and $\hat{\alpha}_{KS}^+$.

Fig. 3.8 compares the precision of the WPDS weights computed by the different methods for each of the example programs. It shows that $\hat{\alpha}_{KS}$ creates strictly more precise weights than KS-reinterpretation for about 10.14% of the WPDS rules (computed as a geometric mean). The “ $\hat{\alpha}_{KS} < \text{KS-reinterp}$ ” column of Fig. 3.8 is particularly interesting in light of the fact that a study of relative precision of abstract transformers created for *individual* instructions via KS-reinterpretation and $\hat{\alpha}_{KS}$ [62, §5.4.1], reported that $\hat{\alpha}_{KS}$ creates strictly more precise transformers than KS-reinterpretation for only about 3.2% of the instructions that occur in the corpus of 19,066

Prog. name	Branches	1-Vocabulary Affine Relations at Branch Points			
		MOS-reinterp < KS-reinterp	KS-reinterp < MOS-reinterp	$\hat{\alpha}_{KS} < \text{KS-reinterp}$	$\hat{\alpha}_{KS}^+ < \hat{\alpha}_{KS}$
write	26	0(0.00%)	0(0.00%)	4(15.38%)	0(0.00%)
finger	48	0(0.00%)	0(0.00%)	14(29.17%)	32(66.67%)
subst	74	1(1.35%)	0(0.00%)	15(20.27%)	0(0.00%)
chkdsk	119	0(0.00%)	0(0.00%)	13(10.92%)	0(0.00%)
convert	161	1(0.62%)	0(0.00%)	49(30.43%)	0(0.00%)
route	243	0(0.00%)	4(1.65%)	63(25.93%)	0(0.00%)
comp	224	0(0.00%)	0(0.00%)	7(3.12%)	0(0.00%)
logoff	306	0(0.00%)	0(0.00%)	91(29.74%)	20(6.54%)
setup	589	0(0.00%)	0(0.00%)	39(6.62%)	0(0.00%)

Figure 3.9: Comparison of the precision of the one-vocabulary affine relations identified to hold at branch points via interprocedural analysis, using weights created using MOS-reinterpretation, KS-reinterpretation, and $\hat{\alpha}_{KS}$. (E.g., KS-reinterp < MOS-reinterp reports the number of branch points at which the KS-reinterp results were more precise than the MOS-reinterp results.)

instructions from §3.5.2. The numbers in Fig. 3.8 differ from that study in two ways: (i) Fig. 3.8 compares the precision of abstract transformers for basic blocks rather than for individual instructions; and (ii) Fig. 3.8 is a comparison for the instructions that appear in specific programs, whereas the corpus of 19,066 instructions used in the per-instruction study from [62, §5.4.1] was created using a randomized instruction generator.

$\hat{\alpha}_{KS}^+$ creates strictly more precise weights than $\hat{\alpha}_{KS}$ for only about 0.1% of the WPDS rules (computed as a geometric mean). Improvements are obtained in only four of the nine programs (finger, chkdsk, logoff, and setup).

Figs. 3.9 and 3.10 answer question (3) posed at the beginning of this section:

For what percentage of program points does $\hat{\alpha}_{KS}$ produce more precise answers than KS-reinterpretation and MOS-reinterpretation?

Prog. name	Procs	2-Vocabulary Procedure Summaries			
		MOS-reinterp < KS-reinterp	KS-reinterp < MOS-reinterp	$\hat{\alpha}_{KS} < \text{KS-reinterp}$	$\hat{\alpha}_{KS}^+ < \hat{\alpha}_{KS}$
write	10	0(0.00%)	0(0.00%)	5(50.00%)	0(0.00%)
finger	18	0(0.00%)	0(0.00%)	10(55.56%)	2(11.11%)
subst	16	0(0.00%)	0(0.00%)	6(37.50%)	0(0.00%)
chkdsk	18	0(0.00%)	0(0.00%)	9(50.00%)	0(0.00%)
convert	38	0(0.00%)	0(0.00%)	8(21.05%)	0(0.00%)
route	40	0(0.00%)	0(0.00%)	18(45.00%)	0(0.00%)
comp	35	1(2.86%)	0(0.00%)	13(37.14%)	0(0.00%)
logoff	46	0(0.00%)	0(0.00%)	14(30.43%)	1(2.17%)
setup	67	0(0.00%)	1(1.49%)	40(59.70%)	0(0.00%)

Figure 3.10: Comparison of the precision of the two-vocabulary affine relations identified to hold at procedure-exit points via interprocedural analysis, using weights created using MOS-reinterpretation, KS-reinterpretation, and $\hat{\alpha}_{KS}$. (E.g., KS-reinterp < MOS-reinterp reports the number of procedure-exit points at which the KS-reinterp results were more precise than the MOS-reinterp results.)

Figs. 3.9 and 3.10 summarize the results obtained from comparing the precision of the affine relations identified via interprocedural analysis using the different weights.⁶

Compared to runs based on either KS-reinterpretation or MOS-reinterpretation, the analysis runs based on $\hat{\alpha}_{KS}$ weights identified more precise affine relations at a substantial number of points (for both one-vocabulary affine relations that hold at branch points—Fig. 3.9, col. 5—and two-vocabulary affine relations that hold at procedure-exit points—Fig. 3.10, col. 5). For one-vocabulary affine relations, the $\hat{\alpha}_{KS}$ analysis results are strictly better than the KS-reinterpretation results at 18.6% of all branch points (computed as a geometric mean). For two-vocabulary

⁶Register `eip` is the x86 instruction pointer. There are some situations that cause the MOS-reinterpretation weights and KS-reinterpretation weights to fail to capture the value of the post-state `eip` value. Therefore, before comparing affine relations, we performed `havoc(eip')`. This adjustment avoids biasing the results merely because of trivial affine relations of the form “`eip' = constant`”.

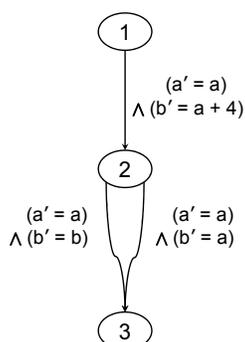


Figure 3.11: Simplified version of an example that caused KS results to be less precise than MOS results, due to compose not distributing over join in the KS domain.

affine relations describing procedure summaries, the $\hat{\alpha}_{\text{KS}}$ analysis results are strictly better than the KS-reinterpretation results at 42% of all procedures (computed as a geometric mean).

For one-vocabulary affine relations, the $\hat{\alpha}_{\text{KS}}^+$ analysis results are strictly better than the $\hat{\alpha}_{\text{KS}}$ results at 7.3% of all branch points (computed as a geometric mean). For two-vocabulary affine relations describing procedure summaries, the $\hat{\alpha}_{\text{KS}}^+$ analysis results are strictly better than the $\hat{\alpha}_{\text{KS}}$ results at 1.4% of all procedures (computed as a geometric mean). However, in both cases improvements are obtained in only two of the nine programs (finger and logoff).

3.5.4 Imprecision Due to Non-Distributivity of KS

Fig. 3.9 shows that in a couple of cases, one in subst and the other in convert, the MOS-reinterpretation results were better than the KS-reinterpretation results. (Although not shown in Fig. 3.9, the MOS-reinterpretation results were also better than the $\hat{\alpha}_{\text{KS}}$ results in the two cases.) We examined these cases, and found that they were an artifact of (i) the evaluation order chosen, and (ii) compose failing to distribute over

join in the KS domain (see §3.4.4).

Fig. 3.11 is a simplified version of the actual transformers in `subst` and `convert` that caused the KS-based analyses to return a less-precise element than the MOS-based analysis. In particular, if the join of the transformers on the two edges from node 2 to node 3 is performed before the composition of the individual $2 \rightarrow 3$ transformers with the $1 \rightarrow 2$ transformer, the combined $2 \rightarrow 3$ KS transformer is “ $a' = a$ ” (i.e., b and b' are unconstrained). The loss of information about b and b' cascades to the $1 \rightarrow 3$ KS transformer, which is also “ $a' = a$ ”. In particular, it fails to contain the conjunct $2^{30}b' = 2^{30}a$, which expresses that the two low-order bits of b' at node 3 are the same as the two low-order bits of a at node 1.

In contrast, in the MOS domain, the combined $2 \rightarrow 3$ transformer is the affine closure of the transformers “ $a' = a \wedge b' = b$ ” and “ $a' = a \wedge b' = a$ ”, which avoids the complete loss of information about b and b' , and hence the $1 \rightarrow 3$ MOS transformer is able to capture the relation “ $a' = a \wedge 2^{30}b' = 2^{30}a$ ”.

3.6 Related Work

3.6.1 Abstract Domains for Affine-Relation Analysis

The original work on affine-relation analysis (ARA) was an intraprocedural ARA algorithm due to [49]. In Karr’s work, a domain element represents a set of points that satisfy affine relations over variables that hold elements of a *field*. Karr’s algorithms are based on linear algebra (i.e., vector spaces).

[75] gave an algorithm for interprocedural ARA, again for vector spaces over a field. Later 2005, 2007, they generalized their techniques to work for modular arithmetic: they introduced the MOS domain, in which an element represents an affine-closed set of affine transformers over variables that hold machine integers, and gave an algorithm for interprocedural

ARA. The algorithms for operations of the MOS domain are based on an extension of linear algebra to modules over a ring.

The version of the KS domain presented in this paper was inspired by, but is somewhat different from, the techniques described in two papers by [52, 53]. Our goals and theirs are similar, namely, to be able to create abstract transformers automatically that are bit-precise, modulo the inherent limitation on precision that stems from having to work with affine-closed sets of values. Compared with their work, we avoid the use of bit-blasting, and work directly with representations of w -bit affine-closed sets. The greatly reduced number of variables that comes from working at word level opened up the possibility of applying our methods to much larger problems, and as discussed in §3.4 and §3.5, we were able to apply our methods to interprocedural program analysis.

As shown in §3.4.2, the algorithm for projection given by [52, §3] does not always find answers that are as precise as the domain is capable of representing. One consequence is that their join algorithm does not always find the least upper bound of its two arguments. In this paper, these issues have been corrected by employing the Howell form of matrices to normalize KS elements (§3.1.2, §3.4.2, and §3.4.4; see also §3.6.2 below).

King and Søndergaard introduced another interesting technique that we did not explore, which is to use affine relations over m -bit numbers, for $m > 1$, to represent sets of 1-bit numbers. To make this approach sensible, their concretization function intersects the “natural” concretization, which yields an affine-closed set of tuples of m -bit numbers, with the set $\{\langle v_1, \dots, v_k \rangle \mid v_1, \dots, v_k \in \{0, 1\}\}$. In essence, this approach restricts the concretization to tuples of 1-bit numbers [53, Defn. 2]. The advantage of the approach is that KS elements over $\mathbb{Z}_{2^m}^k$ can then represent sets of 1-bit numbers that can only be over-approximated using KS elements over \mathbb{Z}_2^k .

Example 3.20. *Suppose that we have three variables $\{v_1, v_2, v_3\}$, and want to represent the set of assignments $\{\langle 001 \rangle, \langle 010 \rangle, \langle 100 \rangle\}$. The best KS element over*

\mathbb{Z}_2^3 is

$$\left[\begin{array}{ccc|c} v_1 & v_2 & v_3 & 1 \\ 1 & 1 & 1 & 1 \end{array} \right], \quad (3.9)$$

which corresponds to the affine relation $v_1 + v_2 + v_3 + 1 = 0$. The set of satisfying assignments is $\{\langle 001 \rangle, \langle 010 \rangle, \langle 100 \rangle, \langle 111 \rangle\}$, which includes the extra tuple $\langle 111 \rangle$.

Now consider what sets can be represented when $m = 2$, so that arithmetic is performed mod 4. In particular, instead of the matrix in Eqn. (3.9) we use the following KS element over \mathbb{Z}_4^3 :

$$\left[\begin{array}{ccc|c} v_1 & v_2 & v_3 & 1 \\ 1 & 1 & 1 & 3 \end{array} \right], \quad (3.10)$$

which corresponds to the affine relation $v_1 + v_2 + v_3 + 3 = 0$. The matrix in Eqn. (3.10) has sixteen satisfying assignments:

$$\begin{array}{cccccccc} \langle 001 \rangle & \langle 010 \rangle & \langle 100 \rangle & \langle 122 \rangle & \langle 212 \rangle & \langle 221 \rangle & \langle 032 \rangle & \langle 023 \rangle \\ \langle 203 \rangle & \langle 230 \rangle & \langle 302 \rangle & \langle 320 \rangle & \langle 113 \rangle & \langle 131 \rangle & \langle 311 \rangle & \langle 333 \rangle \end{array}$$

However, only three of the assignments are in the restricted concretization, namely, the desired set $\{\langle 001 \rangle, \langle 010 \rangle, \langle 100 \rangle\}$. \square

To represent sets of tuples of w -bit numbers, as considered in this paper, the analogous technique would use a y -bit KS domain, $y > w$, in a similar fashion. That is, the “natural” concretization would be intersected with the set $\{\langle v_1, \dots, v_k \rangle \mid v_1, \dots, v_k \in \{0, 1, \dots, 2^w - 1\}\}$. We leave the exploration of these issues for possible future research.

Like some other relational domains, including polyhedra [23, 3, 45] and grids [2], KS/AG fits the dual-representation paradigm of having both a constraint representation (KS) and a generator representation (AG). MOS is based on a generator representation. Whereas many implemen-

tations of domains with a dual representation perform some operations in one representation and other operations in the other representation, converting between representations as necessary, one of the clever aspects of both MOS and KS is that they avoid the need to convert between representations.

[39] describes congruence lattices, where the lattice elements are cosets in the group \mathbb{Z}^k . The generator form for a congruence-lattice element is defined by a point and a basis. The basis is used to describe the coset. The corresponding constraint form for a domain element is a system of Diophantine linear congruence equations. Conversion from the “normalized representation” (generator form) to an equation system is done by an elimination algorithm. The reverse direction is carried out by solving a set of equations. Granger’s normalized representation can be used as a domain for representing affine relations over machine integers. However, Granger’s approach does not have unique normalized representations, because a coset space can have multiple bases. As a result, his method for checking that two domain elements are equal is to check containment in both directions (i.e., to perform two containment checks). Moreover, checking containment is costly because a representation conversion is required: one has to compare the cosets, which involves converting one coset into equational form and then checking if the other coset (in generator form) satisfies the constraints of equational form. In contrast, for the KS domain, one can easily check containment using the KS meet and equality operations:

$$\gamma(X) \subseteq \gamma(Y) \text{ iff } X \sqsubseteq Y \text{ iff } X = (X \sqcap Y).$$

Similarly, containment checking in AG can be performed using the AG join and equality operations. Thus, in both KS and AG, the costly step of converting between generator form and constraint form (or vice versa) is avoided.

Gulwani and Necula introduced the technique of random interpretation (for vector spaces over a field), and applied it to identifying both intraprocedural 2003 and interprocedural 2005 affine relations. The fact that random interpretation involves collecting samples—which are similar to rows of AG elements—suggests that the AG domain might be used as an efficient abstract datatype for storing and manipulating data during random interpretation. Because the AG domain is equivalent to the KS domain (see §3.2), the KS domain would be an alternative abstract datatype for storing and manipulating data during random interpretation.

3.6.2 Howell Form

In contrast with both the Müller-Olm/Seidl work and the King/Søndergaard work, our work takes advantage of the Howell form of matrices. Howell form can be used with each of the domains KS, AG, and MOS defined in §3.1. Because Howell form is canonical for non-empty sets of basis vectors, it provides a way to test pairs of elements for equality of their concretizations—an operation needed by analysis algorithms to determine when a fixed point is reached. In contrast, [78, §2] and [52, Fig. 1], [53, Fig. 2] use “echelon form” (called “triangular form” by King and Søndergaard), which is not canonical.

The algorithms given by Müller-Olm and Seidl avoid computing multiplicative inverses, which are needed to put a matrix in Howell form (line 8 of Alg. 3). However, their preference for algorithms that avoid inverses was originally motivated by the fact that at the time of their original 2005 work they were unaware [77] of Warren’s $O(\log w)$ algorithms [102, §10-15] for computing the inverse of an odd element, and only knew of an $O(w)$ algorithm [76, Lemma 1].

3.6.3 Symbolic Abstraction for Affine-Relation Analysis

[52], 2010 defined the KS domain, and used it to create implementations of best KS transformers for the individual bits of a bit-blasted concrete semantics. They used bit-blasting to express a bit-precise concrete semantics for a statement or basic block. The use of bit-blasting let them track the effect of non-linear bit-twiddling operations, such as shift and xor.

In this paper, we also work with a bit-precise concrete semantics; however, we avoid the need for bit-blasting by working with QFBV formulas expressed in terms of word-level operations; such formulas also capture the precise bit-level semantics of each instruction or basic block. We take advantage of the ability of an SMT solver to decide the satisfiability of such formulas, and use $\hat{\alpha}_{KS}$ to create best word-level transformers.

Prior to our SAS 2011 paper [29], it was not known how to perform $\tilde{\alpha}_{MOS}(\varphi)$ in a non-trivial fashion (other than defining $\tilde{\alpha}_{MOS}$ to be $\lambda f. \top$). The fact that [53, Fig. 2] had been able to devise an algorithm for $\hat{\alpha}_{KS}$ caused us to look more closely at the relationship between MOS and KS. The results presented in §3.3.1 establish that MOS and KS are different, incomparable abstract domains. We were able to give sound interconversion methods (§3.3.2–§3.3.4), and thereby obtained a method for performing $\tilde{\alpha}_{MOS}(\varphi)$ (§3.3.5).

3.7 Chapter Notes

A more detailed discussion of this project can be found in our journal paper [30]. Matt Elder and Thomas Reps developed the improvement to the symbolic abstraction algorithm (§3.4.8). They also designed and implemented the KS reinterpretation. A detailed discussion of KS reinterpretation and its associated implementation can be found in [30, §6]. Junghee Lim and Thomas Reps developed the TSL infrastructure [62] on which our implementation is based on. They also implemented an earlier

version of the MOS domain, which I modified so that Howell matrices are now used as the underlying representation. Junghee Lim, Tycho Andersen, and Thomas Reps helped me in the setup and implementation of the analysis infrastructure used in the experiments.

4 A NEW ABSTRACTION FRAMEWORK FOR AFFINE TRANSFORMERS

This chapter addresses the problem of abstracting a set of affine transformers $\vec{v}' = \vec{v} \cdot C + \vec{d}$, where \vec{v} and \vec{v}' represent the pre-state and post-state, respectively. C is the linear component of the transformation and \vec{d} is a constant vector. For example, $[x' \ y'] = [x \ y] \begin{bmatrix} 1 & 0 \\ 2 & 0 \end{bmatrix} + [10 \ 0]$ denotes the affine transformation $(x' = x + 2y + 10 \wedge y' = 0)$ over variables $\{x, y\}$. We denote an affine transformation by $C : \vec{d}$.

We introduce a framework to harness any base abstract domain \mathcal{B} in an abstract domain of affine transformations. Abstract domains are usually used to define constraints on the variables of a program. In this work, however, abstract domain \mathcal{B} is re-purposed to constrain the elements of C and \vec{d} —thereby defining a set of affine transformers on program states. This framework facilitates intra- and interprocedural analyses to obtain function and loop summaries, as well as to prove program assertions.

We provide analysis techniques to abstract the behavior of the program as a set of affine transformations over bit-vectors. This work is based on the following observation:

Observation 1. *Abstract domains are usually used to define constraints on the variables of a program. However, they can be re-purposed to constrain the elements of $C : \vec{d}$ —thereby defining a set of affine transformers on program states.*

Chapter 3 compared two abstract domains for affine-relation analysis over bitvectors: (i) an affine-closed abstraction of relations over program variables (AG/KS), and (ii) an affine-closed abstraction of affine transformers over program variables (MOS). We observed that the MOS domain can encode two-vocabulary relations that are not affine-closed even though the affine transformers themselves are affine closed. (See §3.3.1 for an

example.) Thus, moving the abstraction from affine relations over program variables to affine relations over affine transformations possibly offers some advantages because it allows some non-affine-closed sets to be representable.

Problem Statement. Our goal is to generalize the ideas used in the MOS domain—in particular, to have an abstraction of *sets of affine transformers*—but to provide a way for a client of the abstract domain to have some control over the performance/precision trade-off. Toward this end, we define a new family of numerical abstract domains, denoted by $\text{ATA}[\mathcal{B}]$. (ATA stands for Affine-Transformers Abstraction.) Following Obs. 1, $\text{ATA}[\mathcal{B}]$ is parameterized by a base numerical abstract domain \mathcal{B} , and allows one to represent a set of affine transformers (or, alternatively, certain disjunctions of transition formulas).

Summary of the Approach. Let the $(k + k^2)$ -tuple $(d_1, d_2, \dots, d_k, c_{11}, c_{12}, \dots, c_{1k}, c_{21}, c_{22}, \dots, c_{kk})$ denote the affine transformation $\bigwedge_{j=1}^k \left(v'_j = \sum_{i=1}^k (c_{ij}v_i) + d_j \right)$, also written as “ $C : \vec{d}$.” The key idea is that we will use $(k + k^2)$ symbolic constants to represent the $(k + k^2)$ coefficients in a transformation of the form $C : \vec{d}$, and use a base abstract domain \mathcal{B} —provided by a client of the framework—to represent *sets* of possible values for these symbolic constants. In particular, \mathcal{B} is an abstract domain for which, for all $b \in \mathcal{B}$, $\gamma(b)$ is a set of $(k + k^2)$ -tuples—each tuple of which provides values for $\{d_i\} \cup \{c_{ij}\}$, and can thus be interpreted as an affine transformation $C : \vec{d}$.

With this approach, a given $b \in \mathcal{B}$ represents the disjunction $\bigvee \{(C : \vec{d}) \in \gamma(b)\}$. When \mathcal{B} is a non-relational domain, each $b \in \mathcal{B}$ constrains the values of $\{d_i\} \cup \{c_{ij}\}$ *independently*. When \mathcal{B} is a relational domain, each $b \in \mathcal{B}$ can impose *intra-component* constraints on the allowed tuples $(d_1, d_2, \dots, d_k, c_{11}, c_{12}, \dots, c_{1k}, c_{21}, c_{22}, \dots, c_{kk})$.

$\text{ATA}[\mathcal{B}]$ generalizes the MOS domain, in the sense that the MOS domain is exactly $\text{ATA}[\text{AG}/\text{KS}]$, where AG/KS is a relational abstract domain that captures affine equalities of the form $\sum_i a_i k_i = b$, where $a_i, b \in \mathbb{Z}_{2^w}$ and \mathbb{Z}_{2^w} is the set of w -bit bitvectors (see §3.1.3 and §3.1.4). For instance, an element in $\text{ATA}[\text{AG}]$ can capture the set of affine transformers “ $x' = k_1 * x + k_1 * y + k_2$, where k_1 is odd, k_2 is even, and k_1 is the coefficient of both x and y .” On the other hand, an element in the abstract domain $\text{ATA}[\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+k^2)}]$, where $\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+k^2)}$ is the abstract domain of $(k+k^2)$ -tuples of intervals over bitvectors, can capture a set of affine transformers such as $x' = k_3 * x + k_4 * y + k_5$, where $k_3 \in [0, 1]$, $k_4 \in [2, 2]$, and $k_5 \in [0, 10]$.

This work addresses a wide variety of issues that arise in defining the $\text{ATA}[\mathcal{B}]$ framework, including describing the abstract-domain operations of $\text{ATA}[\mathcal{B}]$ in terms of the abstract-domain operations available in the base domain \mathcal{B} .

Contributions. The overall contribution of our work is the framework $\text{ATA}[\mathcal{B}]$, for which we present

- methods to perform basic abstract-domain operations, such as equality and join.
- a method to perform abstract composition, which is needed to perform abstract interpretation.
- a faster method to perform abstract composition when the base domain is non-relational.

Organization This chapter is organized as follows: §4.1 introduces the terminology used in the chapter; and presents some needed background material. §4.2 demonstrates the framework with the help of an example. §4.3 formally introduces the parameterized abstract domain $\text{ATA}[\mathcal{B}]$. §4.4 provides discussion and related work. Proofs are given in

4.1 Preliminaries

Matrix addition and multiplication are defined as usual, forming a matrix ring. We denote the transpose of a matrix M by M^t . A *one-vocabulary matrix* is a matrix with $k + 1$ columns. A *two-vocabulary matrix* is a matrix with $2k + 1$ columns. In each case, the “+1” is related to the fact that we capture affine rather than linear relations. I_n denotes the $n \times n$ identity matrix. Given a matrix C , we use $C[i, j]$ to refer to the entry at the i -th column and j -th row of C . Given a vector \vec{d} , we use $\vec{d}[j]$ to refer to the j -th entry in \vec{d} .

4.1.1 Affine Programs

$$\begin{aligned}
 \langle \text{Block} \rangle &:: \text{!} : (\langle \text{Stmt} \rangle ;)^* \langle \text{Next} \rangle \\
 \langle \text{Next} \rangle &:: \text{jump } \text{!}; \\
 &\quad | \text{jump } \langle \text{Cond} \rangle ? \text{!}_1 : \text{!}_2 \\
 \langle \text{Cond} \rangle &:: ? | \langle \text{Expr} \rangle \text{Op} \langle \text{Expr} \rangle \\
 \langle \text{Op} \rangle &:: = | \neq | \geq | \leq \\
 \langle \text{Expr} \rangle &:: c_0 + \sum_{i=1}^k c_i * v_i \\
 \langle \text{Stmt} \rangle &:: v_j := \langle \text{Expr} \rangle \\
 &\quad | v_j := ?
 \end{aligned}$$

We borrow the notion of affine programs from [78]. We restrict our affine programs to consist of a single procedure. The statements are restricted to either affine assignments or non-deterministic assignments. The control-flow instruction consists of either an unconditional jump statement, or a conditional jump with an affine equality, an affine disequality, an affine inequality, or unknown guard condition.

4.1.2 Abstract-Domain Operations

Tab. 4.1 lists the abstract-domain operations needed to generate the program abstraction and perform fixpoint analysis for an affine program. Bottom, equality, and join are standard abstract-domain operations. The widen operation is needed for domains with infinite ascending chains to ensure termination. The two operations of the form $\alpha(\text{Stmt})$ perform

Table 4.1: Abstract-domain operations.

Type	Operation	Description	Type	Operation	Description
\mathcal{A}	\perp	<i>bottom element</i>	\mathcal{A}	$\alpha(v_j := ?)$	<i>abstraction for nondeterministic assignments</i>
bool	$(a_1 == a_2)$	<i>equality</i>	\mathcal{A}	$\alpha(v_j := c_0 + \sum_{i=1}^k c_{ij} * v_i)$	<i>abstraction for affine assignments</i>
\mathcal{A}	$(a_1 \sqcup a_2)$	<i>join</i>	\mathcal{A}	$(a_1 \circ a_2)$	<i>composition</i>
\mathcal{A}	$(a_1 \nabla a_2)$	<i>widen</i>			
\mathcal{A}	Id	<i>identity element</i>			

abstraction on an assignment statement Stmt to generate an abstract transformer. Id is the identity element; which represents the identity transformation ($\bigwedge_{i=1}^k v'_i = v_i$). Finally, the abstract-composition operation $a_1 \circ a_2$ returns a sound overapproximation of the composition of the abstract transformation a_1 with the abstract transformation a_2 .

4.1.3 Relating MOS and AG

Since KS and AG are the same domain, we will use AG to refer to KS/AG domain (See §3.2). There are two ways to relate the MOS and AG domains. One way is to use them as abstractions of two-vocabulary relations and provide (approximate) inter-conversion methods. The other is to use a variant of the AG domain to represent the elements of the MOS domain exactly.

Comparison of MOS and AG elements as abstraction of two-vocabulary relations.

As shown in §3.3.1, the MOS and AG domains are incomparable: some relations are expressible in each domain that are not expressible in the other. Intuitively, the central difference is that MOS is a domain of sets of *functions*, while AG is a domain of *relations*.

Soundly converting an MOS element M to an overapproximating AG element is equivalent to stating two-vocabulary affine constraints satisfied by M §3.3.2.

Table 4.2: Example demonstrating two ways of relating MOS and AG.

MOS element (M)	Overapproximating AG element (A_1)	Reformulation as abstraction over affine transformers (A_2)
$\left\{ \begin{bmatrix} 1 & x & y \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 1 & x & y \\ 0 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix} \right\}$	$\begin{bmatrix} 1 & x & y & x' & y' \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & a_{01} & a_{02} & a_{10} & a_{11} & a_{12} & a_{20} & a_{21} & a_{22} \\ 1 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$

Reformulation of MOS elements as AG elements.

An MOS element $M = \{M_1, M_2, \dots, M_n\}$ represents the set of $(k+1) \times (k+1)$ matrices in the affine closure of the matrices in M . Each matrix can be thought of as a $(k+1) \times (k+1)$ vector, and hence M can be represented by an AG element of size $n \times ((k+1) \times (k+1))$.

Example 4.1. Tab. 4.2 shows the two ways MOS and AG elements can be related. Column 1 shows the MOS element M , which represents the set of matrices in the affine closure of the two $(k+1) \times (k+1)$ matrices, with $k = 2$. The second column gives the AG element A_1 (a matrix with $2k+1$ columns) representing the affine-closed space over $\{x, y, x', y'\}$ satisfied by M . Consequently, $\gamma_{AG}(A_1) \supseteq \gamma_{MOS}(M)$. Column 3 shows the two matrices of M as the $2 \times ((k+1) \times (k+1))$ AG element A_2 . Because A_2 is just a reformulation of M , $\gamma_{AG}(A_2) = \gamma_{MOS}(M)$. \square

4.2 Overview

In this section, we motivate and illustrate the ATA[\mathcal{B}] framework, with the help of several examples. The first two examples illustrate the following principle, which restates Obs. 1 more formally:

Observation 2. Each affine transformation $C : \vec{d}$ in a set of affine transformations involves $(k+1)^2$ coefficients $\in \mathbb{Z}_{2^w} : (1, d_1, d_2, \dots, d_k, 0, c_{11}, c_{12}, \dots, 0, c_{21}, \dots, c_{kk})$.¹ Thus, we may use any abstract domain whose elements concretize

¹ k of the coefficients are always 0, and one coefficient is always 1 (i.e., the first column is always $(1 \mid 0 \ 0 \ \dots \ 0)^t$). For this reason, we really need only $k + k^2$ elements, but we will sometimes refer to $(k+1)^2$ elements for brevity.

to subsets of $\mathbb{Z}_{2^w}^{(k+1)^2}$ as a method for representing a set of affine transformers. \square

Example 4.2. The AG element A_2 in column 3 of Tab. 4.2 illustrates how an AG element with $(k+1)^2$ columns represents the same set of affine transformers as the MOS element M shown in column 1. For instance, the first row of A_2 represents the first matrix in M . \square

Example 4.3. Consider the element $E = ([1, 1], [0, 10], [0, 0], [0, 0], [1, 1], [2, 3], [0, 0], [0, 0], [1, 1])$ of $\mathcal{J}_{\mathbb{Z}_{2^w}}^9$. E can be depicted more mnemonically as the following matrix:

$$\left[\begin{array}{c|cc} & x & y \\ \hline 1 & [1, 1] & [0, 10] & [0, 0] \\ \hline & [0, 0] & [1, 1] & [2, 3] \\ \hline & [0, 0] & [0, 0] & [1, 1] \end{array} \right],$$

where every element in E is an interval ($\mathcal{J}_{\mathbb{Z}_{2^w}}$). E represents the point set $\{(x', y', x, y) : \exists i_1, i_2 \in \mathbb{Z}_{2^w} : x' = x + i_1 \wedge y' = i_2 x + y \wedge 0 \leq i_1 \leq 10 \wedge 2 \leq i_2 \leq 3\}$. \square

Examples 4.2 and 4.3 both exploit Observation 2, but use different abstract domains. Ex. 4.2 uses the AG domain with $(k+1)^2$ columns, whereas Ex. 4.3 uses the domain $\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}$. In particular, an abstract-domain element in our framework $\text{ATA}[\mathcal{B}]$ is a set of affine transformations $\vec{v}' = \vec{v} \cdot C + \vec{d}$, such that the allowed coefficients in the matrix C and the vector \vec{d} are abstracted by a base abstract domain \mathcal{B} .

The remainder of this section shows how different instantiations of Observation 2 allow different properties of a program to be recovered.

Example 4.4. In this example, the variable r of function f is initialized to 0 and conditionally incremented by $2x$ inside a loop with 10 iterations.

```

ENT: int f(int x) {
L0:   int i = 0, r = 0;
L1:   while(i <= 10) {
L2:     if(*)
L3:       r = r + 2*x;
L4:     i = i + 1;
      }
L5:   return r;
      }

```

The exact function summary for function f , denoted by S_f , is $(\exists k. r' = 2kx \wedge 0 \leq k \leq 10)$. Note that S_f expresses two important properties of the function: (i) the return value r' is an even multiple of x , and (ii) the multiplicative factor is contained in an interval. \square

$\mathcal{B} = \mathbf{AG}$ with $(k + 1)^2$ columns: Fig. 4.1(a) shows the abstract transformers generated with the MOS domain.² Each matrix of the form $\left[\begin{array}{c|ccc} 1 & d_1 & d_2 & d_3 \\ \hline 0 & c_{11} & c_{12} & c_{13} \\ 0 & c_{21} & c_{22} & c_{23} \\ 0 & c_{31} & c_{32} & c_{33} \end{array} \right]$ represents the state transformation $(x' = d_1 + c_{11}x + c_{21}i + c_{31}r) \wedge (i' = d_2 + c_{12}x + c_{22}i + c_{32}r) \wedge (r' = d_3 + c_{13}x + c_{23}i + c_{33}r)$.

For instance, the abstract transformer for $L3 \rightarrow L4$ is an MOS-domain element with a single matrix that represents the affine transformation: $(x' = x) \wedge (i' = i) \wedge (r' = 2x + r)$. The edges absent from Fig. 4.1(a), e.g., $L1 \rightarrow L2$, have the identity MOS-domain element.

Edge	Transformer
$L0 \rightarrow L1$	$\left\{ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \right\}$
$L3 \rightarrow L4$	$\left\{ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \right\}$
$L4 \rightarrow L1$	$\left\{ \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \right\}$

(a)

Iteration	Node L1
(i)	$\left\{ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \right\}$
(ii)	$\left\{ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \right\}$
(iii)	$\left\{ \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 & 2 & 0 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \right\}$

(b)

Figure 4.1: Abstract transformers and snapshots in the fixpoint analysis with the MOS domain for Ex. 4.4.

To obtain function summaries, an iterative fixed-point computation needs to be performed. An abstract-domain element a is a *summary* at some program point L , if it describes a two-vocabulary transition relation that overapproximates the (compound) transition relation from the beginning of the function to program point L . Fig. 4.1(b) provides the iteration results for the summary at the program point $L1$. After iteration (i), the result represents $(x' = x) \wedge (i' = 0) \wedge (r' = 0)$. After iteration (ii), it adds the affine transformer $(x' = x) \wedge (i' = 1) \wedge (r' = 2x)$ to the summary. Quiescence is discovered on the third iteration because the affine-closure

²We will continue to refer to the MOS domain directly, rather than “the instantiation of Observation 2 with an AG element containing $(k + 1)^2$ columns” (à la Ex. 4.2).

of the three matrices is the same as the affine-closure of the two matrices after the second iteration. As a result, the function summary that MOS learns, denoted by S_{MOS} , is $\exists k.r' = 2kx$, which is an overapproximation of the exact function summary S_f . Imprecision occurs because the MOS-domain is not able to represent inequality guards. Hence, the summary captures the evenness property, but not the bounds property.

$\mathcal{B} = \mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}$: By using different \mathcal{B} s, an analyzer will be able to recover different properties of a program. Now consider what happens when the program above is analyzed with $\text{ATA}[\mathcal{B}]$ instantiated with the non-relational base domain of environments of intervals $(\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2})$. The identity trans-

formation for the abstract domain $\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]$ is
$$\begin{bmatrix} 1 & | & [0,0] & [0,0] & [0,0] \\ 0 & | & [1,1] & [0,0] & [0,0] \\ 0 & | & [0,0] & [1,1] & [0,0] \\ 0 & | & [0,0] & [0,0] & [1,1] \end{bmatrix}.$$

The bottom element for the abstract domain $\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]$, denoted by

$$\perp_{\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]} \text{ is } \begin{bmatrix} 1 & | & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} \\ 0 & | & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} \\ 0 & | & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} \\ 0 & | & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} & \perp_{\mathcal{J}_{\mathbb{Z}_2^w}} \end{bmatrix}.^3$$

Fig. 4.2 shows the abstract transformers and the fixpoint analysis for the node L1 with the $\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]$ domain. One advantage of using intervals as the base domain is that they can express inequalities. For instance, the abstract transformer for the edge $L1 \rightarrow L2$ specifies the transformation $(x' = x) \wedge (0 \leq i' \leq 10) \wedge (r' = r)$. Consequently, the function summary that $\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]$ learns, denoted by $S_{\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]}$, is $r' = [0, 20]x$. This summary captures the bounds property, but not the evenness property. Notice that, $S_f = S_{\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]} \wedge S_{\text{MOS}}$.

³The abstract domain $\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}$ is the product domain of $(k+1)^2$ interval domains, that is, $\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2} = \mathcal{J}_{\mathbb{Z}_2^w} \times \mathcal{J}_{\mathbb{Z}_2^w} \times \dots \times \mathcal{J}_{\mathbb{Z}_2^w}$. $\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}$ uses smash product to maintain a canonical representation for $\perp_{\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]}$. Thus, if any of the coefficients in an abstract-domain element $b \in \text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]$ is $\perp_{\mathcal{J}_{\mathbb{Z}_2^w}}$, then b is smashed to $\perp_{\text{ATA}[\mathcal{J}_{\mathbb{Z}_2^w}^{(k+1)^2}]}$.

Edge	Transformer
L0 → L1	$\begin{bmatrix} 1 & [0,0] & [0,0] & [0,0] \\ 0 & [1,1] & [0,0] & [0,0] \\ 0 & [0,0] & [0,0] & [0,0] \\ 0 & [0,0] & [0,0] & [0,0] \end{bmatrix}$
	$\begin{bmatrix} 1 & [0,0] & [0,10] & [0,0] \\ 0 & [1,1] & [0,0] & [0,0] \\ 0 & [0,0] & [0,0] & [0,0] \\ 0 & [0,0] & [0,0] & [1,1] \end{bmatrix}$
	$\begin{bmatrix} 0 & [1,1] & [0,0] & [2,2] \\ 0 & [0,0] & [1,1] & [0,0] \\ 0 & [0,0] & [0,0] & [1,1] \\ 1 & [0,0] & [1,1] & [0,0] \end{bmatrix}$
	$\begin{bmatrix} 0 & [1,1] & [0,0] & [0,0] \\ 0 & [0,0] & [1,1] & [0,0] \\ 0 & [0,0] & [0,0] & [1,1] \\ 0 & [0,0] & [0,0] & [1,1] \end{bmatrix}$

(a)

Iteration	Node L1
(i)	$\begin{bmatrix} 1 & [0,0] & [0,0] & [0,0] \\ 0 & [1,1] & [0,0] & [0,0] \\ 0 & [0,0] & [0,0] & [0,0] \\ 0 & [0,0] & [0,0] & [0,0] \end{bmatrix}$
	$\begin{bmatrix} 1 & [0,0] & [0,1] & [0,0] \\ 0 & [1,1] & [0,0] & [0,2] \\ 0 & [0,0] & [0,0] & [0,0] \\ 0 & [0,0] & [0,0] & [0,0] \end{bmatrix}$
	...
	$\begin{bmatrix} 1 & [0,0] & [0,10] & [0,0] \\ 0 & [1,1] & [0,0] & [0,20] \\ 0 & [0,0] & [0,0] & [0,0] \\ 0 & [0,0] & [0,0] & [0,0] \end{bmatrix}$

(b)

Figure 4.2: Abstract transformers and fixpoint analysis with the $\text{ATA}[\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}]$ domain for Ex. 4.4.

Consider the instantiation of the ATA framework with strided-intervals over bitvectors [83], denoted by $\mathcal{S}\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}$. A strided interval represents a set of the form $\{l, l + s, l + 2s, \dots, l + (n - 1)s\}$. Here, l is the beginning of the interval, s is the stride, and n is the interval size. Consequently, $\text{ATA}[\mathcal{S}\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}]$ learns the function summary $\exists k.r' = kx \wedge k = 2[0, 10]$, which captures both the evenness property and the bounds property. Note that a traditional (non-ATA-framework) analysis based on the strided-interval domain alone would not be able to capture the desired summary because the strided-interval domain is non-relational.

Widening concerns. In principle, abstract domains $\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}$ and $\mathcal{S}\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}$ do not need widening operations because the lattice height is finite. However, the height is exponential in the bitwidth w of the program variables, and thus in practice we need widening operations to speed-up the fixpoint iteration. In the presence of widening, neither $\text{ATA}[\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}]$ nor $\text{ATA}[\mathcal{S}\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}]$ will be able to capture the bounds property for Ex. 4.4, because they are missing relational information between the loop counter i and the variable r . However, the reduced product of $\text{ATA}[\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}]$ (or

$\text{ATA}[\mathcal{S}\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}]$ and MOS can learn the exact function summary.

4.3 Affine-Transformer-Abstraction Framework

In this section, we formally introduce the Affine-Transformer-Abstraction framework (ATA) and describe abstract-domain operations for the framework. We also discuss some specific instantiations.

ATA[\mathcal{B}] Definition. Let C be a k -by- k matrix: $[c_{ij}]$, where each c_{ij} is a symbolic constant for the entry at i -th row and j -th column. Let \vec{d} be a k -vector, $[d_i]$, where each d_i is a symbolic constant for the i -th entry in the vector. As mentioned in Chapter 1, an affine transformer, denoted by $C : \vec{d}$, describes the relation $\vec{v}' = \vec{v} \cdot C + \vec{d}$, where \vec{v}' and \vec{v} are row vectors of size k that represent the post-transformation state and the pre-transformation state, respectively, on program variables.

Given a base abstract domain \mathcal{B} , the ATA framework generates a corresponding abstract domain $\text{ATA}[\mathcal{B}]$ whose elements represent a transition relation between the pre-state and the post-state program vocabulary. Each element $a \in \text{ATA}[\mathcal{B}]$ is represented using an element $\text{base}(a) \in \mathcal{B}$, such that:

$$\gamma(a) = \{(\vec{v}, \vec{v}') \mid \exists (C : \vec{d}) \in \gamma(\text{base}(a)) : \vec{v}' = \vec{v} \cdot C + \vec{d}\}.$$

4.3.1 Abstract-Domain Operations for ATA[\mathcal{B}]

In this subsection, we provide all the abstract-domain operations for $\text{ATA}[\mathcal{B}]$, with the exception of abstract composition, which is discussed in §4.3.2.

In the $\text{ATA}[\mathcal{B}]$ framework, the symbolic constants in the base domain \mathcal{B} are denoted by symbols($C : \vec{d}$), where $\text{symbols}(C : \vec{d}) = (d_1, d_2, \dots, d_n, c_{11}, c_{12}, \dots, c_{1k}, c_{21}, c_{22}, \dots, c_{2k}, \dots, c_{kk})$ is the tuple of $k + k^2$

Table 4.3: Base abstract-domain operations.

Type	Operation	Description
\mathcal{B}	\perp	<i>bottom element</i>
\mathcal{B}	\top	<i>top element</i>
bool	$(b_1 == b_2)$	<i>equality</i>
\mathcal{B}	$(b_1 \sqcup b_2)$	<i>join</i>
\mathcal{B}	$(b_1 \nabla b_2)$	<i>widen</i>
\mathcal{B}	$\text{havoc}(b_1, S)$	<i>remove all constraints on symbolic constants in S</i>
\mathcal{B}	$\alpha(\text{ct})$	<i>abstraction for the concrete affine transformer ct, where $\text{ct} \in \text{symbols}(C : \vec{d}) \rightarrow \mathbb{Z}_{2^w}$</i>

symbolic constants in the affine transformation. Tab. 4.3 lists the abstract-domain interface for the base abstract domain \mathcal{B} needed to implement these operations for $\text{ATA}[\mathcal{B}]$. The first five operations in the interface are standard abstract-domain operations. $\text{havoc}(b_1, S)$ takes an element b_1 and a subset $S \subseteq \text{symbols}(C : \vec{d})$ of symbolic constants, and returns an element without any constraints on the symbolic constants in S . The last operation in Tab. 4.3 defines an abstraction for a concrete affine transformer ct . A concrete affine transformer is a mapping from the symbolic constants in the affine transformer to bitvectors of size w . We represent concrete

state ct with the $(k + 1) \times (k + 1)$ matrix:

$$\begin{bmatrix} 1 & \text{ct}(d_1) & \text{ct}(d_2) & \dots & \text{ct}(d_k) \\ 0 & \text{ct}(c_{11}) & \text{ct}(c_{12}) & \dots & \text{ct}(c_{1k}) \\ 0 & \text{ct}(c_{21}) & \text{ct}(c_{22}) & \dots & \text{ct}(c_{2k}) \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \text{ct}(c_{k1}) & \text{ct}(c_{k2}) & \dots & \text{ct}(c_{kk}) \end{bmatrix},$$

where $\text{ct}(s)$ denotes the concrete value in \mathbb{Z}_{2^w} of symbol s in the concrete state ct .

Tab. 4.4 gives the abstract-domain operations for $\text{ATA}[\mathcal{B}]$ in terms of the base abstract-domain operations in \mathcal{B} . The first operation is the \perp element, which is simply defined as $\perp_{\mathcal{B}}$, the bottom element in the base domain. Similarly, equality, join, and widen operations are defined as the equality, join, and widen operations in the base domain. The equality operation is not the exact equality operation; that is, $(a_1 \widetilde{=} a_2)$ can return false, even if $\gamma(a_1) = \gamma(a_2)$. However, the equality operation is sound; that

Table 4.4: Abstract-domain operations for the $\text{ATA}[\mathcal{B}]$ -domain.

Type	Operation	Description
\mathcal{A}	\perp	$\perp_{\mathcal{B}}$
bool	$(a_1 \equiv a_2)$	$(\text{base}(a_1) == \text{base}(a_2))$
\mathcal{A}	$(a_1 \sqcup a_2)$	$(\text{base}(a_1) \sqcup \text{base}(a_2))$
\mathcal{A}	$(a_1 \nabla a_2)$	$(\text{base}(a_1) \nabla \text{base}(a_2))$
\mathcal{A}	$\alpha(v_j := d_j + \sum_{i=1}^k c_{ij} * v_i)$	$\alpha \left(\left[\begin{array}{c ccc} 1 & 0 & d_j & 0 \\ \hline 0 & I_{j-1} & [c_{1j}, c_{2j}, \dots, c_{(j-1)j}]^t & 0 \\ 0 & 0 & c_{jj} & 0 \\ 0 & 0 & [c_{(j+1)j}, c_{(j+2)j}, \dots, c_{kj}]^t & I_{k-j} \end{array} \right] \right)$
\mathcal{A}	$\alpha(v_j :=?)$	$\text{havoc}(\alpha(I_{k+1}), \{d_j, c_{1j}, c_{2j}, \dots, c_{kj}\})$
\mathcal{A}	Id	$\alpha(I_{k+1})$

is, when $(a_1 \equiv a_2)$ returns true, then $\gamma(a_1) = \gamma(a_2)$. The \sqcup operation for the $\text{ATA}[\mathcal{B}]$ is a quasi-join operation [34]. In other words, the least upper bound does not necessarily exist for $\text{ATA}[\mathcal{B}]$, but a sound upper-bound operation \sqcup is available.

The abstraction operation for the affine-assignment statement $\alpha(v_j := d_0 + \sum_{i=1}^k c_{ij} * v_i)$ gives back an $\text{ATA}[\mathcal{B}]$ -element with a single transformer where every variable $v \in V - \{v_j\}$ is left unchanged and the variable v_j is transformed to reflect the assignment by updating the coefficients of the corresponding column. The abstraction operation for the non-deterministic assignment statement $\alpha(v_j :=?)$ gives back an $\text{ATA}[\mathcal{B}]$ -element, such that every variable $v \in V - \{v_j\}$ is left unchanged but the symbolic constant corresponding to the coefficients in the column j of the affine transformation can be any value. This operation is carried out by performing *havoc* on the identity transformation with respect to the set $\{d_j, c_{1j}, c_{2j}, \dots, c_{kj}\}$ of symbolic constants. The identity transformation *Id* is obtained by abstracting the concrete affine transformer *ct* that represents the identity transformer. We provide proofs of soundness for these abstract-domain operations in App. D.

4.3.2 Abstract Composition

We have shown that all the abstract-domain operations for $\text{ATA}[\mathcal{B}]$ can be implemented in terms of abstract-domain operations in \mathcal{B} , with the exception of abstract composition. Let us consider the composition of two abstract values $\alpha, \alpha' \in \text{ATA}[\mathcal{B}]$, representing the two-vocabulary relations $R[\vec{v}; \vec{v}'] = \gamma(\alpha)$ and $R'[\vec{v}'; \vec{v}'] = \gamma(\alpha')$. An abstract operation $\circ^\#$ is a sound abstract-composition operation if, for all $\alpha'' = \alpha' \circ^\# \alpha$, $\gamma(\alpha'') \supseteq \{ (\vec{v}; \vec{v}'') \mid \exists \vec{v}'. R[\vec{v}; \vec{v}'] \wedge R'[\vec{v}'; \vec{v}'] \}$. This condition translates to:

$$\begin{aligned} \gamma(\text{base}(\alpha'')) \supseteq \{ (\vec{v}, \vec{v}'') \mid \exists (C : \vec{d}) \in \gamma(\text{base}(\alpha)), (C' : \vec{d}') \in \gamma(\text{base}(\alpha')), \\ (4.1) \\ (C'' : \vec{d}'') : (\vec{v}'' = \vec{v} \cdot C'' + \vec{d}'') \wedge (C'' = C \cdot C') \\ \wedge (\vec{d}'' = \vec{d} \cdot C' + \vec{d}') \} \end{aligned}$$

The presence of the quadratic components $C \cdot C'$ and $\vec{d} \cdot C'$ makes the implementation of abstract composition non-trivial. One extremely expensive method to implement abstract composition is to enumerate the set of all concrete transformers $(C : \vec{d}) \in \gamma(\text{base}(\alpha))$ and $(C' : \vec{d}') \in \gamma(\text{base}(\alpha'))$, perform matrix multiplication for each pair of concrete transformers, and perform join over all pairs of them. This approach is impractical because the set of all concrete transformers in an abstract value can be very large.

First, we provide a general method to implement abstract composition. Then, we provide methods for abstract composition when the base domain \mathcal{B} has certain properties, like non-relationality and weak convexity. The latter methods are faster, but are only applicable to certain classes of base abstract domains.

General Case.

We present a general method to perform abstract composition by reducing it to the symbolic-abstraction problem. The *symbolic abstraction* of a formula φ in logic \mathbb{L} , denoted by $\hat{\alpha}(\varphi)$, is the best value in \mathcal{B} that over-approximates the set of all concrete affine transformers $(C : \vec{d})$ that satisfy φ [84, 100]. For all $b \in \mathcal{B}$, the *symbolic concretization* of \mathcal{B} , denoted by $\hat{\gamma}(b)$, maps b to a formula $\hat{\gamma}(b) \in \mathbb{L}$ such that b and $\hat{\gamma}(b)$ represent the same set of concrete affine transformers (i.e., $\gamma(b) = \llbracket \hat{\gamma}(b) \rrbracket$). We expect the base domain \mathcal{B} to provide the $\hat{\gamma}$ operation. In our framework, there are slightly different variants of $\hat{\alpha}$ and $\hat{\gamma}$ according to which vocabulary of symbolic constants are involved. For instance, we use $\hat{\gamma}'$ to denote symbolic concretization in terms of the primed symbolic constants $\text{symbols}(C' : \vec{d}')$. Similarly, $\hat{\alpha}''$ denotes symbolic abstraction in terms of the double-primed symbolic constants $\text{symbols}(C'' : \vec{d}'')$. The function *dropPrimes* shifts the vocabulary of symbolic constants by removing the primes from the symbolic constants that an abstract value represents.

We use $\mathbb{L} = \text{QF_BV}$, i.e., quantifier-free bit-vector logic, to express abstract composition symbolically as follows:

$$\begin{aligned} \text{base}(a'') &= \text{dropPrimes}(\hat{\alpha}''(\varphi)), \text{ where} & (4.2) \\ \varphi &= (C'' = C \cdot C') \wedge (\vec{d}'' = \vec{d} \cdot C' + \vec{d}') \\ &\quad \wedge \hat{\gamma}(\text{base}(a)) \wedge \hat{\gamma}'(\text{base}(a')). \end{aligned}$$

(Note that $\hat{\gamma}(\text{base}(a))$ and $\hat{\gamma}'(\text{base}(a'))$ are formulas over $\text{symbols}(C : \vec{d})$ and $\text{symbols}(C' : \vec{d}')$ respectively.) Past literature [84, 100, 30] provides various algorithms to implement symbolic abstraction. Symbolic-abstraction methods are usually slow because they make repeated calls to an SMT solver. Specifically, the symbolic-abstraction algorithms in [84, 100] require $\mathcal{O}(h'')$ calls to SMT, where h'' is the height of the abstract-domain element — i.e., $\text{base}(a'')$ in the lattice \mathcal{B} .

Alg. 5 is a variant of the symbolic-abstraction algorithm from [84]. Alg. 5 needs a method to enumerate a generator set gs for each $b \in \mathcal{B}$. Such a set can easily be obtained from the generator representation of \mathcal{B} . For instance, each row in an AG element is an affine transformer, and a generator set for the AG element is the set of all rows in the AG matrix: the affine combination of the rows generate the concrete affine transformers that the AG element represents. Note that the generator set for an abstract value b is usually much smaller than the set of all affine transformers in b . For the AG domain, the generating set is worst-case polynomial size, whereas the set of all affine transformers is worst-case exponential in the number of variables k .

In Alg. 5, line 3 initializes the value *lower* to the product of each pair of abstract transformers. The product $t \times t'$, where $t = \left[\begin{array}{c|c} 1 & \vec{d} \\ \hline 0 & C \end{array} \right]$ and $t' = \left[\begin{array}{c|c} 1 & \vec{d}' \\ \hline 0 & C' \end{array} \right]$ is $\left[\begin{array}{c|c} 1 & \vec{d} \cdot C' + \vec{d}' \\ \hline 0 & C \cdot C' \end{array} \right]$. Because *lower* is initialized to $\{t \times t' \mid t \in gs_1, t' \in gs_2\}$ rather than \perp , the number of SMT calls in the symbolic abstraction is significantly reduced, compared to the algorithm from [84]. The function *GetModel*, used at line 5, returns the model $M \in \text{symbols}(C'' : \vec{d}'') \rightarrow \mathbb{Z}_{2^w}$ satisfying the formula $(\varphi \wedge \neg \hat{\gamma}(\text{lower}))$ given to the SMT solver at line 4. Thus, the model M is a concrete affine transformer in α'' . The representation function β , used at line 6, maps a singleton model M to the least value in \mathcal{B} that overapproximates $\{M\}$ [84]. While the SMT call at line 4 is satisfiable, the loop keeps improving the value of *lower* by adding the satisfying model M to *lower* via the representation function β and the join operation. When line 4 is unsatisfiable, the loop terminates and returns *lower*. This method is sound because the unsatisfiable call proves that $\varphi \Rightarrow \hat{\gamma}(\text{lower})$. The loop terminates when the height of the base domain \mathcal{B} is finite.

Algorithm 5 Abstract Composition via Symbolic Abstraction

```

1:  $gs_1 \leftarrow \{t_1, t_2, \dots, t_{l_1}\}$   $\triangleright$  where  $\text{base}(a) = \bigsqcup_{i=0}^{l_1} t_i$ 
2:  $gs_2 \leftarrow \{t'_1, t'_2, \dots, t'_{l_2}\}$   $\triangleright$  where  $\text{base}(a') = \bigsqcup_{i=0}^{l_2} t'_i$ 
3:  $lower \leftarrow \{t \times t' \mid t \in gs_1, t' \in gs_2\}$ 
4: while  $r \leftarrow \text{SMTCall}(\varphi \wedge \neg \hat{\gamma}(lower))$  is Sat do
5:    $M \leftarrow \text{GetModel}(r)$ 
6:    $lower \leftarrow lower \sqcup \beta(M)$ 
7: return  $lower$ 

```

Non-relational base domains.

In this section, we present a method to implement abstract composition for $\text{ATA}[\mathcal{B}]$, when \mathcal{B} is non-relational. We focus on the non-relational case separately because it allows us to implement a sound abstract-composition operation efficiently.

Foundation domain. Each element in the non-relational domain \mathcal{B} is a mapping from symbols S to a subset of \mathbb{Z}_{2^w} . We introduce the concept of a *foundation domain*, denoted by $\mathcal{F}_{\mathcal{B}}$, to represent the abstractions of subsets of \mathbb{Z}_{2^w} in the base abstract-domain elements. We can define a non-relational base domain in terms of the foundation domain as follows: $\mathcal{B} \stackrel{\text{def}}{=} S \rightarrow \mathcal{F}_{\mathcal{B}}$. For instance, the non-relational domain of intervals $\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}$ can be represented by $S \rightarrow \mathcal{J}_{\mathbb{Z}_{2^w}}$, where $\mathcal{J}_{\mathbb{Z}_{2^w}}$ represents the interval lattice over \mathbb{Z}_{2^w} , and S is a set of $(k+1)^2$ symbolic constants that represent the coefficients of an affine transformer.

A foundation domain \mathcal{F} is a lattice whose elements concretize to subsets of \mathbb{Z}_{2^w} . Tab. 4.5 present the foundation-domain operations for \mathcal{F} . Bottom, equality, join, widen, and $\alpha(\text{bv})$ are standard abstract-domain operations. The abstract addition and multiplication operations provide a sound reinterpretation of the collecting semantics of concrete addition and multiplication. For instance, with the interval foundation domain, $[0, 7] +^\# [-3, 17] = [-3, 24]$ and $[0, 6] \times^\# [-3, 3] = [-18, 18]$.

Table 4.5: Foundation-domain operations.

Type	Operation	Description	Type	Operation	Description
\mathcal{F}	\perp	<i>empty set</i>	\mathcal{F}	$\alpha(\text{bv})$	<i>abstraction for the bitvector value $\text{bv} \in \mathbb{Z}_{2^w}$</i>
bool	$(f_1 == f_2)$	<i>equality</i>	\mathcal{F}	$(f_1 +^\# f_2)$	<i>abstract addition</i>
\mathcal{F}	$(f_1 \sqcup f_2)$	<i>join</i>	\mathcal{F}	$(f_1 \times^\# f_2)$	<i>abstract multiplication</i>
\mathcal{F}	$(f_1 \nabla f_2)$	<i>widen</i>			

Abstract composition for a non-relational domain is defined as follows:

$$\begin{aligned}
a' \circ_{\text{NR}} a = & \left\{ (\vec{v}, \vec{v}') \mid \exists (C : \vec{d}) : (\vec{v}' = \vec{v} \cdot C + \vec{d}) \wedge b \in (\text{symbols}(C : \vec{d}) \rightarrow \mathcal{F}) \right. \\
& \wedge \left(\bigwedge_{1 \leq i, j \leq k} (b[c_{ij}] = \sum_{1 \leq l \leq k}^\# (\text{base}(a)[c_{il}] \times^\# \text{base}(a')[c_{lj}])) \right) \\
& \left. \wedge \left(\bigwedge_{1 \leq j \leq k} b[d_j] = \sum_{1 \leq l \leq k}^\# (\text{base}(a)[d_l] \times^\# \text{base}(a')[c_{lj}] +^\# \text{base}(a')[d_j]) \right) \right\}.
\end{aligned} \tag{4.3}$$

The term $b[s]$, where $b \in \mathcal{B}$ and $s \in \text{symbols}(C : \vec{d})$, refers to the element in the foundation domain $f \in \mathcal{F}_{\mathcal{B}}$, that corresponds to the symbol s . $\sum_{1 \leq l \leq k}^\#$ is calculated by abstractly adding the k terms indexed by l . Abstract composition for a non-relational domain uses abstract addition and abstract multiplication to soundly overapproximate the quadratic terms occurring in Eqn. (4.1). We provide a proof of the soundness for $a' \circ_{\text{NR}} a$ in App. E.1. The abstract-composition operation requires $\mathcal{O}(k^3)$ abstract-addition operations and $\mathcal{O}(k^3)$ abstract-multiplication operations.

Examples of foundation domains. We now present a few foundation domains that allow to construct the non-relational small-set, interval [19], and strided-interval [83] base domains.

Small sets. $\mathcal{F}_{\text{SS}_n} \stackrel{\text{def}}{=} \{\top\} \cup \{S \mid S \subseteq \mathbb{Z}_{2^w} \wedge |S| \leq n\}$. The join operation is defined by: $(f_1 \sqcup f_2) = \begin{cases} f_1 \cup f_2 & \text{if } |f_1 \cup f_2| \leq n \\ \top & \text{otherwise} \end{cases}$

n denotes the maximum cardinality allowed in the non-top elements of $\mathcal{F}_{\text{SS}_n}$. Other abstract operators, including abstract addition and multiplication, are implemented in a similar manner.

Intervals. $\mathcal{F}_{\mathbb{Z}_{2^w}} \stackrel{\text{def}}{=} \{\perp\} \cup \{[a, b] \mid a, b \in \mathbb{Z}_{2^w}, a \leq b\}$. Most abstract operations are straightforward (See [19] for details). The abstract-addition and abstract-multiplication operations need to be careful about overflows to preserve soundness. For instance,

$$[a_1, b_1] +^\# [a_2, b_2] = \begin{cases} [a_1 + a_2, b_1 + b_2] & \text{if neither } a_1 + a_2 \text{ nor } b_1 + b_2 \\ & \text{overflows} \\ [\min, \max] & \text{otherwise} \end{cases}$$

Strided Interval. $\mathcal{F}_{S\mathbb{Z}_{2^w}} \stackrel{\text{def}}{=} \{\perp\} \cup \{s[a, b] \mid a, b, s \in \mathbb{Z}_{2^w}, a \leq b\}$, where $\gamma(s[a, b]) = \{i \mid a \leq i \leq b, i \equiv a \pmod{s}\}$. (See [83, 91] for the details of the abstract-domain operations.)

Affine-Closed Base Domain.

We discuss the special case when the base domain \mathcal{B} is affine-closed, i.e., $\mathcal{B} = \text{AG}$. The abstract composition is defined as:

$$a' \circ_{\text{AG}} a = a'', \text{ where } \text{base}(a'') = \langle \{t_i \times t'_j \mid 1 \leq i \leq l, 1 \leq j \leq l'\} \rangle \wedge \quad (4.4)$$

$$\text{base}(a) = \langle \{t_1, t_2, \dots, t_l\} \rangle \wedge \text{base}(a') = \langle \{t'_1, t'_2, \dots, t'_{l'}\} \rangle$$

Lemma 5.1 in [78] asserts that the above abstract composition method is sound by linearity of affine-closed abstractions. The abstract composition has time complexity $\mathcal{O}(hh'k^3)$, h (respectively h') is the height of the abstract-domain element $\text{base}(a)$ (or $\text{base}(a')$) in the AG lattice. Because the height of the AG lattice with $(k+1)^2$ columns is $\mathcal{O}(k^2)$, the time complexity for the abstract composition operation translates to $\mathcal{O}(k^7)$. Alg. 5 essentially implements Eqn. (4.4), but makes an extra SMT call to ensure that the result is sound. Because Eqn. (4.4) is sound by linearity for the AG domain, the very first SMT call in the while-loop condition at line 4 in Alg. 5 will be unsatisfiable.

Weakly-Convex Base Domain

We present methods to perform abstract composition when the base domain \mathcal{B} satisfies a property we call *weak convexity*. Base domain \mathcal{B} is *weakly convex* iff

- The abstraction of a single concrete affine transformer is exact: $\gamma(\alpha(t_i)) = \{t_i\}$.
- All abstract-domain elements $b \in \mathcal{B}$ are contained in a convex space over rationals: For any set of concrete affine transformers $\{t_0, t_1, \dots, t_l\}$, such that $b = \bigsqcup_{i=0}^l t_i$, and any $t \in \gamma(b)$:

$$\exists \lambda_1, \lambda_2, \dots, \lambda_l \in \mathbb{Q}. (0 \leq \lambda_1, \lambda_2, \dots, \lambda_l \leq 1) \wedge \sum_{i=0}^l \lambda_i = 1 \wedge \text{cast}_{\mathbb{Q}}(t) = \sum_{i=0}^l \lambda_i \text{cast}_{\mathbb{Q}}(t_i).$$

The $\text{cast}_{\mathbb{Q}}$ function is used to specify the convexity property by moving the point space from bitvectors to rationals. For instance, the expression $\sum_{i=0}^l \lambda_i \cdot \text{cast}_{\mathbb{Q}}(t_i)$ specifies the convex combination of the concrete affine transformers m_i in the rational space $\text{cast}_{\mathbb{Q}}^{(k+k^2)}$.

Any convex abstract domain over rationals, such as polyhedra [23] or octagons [68], can be used to create a weakly-convex domain over bitvectors [96, 94]. Abstract composition for weakly-convex base domains is defined as follows:

$$a' \circ_{\text{WC}} a = a'', \text{ where } \text{base}(a'') = \begin{cases} \{t_i \times t'_j \mid 1 \leq i \leq l, 1 \leq j \leq l'\} & \text{if there are no overflows in any} \\ & \text{matrix multiplication } t_i \times t'_j \\ \top_{\mathcal{B}} & \text{otherwise} \end{cases} \quad (4.5)$$

where $\text{base}(a) = \{t_1, t_2, \dots, t_l\}$ and $\text{base}(a') = \{t'_1, t'_2, \dots, t'_{l'}\}$.

The intuition is that the weak-convexity properties are preserved under matrix multiplication in the absence of overflows. This principle is similar to the linearity argument used to show that abstract composition is sound

when the base domain is affine-closed. (See above for more details.) We provide a proof of the soundness for $a' \circ_{WC} a$ in App. E.2. Similar to the affine-closed case, abstract composition has time complexity $\mathcal{O}(H^2k^3)$, where H is the height of the \mathcal{B} lattice.

Practical concerns. With the exception of the non-relational base domain, the complexity of the abstract-composition algorithms is dependent on the height of the abstract-domain elements involved in the composition, i.e., h and h' . Practical implementations of abstract composition might decide to return \top for abstract composition if the number of matrices to multiply is beyond some threshold, say t , so that the complexity of the abstract composition is $\mathcal{O}(tk^3)$.

4.3.3 Merge Function

Knoop and Steffen [54] extended the Sharir and Pnueli [92] algorithm for interprocedural dataflow analysis to handle local variables. Suppose at a call site CS, procedure P calls procedure Q. The global variables, denoted by \vec{g} , are accessible to Q, but the local variables, denoted by \vec{l} , in P are inaccessible to Q. Thus, the values of local variables after the call site CS come from the values before the call point, and the values of global variables after the call site CS come from the values at the return site in procedure Q. A *merge function* is used to combine the abstract-domain element before the call to Q with the abstract-domain element returned by Q to create the abstract-domain element to use in P after the call to Q has finished.

We assume that in each function, the local variables are initialized to 0. To simplify the discussion, assume that all scopes have the same number of locals, and that each vocabulary \vec{v} consists of subvocabularies \vec{g} and \vec{l} —that is, $\vec{v} = (\vec{g}, \vec{l})$. Suppose that we have two relations, $R[\vec{g}, \vec{l}; \vec{g}', \vec{l}']$ and $R'[\vec{g}, \vec{l}; \vec{g}', \vec{l}']$, each of which is a subset of $\mathbb{Z}_{2^w}^k \times \mathbb{Z}_{2^w}^k$, where R is the

transition relation from the start state of the calling procedure P to the call site CS, and R' is the transition relation from the start state to the return site of the called procedure Q. Operationally, after completing the call at the call site CS, we want $\text{MERGE}(R[\vec{g}, \vec{l}; \vec{g}', \vec{l}'], R'[\vec{g}, \vec{l}; \vec{g}', \vec{l}'])$ to act as a modified relational composition in which R' acts like the identity function on locals, so that \vec{l}' values from R are passed through R' unchanged to become the \vec{l}' values of the result. This semantics can be specified as follows:

$$\begin{aligned} & \text{MERGE}(R[\vec{g}, \vec{l}; \vec{g}', \vec{l}'], R'[\vec{g}, \vec{l}; \vec{g}', \vec{l}'] \quad (4.6) \\ & = \text{REVERTLOCALS}(R'[\vec{g}, \vec{l}; \vec{g}', \vec{l}']) \circ R[\vec{g}, \vec{l}; \vec{g}', \vec{l}'] \end{aligned}$$

We define $\text{REVERTLOCALS}(R'[\vec{g}, \vec{l}; \vec{g}', \vec{l}'])$ as follows:

$$\text{REVERTLOCALS}(R'[\vec{g}, \vec{l}; \vec{g}', \vec{l}']) \stackrel{\text{def}}{=} \{(\vec{g}, \vec{q}, \vec{g}', \vec{q}) \mid R'[\vec{g}, \vec{l}; \vec{g}', \vec{l}']\} \quad (4.7)$$

Recall that the $k + k^2$ symbolic constants in an affine transformation, $\text{symbols}(C: \vec{d})$, can be padded with a one and k zeroes and arranged as follows: $\left[\begin{array}{c|c} 1 & \vec{d} \\ \hline 0 & C \end{array} \right]$. We can partition $\text{symbols}(C: \vec{d})$ into globals and

locals to write the matrix as $\left[\begin{array}{c|cc} 1 & \vec{d}_g & \vec{d}_l \\ \hline 0 & C_{gg} & C_{gl} \\ 0 & C_{lg} & C_{ll} \end{array} \right]$, which represents the affine transformation

$$(\vec{g}' = \vec{g} \cdot C_{gg} + \vec{l} \cdot C_{lg} + \vec{d}_g) \wedge (\vec{l}' = \vec{g} \cdot C_{gl} + \vec{l} \cdot C_{ll} + \vec{d}_l)$$

Let $a, a' \in \text{ATA}[\mathcal{B}]$ be the abstract transformers that represents the relations $R[\vec{g}, \vec{l}; \vec{g}', \vec{l}']$ and $R'[\vec{g}, \vec{l}; \vec{g}', \vec{l}']$, respectively. Then the merge function for a_1 and a_2 is defined as follows:

$$\text{Merge}(a, a') = a'', \text{ where} \quad (4.8)$$

$$\begin{aligned}
\text{base}(a'') &= (b_g \sqcap \text{havoc}(\text{base}(\text{Id}), \text{gsyms})) \circ a \\
b_g &= \text{havoc}(\text{base}(a'), \text{lsyms}), \\
\text{lsyms} &= \text{symbols}(C_{gl}) \cup \text{symbols}(C_{ll}) \cup \\
&\quad \text{symbols}(d_l) \cup \text{symbols}(C_{lg}) \\
\text{gsyms} &= \text{symbols}(C_{gg}) \cup \text{symbols}(d_g)
\end{aligned}$$

lsyms are the symbols in the affine transformation that involve local variables. gsyms are the symbols in the affine transformation that are not in lsyms. The expression $b_{gi} = (b_g \sqcap \text{havoc}(\text{base}(\text{Id}), \text{gsyms}))$ transforms

each affine transformer $\left[\begin{array}{c|c|c} 1 & \vec{d}_g & \vec{d}_l \\ \hline 0 & C_{gg} & C_{gl} \\ \hline 0 & C_{lg} & C_{ll} \end{array} \right] \in \gamma(\text{base}(a))$ to $\left[\begin{array}{c|c|c} 1 & \vec{d}_g & 0 \\ \hline 0 & C_{gg} & 0 \\ \hline 0 & 0 & I \end{array} \right]$. In

this way, b_{gi} ensures that the modifications of the globals at the return point of Q are accounted for, while the locals for P pass through Q unmodified. We provide the proof of soundness for the merge-function definition (Eqn. (4.8)) in App. F.

4.4 Discussion and Related Work

The abstract-domain elements in our framework abstract two-vocabulary relationships arising between the pre-transformation state and post-transformation state. For the sake of simplicity, we assumed that the variable sets in the pre-transformation and post-transformation state are the same, and an affine transformer is represented by a $(k+1) \times (k+1)$ matrix, where k is the number of variables in the pre-transformation state. However, this requirement is not mandatory. We can easily adapt our abstract-domain operations to work on $(k+1) \times (k'+1)$ matrices where k' is the number of variables in the post-transformation state.

The abstract-domain elements in our framework are not necessarily closed under intersection. Consider the two abstract values a_1 and a_2

for the vocabulary $V = \{v_1\}$. Let α_1 represent the affine transformation $v'_1 = 0$ and α_2 represent the identity affine transformation $v'_1 = v_1$. Thus, $\alpha_1 = \alpha\left(\left[\begin{array}{c|c} 1 & 0 \\ 0 & 0 \end{array}\right]\right)$, and $\alpha_2 = \alpha\left(\left[\begin{array}{c|c} 1 & 0 \\ 0 & 1 \end{array}\right]\right)$. The intersection of $\gamma(\alpha_1)$ and $\gamma(\alpha_2)$ is the point $p = (v'_1 = 0, v_1 = 0)$. There does not exist an abstract value in $\text{ATA}[\mathcal{B}]$, that can exactly represent the point p , because any abstract value containing p must contain at least one affine transformer of the form $v'_1 = v_1 \cdot c$, and thus must contain all points of the form $(v'_1 = t \cdot c, v_1 = t)$, where $t \in \mathbb{Z}_{2^w}$. As a consequence, there does not exist a Galois connection between $\text{ATA}[\mathcal{B}]$ and the concrete domain \mathcal{C} of all two-vocabulary relations $R[V; V']$, which implies that there does not exist a best abstraction for a set of concrete points. For instance, consider the abstraction of the guard statement $S_G = \{v_1 \leq 10\}$, with the $\text{ATA}[\mathcal{J}_{\mathbb{Z}_{2^w}}^{(k+1)^2}]$ domain. Consider $\alpha_3 = \left[\begin{array}{c|c} 1 & [0, 10] \\ 0 & [0, 0] \end{array}\right]$ and $\alpha_4 = \left[\begin{array}{c|c} 1 & [0, 0] \\ 0 & [1, 1] \end{array}\right]$. α_3 specifies the guard constraint $0 \leq v'_1 \leq 10$, while α_4 is the identity transformation $v'_1 = v_1$. Note that these abstract values are incomparable and can be used to represent the abstract transformer for S_G . Furthermore, $\alpha_3 \sqcap \alpha_4$ does not exist. Thus, an analysis has to settle for either α_3 or α_4 . (In §4.2, we used an abstract transformer similar to α_3 for the guard in the while statement in Ex. 4.4. Using an identity transfer for the guard statement would not have been useful to capture the desired bounds constraint.)

The ATA constructor preserves finiteness; that is, if the base domain \mathcal{B} is finite, then the domain $\text{ATA}[\mathcal{B}]$ is finite as well.

It is also possible to use the ATA constructor to infer affine transformations over rationals or reals. In these cases, the symbolic-composition methods for weakly-convex base domains (see §4.3.2) will carry over to affine transformations over rationals or reals for convex base domains (e.g., polyhedra) with only slight modifications. For instance, abstract composition for convex base domains over rationals or reals is defined as

follows:

$$\begin{aligned} a' \circ a = a'', \text{ where } \text{base}(a'') &= \{t_i \times t'_j \mid 1 \leq i \leq l, 1 \leq j \leq l'\} \\ \text{where } \text{base}(a) &= \{t_1, t_2, \dots, t_l\} \text{ and } \text{base}(a') = \{t'_1, t'_2, \dots, t'_{l'}\}. \end{aligned}$$

Chen et al. [18] devised the *interval-polyhedra* domain which can express constraints of the form $\sum_k [a_k, b_k]x_k \leq c$ over rationals. Interval polyhedra are more expressive than classic convex polyhedra, and thus can express certain non-convex properties. Abstract-domain operations for interval polyhedra are constructed by linear programming and interval Fourier-Motzkin elimination. The domain has similarities to the $\text{ATA}[\mathbb{J}_{\mathbb{Z}^{2w}}^{(k+1)^2}]$ domain because the coefficients in the abstract values are intervals.

Miné [69] introduced *weakly relational domains*, which are a parameterized family of relational domains, parameterized by a non-relational base abstract domain. They can express constraints of the form $(v_j - v_i) \in \mathcal{F}$, where \mathcal{F} is an abstraction over $\mathcal{P}(\mathbb{Z})$. Similar to $\text{ATA}[\mathcal{B}]$, Miné's framework requires the base non-relational domain to provide abstract-addition and abstract-unary-minus operations. These operations are used to propagate information between constraints via a closure operation that is similar to finding shortest paths.

Sankaranarayanan et al. [89] introduced a domain based on template constraint matrices (TCMs) that is less powerful than polyhedra, but more general than intervals and octagons. Their analysis discovers linear-inequality invariants using polyhedra with a predefined fixed shape. The predefined shape is given by the client in the form of a template matrix. Our approach is similar because an affine transformer with symbolic constants can be seen as a template. However, the approaches differ because Sankaranarayanan et al. use an LP solver to find values for template parameters, whereas we use operations and values from an abstract domain to find and represent a set of allowed values for template parameters.

An abstract-domain element in $\text{ATA}[\mathcal{B}]$ can be seen as an abstraction over sets of functions: $\mathbb{Z}_{2^w}^k \rightarrow \mathbb{Z}_{2^w}^k$. Jeannet et al. [46] provide a theoretical treatment of the relational abstraction of functions. They describing existing and new methods of abstracting functions of signature: $D_1 \rightarrow D_2$, resulting in a family of relational abstract domains. $\text{ATA}[\mathcal{B}]$ is not captured by their framework of functional abstractions.

4.5 Chapter Notes

Thomas Reps supervised me in the writing of the paper, and provided some key insights for this work. Jason Breck provided useful comments about this work.

5 AN ABSTRACT DOMAIN FOR BIT-VECTOR INEQUALITIES

The chapter describes the design and implementation of a new abstract domain, called the *Bit-Vector Inequality* (\mathcal{BVJ}) domain, that is capable of capturing certain inequalities over bit-vector-valued variables. We also consider some variants of the \mathcal{BVJ} domain.

Key Insight: The View-Product Combinator. The key insight used to design \mathcal{BVJ} domain (and its variants) involves a new domain combinator (denoted by \mathcal{V}), called the *view-product combinator*. \mathcal{V} constructs a reduced product of two domains [21], using shared *view-variables* to communicate information between the domains. The following example illustrates the concept of a view-variable:

Example 5.1. Consider the equality constraint $H := x + 2\text{mm}[y + 2] = 4$, where x and y are bit-vector variables and mm is the memory map. The term $\text{mm}[e]$ denotes the contents of mm at address e . Let $\mathcal{E}_{\mathbb{Z}_{2^w}}$ denote any of the abstract domains of relational affine equalities over bit-vector variables (See Chapters 3 and 4). H cannot be expressed using $\mathcal{E}_{\mathbb{Z}_{2^w}}$ alone. However, the formulas $x + 2\text{mm}[y + 2] = 4$ and $u = \text{mm}[y + 2] \wedge x + 2u = 4$ are equisatisfiable. The variable u is called a view-variable; the constraint $u = \text{mm}[y + 2]$ is the view-constraint associated with u .

The equality $x + 2u = 4$ can be expressed using $\mathcal{E}_{\mathbb{Z}_{2^w}}$; therefore, what we require is a second abstract domain capable of expressing invariants that involve memory accesses, such as $u = \text{mm}[y + 2]$. This need motivates the bit-vector memory domain \mathcal{M} that we introduce in §5.3. The constraint H can then be expressed by (i) introducing view-variable u , (ii) representing u 's view-constraint in \mathcal{M} , (iii) extending $\mathcal{E}_{\mathbb{Z}_{2^w}}$ with u , and (iv) using \mathcal{M} and $\mathcal{E}_{\mathbb{Z}_{2^w}}$ together. \square

The *Bit-Vector Memory-Equality Domain* $\mathcal{BVM}\mathcal{E}$, a domain of bit-vector affine-equalities over variables and memory-values, is created by applying the view-product combinator \mathcal{V} to the bit-vector memory domain (§5.3) and the bit-vector equality domain. The *Bit-Vector Inequality Domain* $\mathcal{BV}\mathcal{J}$, a domain of bit-vector affine-inequalities over variables, is created by applying \mathcal{V} to the bit-vector equality domain and a bit-vector interval domain. The *Bit-Vector Memory-Inequality Domain* $\mathcal{BVM}\mathcal{J}$, a domain of relational bit-vector affine-inequalities over variables and memory, is then created by applying \mathcal{V} to the $\mathcal{BVM}\mathcal{E}$ domain and the bit-vector interval domain. The latter construction illustrates that \mathcal{V} composes: the $\mathcal{BVM}\mathcal{J}$ domain is created via two applications of \mathcal{V} .

The design of the view-product combinator \mathcal{V} was inspired by the Subpolyhedra domain [58] (SubPoly), a domain for inferring relational linear inequalities over rationals. SubPoly is constructed as a reduced product of an affine-equality domain \mathcal{K} over rationals [49], and an interval domain \mathcal{J} over rationals [19] in which *slack variables* are used to communicate between the two domains. Such a design enables SubPoly to be as expressive as Polyhedra, but more scalable. \mathcal{V} provides a generalization of the construction used in SubPoly. In fact, SubPoly can be constructed by applying \mathcal{V} to \mathcal{K} and \mathcal{J} (see Eqn. (5.1) in §5.7).

Enabling Technology: Automatic Synthesis of Best Abstract Operations. Using a symbolic abstraction method §2.3.2, we give a procedure for *synthesizing best abstract operations* for general reduced-product domains. In particular, we use this framework to ensure that the transformers for the reduced-product domains constructed via \mathcal{V} are sound and precise, thereby guaranteeing that analysis results will be a conservative over-approximation of the concrete semantics.

Contributions.

The contributions of the chapter are:

- The bit-vector memory domain, a non-relational memory domain capable of expressing invariants involving memory accesses (§5.3).
- The view-product combinator \mathcal{V} , a general procedure to construct more expressive domains (§5.4).
- Three domains for machine-code analysis constructed using \mathcal{V} :
 - The bit-vector memory-equality domain $\mathcal{BVM}\mathcal{E}$, which captures equality relations among bit-vector variables and memory (Defn. 5.6).
 - The bit-vector inequality domain $\mathcal{BV}\mathcal{I}$, which captures inequality relations among bit-vector variables (Defn. 5.7).
 - The bit-vector memory-inequality domain $\mathcal{BVM}\mathcal{I}$, which captures inequality relations among bit-vector variables and memory (Defn. 5.8).
- A procedure for synthesizing best abstract operations for reduced products of domains that meet certain requirements (§5.5).
- Experimental results that illustrate the effectiveness of the $\mathcal{BV}\mathcal{I}$ domain applied to machine-code analysis (§5.6). On average (geometric mean), our $\mathcal{BV}\mathcal{I}$ -based analysis is about 3.5 times slower than an affine-equality-based analysis, while finding improved (more-precise) invariants at 29% of the branch points.

§5.1 provides an overview of our solution. §5.2 defines terminology. §5.7 describes related work.

5.1 Overview

In this section, we illustrate the design of the bit-vector memory-inequality domain. Consider the Intel x86 machine-code snippet shown in Fig. 5.1. Pseudo-code for each instruction is shown at the right-hand side of each

```

mov  ecx, [ebp-4]    //ecx = mm[ebp-4]
add  ecx, eax       //ecx = ecx + eax
cmp  ecx, 10d       //if ecx >_u 10
ja   L              //  goto L
xor  ecx, ecx       // ecx = 0

```

$0 \leq \text{mm}[\text{ebp} - 4] + \text{eax} \leq 10$

L:

Figure 5.1: Example snippet of Intel x86 machine code.

instruction. Note that the **ja** instruction, “jump if above, unsigned”, treats **ecx** as unsigned. Thus, control reaches the **xor** instruction only if the value in **ecx** is greater than or equal to 0, and less than or equal to 10. Let $R := \{\text{eax}, \text{ecx}, \text{ebp}\}$ denote the set of register variables, and mm denote the memory map.

The highlighted text in Fig. 5.1 states the invariant $H := 0 \leq \text{mm}[\text{ebp} - 4] + \text{eax} \leq 10$ that holds after the **xor** instruction. Again, let $\mathcal{E}_{\mathbb{Z}_{2^w}}$ denote any of the abstract domains of relational affine equalities over bit-vector variables. The invariant H cannot be represented using $\mathcal{E}_{\mathbb{Z}_{2^w}}$, because (i) H involves an invariant about a value in memory mm , and (ii) H is an inequality. To handle memory, we introduce the bit-vector memory domain \mathcal{M} , which is capable of expressing certain constraints on memory accesses. In particular, we can use \mathcal{M} to express the constraint $u = \text{mm}[\text{ebp} - 4]$, where u is a fresh variable. We call $u = \text{mm}[\text{ebp} - 4]$ a *view-constraint* for u . H can be written as the equisatisfiable formula $H_m := u = \text{mm}[\text{ebp} - 4] \wedge 0 \leq u + \text{eax} \leq 10$.

Notice there are no memory accesses in the constraint $0 \leq u + \text{eax} \leq 10$. The inequality $0 \leq u + \text{eax} \leq 10$ and $u + \text{eax} = s \wedge 0 \leq s \leq 10$ are equisatisfiable, where s is a fresh variable. Similar to u , s is a view-variable with $u + \text{eax} = s$ is a view-constraint for s . Thus, H_m can be rewritten as $H_{m_i} := u = \text{mm}[\text{ebp} - 4] \wedge u + \text{eax} = s \wedge 0 \leq s \leq 10$. Furthermore, $u + \text{eax} = s$ can be expressed in $\mathcal{E}_{\mathbb{Z}_{2^w}}$, and $0 \leq s \leq 10$ can be expressed

using a bit-vector interval domain $\mathcal{J}_{\mathbb{Z}_{2^w}}$. Thus, by introducing the view-variables u and s , the invariant H can be expressed using \mathcal{M} , $\mathcal{E}_{\mathbb{Z}_{2^w}}$, and $\mathcal{J}_{\mathbb{Z}_{2^w}}$ together. The following derivation illustrates the above decomposition of the invariant H :

$$\text{MEMORY} \frac{u = \text{mm}[\text{ebp} - 4] \quad \frac{u + \text{eax} = s \quad s \in [0, 10]}{0 \leq u + \text{eax} \leq 10} \text{INEQUALITY}}{0 \leq \text{mm}[\text{ebp} - 4] + \text{eax} \leq 10}$$

Note that the view-constraints $u = \text{mm}[\text{ebp} - 4]$ and $s = u + \text{eax}$ do not directly constrain the values of R and mm ; they only constrain the view-variables u and s . The shared view-variables are used to exchange information among the various domains. In particular, the view-variable u is used to exchange information between memory domain \mathcal{M} and equality domain $\mathcal{E}_{\mathbb{Z}_{2^w}}$, and the view-variable s is used to exchange information between $\mathcal{E}_{\mathbb{Z}_{2^w}}$ and interval domain $\mathcal{J}_{\mathbb{Z}_{2^w}}$.

5.2 Terminology

For a somewhat technical reason, we introduce the device of an *abstract-domain constructor*. Given an abstract-domain family \mathcal{A} and vocabularies V_1 and V_2 , an abstract-domain constructor for \mathcal{A} , denoted by $\mathcal{C}_{\mathcal{A}}(V_1, V_2)$, constructs $\mathcal{A}[V_3]$, where $V_1 \subseteq V_3 \subseteq V_1 \uplus V_2$. In particular, $\mathcal{C}_{\mathcal{A}}$ is free to decide what subset of V_2 to use when constructing $\mathcal{A}[V_3]$. (In our applications, the abstract-domain constructors either use all of V_2 or none of V_2 .)

Let $A \in \mathcal{A}[V]$; we denote by $A \downarrow_{V_1}$ the value obtained by projecting A onto the vocabulary $V_1 \subseteq V$. We use \oplus to denote a *vocabulary-extension operator* over domains; in particular, given domain $\mathcal{A}[V_1]$ and vocabulary V_2 , $\mathcal{A}[V_1] \oplus V_2 = \mathcal{A}[V_1 \uplus V_2]$.

Let $\mathcal{G}_1 = \mathcal{C} \xleftrightarrow[\alpha_1]{\gamma_1} \mathcal{A}_1$ and $\mathcal{G}_2 = \mathcal{C} \xleftrightarrow[\alpha_2]{\gamma_2} \mathcal{A}_2$ be two Galois connections.

We use $\mathcal{A}_1[V_1] \star \mathcal{A}_2[V_2]$ to denote the reduced product of the domains [21, §10.1], and $\langle A_1; A_2 \rangle$ to denote an element of $\mathcal{A}_1[V_1] \star \mathcal{A}_2[V_2]$.

5.3 Base Domains

Let mm denote a memory map, and $\text{mm}[t]$ denote the w -bit value at memory address t . For instance, for Intel x86 machine code, $\text{mm}[t]$ denotes (the little-endian interpretation of) the four bytes in memory pointed to by term t (cf. Ex. 1.3).

Bit-Vector Memory Domain (\mathcal{M}). Domain \mathcal{M} can capture a limited class of invariants involving memory accesses. In particular, the domain is capable of capturing the constraint that the value of a variable v equals the value of the memory mm at address e : $v = \text{mm}[e]$.

Definition 5.2. *Given variables P and memory map mm , an element M of the **bit-vector memory domain** $\mathcal{M}[(\text{mm}, P)]$ is either (i) \perp , or (ii) a set of constraints, where each constraint C_i is of the form $v_i = \text{mm}[\sum_j a_{ij}v_j + b_i]$, $a_{ij}, b_i \in \mathbb{Z}_{2^w}$, $v_i, v_j \in P$. The concretization of M is*

$$\gamma(M) = \{(\text{mm}, \vec{v}) \mid \bigwedge_{C_i \in M} (\text{mm}, \vec{v}) \models C_i\}.$$

□

The join and meet operation are defined as intersection and union of constraints, respectively. The join operation $\sqcup_{\mathcal{M}}$ is: $M_1 \sqcup_{\mathcal{M}} M_2 \stackrel{\text{def}}{=} \{C \mid C \in M_1 \text{ and } C \in M_2\}$. The meet operation $\sqcap_{\mathcal{M}}$ is: $M_1 \sqcap_{\mathcal{M}} M_2 \stackrel{\text{def}}{=} \{C \mid C \in M_1 \text{ or } C \in M_2\}$. The domain constructor for the bit-vector memory domain is defined as $\mathfrak{C}_{\mathcal{M}}(V_1, V_2) \stackrel{\text{def}}{=} \mathcal{M}[V_1 \uplus V_2]$.

Bit-Vector Equality Domain ($\mathcal{E}_{\mathbb{Z}_{2^w}}$). Domain $\mathcal{E}_{\mathbb{Z}_{2^w}} = \text{KS}$ can capture relational affine equalities over bit-vectors.

Bit-Vector Interval Domain ($\mathcal{J}_{\mathbb{Z}_2^w}$). Domain $\mathcal{J}_{\mathbb{Z}_2^w}$ can capture non-relational bit-vector interval constraints.

Definition 5.3. Given variables V , an element I of the *bit-vector interval domain* $\mathcal{J}_{\mathbb{Z}_2^w}[V]$ is either (i) \perp , or (ii) a set of interval constraints, where each constraint C_i is of the form $l_i \leq v_i \leq u_i$, $l_i, u_i \in \mathbb{Z}_2^w, v_i \in V$. The concretization of I is

$$\gamma(I) = \{\vec{v} \mid \bigwedge_{C_i \in I} \vec{v} \models C_i\}.$$

□

The domain constructor for the bit-vector interval domain is defined as $\mathfrak{C}_{\mathcal{J}}(V_1, V_2) \stackrel{\text{def}}{=} \mathcal{J}_{\mathbb{Z}_2^w}[V_1]$.

5.4 The View-Product Combinator

In this section, we define the view-product combinator \mathcal{V} , and construct three domains using different applications of \mathcal{V} . \mathcal{V} constructs a reduced-product of domains \mathcal{A}_1 and \mathcal{A}_2 so that a set of shared view-variables communicates information between the domains. \mathcal{V} makes use of two principles:

Principle 1: View-variables are constrained by view-constraints. Consider abstract domain $\mathcal{A}_1[V_1]$, and let $V_1 \cap V_2 = \emptyset$. The variables in V_2 are the *view-variables*. \mathcal{V} will use the domain $\mathcal{A}_1[V_1] \oplus V_2 = \mathcal{A}_1[V_1 \uplus V_2]$ (cf. §5.2). A *view-constraint* C for V_2 is an element of $\mathcal{A}[V_1 \uplus V_2]$ that serves to constrain the variable set V_2 .

Principle 2: A view-constraint does not constrain the values of variables in V_1 . Given abstract domain $\mathcal{A}[V_1 \uplus V_2]$, an *acceptable view-constraint* C for V_2 is an abstract value $C \in \mathcal{A}[V_1 \uplus V_2]$ such that $C \downarrow_{V_1} = \top$.

We are now in a position to describe how \mathcal{V} works.

Arguments. \mathcal{V} takes four arguments:

- $\mathcal{A}_1[V_1]$, an abstract domain defined over vocabulary V_1 ,
- V_2 , a vocabulary such that $V_1 \cap V_2 = \emptyset$,
- $C \in \mathcal{A}_1[V_1 \uplus V_2]$, a view-constraint for V_2 , and
- $\mathfrak{C}_{\mathcal{A}_2}$, an abstract-domain constructor for abstract domain \mathcal{A}_2 ,

Enforcement of view-constraint C . The view-constraint C constrains the variables in V_2 . Therefore, \mathcal{V} should only consider those elements $\mathcal{A}[V_1 \uplus V_2]$ that satisfy C ; that is, only elements $A \in \mathcal{A}[V_1 \uplus V_2]$ such that $A \sqsubseteq_{\mathcal{A}[V_1 \uplus V_2]} C$ holds. Put another way, the view-constraint C can be seen as an *integrity constraint*. The next definition formalizes this notion (for a general integrity constraint $D \in \mathcal{A}[V]$):

Definition 5.4. *Given an abstract domain $\mathcal{A}[V]$ and an element $D \in \mathcal{A}[V]$, the **abstract domain $\mathcal{A}[V]$ modulo D** , denoted by $\mathcal{A}[V] \mid_D$, is an abstract domain $\mathcal{A}'[V]$ that contains exactly the elements $D \sqcap A$, $A \in \mathcal{A}$; that is,*

$$\mathcal{A}[V] \mid_D \stackrel{\text{def}}{=} \{D \sqcap A \mid D \in \mathcal{A}[V]\}$$

□

Using the notation from Defn. 5.4, the domain that only contains elements that satisfy view-constraint C is $\mathcal{A}_1[V_1 \uplus V_2] \mid_C$.

Reduced product. Let $\mathcal{A}'_1[V_1 \uplus V_2]$ be $(\mathcal{A}_1[V_1] \oplus V_2) \mid_C$, where C is the view-constraint for view-variables V_2 . All that is left is to perform a reduced product of \mathcal{A}'_1 and \mathcal{A}_2 . However, for \mathcal{A}'_1 and \mathcal{A}_2 to be able to exchange information, the vocabulary of \mathcal{A}_2 should include at least the view-variables V_2 . This condition is satisfied by the domain $\mathfrak{C}_{\mathcal{A}_2}(V_2, V_1)$ constructed using the abstract-domain constructor for \mathcal{A}_2 supplied as the fourth argument to \mathcal{V} (cf. §5.2).

Summing up, the reduced product performed by \mathcal{V} creates $\mathcal{A}_3[\mathbb{V}_1 \uplus \mathbb{V}_2] := \mathcal{A}'_1[\mathbb{V}_1 \uplus \mathbb{V}_2] \star \mathfrak{C}_{\mathcal{A}_2}(\mathbb{V}_2, \mathbb{V}_1)$.

Definition 5.5 (View-Product Combinator (\mathcal{V})). *Given*

- $\mathcal{A}_1[\mathbb{V}_1]$, an abstract domain defined over vocabulary \mathbb{V}_1 ,
- \mathbb{V}_2 , a vocabulary such that $\mathbb{V}_1 \cap \mathbb{V}_2 = \emptyset$,
- $\mathbb{C} \in \mathcal{A}_1[\mathbb{V}_1 \uplus \mathbb{V}_2]$, a view-constraint for \mathbb{V}_2 , and
- $\mathfrak{C}_{\mathcal{A}_2}$, an abstract-domain constructor for abstract domain \mathcal{A}_2 ,

the **view-product combinator** $\mathcal{V}[\mathcal{A}_1[\mathbb{V}_1], \mathbb{V}_2, \mathbb{C}, \mathfrak{C}_{\mathcal{A}_2}]$ constructs a domain $\mathcal{A}_3[\mathbb{V}_1 \uplus \mathbb{V}_2]$ such that

$$\mathcal{A}_3[\mathbb{V}_1 \uplus \mathbb{V}_2] \stackrel{\text{def}}{=} (\mathcal{A}_1[\mathbb{V}_1] \oplus \mathbb{V}_2) \upharpoonright_{\mathbb{C}} \star \mathfrak{C}_{\mathcal{A}_2}(\mathbb{V}_2, \mathbb{V}_1)$$

□

Instantiations. We now describe three applications of \mathcal{V} that use the bit-vector domains defined in §5.3.

\mathcal{V} applied to the bit-vector memory domain \mathcal{M} and the bit-vector equality domain $\mathcal{E}_{\mathbb{Z}_2^w}$ constructs the *Bit-Vector Memory-Equality Domain* $\mathcal{BVM}\mathcal{E}$, a domain of bit-vector affine-equalities over variables and memory-values.

Definition 5.6 (Bit-Vector Memory-Equality Domain ($\mathcal{BVM}\mathcal{E}$)).

$$\mathcal{BVM}\mathcal{E}[(\mathbb{m}, \mathbb{P} \uplus \mathbb{U})] \stackrel{\text{def}}{=} \mathcal{V}[\mathcal{M}[(\mathbb{m}, \mathbb{P})], \mathbb{U}, \mathbb{C}_{\mathbb{m}}, \mathfrak{C}_{\mathcal{E}_{\mathbb{Z}_2^w}}],$$

where

- $\mathcal{M}[(\mathbb{m}, \mathbb{P})]$, the bit-vector memory domain over memory map \mathbb{m} and variables \mathbb{P} (cf. Defn. 5.2),
- \mathbb{U} , a vocabulary of variables such that $\mathbb{U} \cap \mathbb{P} = \emptyset$,
- $\mathbb{C}_{\mathbb{m}} \in \mathcal{M}[(\mathbb{m}, \mathbb{P} \uplus \mathbb{U})]$, a view-constraint for \mathbb{U} , and
- $\mathfrak{C}_{\mathcal{E}_{\mathbb{Z}_2^w}}(\mathbb{V}_1, \mathbb{V}_2) \stackrel{\text{def}}{=} \mathcal{E}_{\mathbb{Z}_2^w}[\mathbb{V}_1 \uplus \mathbb{V}_2]$.

□

\mathcal{V} applied to the bit-vector equality domain $\mathcal{E}_{\mathbb{Z}_2^w}$ and the bit-vector interval domain $\mathcal{J}_{\mathbb{Z}_2^w}$ constructs the *Bit-Vector Inequality Domain* \mathcal{BVJ} , a domain of bit-vector affine-equalities over variables.

Definition 5.7 (The Bit-Vector Inequality Domain (\mathcal{BVJ})).

$$\mathcal{BVJ}[\mathbf{P} \uplus \mathbf{S}] \stackrel{\text{def}}{=} \mathcal{V}[\mathcal{E}_{\mathbb{Z}_2^w}[\mathbf{P}], \mathbf{S}, \mathbf{C}_s, \mathfrak{C}_J],$$

where

- $\mathcal{E}_{\mathbb{Z}_2^w}[\mathbf{P}]$, the bit-vector equality domain over variables \mathbf{P} ,
- \mathbf{S} , a vocabulary of variables such that $\mathbf{S} \cap \mathbf{P} = \emptyset$,
- $\mathbf{C}_s \in \mathcal{E}_{\mathbb{Z}_2^w}[\mathbf{P} \uplus \mathbf{S}]$, a view-constraint for \mathbf{S} , and
- $\mathfrak{C}_J(\mathbf{V}_1, \mathbf{V}_2) \stackrel{\text{def}}{=} \mathcal{J}_{\mathbb{Z}_2^w}[\mathbf{V}_1]$.

□

The combinator \mathcal{V} applied to the bit-vector memory-equality domain $\mathcal{BVM}\mathcal{E}$, and a bit-vector interval domain $\mathcal{J}_{\mathbb{Z}_2^w}$ constructs the *Bit-Vector Memory-Inequality Domain* \mathcal{BVMJ} , a domain of relational bit-vector affine-inequalities over variables and memory.

Definition 5.8 (Bit-Vector Memory-Inequality Domain (\mathcal{BVMJ})).

$$\mathcal{BVMJ}[(\mathbf{mm}, \mathbf{P} \uplus \mathbf{U} \uplus \mathbf{S})] \stackrel{\text{def}}{=} \mathcal{V}[\mathcal{BVM}\mathcal{E}[(\mathbf{mm}, \mathbf{P} \uplus \mathbf{U})], \mathbf{S}, \mathbf{C}_s, \mathfrak{C}_J]$$

where

- $\mathcal{BVM}\mathcal{E}[(\mathbf{mm}, \mathbf{P} \uplus \mathbf{U})]$, the bit-vector memory-equality domain over memory map \mathbf{mm} and variables $\mathbf{P} \uplus \mathbf{U}$ (cf. Defn. 5.6),
- \mathbf{S} , a vocabulary of variables such that $\mathbf{S} \cap (\mathbf{P} \uplus \mathbf{U}) = \emptyset$,
- $\mathbf{C}_s \in \mathcal{BVM}\mathcal{E}[(\mathbf{mm}, \mathbf{P} \uplus \mathbf{U})]$, a view-constraint for \mathbf{S} , and
- $\mathfrak{C}_J(\mathbf{V}_1, \mathbf{V}_2) \stackrel{\text{def}}{=} \mathcal{J}_{\mathbb{Z}_2^w}[\mathbf{V}_1]$.

□

5.5 Synthesizing Abstract Operations for Reduced-Product Domains

In this section, we first discuss a method for automatically synthesizing abstract operations for a general reduced-product domain $\mathcal{A}_3[V_1 \uplus V_2] := \mathcal{A}_1[V_1] \star \mathcal{A}_2[V_2]$ using the abstract operations of \mathcal{A}_1 and \mathcal{A}_2 , which themselves may be automatically synthesized. We then discuss some pragmatic choices that we made for the reduced-product domains that the view-product combinator \mathcal{V} creates.

Semantic Reduction. The semantic reduction of $\langle A_1; A_2 \rangle \in \mathcal{A}_1 \star \mathcal{A}_2 = \mathcal{A}_3$ can be computed as $\hat{\alpha}_3(\psi)$, where ψ is $\hat{\gamma}_1(A_1) \wedge \hat{\gamma}_2(A_2)$. Computing the semantic reduction in this way can be computationally expensive. Instead we assume there exists a *weak semantic-reduction* operator Reduce . Using this weak semantic-reduction operator we can define the other abstract operations for the reduced-product domain. Furthermore, because it can be expensive to determine whether $\langle A_1; A_2 \rangle \sqsubseteq \langle A'_1; A'_2 \rangle$ holds, we define a weaker approximation order \sqsubseteq :

Definition 5.9. *The weak approximation order \sqsubseteq is defined as follows: $\langle A'_1; A'_2 \rangle \sqsubseteq \langle A''_1; A''_2 \rangle$ if and only if $A'_1 \sqsubseteq_{\mathcal{A}_1} A''_1$ and $A'_2 \sqsubseteq_{\mathcal{A}_2} A''_2$. \square*

It is easy to show that if $\langle A'_1; A'_2 \rangle \sqsubseteq \langle A''_1; A''_2 \rangle$, then $\langle A'_1; A'_2 \rangle \sqsubseteq \langle A''_1; A''_2 \rangle$, though the converse may not always hold.

We define *quasi-join*, an approximation to the join operator for reduced-product domain, which is not guaranteed to return the least-upper bound, but is sound and simple.

Definition 5.10. *The quasi-join operator, denoted by \sqcup , is defined as follows: $\langle A'_1; A'_2 \rangle \sqcup \langle A''_1; A''_2 \rangle \stackrel{\text{def}}{=} \text{Reduce}(\langle A'_1 \sqcup_{\mathcal{A}_1} A''_1; A'_2 \sqcup_{\mathcal{A}_2} A''_2 \rangle)$. \square*

Theorem 5.11. [Soundness of Quasi-Join] *Let $\langle A_1; A_2 \rangle = \langle A'_1; A'_2 \rangle \sqcup \langle A''_1; A''_2 \rangle$. Then $\gamma(\langle A_1; A_2 \rangle) \supseteq \gamma(\langle A'_1; A'_2 \rangle) \cup \gamma(\langle A''_1; A''_2 \rangle)$. \square*

The *quasi-meet* operation is similar in flavor to quasi-join:

Definition 5.12. The *quasi-meet* operator, denoted by \sqcap , is defined as follows:

$$\langle A'_1; A'_2 \rangle \sqcap \langle A''_1; A''_2 \rangle \stackrel{\text{def}}{=} \text{Reduce}(\langle A'_1 \sqcap_{\mathcal{A}_1} A''_1; A'_2 \sqcap_{\mathcal{A}_2} A''_2 \rangle). \quad \square$$

Theorem 5.13. [Soundness of Quasi-Meet] Let $\langle A_1; A_2 \rangle = \langle A'_1; A'_2 \rangle \sqcap \langle A''_1; A''_2 \rangle$. Then $\gamma(\langle A_1; A_2 \rangle) \supseteq \gamma(\langle A'_1; A'_2 \rangle) \cap \gamma(\langle A''_1; A''_2 \rangle)$. \square

We now define weak semantic-reduction operators for each of the pairs of domains used in our constructions. The algorithms for these specific domains have not been stated explicitly in the literature, and are stated here for completeness. (Previous work tackled the rational-arithmetic variants of these domains.)

Algorithm 6 Algorithm for weak semantic-reduction for $\mathcal{M}[V] \star \mathcal{E}_{\mathbb{Z}_2^w}[V]$.

```

1: for constraint  $v_1 = \text{mm}[e_1] \in M$  do
2:   for constraint  $v_2 = \text{mm}[e_2] \in M$  do
3:     if  $E \sqsubseteq_{\mathcal{E}_{\mathbb{Z}_2^w}} \{e_1 = e_2\}$  then
4:        $E \leftarrow E \sqcap_{\mathcal{E}_{\mathbb{Z}_2^w}} \{v_1 = v_2\}$ 
5: return  $\langle M; E \rangle$ 

```

Weak Reduce for $\mathcal{M} \star \mathcal{E}_{\mathbb{Z}_2^w}$. Alg. 6 computes the weak semantic-reduction for the bit-vector memory domain and bit-vector equality domain. Given $\langle M; E \rangle \in \mathcal{M}[(\text{mm}, V)] \star \mathcal{E}_{\mathbb{Z}_2^w}[V]$, $\text{Reduce}(\langle M; E \rangle)$ infers further equalities among V . The key insight is to model the memory map mm as an uninterpreted function; that is, $t_1 = t_2$ implies $\text{mm}[t_1] = \text{mm}[t_2]$. We are effectively approximating the theory of arrays using the theory of uninterpreted functions with equality (EUF) [14], thereby ensuring that the reduction operation is efficient. As shown in lines 3 and 4, if $e_1 = e_2$ then the algorithm infers that $v_1 = v_2$. The following example illustrates the working of Alg. 6.

Example 5.14. Let V be $\{x, y, z, u_1, u_2\}$. Consider $\langle M; E \rangle := \langle u_1 = \text{mm}[x], u_2 = \text{mm}[y + 8]; x = y + z - 2, z = 10 \rangle$. From E we can infer that $x = y + 8$. Thus, $\text{mm}[x] = \text{mm}[y + 8]$, and we can infer that $u_1 = u_2$. Thus, E can be updated to $E \sqcap_{\mathcal{E}_{\mathbb{Z}_{2^w}}} \{u_1 = u_2\}$.

No further reduction is possible; thus, $\text{Reduce}_{\mathcal{M} \star \mathcal{E}_{\mathbb{Z}_{2^w}}}(\langle M; E \rangle) = \langle u_1 = \text{mm}[x], u_2 = \text{mm}[y + 8]; u_1 = u_2, x = y + z - 2, z = 10 \rangle$. \square

The following example shows a case when $\text{Reduce}_{\mathcal{M} \star \mathcal{E}_{\mathbb{Z}_{2^w}}}$ fails to find the most precise answer.

Example 5.15. Consider the situation when analyzing Intel x86 machine code. The bit-width w is 32. The memory map mm is a map from 32-bit bit-vector to 8-bit bit-vectors, and the addressing mode is little-endian (cf. Ex. 1.3). Let V be $\{r, u_1, u_2, u_3\}$. Consider $\langle M; E \rangle := \langle u_1 = \text{mm}[r], u_2 = \text{mm}[r + 2], u_3 = \text{mm}[r + 4]; u_1 = 0, u_3 = 0; \top \rangle$. Because neither $r = r + 2$ nor $r = r + 4$ hold, Alg. 6 cannot infer any further equalities among u_1, u_2 , and u_3 . Thus, $\text{Reduce}_{\mathcal{M} \star \mathcal{E}_{\mathbb{Z}_{2^w}}}(\langle M; E \rangle) = \langle M; E \rangle$.

However, $\langle M; E \rangle$ is not the best reduction possible. Because $u_1 = 0$ and $u_3 = 0$, we can infer that $\text{mm}[r]$ and $\text{mm}[r + 4]$ are 0. Thus, the bytes at addresses r through $r + 7$ are all 0, because we assumed little-endian addressing. Consequently, $\text{mm}[r + 2] = 0$, and $u_2 = 0$. Thus, the best reduction for $\langle M; E \rangle$ is $\langle u_1 = \text{mm}[r], u_2 = \text{mm}[r + 2], u_3 = \text{mm}[r + 4]; u_1 = 0, u_2 = 0, u_3 = 0 \rangle$. The reason the $\text{Reduce}_{\mathcal{M} \star \mathcal{E}_{\mathbb{Z}_{2^w}}}$ algorithm was unable to deduce this was because the algorithm treats the memory map as an uninterpreted function. Thus, though sound, Alg. 6 is not always able to deduce the best possible reduced value. \square

Theorem 5.16. [Soundness of Alg. 6] Let $\langle M'; E' \rangle = \text{Reduce}_{\mathcal{M} \star \mathcal{E}_{\mathbb{Z}_{2^w}}}(\langle M; E \rangle)$. Then $\gamma(\langle M'; E' \rangle) = \gamma(\langle M; E \rangle)$, and $\langle M'; E' \rangle \sqsubseteq \langle M; E \rangle$. \square

Weak Reduce for $\mathcal{E}_{\mathbb{Z}_{2^w}} \star \mathcal{J}_{\mathbb{Z}_{2^w}}$. Alg. 7 computes a weak semantic-reduction for the bit-vector equality domain and the bit-vector interval domain.

Algorithm 7 Algorithm for weak semantic-reduction for $\mathcal{E}_{\mathbb{Z}_{2^w}}[P \uplus S] \star \mathcal{J}_{\mathbb{Z}_{2^w}}[S]$.

```

1:  $E' \leftarrow E \downarrow_S$ 
2: for  $s \in S$  do
3:   for constraint  $(s = c + \sum_{t \in S \wedge t \neq s} \alpha_t t) \in E'$  do
4:      $\iota \leftarrow \llbracket c + \sum_{t \in S \wedge t \neq s} \alpha_t t \rrbracket(I)$ 
5:      $I \leftarrow I \sqcap_{\mathcal{J}_{\mathbb{Z}_{2^w}}} \top[s \in \iota]$ 
6: return  $\langle E; I \rangle$ 

```

Given $\langle E; I \rangle \in \mathcal{E}_{\mathbb{Z}_{2^w}}[P \uplus S] \star \mathcal{J}_{\mathbb{Z}_{2^w}}[S]$, $\text{Reduce}_{\mathcal{E}_{\mathbb{Z}_{2^w}} \star \mathcal{J}_{\mathbb{Z}_{2^w}}}(\langle E; I \rangle)$ infers tighter interval bounds on the variables in S . Reduction is performed by first projecting E onto S to determine the affine relations E' that hold among the variables in S (line 1). These affine relations are used in lines 2–5 to infer tighter intervals for the variables in S . “ $\llbracket expr \rrbracket$ ” denotes the evaluation of expression $expr$ over interval domain $\mathcal{J}_{\mathbb{Z}_{2^w}}$ via interval arithmetic. New interval constraints are identified by evaluating expressions of the form $c + \sum_{t \in S \wedge t \neq s} \alpha_t t$ over $\mathcal{J}_{\mathbb{Z}_{2^w}}$ using the current interval value I (line 4). These constraints are then incorporated into I via meet (line 5).

Example 5.17. Let $P := \{x, y\}$, and $S := \{s_1, s_2\}$, where the bit-width of the bit-vector is 4. Consider $\langle E; I \rangle := \langle s_1 = 2x + 2y, s_2 = x + y; s_1 \in [4, 9], s_2 \in [3, 5] \rangle$. By projecting E onto the variables $S := \{s_1, s_2\}$, we obtain $E' = \{s_1 = 2s_2\}$ on line 1. On line 4, we have $\iota = \llbracket 2s_2 \rrbracket(I) = [6, 10]$. Using this equation, I is updated on line 5; that is, $I = \{s_1 \in [4, 9], s_2 \in [3, 5]\} \sqcap \{s_1 \in [6, 10], s_2 \in \top\} = \{s_1 \in [6, 9], s_2 \in [3, 5]\}$. No further reduction is possible.

Thus, $\text{Reduce}_{\mathcal{E}_{\mathbb{Z}_{2^w}} \star \mathcal{J}_{\mathbb{Z}_{2^w}}}(\langle E; I \rangle) = \langle s_1 = 2x + 2y, s_2 = x + y; s_1 \in [6, 9], s_2 \in [3, 5] \rangle$. Note that Alg. 7 is not guaranteed to deduce the best possible interval constraints. The best possible reduction of $\langle E; I \rangle$ is $\langle s_1 = 2x + 2y, s_2 = x + y; s_1 \in [6, 8], s_2 \in [3, 4] \rangle$. \square

Theorem 5.18. [Soundness of Alg. 7] Let $\langle E'; I' \rangle = \text{Reduce}_{\mathcal{E}_{\mathbb{Z}_{2^w}} \star \mathcal{J}_{\mathbb{Z}_{2^w}}}(\langle E; I \rangle)$. Then $\gamma(\langle E'; I' \rangle) = \gamma(\langle E; I \rangle)$, and $\langle E'; I' \rangle \sqsubseteq \langle E; I \rangle$. \square

Program Name	Measures of Size				Performance (sec.)		Precision ($\mathcal{BVJ} \sqsubset \mathcal{E}_{\mathbb{Z}_2^w}$)	
	Instrs	CFGs	BBs	Branches	$\mathcal{E}_{\mathbb{Z}_2^w}$	\mathcal{BVJ}	Control-Point Invariants	Procedure Summaries
finger	532	18	298	48	185	639	24/48	3/18
subst	1093	16	609	74	303	1151	11/74	3/16
label	1167	16	573	103	236	986	24/103	2/16
chkdsk	1468	18	787	119	631	1675	12/119	3/18
convert	1927	38	1013	161	441	1744	101/161	0/38
route	1982	40	931	243	749	2497	78/243	2/39
comp	2377	35	1261	224	849	2740	22/224	0/35
logoff	2470	46	1145	306	861	3253	124/306	13/46

Table 5.1: Machine-code analysis using \mathcal{BVJ} . Columns 6–9 show the times (in seconds) for the $\mathcal{E}_{\mathbb{Z}_2^w}$ -based analysis, and for the \mathcal{BVJ} -based analysis; and the degree of improvement in precision measured as the number of control points at which \mathcal{BVJ} -based analysis gave more precise invariants compared to $\mathcal{E}_{\mathbb{Z}_2^w}$ -based analysis, and the number of procedures for which \mathcal{BVJ} -based analysis gave more precise summaries compared to $\mathcal{E}_{\mathbb{Z}_2^w}$ -based analysis.

5.6 Experimental Evaluation

In this section, we compare the performance and precision of the bit-vector equality domain $\mathcal{E}_{\mathbb{Z}_2^w}$ with that of the bit-vector inequality domain \mathcal{BVJ} . The abstract transformers for the \mathcal{BVJ} domain were synthesized using the approach given in Section 6. The weak semantic-reduction operator described in Alg. 7 was used in the implementation of the compose operation needed for interprocedural analysis.

Experimental Setup. We analyzed a corpus of Windows utilities using the WALi [50] system for weighted pushdown systems (WPDSs). Tab. 5.1 lists several size parameters of the examples (number of instructions, procedures, basic blocks, and branches).¹ The weight on each WPDS rule encodes the abstract transformer for a basic block B of the program, including a jump or branch to a successor block. A formula φ_B is created

¹Due to the high cost of the WPDS construction, all analyses excluded the code for libraries. Because register `eax` holds the return value from a call, library functions were modeled approximately (albeit unsoundly, in general) by “`havoc(eax)`”.

that captures the concrete semantics of B , and then the weight for B is obtained by performing $\hat{\alpha}(\varphi_B)$. We used EWPDS merge functions [56] to preserve caller-save and callee-save registers across call sites. The post query used the FWPDS algorithm [55].

View-selection Heuristic. Given the set of machine registers P , the view-constraints C_s were computed as $C_s := \bigcup_{r_i \in P} \{s_{i1} = r_i, s_{i2} = r_i + 2^{31}\}$. In particular, the view-variable s_{i2} allows us to keep track of whether r_i , when treated as a signed value, is less-than or equal 0.

Performance. Columns 6 and 7 of Tab. 5.1 list the time taken, in seconds, for $\mathcal{E}_{\mathbb{Z}_{2^w}}$ -based analysis, and the \mathcal{BVJ} -based analysis. On average (geometric mean), \mathcal{BVJ} -based analysis is about 3.5 times slower than $\mathcal{E}_{\mathbb{Z}_{2^w}}$ -based analysis.

Precision. We compare the procedure summaries, and the invariants for each control point—i.e., the point just before a branch instruction. Column 8 lists the number of control points at which \mathcal{BVJ} -based analysis gave more precise invariants compared to $\mathcal{E}_{\mathbb{Z}_{2^w}}$ -based analysis. \mathcal{BVJ} -based analysis gives more precise invariants at up to 63% of control points, and, on average, \mathcal{BVJ} -based analysis improves precision for 29% of control points. Column 9 lists the number of procedures for which \mathcal{BVJ} -based analysis gave more precise summaries compared to $\mathcal{E}_{\mathbb{Z}_{2^w}}$ -based analysis. \mathcal{BVJ} -based analysis gives more precise summaries for up to 17% of procedures, and, on average \mathcal{BVJ} -based analysis gave better summaries for 9.3% of procedures.

5.7 Related Work

Other work on identifying bit-vector-inequality invariants includes Brauer and King [11, 12] and Masdupuy [66]. Masdupuy proposed a relational

abstract domain of interval congruences on rationals. One limitation of his machinery is that the domain represents diagonal grids of parallelepipeds, where the dimension of each parallelepiped equals the number of variables tracked (say n). In our work, we can have any number of view-variables, which means that the point-spaces represented can be constrained by more than n constraints.

Brauer and King employ bit-blasting to synthesize abstract transformers for the interval and octagon [68] domains. One of their papers uses universal-quantifier elimination on Boolean formulas [11]; the other avoids quantifier elimination [12]. Compared with their work, we avoid the use of bit-blasting and work directly with representations of sets of w -bit bit-vectors. The greatly reduced number of variables that comes from working at word level opens up the possibility of applying our methods to much larger problems; as discussed in §5.6, we were able to apply our methods to interprocedural program analysis. The equality domain $\mathcal{E}_{\mathbb{Z}_2^w}$ that we work with can capture relations on an arbitrary number of variables, and thus so can the domains that we construct using \mathcal{V} . Compared with octagons, which are limited to two variables and coefficients of ± 1 , the advantage is that our domains can express more interesting invariants and procedure summaries. In particular, Octagon-based summaries would be limited to one pre-state variable and one post-state variable.

The view-product combinator \mathcal{V} is a compositional generalization of the construction used in SubPoly. SubPoly is constructed as a reduced product of an affine-equality domain \mathcal{K} over rationals [49], and an interval domain \mathcal{J} over rationals [19] with slack view-variables S to communicate between the two domains. SubPoly can be constructed by applying \mathcal{V} to \mathcal{K} and \mathcal{J} , as follows:

$$\text{SubPoly}[P \uplus S] \stackrel{\text{def}}{=} \mathcal{V}[\mathcal{K}[P], S, C_s, \mathfrak{J}], \quad (5.1)$$

where C_s is the view-constraint for S and $\mathfrak{J}(V_1, V_2) \stackrel{\text{def}}{=} \mathcal{J}[V_1]$.

When an analysis system works with two or more reasoning techniques, there is often an opportunity to share information to improve the precision of both. The principle is found in the classic papers of Cousot and Cousot [21] and Nelson and Oppen [80]. In practice, there are a range of choices as to what might be shared, and our work represents one point in that design space. The algorithms for weak semantic-reduction (Alg. 6 and 7) adapt techniques for theory combination that have been used in Satisfiability Modulo Theory (SMT) solvers [27, 33].

The work of Chang and Leino [16] is similar in spirit to ours. They developed a technique for extending the properties representable by a given abstract domain from schemas over variables to schemas over terms. To orchestrate the communication of information between domains, they designed the congruence-closure abstract domain, which introduces variables to stand for subexpressions that are alien to a base domain; to the base domain, these expressions are just variables. Their scheme for propagating information between domains is mediated by the e-graph of the congruence-closure domain. In contrast, our method can make use of past work on synthesizing best abstract operations [84, 100] to propagate information between domains. As discussed in §5.5, we also employ less precise, more pragmatic procedures that use information in one domain to iteratively refine information in another domain. Cousot et al. [22] have recently studied the iterated pairwise exchange of observations between components as a way to compute an overapproximation of a reduced product.

5.8 Chapter Notes

Aditya Thakur and Thomas Reps supervised me in this work. Aditya Thakur was instrumental in providing the key insights and helped me in the implementation of the view-combinator and reduced-product framework.

6 SOUND BIT-PRECISE NUMERICAL DOMAINS

FRAMEWORK FOR INEQUALITIES

This chapter tackles the challenges of implementing a bit-precise relational domain capable of expressing program invariants. The chapter describes the design and implementation of a new framework for abstract domains, called the *Bit-Vector-Sound Finite-Disjunctive (BVSFD)* domains, which are capable of capturing useful program invariants such as inequalities over bit-vector-valued variables.

The *BVJ* domain in Chapter 5 can handle certain kind of bit-vector inequalities, but it needs the client to provide a template for inequalities. Moreover, the domain is incapable of expressing simple inequalities of form $x \leq y$, because they have variables on both side of the inequality.

Simon et al. [96] introduced sound wrap-around to ensure that the polyhedral domain is sound over bit-vectors. The wrap-around operation is called selectively on the abstract-domain elements while calculating the fixpoint. The operation is called selectively to preserve precision while not compromising on soundness with respect to the concrete semantics. The Verasco static analyzer [48] provides a bit-precise parameterized framework for abstract domains using the wrap-around operation. This approach has two disadvantages:

- The wrap-around operation almost always loses information due to calls on join: the convex hull of the elements that did not overflow with those that did usually does not satisfy many inequality constraints.
- They do not show how to create abstract transformers automatically.

We introduce a class of abstract domains, called $BVS(\mathcal{A})$, that is sound with respect to bitvectors whenever \mathcal{A} is sound with respect to mathematical integers. The \mathcal{A} domain can be any numerical abstract domain. For example, it can be the polyhedral domain, which can represent useful

program invariants as inequalities. We also describe how to create abstract transformers for $BVS(\mathcal{A})$ that are sound with respect to bitvectors. For $v \subseteq \text{Var}$ and $av \in \mathcal{A}$, we denote the result by $WRAP_v(av)$; the operation performs wraparound on av for variables in v . We give an algorithm for $WRAP_v(av)$ that works for any relational abstract domain (see §6.3.1). We use a finite number of disjunction of \mathcal{A} elements—captured in the domain $\mathcal{FD}_d(\mathcal{A})$ —to help retain precision. The finite disjunctive domain is parameterized by the maximum number of disjunctions allowed in the domain (referred to as d). Note that $d=1$ is the same as convex polyhedra with wrap-around [96].

Problem statement.

Given a relational numeric domain over integers, capable of expressing inequalities, (i) provide an automatic method to create a relational abstract domain that can capture inequalities over bit-vector-valued variables; (ii) create sound bit-precise abstract transformers; and (iii) use them to identify inequality invariants over a set of program variables.

Related work and contributions. Our work incorporates a number of ideas known from the literature, including

- the use of relational abstract domains [23, 71, 58, 69, 97, 82] that are sound over mathematical integers and capable of expressing inequalities.
- the use of a wrap-around operation [96, 13, 3] to ensure that the abstraction is sound with respect to the concrete semantics of the bitvector operations.
- the use of finite disjunctions [88, 4, 37] over abstract domains to obtain more precision.

- the use of instruction reinterpretation [47, 79, 81, 64, 60, 30] to obtain an abstract transformer automatically for an edge from a basic block to its successor.

Our contribution is that we put all of these to work together in a parameterized framework, along with a mechanism to increase precision by performing wrap-around on abstract values lazily.

- We propose a framework for abstract domains, called $BVSFD_d(\mathcal{A})$, to express bit-precise relational invariants by performing wrap-around over abstract domain \mathcal{A} and using disjunctions to retain precision. This abstract domain is parameterized by a positive value d , which provides the maximum number of disjunctions that the abstract domain can make use of.
- We provide a generic technique via reinterpretation to create the abstract transformer for the path through a basic block to a given successor, such that the transformer incorporates lazy wrap-around.
- We present experiments to show how the performance and precision of $BVSFD_d$ analysis changes with the tunable parameter d .

§6.1 introduces the terminology and notation used in the rest of the chapter. §6.2 demonstrates our framework with the help of an example. §6.3 introduces the $BVSFD_d$ abstract-domain framework, and formalizes abstract-transformer generation for the framework. §6.4 presents experimental results.

6.1 Terminology

$\mathcal{A}[V]$ denotes the specific instance of \mathcal{A} that is defined over vocabulary V , where V is a tuple of variables (v_1, v_2, \dots, v_k) . Each variable v_i also has an associated size in bits, denoted by $s(v_i)$. The domain \mathcal{C} is the powerset of the set of concrete states.

6.1.1 Concretization.

Given an abstract value $A \in \mathcal{A}[V]$, where V consists of n variables (v_1, v_2, \dots, v_k) , the concretization of A , denoted by $\gamma_{\mathcal{A}[V]}(A)$, is the set of concrete states covered by A .

A concrete state σ is a mapping from variables to their concrete values, $\sigma : V \rightarrow \prod_{v \in V} BV^{s(v)}$, where $s(v)$ is the size of variable v in bits and BV^b is a bitvector with b bits.

$$\gamma_{\mathcal{A}[V]}(A) = \bigcup_{(a_1, a_2, \dots, a_k) \in A} \mu_V(bv_1, bv_2, \dots, bv_k), \text{ where } bv_i = a_i \% 2^{s(v_i)} \text{ for } i \in 0..n$$

$\mu_V(bv_1, bv_2, \dots, bv_k)$ takes a tuple of bitvectors corresponding to vocabulary V and returns all concrete stores where each variable v_i in V has the value bv_i .

6.1.2 Wrap-around operation.

The wrap-around operation, denoted by $WRAP_{V'}^{ty}(A)$, takes an abstract-domain value $A \in \mathcal{A}[V]$, a subset V' of the vocabulary V , and the desired type ty for vocabulary V' . It returns an abstract value A' such that the wrap-around behavior of the points in A is soundly captured. This operation is performed by displacing the concrete values in $\gamma_{\mathcal{A}[V]}(A)$ that are outside the bitvector range for ty in vocabulary V' to the correct bitvector range by appropriate linear transformations.

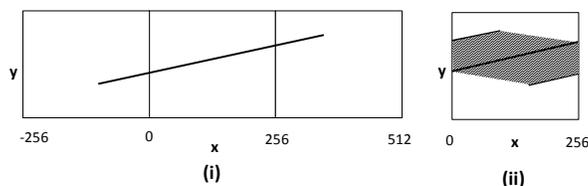


Figure 6.1: Wrap-around on variable x , treated as an unsigned char.

For example, in Fig. 6.1, the result of calling wrap-around on the line

$$\begin{aligned}
\langle \text{ELang} \rangle &:: (\text{Block})^* \\
\langle \text{Block} \rangle &:: \text{!} : (\langle \text{Stmt} \rangle ;)^* \langle \text{Next} \rangle \\
\langle \text{Next} \rangle &:: \mathbf{jump \text{!}} ; \\
&| \mathbf{if } v \langle \text{Op} \rangle_{\langle \text{Type} \rangle} \langle \text{Expr} \rangle \mathbf{ then jump \text{!}} ; \langle \text{Next} \rangle \\
\langle \text{Op} \rangle &:: < | \leq | = | \neq | \geq | > \\
\langle \text{Expr} \rangle &:: n | n * v + \langle \text{Expr} \rangle \\
\langle \text{Stmt} \rangle &:: v = \langle \text{Expr} \rangle \\
&| v : \langle \text{Type} \rangle = v : \langle \text{Type} \rangle \\
\langle \text{Type} \rangle &:: (\mathbf{uint} | \mathbf{int}) \langle \text{Size} \rangle \\
\langle \text{Size} \rangle &:: \mathbf{1} | \mathbf{2} | \mathbf{4} | \mathbf{8}
\end{aligned}$$

in (i) for variable x leads to an abstract value that is the abstraction of the points in the three line segments in (ii). For the abstract domain of polyhedra, that abstraction is the shaded area in (ii).

6.1.3 Soundness.

An abstract value $A \in \mathcal{A}[V]$ is sound with respect to a set of concrete values $C \in \mathcal{C}$, if $\gamma_{\mathcal{A}[V]} \supseteq C$.

6.1.4 $\mathcal{L}(\text{ELang})$: A Concrete language featuring finite integer arithmetic.

We borrow the simple language featuring finite-integer arithmetic from § 2.1 of [96] (with minor syntactic changes). An ELang program is a sequence of basic blocks with execution starting from the first block. Each basic block consists of a sequence of statements and a list of control-flow instructions.

The statements are restricted to an assignment of a linear expression or a cast operation. The control-flow instruction consists of either a jump statement, or a conditional that is followed by more control-flow instructions. The assignment and condition instructions expect the variable and the expression involved to have the same type.

6.2 Overview

In this section, we motivate and illustrate the design of our analysis using the *BVSFD* domain.

Example 6.1. *This example illustrates a function f that takes two 32-bit integers x and y at different rates, and resets their values to zero in case y is negative or overflows to a negative value. The function summary that we would like to obtain states that the relationship $x' \leq y'$ holds. Here, the unprimed and primed variables denote the pre-state vocabulary variables and the post-state vocabulary variables, respectively.*

```

L0: f(int x, int y) {
L1:   assume(x<=y)
L2:   while(*) {
L3:     if(*)
L4:       x=x+1, y=y+1
L5:     y=y+1
L6:     if(y<=0)
L7:       x=0, y=0
L8:   }
END: }

```

This example illustrates that merely detecting overflow would not be useful to assert the $x \leq y$ relationship at the end of the function. \square

6.2.1 Creation of Abstract Transformers

Consider the analysis for Ex. 6.1 with the abstract domain $BVSFD_2(\mathcal{OC}\mathcal{T})$. The first step involves constructing the abstraction of the concrete operations in the program as abstract transformers.

For instance, the abstract transformer for the concrete operations starting from node L0 and ending at node L2, denoted by $\tau_{L0 \rightarrow L2}^\sharp$, is defined as

$$\{m \leq x, y \leq M \wedge x' = x \wedge y' = y \wedge x' \leq y'\},$$

where m and M represent the minimum and maximum values for a signed 32-bit integer, respectively. The constraints $\{m \leq x, y \leq M\}$ are the bounding constraints on the pre-state vocabulary that are added because L0 is the entry point of the function, and the function expects three 32-bit signed values x and y as input. The equality constraints $\{x' = x, y' = y\}$ specify that the variables x and y are unchanged. Finally, the constraint $\{x' \leq y'\}$ is added as a consequence of the assume call.

Now consider other concrete transformations, such as $L4 \rightarrow L5$ and $L5 \rightarrow L6$. For the transformation $L4 \rightarrow L5$, the values for x' and y' might overflow because of the increment operations at L4. Consequently, the value of the incoming variable y in the transformation $L5 \rightarrow L6$ might have overflowed as well. There are two ways to design the abstract transformer to deal with these kind of scenarios: 1) a naive eager approach, 2) a lazy approach.

Eager Abstract Transformers.

In the naive eager approach, the abstract transformers are created such that the pre-state vocabulary is always bounded as per the type requirements. For this example, that would mean that the pre-state vocabulary variables x and y are bounded in the range $[m, M]$. Consequently, the abstract transformers for $L4 \rightarrow L5$ and $L5 \rightarrow L6$ are:

- $\tau_{L4 \rightarrow L5}^{\#E} = \{m \leq x, y \leq M \wedge x' = x + 1 \wedge y' = y + 1\}$
- $\tau_{L5 \rightarrow L6}^{\#E} = \{m \leq x, y \leq M \wedge x' = x \wedge y' = y + 1\}$.

Because the eager approach expects the pre-state vocabulary to be bounded, an abstract-composition operation $\alpha_1 \circ \alpha_2$, where α_1 and α_2 are abstract transformers, needs to call the *WRAP* operation (§6.1.2) for the entire post-state vocabulary of α_2 , for correctness. For instance, let $\alpha_1 = \{m \leq u \leq M \wedge u' = u\}$ and $\alpha_2 = \{m \leq u \leq M \wedge u' = M + 1\}$. The abstract transformer α_1 preserves u and α_2 changes the value of u' to $M + 1$. The composition of these operations matches the pre-state vocabulary of α_1

with the post-state vocabulary of α_2 , by renaming them to the same temporary variables and performing a meet. For this particular example, it will perform $\{m \leq u'' \leq M \wedge u' = u''\} \sqcap WRAP_{u''}(\{m \leq u \leq M \wedge u'' = M+1\})$, where it has matched the pre-state vocabulary variable u of α_1 with the post-state vocabulary variable u' of α_2 , by renaming them both to a temporary variable u'' . Note that in the absence of the *WRAP* operation on the post-state vocabulary of α_2 , the meet operation above will return the empty element \perp . This result would be unsound because the value of u' in α_2 should have overflowed to m .

Now consider the composition $\tau_{L5 \rightarrow L6}^{\#E} \circ \tau_{L4 \rightarrow L5}^{\#E}$. After matching, composition will perform the meet of:

- $\{m \leq x'', y'' \leq M \wedge x' = x'' \wedge y' = y'' + 1\}$
- $WRAP_{\{x'', y''\}}\{m \leq x, y \leq M \wedge x'' = x + 1 \wedge y'' = y + 1\}$

The result of *WRAP* will be a join of four values. The four values are the combinations of cases where x'' and y'' might or might not overflow. As a result, the final composition will give an abstract transformer that overapproximate those four values. For $BVSFD_2(\text{OCT})$, it will result in a loss of precision because it cannot express the disjunction of these four values precisely.

Lazy Abstract Transformers.

The eager approach to creating abstract transformers forces a call to *WRAP* at each compose operation. The lazy approach can avoid unnecessary calls to *WRAP* by not adding any bounding constraints to the pre-state vocabulary. The abstract transformer $\tau_{L4 \rightarrow L5}^{\#L}$ is defined as $\{x' = x + 1 \wedge y' = y + 1\}$ and $\tau_{L5 \rightarrow L6}^{\#L}$ is defined as $\{x' = x \wedge y' = y + 1\}$. The abstract transformer $\tau_{L4 \rightarrow L5}^{\#L}$ is sound, because the concretization of the abstract transformer, denoted by $\gamma(\tau_{L4 \rightarrow L5}^{\#L})$, overapproximates the collecting concrete semantics for $L4 \rightarrow L5$ (see proposition 1 in [96]). A similar argument can be made for the abstract transformer $\tau_{L5 \rightarrow L6}^{\#L}$. The composition $\tau_{L5 \rightarrow L6}^{\#L} \circ \tau_{L4 \rightarrow L5}^{\#L}$

gives $\{x' = x + 1 \wedge y' = y + 2\}$. Thus, the lazy abstract transformer can retain precision by avoiding unnecessary calls to *WRAP*. However, one cannot avoid calling *WRAP* for every kind of abstract transformer and still maintain soundness. Consider the abstract transformer $\tau_{L5 \rightarrow L7}^{\#L}$, which is similar to $\tau_{L5 \rightarrow L6}^{\#L}$, but must additionally handle the branch condition for $L6 \rightarrow L7$. Defining it in a similar vein as $L4 \rightarrow L5$ will result in $\{x' = x \wedge y' = y + 1 \wedge y' \leq 0\}$. While the first three constraints are sound, the fourth constraint representing the branch condition, is unsound with respect to the concrete semantics. The reason is that y' might overflow to a negative value, in which case the condition evaluates to true and the branch to $L7$ is taken. However, the abstract transformer does not capture that behavior and is, therefore, unsound with respect to the concrete semantics. To achieve soundness in the presence of a branch condition, the following steps are performed for each variable v' involved in a branch condition:

- A backward dependency analysis is performed to find the subset V_b of the pre-state vocabulary on which v' depends. For the edge $L5 \rightarrow L6$, only y' is involved in a branch condition. The backward-dependency analysis yields $V_b = \{y\}$, because the only pre-state vocabulary variable that y' depends on is y .
- The bounding constraints for the vocabulary V_b are added to the abstract transformer. For example, after adding bounding constraints, $\tau_{L5 \rightarrow L6}^{\#L}$ becomes $\{m \leq y \leq M \wedge x' = x \wedge y' = y + 1 \wedge y' \leq 0\}$.
- The wrap operation is called on the abstract transformer for the variable v' . For the edge $L5 \rightarrow L6$, this step will soundly set the abstract transformer to the disjunction of
 - $\{m \leq y \leq M - 1 \wedge x' = x \wedge y' = y + 1 \wedge y' \leq 0\}$
 - $\{y = M - 1 \wedge x' = x \wedge y' = m\}$

Note that the abstract domain $BVSFD_2(\mathcal{OC}\mathcal{T})$ can precisely express the above disjunction because the number of disjunctions is ≤ 2 .

Our analysis uses the lazy approach to create abstract transformers, because it can provide more precise function summaries. Fig. 6.2 illustrates the lazy abstract transformers generated for Ex. 6.1.

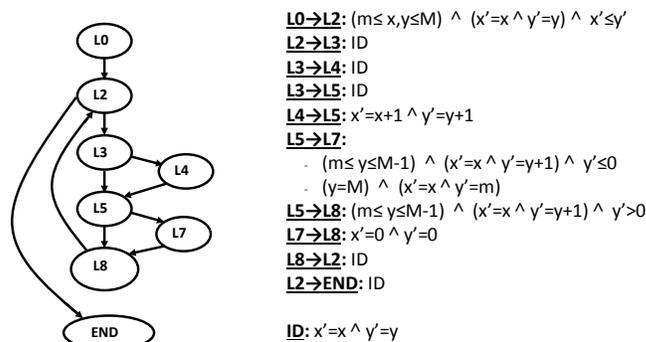


Figure 6.2: Lazy abstract transformers with the $BVSFD_2(\mathcal{POLY})$ domain for Ex. 6.1. ID refers to the identity transformation.

6.2.2 Fixed-point computation

To obtain function summaries, an iterative fixed-point computation needs to be performed. Tab. 6.1 provides some snapshots of the fixpoint analysis with the $BVSFD_2(\mathcal{OCT})$ domain for Ex. 6.1.

To simplify the discussion, we focus only on three program points: L2, L5, and L7. Each row in the table shows the intermediate value of path summaries from L0 to each of the three program points. Quiescence is discovered during the fifth iteration. The abstract value in row (i) and column L2 shows the intermediate path summary for L2 calculated after one iteration of the analysis. It states that the pre-state vocabulary variables x and y are bounded and neither of them has been modified, because at this iteration the analysis has not considered the paths that go through the loop. At row (i) and column L5, the domain precisely captures the disjunction of two paths arising at the conditional at L3. The abstract

Table 6.1: Snapshots in the fixed-point analysis for Ex. 6.1 using the $BVSFD_2(\mathcal{OC}\mathcal{T})$ domain. B_{v_1, v_2, \dots, v_k} are the bounding constraints for the variables v_1, v_2, \dots, v_k .

Node	L2	L5	L7
(i)	<ul style="list-style-type: none"> $B_{x,y} \wedge (x' = x) \wedge (y' = y)$ 	<ul style="list-style-type: none"> $B_{x,y} \wedge (x' = x) \wedge (y' = y)$ $B_{x,y} \wedge (x' = x + 1) \wedge (y' = y + 1)$ 	<ul style="list-style-type: none"> $B_{x,y} \wedge (m \leq y' \leq 0) \wedge (x \leq x' \leq x + 1) \wedge (y \leq y' \leq y + 2) \wedge (x' \leq y')$ $B_{x,y} \wedge (m \leq y' \leq m + 1) \wedge (x \leq x' \leq x + 1)$
(ii)	<ul style="list-style-type: none"> $B_{x,y,y'} \wedge (x \leq x' \leq x + 1) \wedge (y \leq y' \leq y + 2) \wedge (x' \leq y')$ $B_{x,y} \wedge (x' = 0) \wedge (y' = 0)$ 	<ul style="list-style-type: none"> $B_{x,y} \wedge (x \leq x' \leq x + 2) \wedge (y \leq y' \leq y + 3) \wedge (x' \leq y')$ $B_{x,y} \wedge (0 \leq x' \leq 1) \wedge (y' = x')$ 	<ul style="list-style-type: none"> $B_{x,y} \wedge (m \leq y' \leq 0) \wedge (x \leq x' \leq x + 2) \wedge (y \leq y' \leq y + 4) \wedge (x' \leq y')$ $B_{x,y} \wedge (m \leq y' \leq m + 3) \wedge (x \leq x' \leq x + 2)$
(iii)	<ul style="list-style-type: none"> $B_{x,y,y'} \wedge (x \leq x' \leq x + 2) \wedge (y \leq y' \leq y + 4) \wedge (x' \leq y')$ $B_{x,y} \wedge (0 \leq x' \leq 1) \wedge (0 \leq y' \leq 2) \wedge (x' \leq y')$ 	<ul style="list-style-type: none"> $B_{x,y} \wedge (x \leq x' \leq x + 3) \wedge (y \leq y' \leq y + 5) \wedge (x' \leq y')$ $B_{x,y} \wedge (0 \leq x' \leq 2) \wedge (0 \leq y' \leq 4) \wedge (x' \leq y')$ 	<ul style="list-style-type: none"> $B_{x,y} \wedge (m \leq y' \leq 0) \wedge (x \leq x' \leq x + 3) \wedge (y \leq y' \leq y + 6) \wedge (x' \leq y')$ $B_{x,y} \wedge (m \leq y' \leq m + 5) \wedge (x \leq x' \leq x + 3)$
(iv)	<ul style="list-style-type: none"> $B_{x,y,y'} \wedge (x \leq x') \wedge (y \leq y') \wedge (x' \leq y')$ $B_{x,y,y'} \wedge (0 \leq x') \wedge (x' \leq y')$ 	<ul style="list-style-type: none"> $B_{x,y} \wedge (x \leq x') \wedge (y \leq y') \wedge (x' \leq y')$ $B_{x,y} \wedge (0 \leq x') \wedge (x' \leq y')$ 	<ul style="list-style-type: none"> $B_{x,y} \wedge (m \leq y' \leq 0) \wedge (x \leq x') \wedge (y \leq y') \wedge (x' \leq y')$ $B_{x,y,y'} \wedge (m \leq y' \leq 0)$

value at row (i), column L7 is obtained as the composition of the abstract transformer $L5 \rightarrow L7$ with the path summary at L5 in row (i). The abstract-composition operations for abstract values with disjunctions performs abstract composition for all pairs of abstract transformers in the arguments, and then does a join on that set of values. To obtain the abstract transformer for L7 in row (i), it computes the join of the following values:

1. $B_x \wedge (m \leq y \leq M - 1) \wedge (x' = x) \wedge (y' = y + 1) \wedge (y' \leq 0)$
2. $B_x \wedge (m \leq y \leq M - 2) \wedge (x' = x + 1) \wedge (y' = y + 2) \wedge (y' \leq 0)$
3. $B_x \wedge (y = M) \wedge (x' = x) \wedge (y' = m)$

$$4. B_x \wedge (M \leq y \leq M - 1) \wedge (x' = x) \wedge (m \leq y' \leq m + 1)$$

Our abstract-domain framework uses a distance heuristic (see §6.3.1) to merge abstract values that are closest to each other. For this particular case, the abstract transformers (1) and (2) describe the scenarios where y' does not overflow, and are merged to give the first disjunct of row (i), column L7. Similarly, the abstract transformers (3) and (4) describe the scenarios where y' overflows, and are merged to give the second disjunct of row (i), column L7.

In the second iteration, shown in row (ii), the first disjunct for L2 is the join of the effect of first iteration of the loop, where x and y are incremented, with the old value, where x and y are unchanged. Additionally, the second disjunct in row (ii), column L2 captures the case where both x and y are set to 0 at program point L7. Iteration (iii) proceeds in a similar manner, and finally the value at L2 saturates due to widening. The value of L2 at iteration (iv) is propagated to the end of the function to give the following function summary:

- $B_{x,y,y'} \wedge (x \leq x') \wedge (y \leq y') \wedge (x' \leq y')$
- $B_{x,y,y'} \wedge (0 \leq x') \wedge (x' \leq y')$

Thus, the function summary enables us to establish that $x' \leq y'$ is true at the end of the function.

6.3 The *BVSFD* Abstract-Domain Framework

In this section, we present the intuition and formalism behind the design and implementation of the *BVSFD* abstract-domain framework.

6.3.1 Abstract-Domain Constructors

BVSFD uses of the following abstract-domain constructors:

- **Bit-Vector-Sound Constructor:** This constructor, denoted by $BVS[\mathcal{A}]$, takes an arbitrary abstract domain and constructs a bit-precise ver-

Algorithm 8 Wrap for a single variable

```

1: function WRAP(a, v, ty)
2:   if a is ⊥ then
3:     return ⊥
4:   (m, M) ← Range(ty)
5:   s ← (M - m) + 1
6:   [l, u] ← GetBounds(a, v)
7:   if l ≠ -∞ ∧ u ≠ ∞ then
8:     ⟨ql, qu⟩ ← ⟨[(l - m) / s], [(u - m) / s]⟩
9:   b ← C(m ≤ v) ∩ C(v ≤ M)
10:  if l = -∞ ∨ u = ∞ ∨ (qu - ql) > t then
11:    return RM{v}(a) ∩ b
12:  else
13:    return ⋃q ∈ {ql, qu} ((a ▷ v := v - qs) ∩ b)

```

Type	Operation	Description
\mathcal{A}	\top	top element
\mathcal{A}	\perp	bottom element
bool	$(a_1 == a_2)$	equality
\mathcal{A}	$(a_1 \sqcap a_2)$	meet
\mathcal{A}	$(a_1 \sqcup a_2)$	join
\mathcal{A}	$(a_1 \nabla a_2)$	widen
\mathcal{A}	$\pi_W(a)$	project on vocabulary W
\mathcal{A}	$RM_W(a)$	remove vocabulary W
\mathcal{A}	$\rho(a_1, v_1, v_2)$	rename variable v_1 to v_2
\mathcal{A}	$C(l e_1 \text{ op } l e_2)$	construct abstract value
set[\mathcal{A}]	$WRAP_W^{ty}(a_1)$	wrap vocabulary W
\mathcal{D}	$\mathcal{D}(a_1, a_2)$	distance

Figure 6.3: Abstract-domain interface for \mathcal{A} .

sion of the domain that is sound with respect to the concrete semantics. It needs the base domain \mathcal{A} to provide a *WRAP* operator.

- **Finite-Disjunctive Constructor:** This constructor, denoted by $\mathcal{FD}_d[\mathcal{A}]$, takes an abstract domain \mathcal{A} and a parameter d , and constructs a finite-disjunctive version of the domain, where the number of disjunctions in any abstract value should not exceed d . This constructor uses a distance measure, denoted by \mathcal{D} , to determine which disjuncts are combined when the number of disjunctions exceeds d .

The $BVSFD_d[\mathcal{A}]$ domain is constructed as $BVS[\mathcal{FD}_d[\mathcal{A}]]$. Fig. 6.3 shows the interface that the base abstract domain \mathcal{A} needs to provide to instantiate the $BVSFD_d[\mathcal{A}]$ framework.

The first seven operations are standard abstract-domain operations. The remove-vocabulary operation $RM_W(\mathcal{A})$, can be implemented as $\pi_{V-W}(\mathcal{A})$, where V is the full vocabulary. The rename operation $\rho(\mathcal{A}, v_1, v_2)$ can be easily implemented in most abstract-domain implementations through simple variable renaming and/or variable-order permutation. The construct operation, denoted by C , constructs an abstract value from the linear constraint $l e_1 = l e_2$, where $l e_1$ and $l e_2$ are linear ex-

pressions, and operation $op \in \{=, \leq, \geq\}$. This operation is available for any numeric abstract domain that can capture linear constraints. If the domain cannot express a specific type of linear constraint (for instance, the octagon domain cannot express linear constraints with more than two variables), it can safely return \top . The *WRAP* operation is similar to the wrap operation in [96], except that it returns a set of abstract-domain values, whose disjunction correctly captures the wrap-around behavior. The *WRAP* operation from [96] is modified to return a *set* of abstract-domain values by placing values in a set instead of calling *join*. Alg. 6.3.1 shows how wrap is performed for a single variable. It takes the abstract value a and perform wrap-around on variable v , treated as type ty . Line 4 obtains the range for a type, and line 5 calculates the size of that range. Line 6 obtains the range of v in abstract value a . This operation can be implemented by projecting a on v and reading the resultant interval. Lines 7-8 calculate the range of the quadrants for the variable v . Line 9 computes the bounding constraints on v , treated as type ty . Line 10 compares the number of quadrants to a threshold t . If the number of quadrants exceeds t , the result is computed by removing constraints on v in a using the *RM* operation, and adding the bounding constraints to the final result. Otherwise, for each quadrant, the appropriate value is computed by displacing the quadrants to the correct range. The displacing of the abstract value a for the quadrant q , denoted by $a \triangleright v := v - qs$, is implemented as $RM_{\{u\}}(\rho(a, v, u) \sqcap \mathcal{C}(v = u - qs))$. We used $t = 16$ in the experiments reported in §6.4.

We implement $\mathcal{D}(a_1, a_2)$ by converting a_1 and a_2 into the strongest boxes b_1 and b_2 that overapproximate a_1 and a_2 , and computing the distance between b_1 and b_2 . A box is essentially a conjunction of intervals on each variable in the vocabulary. We measure the distance between two boxes as a tuple (d_1, d_2) , where d_1 is the number of incompatible intervals, and d_2 is the sum of the distances between *compatible* intervals. Two intervals are considered to be *incompatible* if one is unbounded in

a direction that the other one is not. For example, intervals $[0, \infty]$ and $[-7, \infty]$ are compatible, but $[0, 17]$ and $[-7, \infty]$, and $[-\infty, 12]$ and $[-7, \infty]$ are not. The *distance* between two compatible intervals is 0 if their intersection is non-empty; otherwise, it is the difference of the lower bound of the higher interval and the upper bound of the lower interval. For example, the distance between $[0, 11]$ and $[17, 21]$ is $(17 - 11) = 6$. Given two distances $d = (d_1, d_2)$ and $d' = (d'_1, d'_2)$, $d > d'$ iff either (i) $d_1 > d'_1$, or (ii) $d_1 = d'_1 \wedge d_2 > d'_2$. If the number of disjunctions in an abstract value exceeds parameter d , the abstract-domain constructor $\mathcal{FD}_d[\mathcal{A}]$ merges (using join) the pair of abstract-domain elements that are closest as measured by the distance measure.

6.3.2 Abstract Transformers

In this section, we describe how the abstract transformers are generated using reinterpretation (§2.3.1). The reinterpretation consists of a domain of abstract transformers $BVSFD_d[\mathcal{A}[V; V']]$, a domain of abstract integers $BVSFD_d^{\text{INT}}[\mathcal{A}[t; V]]$, and operations to lookup a variable's value in the post-state of an abstract transformer and to create an updated version of a given abstract transformer [30]. Here, V denotes the pre-state vocabulary variables, V' denotes the post-state vocabulary variables, and t denotes a temporary variable not in V or V' . Given blocks $B : [l : s_1; \dots; s_l; \text{next}]$ and $B' : [l' : t'; \dots; t_l; \text{next}]$ in an ELang program (see §6.1.4), where B' is a successor of B , reinterpretation of B can provide an abstract transformer for the transformation that starts from the first instruction in B and ends in the first instruction in B' , denoted by $B \rightarrow B'$.

Rule 1 in Fig. 6.4 specifies how abstract-transformer evaluation for basic-block pairs feeds into abstract-transformer evaluation on a sequence of statements. The evaluation on a sequence of statements starts with the identity abstract transformer, denoted by id . Rule 2 states that the abstract transformer for a sequence of instruction can be broken down into an

Basic Block:

$$\llbracket B \rightarrow B' \rrbracket_{\text{Block}}^\sharp = \llbracket [s_1; \dots; s_l; \text{next}] \rrbracket_{\text{Seq}}^\sharp(\text{id}, l') \quad (6.1)$$

$$\llbracket [s_1; \dots; s_l; \text{next}] \rrbracket_{\text{Seq}}^\sharp(a, l') = \llbracket [s_2; \dots; s_l; \text{next}] \rrbracket_{\text{Seq}}^\sharp(\llbracket [s_1] \rrbracket_{\text{Stmnt}}^\sharp a, l') \quad (6.2)$$

Control Flow:

$$\llbracket [\text{next}] \rrbracket_{\text{Seq}}^\sharp(a, l') = \llbracket [\text{next}] \rrbracket_{\text{Next}}^\sharp(a, l') \quad (6.3)$$

$$\llbracket \text{jump } l'' \rrbracket_{\text{Next}}^\sharp(a, l') = \text{if } l'' \text{ is } l' \text{ then } a \text{ else } \perp \quad (6.4)$$

$$\llbracket \text{if } v \text{ op}_{\text{type}} \text{ exp then jump } l''; \text{next} \rrbracket_{\text{Next}}^\sharp(a, l') = \quad (6.5)$$

if l'' **is** l' **then** $a \sqcap p$ **else** $a \sqcap n$, **where**

$$p = v_{\text{Int}} \text{ op}_{\text{type}} \text{ exp}_{\text{Int}}, n = v_{\text{Int}} \text{ !op } \text{ exp}_{\text{Int}},$$

$$v_{\text{Int}} = \text{LazyWrap}^{\text{type}}(\llbracket v \rrbracket_{\text{Expr}}^\sharp a), \text{exp}_{\text{Int}} = \text{LazyWrap}^{\text{type}}(\llbracket \text{exp} \rrbracket_{\text{Expr}}^\sharp a)$$

Assignments:

$$\llbracket v = \text{exp} \rrbracket_{\text{Stmnt}}^\sharp a = \text{update}(a, v', \llbracket \text{exp} \rrbracket_{\text{Expr}}^\sharp a) \quad (6.6)$$

$$\llbracket v_1 : t_1 = v_2 : t_2 \rrbracket_{\text{Stmnt}}^\sharp a = \text{if } s(t_1) \leq s(t_2) \quad (6.7)$$

then $\text{update}(a, v'_1, \llbracket v_2 \rrbracket_{\text{Expr}}^\sharp a)$

else $\text{update}(a, v'_1, \text{LazyWrap}^{t_1}(\llbracket v_2 \rrbracket_{\text{Expr}}^\sharp a))$

Expressions:

$$\llbracket n \rrbracket_{\text{Expr}}^\sharp a = \text{const_int}(n) \quad (6.8)$$

$$\llbracket v \rrbracket_{\text{Expr}}^\sharp a = \text{lookup}(v', a) \quad (6.9)$$

$$\llbracket \text{exp}_1 * \text{exp}_2 \rrbracket_{\text{Expr}}^\sharp a = \text{mult}(\llbracket \text{exp}_1 \rrbracket_{\text{Expr}}^\sharp a, \llbracket \text{exp}_2 \rrbracket_{\text{Expr}}^\sharp a) \quad (6.10)$$

$$\llbracket \text{exp}_1 + \text{exp}_2 \rrbracket_{\text{Expr}}^\sharp a = \text{add}(\llbracket \text{exp}_1 \rrbracket_{\text{Expr}}^\sharp a, \llbracket \text{exp}_2 \rrbracket_{\text{Expr}}^\sharp a) \quad (6.11)$$

Abstract Integers:

$$\text{const_int}(n) = \mathcal{C}(t = n) \quad (6.12)$$

$$\text{get_const}(i) = \text{if } \pi_{\{t\}}(i) \text{ is } \{t = n\} \text{ then } (\text{true}, n) \text{ else } (\text{false}, 0), \quad (6.13)$$

$$\text{mult}(i_1, i_2) = \text{let } (b, n) = \text{get_const}(i_1) \text{ in} \quad (6.14)$$

(if b **then** $\text{RM}_{\{t', t''\}}(\rho(i_1, t, t') \sqcap \mathcal{C}(t = n * t'))$ **else** \top)

$$\text{add}(i_1, i_2) = \text{RM}_{\{t', t''\}}(\rho(i_1, t, t') \sqcap \rho(i_2, t, t'') \sqcap \mathcal{C}(t = t' + t'')) \quad (6.15)$$

$$i_1 \text{ op } i_2 = \text{RM}_{\{t', t''\}}(\rho(i_1, t, t') \sqcap \rho(i_2, t, t'') \sqcap \mathcal{C}(t' \text{ op } t'')), \quad (6.16)$$

where $\text{op} \in \{=, \leq, \geq\}$

$$i_1 > i_2 = \text{RM}_{\{t', t''\}}(\rho(i_1, t, t') \sqcap \rho(i_2, t, t'') \sqcap \mathcal{C}(t' \geq t'' + 1)) \quad (6.17)$$

$$i_1 < i_2 = \text{RM}_{\{t', t''\}}(\rho(i_1, t, t') \sqcap \rho(i_2, t, t'') \sqcap \mathcal{C}(t' \leq t'' - 1)) \quad (6.18)$$

$$i_1 \neq i_2 = (i_1 < i_2) \sqcup (i_1 > i_2) \quad (6.19)$$

$$\text{LazyWrap}^{\text{type}}(i) = \text{WRAP}_{\{i\}}^{\text{type}}(i \sqcap \mathcal{C}(B_W)), \text{where } W = \text{DependentVoc}_t(i) \quad (6.20)$$

Variable lookup and update:

$$\text{lookup}(a, v') = \pi_{V \cup \{t\}}(a \sqcap \{t = v'\}) \quad (6.21)$$

$$\text{update}(a, v', i) = \text{RM}_{\{v'\}}(a) \sqcap \text{RM}_{\{t\}}(i \sqcap \{v' = t\}) \quad (6.22)$$

Figure 6.4: Reinterpretation semantics for $\mathcal{L}(\text{ELANG})$.

abstract transformer for a smaller sequence of instruction, by recursively performing statement-level abstract interpretation $\llbracket \cdot \rrbracket_{\text{Stmnt}}^\#$ on the first instruction in the sequence. In this rule and subsequent $\llbracket \cdot \rrbracket_{\text{Next}}^\#$ and $\llbracket \cdot \rrbracket_{\text{Stmnt}}^\#$ rules, “ a ” denotes the intermediate abstract transformer value. It starts as id at the beginning of the instruction sequence, and gets updated or accessed by assignment and control-flow statements in the sequence.

Rules 3,4, and 5 handle control-flow statements. Rule 3 delegates the responsibility of executing the last instruction in the statement sequence to $\llbracket \cdot \rrbracket_{\text{Next}}^\#$. Rule 4 deals with unconditional-jump instructions. The label is checked against a goal label and either \perp or the current transformer a is returned accordingly. Rule 5 handles conditional branching. It conjoins the input transformer with p in the true case and n in the false case. p and n are calculated by performing abstract versions of op and $!\text{op}$, respectively, on the sound abstract integers corresponding to v_{Int} and v_{Exp} . Here, $!\text{op}$ denotes the negation of the op symbol. For example, negation of \leq is $>$. The sound version of an abstract integer is created by calling $\text{LazyWrap}^{\text{type}}$ (see rule 20). This function is the key component behind lazy abstract-transformer generation (see §6.2.1). In our implementation, we compute $\text{DependentVoc}_t(i)$ by looking at the constraints in i and returning the vocabulary subset $W \subseteq V$ that depend on t . B_W refers to the bounding constraints on the variables in vocabulary W .

Rules 6 and 7 handle assignment statements. Assignment to a linear expression merely performs a post-state-vocabulary update on the current abstract transformer “ a .” Note that this rule does not call WRAP even though the result of computing exp can go out of bounds. Rule 7 handles the cast operation. $s(t_1)$ and $s(t_2)$ gets the size for the types t_1 and t_2 , respectively. For a downcast operation, it performs simple update. In the case of upcast, $\text{LazyWrap}^{\text{type}}$ is called to preserve soundness (see Section 6 of [96]).

Rules 8, 9, 10, and 11 handle reinterpretation of expressions. Rules

8, 10, and 11 delegate computation to the corresponding abstract-integer operations. Rule 9 performs a variable lookup in the current value of abstract transformer “ α .”

Rules 12 to 20 deal with the operations on abstract integers in $BVSFD_d^{\text{INT}}[A[t; V]]$. Rule 12 constructs an abstract integer from a constant. Rule 13 finds out if a variable is a constant. This operation is used by abstract multiplication (Rule 14) to determine if the multiplication of two abstract integers is linear or not. Rules 14-18 use vocabulary-removal (RM) and variable-rename operations (ρ) to ensure that the vocabulary of the output is $\{t\} \cup V$.

Rules 21 and 22 are variable lookup and update operations. Lookup takes an abstract transformer $\alpha \in BVSFD_d$ and a variable $v' \in V'$, and returns the abstract integer $i \in BVSFD_d^{\text{INT}}$ such that the relationship of t with V in i is the same as the relationship of v' with V in α . The variable-update operation works in the opposite direction. Update takes an abstract transformer $\alpha \in BVSFD_d$, a variable $v' \in V'$, and $i \in BVSFD_d^{\text{INT}}$, and returns $\alpha' \in BVSFD_d$ such that the relationship of v' with V in α' is the same as the relationship of t with V in i , and all the other relationships that do not involve v' remain the same.

6.4 Experimental Evaluation

In this section, we compare the performance and precision of the bit-precise disjunctive-inequality domain $BVSFD_d$ for different values of d . We perform this comparison for the base domains of octagons and polyhedra. The abstract transformers for the $BVSFD$ domain were automatically synthesized for each path through a basic block to one of its successor by using reinterpretation (see §6.3.2). We also perform array-out-of-bounds checking to quantify the usefulness of the precision gain for different values of d . The experiments were designed to shed light on the following

questions:

1. How much does the performance of the analysis degrade as d is increased?
2. How much does the precision of the analysis increase as d is increased?
3. What is the value of d beyond which no further precision is gained?
4. What is the effect of adding sound bit-precise handling of variables on performance and precision?

6.4.1 Experimental Setup

Given a C file, we first create the corresponding LLVM bitcode [57, 63]. We then feed the bitcode to our solver, which uses the WALi [50] system to create a Weighted Pushdown System (WPDS) corresponding to the LLVM CFG. The transitions in the WPDS correspond to CFG edges from a basic block to one of its successors. The semiring weights on the edges are abstract transformers in the $BVSFD_a$ abstract domain. We then perform interprocedural analysis by performing post^* followed by the path-summary operation [85] to obtain overapproximating function summaries and an overapproximation of the reachable states at all branch points. We used EWPDS merge functions [56] to preserve local variables across call sites. We used the Pointset_Powerset [4] framework in the Parma Polyhedra Library [3, 82] to implement the $\mathcal{FD}_a[A]$ constructor. For each example used in the experiments, we use a timeout of 200 seconds.

6.4.2 Assertion Checking

For this set of experiments, we picked the subset of the SVCOMP [7] loop benchmarks for which all assertions hold. Because our analyzer is sound, we are interested in the percentage of true assertions that it can verify. Tab. 6.2 provides information about the benchmarks that we used. We

Table 6.2: Information about the loop benchmarks containing true assertions, a subset of the SVCOMP benchmarks.

Benchmark	examples	instructions	assertions
loop-invgen	18	2373	90
loop-lit	15	1173	16
loops	34	3974	32
loop-acceleration	19	1001	19
total	86	8521	158

performed the analysis on these examples and performed assertion checking by checking whether the program points corresponding to assertion failures had the bottom abstract state.

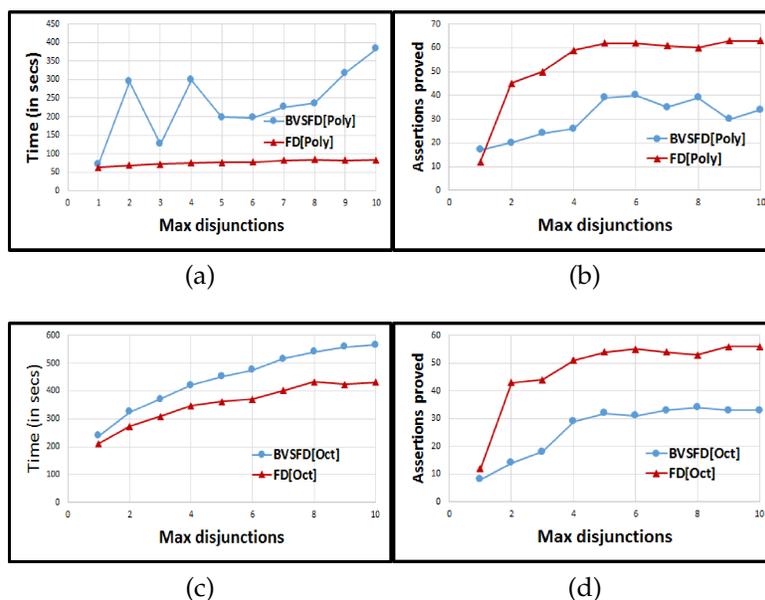


Figure 6.5: Precision and performance numbers for SV-COMP loop benchmarks.

Fig. 6.5a and Fig. 6.5b show the performance and precision numbers, respectively, for the loop SVCOMP benchmarks, with $POLY$ as the base domain. The results answer the experimental questions as follows:

1. With two exceptions, at $d = 2$ and $d = 4$, the performance steadily decreases as the number of maximum allowed disjunctions d is

increased. The analysis times for $d = 2$ and $d = 4$ do not fit the trend because one example times out for $d = 2$ or $d = 4$, but does not time out for $d = 3$ or $d = 5$. This behavior can be attributed to the non-monotonic behaviors of the finite-disjunctive join and widening operations.

2. The precision, measured as the number of proved assertions, increases from $d = 1$ to $d = 6$. From $d = 7$ onwards the change in precision is haphazard.
3. The analysis achieves the best precision at $d = 6$, where it proves 40 out of 157 assertions.
4. The sound analysis using $BVSFD_d[\mathcal{POLY}]$ is 1.1-4.6 times slower than the unsound analysis using $\mathcal{FD}_d[\mathcal{POLY}]$, and is able to prove 44-142% of the assertions obtained with $\mathcal{FD}_d[\mathcal{POLY}]$.

Fig. 6.5c and Fig. 6.5d show the performance and precision numbers, respectively, for the loop SVCOMP benchmarks, with \mathcal{OCT} as the base domain. The results answer the experimental questions as follows:

1. The performance steadily decreases with increase in d .
2. The precision, measured as the number of proved assertions, increases from $d = 1$ to $d = 5$. From $d = 5$ onwards the precision is essentially unchanged.
3. The analysis achieves the best precision at $d = 8$, where it proves 34 out of 157 assertions.
4. The sound analysis using $BVSFD_d[\mathcal{OCT}]$ is 1.1-1.3 times slower than the unsound analysis using $\mathcal{FD}_d[\mathcal{OCT}]$, and is able to prove 33-67% of the assertions obtained with $\mathcal{FD}_d[\mathcal{OCT}]$.

In our experiments, we found the $BVSFD_d[\mathcal{OCT}]$ -based analysis to be 1-3.3x slower than $BVSFD_d[\mathcal{POLY}]$. This slowdown occurs because the maximum vocabulary size in abstract transformers is ≤ 12 , and abstract operations for octagons are slower than that of polyhedra for such a small vocabulary size.

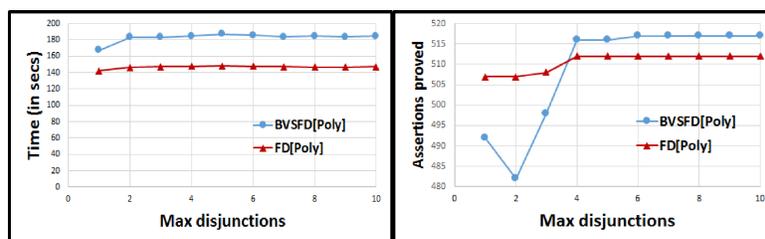


Figure 6.6: Precision and performance numbers for SV-COMP array benchmarks with \mathcal{POLY} as the base domain.

6.4.3 Array-Bounds Checking

We perform array-bound checking using invariants from the $BVSFD_d$ analysis. For each array access and update we create an error state that is reached when an array bound is violated. These array-bounds checks are verified by checking if the path summaries at the error states are \perp . There are 88 examples in the benchmark, with a total of 14,742 instructions and 598 array-bounds checks. Fig. 6.6 lists the number of array-bound checks proven for each application, for different values of d , for the SVCOMP array benchmarks.

The results answer the experimental questions as follows:

1. The performance of the analysis increases by 9% from $d = 1$ to $d = 2$. After $d = 2$, the performance stabilizes.
2. With one exception at $d = 2$, the precision—measured as the number of array-bounds checks proven—increases from $d = 1$ to $d = 4$. From $d = 4$ onwards the precision is essentially unchanged.
3. The analysis achieves the best precision at $d = 4$, where it proves 515 out of the 598 array-bounds checks.
4. The sound analysis using $BVSFD_d[\mathcal{POLY}]$ is 1.18-1.26 times slower than the unsound analysis using $\mathcal{FD}_d[\mathcal{POLY}]$, and is able to prove 95-101% of the array-bound checks obtained with $\mathcal{FD}_d[\mathcal{POLY}]$.

6.5 Chapter Notes

Thomas Reps supervised me in the writing, and provided useful comments for this work.

7 CONCLUSION AND FUTURE WORK

This dissertation described novel abstract domains and abstract-domain frameworks that can soundly provide bit-vector-sound program invariants, that can be used to verify assertions and provide loop and function summaries. This dissertation i) compared two abstract domains for bit-vector-sound equalities, ii) extended one of the domains to provide a novel framework of bit-vector-sound equalities, iii) introduced a bit-vector-sound inequality domain capable of expressing memory variables, and iv) introduced a bit-vector-sound numerical inequality domain framework. The thesis makes the following contributions:

- **Bit-Vector Precise Equality Abstract Domains**
 - **Abstract Domains of Affine Relations:** An affine relation is a *linear-equality* constraint over a given set of variables that hold machine integers. In this work, we compare the MOS and KS abstract domains, along with several variants. We show that MOS and KS are, in general, *incomparable* and give sound interconversion methods for KS and MOS. We introduce a third domain for representing affine relations, called AG, which stands for *affine generators*. Furthermore, we present an experimental study comparing the precision and performance of analyses with the KS and MOS domains.
 - **Abstraction Framework for Affine Transformers:** In this work, we define the Affine-Transformers Abstraction Framework, which represents a new family of numerical abstract domains. This framework is parameterized by a base numerical abstract domain, and allows one to represent a set of affine transformers (or, alternatively, certain disjunctions of transition formulas). Specifically, this framework is a generalization of the MOS domain. The choice of the base abstract domain allows the client

to have some control over the performance/precision trade-off.

- **Bit-Vector Precise Inequality Abstract Domains**

- **An Abstract Domain for Bit-vector Inequalities:** This work describes the design and implementation of a new abstract domain, called the *Bit-Vector Inequality* domain, which is capable of capturing certain inequalities over bit-vector-valued variables (which represent a program’s registers and/or its memory variables). This domain tracks properties of the values of selected registers and portions of memory via *views*, and provides automatic heuristics to gather equality and inequality views from the program. Furthermore, experiments are provided to show the usefulness of the Bit-Vector Inequality domain.
- **Sound Bit-Precise Numerical Domains Framework for Inequalities:** This work introduces a class of abstract domains, parameterized on a base domain, that is sound with respect to bitvectors whenever the base domain is sound with respect to mathematical integers. The base domain can be any numerical abstract domain. We also describe how to create abstract transformers for this framework that incorporate lazy wrap-around to achieve more precision, without sacrificing soundness with respect to bitvectors. We use a finite number of disjunctions of base-domain elements to help retain precision. Furthermore, we present experiments to empirically demonstrate the usefulness of the framework.

In the remainder of this chapter, we present conclusion for the thesis and future directions for the bit-vector-precise equality domains and bit-vector-precise inequality domains, respectively, described in this dissertation.

7.1 Bit-vector-precise Equality Domains

In Chapter 3, we considered MOS and KS abstract domains, along with several variants, and studied how they relate to each other. We showed that MOS and KS are, in general, *incomparable*, and introduced a third domain for representing affine relations, called AG, which stands for *affine generators*, and showed that it is the generator representation of the KS domain.

In Chapter 4, we provided analysis techniques to abstract the behavior of the program as a set of affine transformations over bit-vectors. While the relationship between MOS and KS/AG is interesting, we realized that there is a general framework that can allow one to build the MOS domain from the KS/AG domain. Thus, we generalized the ideas used in the MOS domain—in particular, to have an abstraction of *sets of affine transformers*—but to provide a way for a client of the abstract domain to have some control over the performance/precision trade-off. Toward this end, we defined a new family of numerical abstract domains, denoted by $\text{ATA}[\mathcal{B}]$. (ATA stands for Affine-Transformers Abstraction.) Following our observation, $\text{ATA}[\mathcal{B}]$ is parameterized by a base numerical abstract domain \mathcal{B} , and allows one to represent a set of affine transformers (or, alternatively, certain disjunctions of transition formulas). Specifically, MOS can be defined as $\text{ATA}[\text{KS/AG}]$.

Future Directions. One obvious opportunity for future work is to provide an implementation of the ATA framework, and compare the performance and precision of this framework for different base domains \mathcal{B} . The main challenge in this work would be to figure out the best way of generating abstract transformers for guard statements. The best method to use is not obvious because there are multiple choices for implementing guards and depending on the assertion, one choice might work better than the other. A second direction would be to extend this generic framework so

that it can be applied for sound program analysis for floating-point numbers. There has been a lot of work on tracking floating-point arithmetic soundly in an abstract domain [70, 72, 17]. Once the generic abstract-transformation-abstraction framework is implemented for floating-point numbers, these base domains can be used to provide new sound abstract domains capable of expressing a certain class of disjunctions of affine transformers over floating-point numbers. The biggest challenge is to provide floating-point-sound — yet precise — methods to perform abstract composition.

7.2 Bit-vector-precise Inequality Domains

In Chapter 5, we expanded the set of techniques available for abstract interpretation and model checking of machine code. This work described the design and implementation of a new abstract domain, called the *Bit-Vector Inequality* domain, which is capable of capturing certain inequalities over bit-vector-valued variables (which represent a program’s registers and/or its memory variables). This domain tracks properties of the values of selected registers and portions of memory via *views*, and provides automatic heuristics to gather equality and inequality views from the program. Furthermore, experiments are provided to show the usefulness of the Bit-Vector Inequality domain.

Chapter 6 described a second approach to designing and implementing a bit-precise relational domain capable of expressing inequality invariants. This work presents the design and implementation of a new framework for abstract domains, called the *Bit-Vector-Sound Finite-Disjunctive (BVSEFD)* domains, which are capable of capturing useful program invariants such as inequalities over bit-vector-valued variables. We introduced a class of abstract domains, called $BVS(\mathcal{A})$, that is sound with respect to bitvectors whenever \mathcal{A} is sound with respect to mathematical integers. The \mathcal{A} do-

main can be any numerical abstract domain. We use a finite number of disjunctions of \mathcal{A} elements—captured in the domain $\mathcal{FD}_d(\mathcal{A})$ —to help retain precision, where d is the maximum number of disjunctions that the abstract domain can make use of. We also described a generic technique to create abstract transformers for $BVS(\mathcal{A})$ that are sound with respect to bitvectors. The abstract transformers incorporate a lazy wrap-around mechanism to achieve more precision. Finally, we presented experiments to show how the performance and precision of $BVSFD_d$ analysis changes with the tunable parameter d .

Future Directions One interesting future research direction is improving the lazy wrap-around method to improve precision. Intuitively, the wrap-around operation performs join operations that leads to loss of precision. Lazy wrap-around aims to improve the precision by only calling wrap-around procedure at guard and cast statements. While calling wrap-around at these points is essential for correctness, the result of a wrap-around can be modified to a value that requires fewer disjunctions to express by using a dual of the wrap-around operation, thus leading to improved precision. The challenge is in implementing the dual operation efficiently.

§2.3.2 introduced symbolic abstraction as one of the automatic methods to create abstract transformers. However, this method is not applicable to the polyhedral domain because the polyhedral domain has both infinite ascending and descending chains. The symbolic abstraction method is desirable because the concrete semantics of bit-twiddling operations can be precisely described in Quantifier-Free Bit-Vector (QFBV) logic. Recent improvements in SMT optimization techniques [59, 8] provides insights in designing a symbolic-abstraction operation for the polyhedral domain. Moreover, with the help of the dual of wrap-around operation, sound and precise abstract transformers can be created for $BVSFD$ domain.

A DOMAIN CONVERSIONS

A.1 Soundness of MOS to AG transformation

Thm. 3.7 states that the transformation from MOS to AG given in §3.3.2 is sound. Suppose that \mathcal{B} is an MOS element such that, for every $B \in \mathcal{B}$, $B = \left[\begin{array}{c|c} 1 & c_B \\ \hline 0 & M_B \end{array} \right]$ for some $c_B \in \mathbb{Z}_{2^w}^{1 \times k}$ and $M_B \in \mathbb{Z}_{2^w}^{k \times k}$. Define $G_B = \left[\begin{array}{c|cc} 1 & 0 & c_B \\ \hline 0 & I & M_B \end{array} \right]$ and $G = \bigsqcup_{AG} \{G_B \mid B \in \mathcal{B}\}$. Then, $\gamma_{MOS}(\mathcal{B}) \subseteq \gamma_{AG}(G)$.

Proof. First, recall that for any two AG elements E and F , $E \sqcup_{AG} F$ equals $\text{HOWELLIZE} \left(\left[\begin{array}{c} E \\ F \end{array} \right] \right)$. Because HOWELLIZE does not change the row space of a matrix, $\gamma_{AG}(E \sqcup_{AG} F)$ equals $\gamma_{AG} \left(\left[\begin{array}{c} E \\ F \end{array} \right] \right)$. By the definition of G , we know that $\gamma_{AG}(G) = \gamma_{AG}(\mathbb{G})$, where \mathbb{G} is all of the matrices G_B stacked vertically. Therefore, to show that $\gamma_{MOS}(\mathcal{B}) \subseteq \gamma_{AG}(G)$, we show that $\gamma_{MOS}(\mathcal{B}) \subseteq \gamma_{AG}(\mathbb{G})$.

Suppose that $(\vec{w}, \vec{v}') \in \gamma_{MOS}(\mathcal{B})$. Then, for some vector \vec{w} ,

$$\left[\begin{array}{c|c} 1 & \vec{v} \end{array} \right] \left(\sum_{B \in \mathcal{B}} w_B B \right) = \left[\begin{array}{c|c} 1 & \vec{v}' \end{array} \right].$$

If we break this equation apart, we see that

$$\sum_{B \in \mathcal{B}} w_B = 1 \quad \text{and} \quad \sum_{B \in \mathcal{B}} w_B c_B + \vec{v} \left(\sum_{B \in \mathcal{B}} w_B M_B \right) = \vec{v}'.$$

Let \otimes denote Kronecker product. Now consider the following product, which uses $(\vec{w} \otimes \left[\begin{array}{c|c} 1 & \vec{v} \end{array} \right])$ as a vector of coefficients for the rows of \mathbb{G} :

$$\begin{aligned} (\vec{w} \otimes \left[\begin{array}{c|c} 1 & \vec{v} \end{array} \right]) \mathbb{G} &= \sum_{B \in \mathcal{B}} \left[\begin{array}{cc} w_B | w_B \vec{v} I & w_B c_B + w_B \vec{v} M_B \end{array} \right] \\ &= \left[\begin{array}{c|c} \sum_{B \in \mathcal{B}} w_B | \vec{v} \left(\sum_{B \in \mathcal{B}} w_B \right) & \sum_{B \in \mathcal{B}} w_B c_B + \vec{v} \left(\sum_{B \in \mathcal{B}} w_B M_B \right) \end{array} \right] \\ &= \left[\begin{array}{c|c} 1 | \vec{v} & \vec{v}' \end{array} \right]. \end{aligned}$$

Thus, $\left[1 \mid \vec{v} \quad \vec{v}' \right]$ is a linear combination of the rows of G , and so $(x, x') \in \gamma_{AG}(G)$. Therefore, $\gamma_{MOS}(\mathcal{B}) \subseteq \gamma_{AG}(G)$. \square

A.2 Soundness of KS Without Pre-State Guards to MOS transformation

When $G = \left[\begin{array}{c|cc} 1 & 0 & b \\ \hline 0 & I & M \\ 0 & 0 & R \end{array} \right]$, then $\gamma_{AG}(G) = \gamma_{MOS}(\text{SHATTER}(G))$.

Proof.

$$\begin{aligned} (x, x') \in \gamma_{AG}(G) &\Leftrightarrow \exists v: \left[1 \mid x \quad v \right] G = \left[1 \mid x \quad x' \right] \\ &\Leftrightarrow \exists v: b + xM + vR = x' \\ &\Leftrightarrow \exists v: \left[1 \mid x \right] \left(\left[\begin{array}{c|c} 1 & b \\ \hline 0 & M \end{array} \right] + \sum_i v_i \left[\begin{array}{c|c} 0 & R_i \\ \hline 0 & 0 \end{array} \right] \right) = \left[1 \mid x' \right] \\ &\Leftrightarrow (x, x') \in \gamma_{MOS}(\text{SHATTER}(G)) \end{aligned}$$

\square

A.3 Soundness of KS Without Pre-State Guards to MOS transformation

Lemma A.1. *Suppose that M and N are square matrices of equal dimension such that*

1. M has only ones and zeroes on its diagonal,
2. if $M_{i,i} = 1$, then $M_{h,i} = 0$ for all $h \neq i$, and
3. if $M_{i,i} = 0$, then $N_{i,h} = 0$ for all h .

Then, $MN = N$.

Proof. We know $(MN)_{i,j} = \sum_h M_{i,h}N_{h,j}$. By Items 2 and 3, if $h \neq i$ then either $M_{i,h} = 0$ or $N_{h,j} = 0$, so $(MN)_{i,j} = M_{i,i}N_{i,j}$. If $M_{i,i} = 0$, then by Item 2, $N_{i,j} = 0$; otherwise, $M_{i,i} = 1$. In either case, $(MN)_{i,j} = N_{i,j}$, as we require. \square

Lemma A.2. When $G = \left[\begin{array}{c|cc} 1 & a & b \\ 0 & J & M \\ 0 & 0 & R \end{array} \right]$, such that $\left[\begin{array}{c|c} 1 & a \\ 0 & J \end{array} \right]$ and $\left[\begin{array}{c|c} 1 & b \\ 0 & M \end{array} \right]$ satisfy the conditions of Lem. A.1, then $\gamma_{AG}(G) \subseteq \gamma_{MOS}(\text{SHATTER}(G))$.

Proof.

$$\begin{aligned} (x, x') \in \gamma_{AG}(G) &\Rightarrow \exists v, v': \left[\begin{array}{c|cc} 1 & v & v' \end{array} \right] G = \left[\begin{array}{c|cc} 1 & x & x' \end{array} \right] \\ &\Rightarrow \exists v, v': \left[\begin{array}{c|c} 1 & v \end{array} \right] \left[\begin{array}{c|c} 1 & a \\ 0 & J \end{array} \right] = \left[\begin{array}{c|c} 1 & x \end{array} \right] \\ &\quad \wedge \left[\begin{array}{c|c} 1 & v \end{array} \right] \left[\begin{array}{c|c} 1 & b \\ 0 & M \end{array} \right] + v' \left[\begin{array}{c|c} 0 & R \end{array} \right] = \left[\begin{array}{c|c} 1 & x' \end{array} \right] \end{aligned}$$

By Lem. A.1, $\left[\begin{array}{c|c} 1 & v \end{array} \right] \left[\begin{array}{c|c} 1 & b \\ 0 & M \end{array} \right] = \left[\begin{array}{c|c} 1 & v \end{array} \right] \left[\begin{array}{c|c} 1 & a \\ 0 & J \end{array} \right] \left[\begin{array}{c|c} 1 & b \\ 0 & M \end{array} \right] = \left[\begin{array}{c|c} 1 & x \end{array} \right] \left[\begin{array}{c|c} 1 & b \\ 0 & M \end{array} \right]$, so

$$\begin{aligned} (x, x') \in \gamma_{AG}(G) &\Rightarrow \exists v': \left[\begin{array}{c|c} 1 & x \end{array} \right] \left[\begin{array}{c|c} 1 & b \\ 0 & M \end{array} \right] + v' \left[\begin{array}{c|c} 0 & R \end{array} \right] = \left[\begin{array}{c|c} 1 & x' \end{array} \right] \\ &\Rightarrow \exists v': \left[\begin{array}{c|c} 1 & x \end{array} \right] \left(\left[\begin{array}{c|c} 1 & b \\ 0 & M \end{array} \right] + \sum_i v'_i \left[\begin{array}{c|c} 0 & R_i \\ 0 & 0 \end{array} \right] \right) = \left[\begin{array}{c|c} 1 & x' \end{array} \right] \\ &\Rightarrow (x, x') \in \gamma_{MOS}(\text{SHATTER}(G)) \end{aligned}$$

\square

Now we prove Thm. 3.11. For $G \in AG$, $\gamma_{AG}(G) \subseteq$

$\gamma_{MOS}(\text{SHATTER}(\text{MAKEEXPLICIT}(G)))$.

Proof. Without loss of generality, assume that G has $2k + 1$ columns and is in Howell form.

$\text{MAKEEXPLICIT}(G)$ consists of two loops. In the first loop, every row r with leading index $i \leq k + 1$ for which the rank of the leading value is greater than 0 is generalized by creating from r a row s , which is added to G , such that s 's leading index is also i , but its leading value is 1. Consequently, after the call on $\text{HOWELLIZE}(G)$ in line 8 of MAKEEXPLICIT , the leading value of the row with leading index i is 1.

In the second loop, the matrix is expanded by all-zero rows so that any row with leading index $i \leq k + 1$ is placed in row i .

Thus, for any AG element G , we can decompose $\text{MAKEEXPLICIT}(G)$ into the matrix $\left[\begin{array}{c|cc} 1 & c & b \\ \hline 0 & J & M \\ 0 & 0 & R \end{array} \right]$, where $c, b \in \mathbb{Z}_{2^w}^{1 \times k}$; $J, M \in \mathbb{Z}_{2^w}^{k \times k}$; and $R \in \mathbb{Z}_{2^w}^{r \times k}$ for some $r \leq k$. Moreover, we know that

- J is upper-triangular,
- J has only ones and zeroes on its diagonal,
- if $J_{j,j} = 1$, then column j of J is zero everywhere else, and
- if $J_{j,j} = 0$, then row j of J and row j of M are all zeroes.

By these properties, Lem. A.2 holds, so we know that

$$\gamma_{AG}(G) \subseteq \gamma_{MOS}(\text{SHATTER}(\text{MAKEEXPLICIT}(G))).$$

□

B HOWELL PROPERTIES

Definition B.1. Two module spaces R and S are **perpendicular** (denoted by $R \perp S$) if

1. $r \in R \wedge s \in S \Rightarrow rs^t = 0$,
2. $(\forall r \in R: rs^t = 0) \Rightarrow s \in S$, and
3. $(\forall s \in S: rs^t = 0) \Rightarrow r \in R$.

□

Lemma B.2. If $R \perp S$ and $R \perp S'$, then $S = S'$.

Lemma B.3. For any matrix M , $\text{row } M \perp \text{null}^t M$.

These are standard facts in linear algebra; their standard proofs essentially carry over for module spaces.

Lemma B.4. If $R \perp R'$ and $S \perp S'$, then $R + S \perp R' \cap S'$.

Proof. Pick G_R and G_S so that $\text{row } G_R = R$ and $\text{row } G_S = S$. Because the rows of a matrix are linear generators of its row space,

$$R + S = \text{row} \begin{bmatrix} G_R \\ G_S \end{bmatrix}, \quad \text{so, by Lem. B.3,} \quad R + S \perp \text{null}^t \begin{bmatrix} G_R \\ G_S \end{bmatrix}.$$

Because each row of a matrix acts as a constraint on its null space,

$$R + S \perp \text{null}^t G_R \cap \text{null}^t G_S.$$

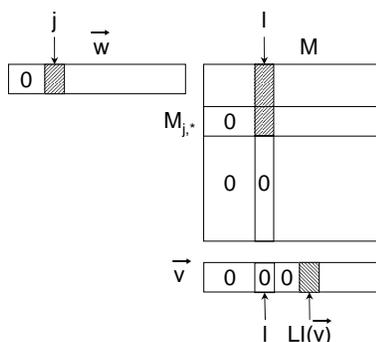
By Lem. B.3 again, we know that $\text{row } G_R \perp \text{null}^t G_R = R \perp R'$, so $\text{null}^t G_R = R'$ by Lem. B.2. Similarly, $\text{null}^t G_S = S'$. Thus, $R + S \perp R' \cap S'$. □

Note. Recall from §3.1 that $[M]_i$ is the matrix that consists of all rows of M whose leading index is i or greater. For any row r , define $LI(r)$ to be

the leading index of r . Define e_i to be the vector with 1 at index i and 0 everywhere else.

Theorem B.5. *If matrix M is in Howell form, and $x \in \text{row } M$, then $x \in \text{row}([M]_{LI(x)})$.*

Proof. Pick v so that $x = vM$, let $j \stackrel{\text{def}}{=} LI(v)$, and let $\ell \stackrel{\text{def}}{=} LI(M_{j,*})$. If $\ell \geq LI(x)$, then we already know that $x \in \text{row}([M]_{LI(x)})$. Otherwise, assume $\ell < LI(x)$. Under these conditions, as depicted in the diagram below,



- $(vM)_\ell = 0$, because $LI(vM) = LI(x) > \ell$,
- $M_{h,\ell} = 0$ for any $h > j$, by Rule 1 of Defn. 3.1, and
- $v_h = 0$ for any $h < j$, because $j = LI(v)$.

Therefore, $0 = (vM)_\ell = \sum_h v_h M_{h,\ell} = v_j M_{j,\ell}$. Thus, because $j = LI(v)$, we know that $LI(v_j M_{j,*})$ is strictly greater than $\ell = LI(M_{j,*})$.

Because multiplication by invertible values can never change nonzero values to zero, we have $LI(v_j M_{j,*}) = LI(2^{\text{rank}(v_j)} M_{j,*})$. Thus, by Rule 4 of Defn. 3.1, we know that $v_j M_{j,*}$ can be stated as a linear combination of rows $j+1$ and greater. That is, $v_j M_{j,*} \in \text{row}([M]_{j+1})$, or equivalently, $v_j M_{j,*} = uM$ with $LI(u) \geq j+1$. We can thus construct $v' = v - v_j e_j + u$ for which $x = v'M$ and $LI(v') \geq j+1$.

By employing this construction iteratively for increasing values of j , we can construct $x = yM$ with $LI(M_{LI(y),*}) \geq LI(x)$. Consequently, x can

be stated as a linear combination of rows with leading indexes $LI(x)$ or greater; i.e., $x \in \text{row}([M]_{LI(x)})$. \square

Now we prove Thm. 3.14. *Suppose that M has c columns. If matrix M is in Howell form, $x \in \text{null}^t M$ if and only if $\forall i : \forall y_1, \dots, y_{i-1} : \begin{bmatrix} y_1 & \cdots & y_{i-1} & x_i & \cdots & x_c \end{bmatrix} \in \text{null}^t([M]_i)$.*

Proof. We know that $\text{row } M \perp \text{null}^t M$, and that $\text{row}([M]_i) \perp \text{null}^t([M]_i)$. Let E_i be the module space generated by $\{e_j \mid j < i\}$, and let F_i be the module space generated by $\{e_j \mid j \geq i\}$. Clearly, $E_i \perp F_i$. By Thm. B.5, we have that $\text{row}([M]_i) = \text{row } M \cap F_i$. Thus,

$$\text{null}^t([M]_i) \perp \text{row } M \cap F_i.$$

By Lem. B.4, we therefore have

$$\text{null}^t([M]_i) = \text{null}^t M + E_i, \tag{B.1}$$

Because $(\text{null}^t M + E_i)$ is the set $\{x + y \mid x \in \text{null}^t M \wedge \forall h \geq i: y_h = 0\}$, Eqn. (B.1) is an equivalent way of stating the property to be shown. \square

C CORRECTNESS OF KS JOIN

We present the proof of Thm. 3.17 in this appendix. *If Y and Z are both $N+1$ -column KS matrices, and $\gamma_{\text{KS}}(Y)$ and $\gamma_{\text{KS}}(Z)$ are both non-empty sets, then $Y \sqcup Z$ is the projection of $\begin{bmatrix} -Y & Y \\ Z & 0 \end{bmatrix}$ onto its right-most $N+1$ columns.*

Proof. $\gamma_{\text{KS}}(Y \sqcup Z)$ is the affine closure of $\gamma_{\text{KS}}(Y) \cup \gamma_{\text{KS}}(Z)$. Thus, we need to show that, for all $\vec{v} \in \mathbb{Z}_{2^w}^N$,

$$\exists \vec{u} \in \mathbb{Z}_{2^w}^N, \sigma \in \mathbb{Z}_{2^w}: \begin{bmatrix} -Y & Y \\ Z & 0 \end{bmatrix} \begin{bmatrix} \vec{u} \\ \sigma \\ \vec{v} \\ 1 \end{bmatrix} = 0$$

$\Leftrightarrow \vec{v}$ is an affine combination of values in $\gamma_{\text{KS}}(Y)$ and $\gamma_{\text{KS}}(Z)$.

Recall that an affine combination is a linear combination whose coefficients sum to 1.

Proof of the “if” direction: Fix a particular $\vec{v} \in \mathbb{Z}_{2^w}^N$, and suppose that we have specific values for $\lambda \in \mathbb{Z}_{2^w}$, $\vec{y} \in \gamma_{\text{KS}}(Y)$, and $\vec{z} \in \gamma_{\text{KS}}(Z)$, such that $\vec{v} = \lambda \vec{y} + \vec{z}(1 - \lambda)$. Pick $\sigma = 1 - \lambda$, and $\vec{u} = (1 - \lambda) \vec{z}$. Then,

$$\begin{aligned} & \begin{bmatrix} -Y & Y \\ Z & 0 \end{bmatrix} \begin{bmatrix} (1 - \lambda) \vec{z} \\ 1 - \lambda \\ \vec{v} \\ 1 \end{bmatrix} = 0 \\ \Leftrightarrow & -Y \begin{bmatrix} (1 - \lambda) \vec{z} \\ 1 - \lambda \end{bmatrix} + Y \begin{bmatrix} \vec{v} \\ 1 \end{bmatrix} = 0 \text{ and } Z \begin{bmatrix} (1 - \lambda) \vec{z} \\ 1 - \lambda \end{bmatrix} = 0 \\ \Leftrightarrow & Y \begin{bmatrix} -(1 - \lambda) \vec{z} + \vec{v} \\ \lambda \end{bmatrix} = 0 \text{ and } (1 - \lambda) Z \begin{bmatrix} \vec{z} \\ 1 \end{bmatrix} = 0 \\ \Leftrightarrow & \lambda Y \begin{bmatrix} \vec{y} \\ 1 \end{bmatrix} = 0 \text{ and } (1 - \lambda) Z \begin{bmatrix} \vec{z} \\ 1 \end{bmatrix} = 0. \end{aligned}$$

These last equations are true because $\vec{y} \in \gamma_{\text{KS}}(Y)$ and $\vec{z} \in \gamma_{\text{KS}}(Z)$. Thus, if \vec{v} is in the affine closure of $\gamma_{\text{KS}}(Y) \cup \gamma_{\text{KS}}(Z)$, then $\begin{bmatrix} \vec{v} \\ 1 \end{bmatrix}$ is in the null space of the projected matrix.

Proof of the “only if” direction: Suppose that x is in the null space of the projected matrix. Pick $\vec{u} \in \mathbb{Z}_{2^w}^N$ and $\sigma \in \mathbb{Z}_{2^w}$ such that

$$\begin{bmatrix} -Y & Y \\ Z & 0 \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \sigma \\ \vec{v} \\ 1 \end{bmatrix} = 0.$$

We must show that \vec{v} is in the affine closure of $\gamma_{\text{KS}}(Y) \cup \gamma_{\text{KS}}(Z)$.

Immediately, we know that $Z \begin{bmatrix} \vec{u} \\ \sigma \end{bmatrix} = 0$ and $Y \begin{bmatrix} \vec{v} - \vec{u} \\ 1 - \sigma \end{bmatrix} = 0$. Because $\gamma_{\text{KS}}(Y)$ and $\gamma_{\text{KS}}(Z)$ are nonempty, we can select an arbitrary $\vec{y}_0 \in \gamma_{\text{KS}}(Y)$ and $\vec{z}_0 \in \gamma_{\text{KS}}(Z)$. Thus,

$$0 = Y \begin{bmatrix} \vec{v} - \vec{u} \\ 1 - \sigma \end{bmatrix} + \sigma Y \begin{bmatrix} \vec{y}_0 \\ 1 \end{bmatrix} = Y \begin{bmatrix} \vec{v} - \vec{u} + \sigma \vec{y}_0 \\ 1 \end{bmatrix}, \text{ and}$$

$$0 = Z \begin{bmatrix} \vec{u} \\ \sigma \end{bmatrix} + (1 - \sigma) Z \begin{bmatrix} \vec{z}_0 \\ 1 \end{bmatrix} = Z \begin{bmatrix} \vec{u} + (1 - \sigma) \vec{z}_0 \\ 1 \end{bmatrix}.$$

Now define \vec{y} and \vec{z} to be the values that we have just shown to be in $\gamma_{\text{KS}}(Y)$ and $\gamma_{\text{KS}}(Z)$:

$$\vec{y} \stackrel{\text{def}}{=} \vec{v} - \vec{u} + \sigma \vec{y}_0 \text{ and } \vec{z} \stackrel{\text{def}}{=} \vec{u} + (1 - \sigma) \vec{z}_0.$$

If we solve for \vec{v} and eliminate \vec{u} in these equations, we get:

$$\vec{v} = \vec{y} - \sigma \vec{y}_0 + \vec{z} + (\sigma - 1) \vec{z}_0.$$

Because $\vec{y}, \vec{y}_0 \in \gamma_{\text{KS}}(Y)$, $\vec{z}, \vec{z}_0 \in \gamma_{\text{KS}}(Z)$, and $1 - \sigma + 1 + (\sigma - 1) = 1$, we

have now stated \vec{v} as an affine combination of values in $\gamma_{\text{KS}}(Y)$ and $\gamma_{\text{KS}}(Z)$, as required. \square

D SOUNDNESS OF THE ABSTRACT-DOMAIN OPERATIONS FOR AFFINE-TRANSFORMERS-ABSTRACTION FRAMEWORK

In this section, we show that the abstract-domain operations for the ATA[\mathcal{B}] framework are sound with respect to the concrete semantics of the programming language.

Lemma D.1. *The bottom element represents the empty set.*

$$\gamma(\perp) = \emptyset \tag{D.1}$$

Proof.

$$\begin{aligned} \gamma(\perp) &= \{(\vec{v}, \vec{v}') : \vec{v}' = \vec{v} \cdot C + \vec{d} \wedge (C : \vec{d}) \in \gamma(\perp_{\mathcal{B}})\} \\ \Rightarrow \gamma(\perp) &= \{(\vec{v}, \vec{v}') : \vec{v}' = \vec{v} \cdot C + \vec{d} \wedge (C : \vec{d}) \in \emptyset\} \\ \Rightarrow \gamma(\perp) &= \emptyset \end{aligned}$$

□

Lemma D.2. *The equality operation is sound.*

$$(\mathbf{a}_1 \widetilde{=} \mathbf{a}_2) \Rightarrow (\gamma(\mathbf{a}_1) == \gamma(\mathbf{a}_2)) \tag{D.2}$$

Proof. We will prove Lemma D.2 by contradiction. Assume that $\mathbf{a}_1 \widetilde{=} \mathbf{a}_2$ but $\gamma(\mathbf{a}_1) \neq \gamma(\mathbf{a}_2)$. Without loss of generality, we can assume that there exists (\vec{v}, \vec{v}') such that:

$$\begin{aligned} &(\vec{v}, \vec{v}') \in \gamma(\text{base}(\mathbf{a}_1)) \wedge (\vec{v}, \vec{v}') \notin \gamma(\text{base}(\mathbf{a}_2)) \\ \Rightarrow &\vec{v}' = \vec{v} \cdot C + \vec{d} \wedge (C : \vec{d}) \in \gamma(\text{base}(\mathbf{a}_1)) \wedge (C : \vec{d}) \notin \gamma(\text{base}(\mathbf{a}_2)) \\ \Rightarrow &\exists \mathbf{b}. \mathbf{b} \in \gamma(\text{base}(\mathbf{a}_1)) \wedge \mathbf{b} \notin \gamma(\text{base}(\mathbf{a}_2)) \\ \Rightarrow &\gamma(\text{base}(\mathbf{a}_1)) \neq \gamma(\text{base}(\mathbf{a}_2)) \end{aligned}$$

$\Rightarrow \text{base}(a_1) \neq \text{base}(a_2)$ (by soundness of equality on \mathcal{B})

$\Rightarrow a_1 \neq a_2$ (Contradiction!)

□

Lemma D.3. *The join operation is sound.*

$$\gamma(a_1 \sqcup a_2) \supseteq \gamma(a_1) \cup \gamma(a_2) \quad (\text{D.3})$$

Proof. Assume that Lemma D.3 is incorrect. Then there exists $(\vec{v}, \vec{v}') \notin \gamma(a_1 \sqcup a_2)$ such that:

$$\begin{aligned} & (\vec{v}, \vec{v}') \in \gamma(a_1) \cup \gamma(a_2) \\ \Rightarrow & \vec{v}' = \vec{v} \cdot C + \vec{d} \wedge (C : \vec{d}) \in \gamma(\text{base}(a_1)) \\ & \text{(Without loss of generality.)} \\ \Rightarrow & \vec{v}' = \vec{v} \cdot C + \vec{d} \wedge (C : \vec{d}) \in \gamma(\text{base}(a_1) \sqcup \text{base}(a_2)) \\ \Rightarrow & (\vec{v}, \vec{v}') \in \gamma(a_1 \sqcup a_2) \text{ (Contradiction!)} \end{aligned}$$

□

The soundness of widening, statement abstractions, and identity function are easy to prove, and follow similar reasoning.

E SOUNDNESS OF ABSTRACT COMPOSITION FOR AFFINE-TRANSFORMERS-ABSTRACTION FRAMEWORK

In this section, we show that the abstract-composition operations defined in §4.3.2 are sound. From Eqn. (4.1), an abstract composition $a_e'' = a' \circ a$ is exact iff:

$$\gamma(a_e'') = \{(\vec{v}, \vec{v}') \mid \exists(C : \vec{d}) \in \gamma(\text{base}(a)), (C' : \vec{d}') \in \gamma(\text{base}(a')), (C'' : \vec{d}'') : \\ (\vec{v}' = \vec{v} \cdot C'' + \vec{d}'') \wedge (C'' = C \cdot C') \wedge (\vec{d}'' = \vec{d} \cdot C' + \vec{d}')\}$$

E.1 Non-Relational Base Domain

In this section, we show that the fast abstract composition for $\text{ATA}[\mathcal{B}]$ (Eqn. (4.3)), when \mathcal{B} is non-relational, is sound. Remember that any non-relational domain can be formulated as follows: $\mathcal{B} \stackrel{\text{def}}{=} \text{symbols}(C : \vec{d}) \rightarrow \mathcal{F}_{\mathcal{B}}$. The term $b[s]$, where $b \in \mathcal{B}$ and $s \in \text{symbols}(C : \vec{d})$, refers to the element in the foundation domain $f \in \mathcal{F}_{\mathcal{B}}$ corresponding to the symbol s .

Axiom 1. *Abstract addition is sound for $\mathcal{F}_{\mathcal{B}}$.*

$$e_1 \in \gamma(f_1) \wedge e_2 \in \gamma(f_2) \Rightarrow e_1 + e_2 \in \gamma(f_1 +^\# f_2) \quad (\text{E.1})$$

Axiom 2. *Abstract multiplication is sound for $\mathcal{F}_{\mathcal{B}}$.*

$$e_1 \in \gamma(f_1) \wedge e_2 \in \gamma(f_2) \Rightarrow e_1 \times e_2 \in \gamma(f_1 \times^\# f_2) \quad (\text{E.2})$$

Theorem E.1.

$$\gamma(a_e'') \subseteq \gamma(a' \circ_{\text{NR}} a). \quad (\text{E.3})$$

Proof. We will prove Thm. E.1 by contradiction. Consider a model $m = (\vec{v}, \vec{v}'')$, such that $m \in \gamma(a_e'')$ and $m \notin \gamma(a' \circ_{NR} a)$. We will show that such a model cannot exist.

$$\begin{aligned}
& m \in \gamma(a_e'') \\
\Leftrightarrow & \exists (C : \vec{d}) \in \gamma(\text{base}(a)), (C' : \vec{d}') \in \gamma(\text{base}(a')), (C'' : \vec{d}'') : \\
& (\vec{v}' = \vec{v} \cdot C'' + \vec{d}'') \wedge (C'' = C \cdot C') \wedge (\vec{d}'' = \vec{d} \cdot C' + \vec{d}') \\
\Leftrightarrow & \exists (C : \vec{d}), (C' : \vec{d}'), (C'' : \vec{d}'') : (\vec{v}' = \vec{v} \cdot C'' + \vec{d}'') \\
& \wedge \left(\bigwedge_{1 \leq i, j \leq k} (C''[i, j] = \sum_{1 \leq l \leq k} (C[i, l] \times C'[l, j])) \right) \\
& \wedge \left(\bigwedge_{1 \leq j \leq k} (\vec{d}''[j] = (\sum_{1 \leq l \leq k} (\vec{d}[l] \times C'[l, j]')) + \vec{d}'[j]) \right) \\
& \wedge \left(\bigwedge_{1 \leq i, j \leq k} C[i, j] \in \gamma(\text{base}(a)[c_{ij}]) \right) \wedge \left(\bigwedge_{1 \leq j \leq k} \vec{d}[j] \in \gamma(\text{base}(a)[d_j]) \right) \\
& \wedge \left(\bigwedge_{1 \leq i, j \leq k} C'[i, j] \in \gamma(\text{base}(a')[c_{ij}]) \right) \wedge \left(\bigwedge_{1 \leq j \leq k} \vec{d}'[j] \in \gamma(\text{base}(a')[d_j]) \right) \\
\Rightarrow & \exists (C'' : \vec{d}'') : (\vec{v}' = \vec{v} \cdot C'' + \vec{d}'') \\
& \wedge \left(\bigwedge_{1 \leq i, j \leq k} (C''[i, j] \in \sum_{1 \leq l \leq k}^{\#} (\text{base}(a)[c_{il}] \times^{\#} \text{base}(a')[c_{lj}])) \right) \\
& \wedge \left(\bigwedge_{1 \leq j \leq k} (\vec{d}''[j] \in \sum_{1 \leq l \leq k}^{\#} (\text{base}(a)[d_l] \times^{\#} \text{base}(a')[c_{lj}] +^{\#} \text{base}(a')[d_j])) \right)
\end{aligned}$$

(by application of axioms 1 and 2 to the expressions

$$\begin{aligned}
& \sum_{1 \leq l \leq k} (C[i, l] \times C'[l, j]) \text{ and } \sum_{1 \leq l \leq k} (\vec{d}[l] \times C'[l, j]') + \vec{d}'[j]) \\
\Leftrightarrow & \exists (C'' : \vec{d}'') : (\vec{v}' = \vec{v} \cdot C'' + \vec{d}'') \wedge \mathbf{b} \in (\text{symbols}(C'' : \vec{d}'') \rightarrow \mathcal{F}) \\
& \wedge \left(\bigwedge_{1 \leq i, j \leq k} \left(\mathbf{b}[c_{ij}'] = \sum_{1 \leq l \leq k}^{\#} (\text{base}(a)[c_{il}] \times^{\#} \text{base}(a')[c_{lj}]) \right) \right)
\end{aligned}$$

$$\wedge \left(\bigwedge_{1 \leq j \leq k} b[d_j''] = \sum_{1 \leq l \leq k}^{\#} (\text{base}(a)[d_l] \times^{\#} \text{base}(a')[c_{lj}]) +^{\#} \text{base}(a')[d_j] \right)$$

$\Leftrightarrow m \in \gamma(a' \circ_{\text{NR}} a)$ (by Eqn. (4.3))

□

E.2 Weakly-Convex Base Domain

In this section, we present a proof of soundness of abstract composition for weakly-convex base domains, denoted by $a' \circ_{\text{WC}} a$ (Eqn. (4.5)).

We present some useful axioms and lemmas before presenting the soundness theorem and its proof. Let $\min_{\mathbb{Z}_{2^w}}$ and $\max_{\mathbb{Z}_{2^w}}$ be the minimum and maximum bitvector values in \mathbb{Z}_{2^w} . Let $\min_{\mathbb{Q}} = \min_{\mathbb{Z}_{2^w}}$ and $\max_{\mathbb{Q}} = \max_{\mathbb{Z}_{2^w}}$.

Axiom 3. $\text{cast}_{\mathbb{Q}}$ is distributive over bitvector addition in the absence of overflows: that is, if $\min_{\mathbb{Q}} \leq \text{cast}_{\mathbb{Q}}(bv_1 + bv_2) \leq \max_{\mathbb{Q}}$, where $bv_1, bv_2 \in \mathbb{Z}_{2^w}$, then

$$\text{cast}_{\mathbb{Q}}(bv_1) + \text{cast}_{\mathbb{Q}}(bv_2) = \text{cast}_{\mathbb{Q}}(bv_1 + bv_2) \quad (\text{E.4})$$

Axiom 4. $\text{cast}_{\mathbb{Q}}$ is distributive over bitvector multiplication in the absence of overflows: that is, if $\min_{\mathbb{Q}} \leq bv_1 \cdot bv_2 \leq \max_{\mathbb{Q}}$, where $bv_1, bv_2 \in \mathbb{Z}_{2^w}$, then

$$\text{cast}_{\mathbb{Q}}(bv_1) \cdot \text{cast}_{\mathbb{Q}}(bv_2) = \text{cast}_{\mathbb{Q}}(bv_1 \cdot bv_2) \quad (\text{E.5})$$

Lemma E.2. $\text{cast}_{\mathbb{Q}}$ is distributive over matrix multiplication for bitvectors, if there are no overflows in the matrix multiplication. That is, for $n \times n$ matrices M and M' where $\forall_{1 \leq i, j \leq n} M[i, j], M'[i, j] \in \mathbb{Z}_{2^w}$,

$$\text{cast}_{\mathbb{Q}}(M) \times \text{cast}_{\mathbb{Q}}(M') = \text{cast}_{\mathbb{Q}}(M \times M'). \quad (\text{E.6})$$

Proof. Let $M'' = \text{cast}_{\mathbb{Q}}(M) \times \text{cast}_{\mathbb{Q}}(M')$. Then,

$$\begin{aligned} \forall_{1 \leq i, j \leq n} : M''[i, j] &= \sum_{1 \leq l \leq n} \text{cast}_{\mathbb{Q}}(M[i, l]) \cdot \text{cast}_{\mathbb{Q}}(M[l, j]) \\ \Rightarrow \forall_{1 \leq i, j \leq n} : M''[i, j] &= \sum_{1 \leq l \leq n} \text{cast}_{\mathbb{Q}}(M[i, l] \cdot M[l, j]) \text{ (by Axiom 4)} \\ \Rightarrow \forall_{1 \leq i, j \leq n} : M''[i, j] &= \text{cast}_{\mathbb{Q}}(\sum_{1 \leq l \leq n} M[i, l] \cdot M[l, j]) \text{ (by Axiom 3)} \\ \Rightarrow \text{cast}_{\mathbb{Q}}(M) \times \text{cast}_{\mathbb{Q}}(M') &= \text{cast}_{\mathbb{Q}}(M \times M') \end{aligned}$$

□

Lemma E.3. *A convex combination of a set of rationals is inside bitvector boundaries if each of the rational values in the set is inside bitvector boundaries. Given any $1 \leq \lambda_1, \lambda_2, \dots, \lambda_l \leq 1$, such that $(\sum_{i=1}^l \lambda_i = 1)$.*

$$\min_{\mathbb{Q}} \leq q_1, q_2, \dots, q_l \leq \max_{\mathbb{Q}} \Rightarrow \min_{\mathbb{Q}} \leq \sum_{i=1}^l \lambda_i q_i \leq \max_{\mathbb{Q}} \quad (\text{E.7})$$

Proof. Suppose $\min_{\mathbb{Q}} > \sum_{i=1}^l \lambda_i q_i$.

$$\begin{aligned} \min_{\mathbb{Q}} > \sum_{i=1}^l \lambda_i q_i &\Rightarrow (\min_{\mathbb{Q}} > \sum_{i=1}^l \lambda_i \min_{\mathbb{Q}}) \text{ (because } \min_{\mathbb{Q}} \leq q_1, q_2, \dots, q_l) \\ \Leftrightarrow \min_{\mathbb{Q}} > \sum_{i=1}^l \lambda_i \min_{\mathbb{Q}} &\Rightarrow (\min_{\mathbb{Q}} > \min_{\mathbb{Q}}) \text{ (because } (\sum_{i=1}^l \lambda_i = 1).) \\ \Leftrightarrow \text{false} \end{aligned}$$

Consequently, $\min_{\mathbb{Q}} \leq \sum_{i=1}^l \lambda_i q_i$. The other inequality $\sum_{i=1}^l \lambda_i q_i \leq \max_{\mathbb{Q}}$ can be proved in a similar fashion. □

Lemma E.4. *There are no overflows in a matrix multiplication of a convex combination of matrices, if there are no overflows in the matrix multiplications of the underlying matrices involved in the convex combination.*

$\forall_{1 \leq i \leq l, 1 \leq j \leq l'} : (t_i \times t'_j)$ does not overflow $\Rightarrow (t \times t')$ does not overflow, where

$$\left(\text{cast}_{\mathbb{Q}}(t) = \sum_{i=1}^l \lambda_i \text{cast}_{\mathbb{Q}}(t_i) \right) \wedge \bigwedge_{i=1}^l (1 \leq \lambda_i \leq 1) \wedge \sum_{i=1}^l \lambda_i = 1, \quad (\text{E.8})$$

$$\text{and, } \left(\text{cast}_{\mathbb{Q}}(t') = \sum_{j=1}^{l'} \lambda'_j \text{cast}_{\mathbb{Q}}(t'_j) \right) \wedge \bigwedge_{j=1}^{l'} (1 \leq \lambda'_j \leq 1) \wedge \sum_{j=1}^{l'} \lambda'_j = 1. \quad (\text{E.9})$$

Proof. Because $(t_i \times t'_j)$ does not overflow, we know that each entry in the computation of the matrix multiplication does not overflow:

$$\forall_{1 \leq p, q \leq o} : \min_{\mathbb{Q}} \leq \sum_{n=1}^o \text{cast}_{\mathbb{Q}}(t_i[p, n]) \cdot \text{cast}_{\mathbb{Q}}(t'_j[n, q]) \leq \max_{\mathbb{Q}} \quad (\text{E.10})$$

where t_i and t'_j are $o \times o$ matrices.

Suppose that $l'' = l \cdot l'$ and $\bigwedge_{m=1}^{l''} (1 \leq (\lambda''_m = \lambda_{\lfloor m/l' \rfloor} \cdot \lambda'_{(m-1)\%l'+1}) \leq 1)$.

Then,

$\sum_{m=1}^{l''} \lambda''_m = \sum_{i=1}^l \lambda_i \cdot \sum_{j=1}^{l'} \lambda'_j = 1 \cdot 1 = 1$. Then, by applying Lem. E.3 to Eqn. (E.10), we get for all $1 \leq p, q \leq o$:

$$\begin{aligned} \min_{\mathbb{Q}} &\leq \sum_{m=1}^{l''} \lambda''_m (\sum_{n=1}^o \text{cast}_{\mathbb{Q}}(t_{\lfloor m/l' \rfloor}[p, n]) \cdot \text{cast}_{\mathbb{Q}}(t'_{m\%l'+1}[n, q])) \leq \max_{\mathbb{Q}} \\ \Leftrightarrow \min_{\mathbb{Q}} &\leq \sum_{i=1}^l \sum_{j=1}^{l'} \lambda_i \lambda'_j (\sum_{n=1}^o \text{cast}_{\mathbb{Q}}(t_i[p, n]) \cdot \text{cast}_{\mathbb{Q}}(t'_j[n, q])) \leq \max_{\mathbb{Q}} \\ \Leftrightarrow \min_{\mathbb{Q}} &\leq (\sum_{n=1}^o (\sum_{i=1}^l \lambda_i \text{cast}_{\mathbb{Q}}(t_i[p, n]))) \cdot (\sum_{j=1}^{l'} \lambda'_j \text{cast}_{\mathbb{Q}}(t'_j[n, q])) \leq \max_{\mathbb{Q}} \\ &\text{(by distributivity of multiplication over addition for rationals)} \\ \Leftrightarrow \min_{\mathbb{Q}} &\leq (\sum_{n=1}^o \text{cast}_{\mathbb{Q}}(t[p, n]) \cdot \text{cast}_{\mathbb{Q}}(t'_j[n, q])) \leq \max_{\mathbb{Q}} \\ &\text{(by the definition of } \text{cast}_{\mathbb{Q}}(t) \text{ and } \text{cast}_{\mathbb{Q}}(t') \text{ in Eqn. (E.8) and Eqn. (E.9))} \end{aligned}$$

Hence $(t \times t')$ does not overflow. \square

Theorem E.5.

$$\gamma(a''_e) \subseteq \gamma(a' \circ_{\text{WC}} a). \quad (\text{E.11})$$

Proof. Consider any model $m = (\vec{v}, \vec{v}'')$, such that $m \in \gamma(\mathbf{a}_e'')$. To prove Thm. E.1, we need to show that $m \in \gamma(\mathbf{a}' \circ_{\text{WC}} \mathbf{a})$. Eqn. (4.5) defines $\mathbf{a}'' = \mathbf{a}' \circ_{\text{WC}} \mathbf{a}$ as follows

$$\text{base}(\mathbf{a}'') = \begin{cases} \{t_i \times t'_j \mid 1 \leq i \leq l, 1 \leq j \leq l'\} & \text{if there are no overflows in any} \\ & \text{matrix multiplication } t_i \times t'_j \\ \top_{\mathcal{B}} & \text{otherwise} \end{cases}$$

where $\text{base}(\mathbf{a}) = \{t_1, t_2, \dots, t_l\}$ and $\text{base}(\mathbf{a}') = \{t'_1, t'_2, \dots, t'_{l'}\}$.

(E.12)

We know that for $m = (\vec{v}, \vec{v}')$,

$$\begin{aligned} \exists (C : \vec{d}) \in \gamma(\text{base}(\mathbf{a})), (C' : \vec{d}') \in \gamma(\text{base}(\mathbf{a}')), (C'' : \vec{d}'') : \quad & \text{(E.13)} \\ (\vec{v}' = \vec{v} \cdot C'' + \vec{d}'') \wedge (C'' = C \cdot C') \wedge (\vec{d}'' = \vec{d} \cdot C' + \vec{d}'). \end{aligned}$$

By the properties of weakly-convex domains (see §4.3.2), we know that

$$\text{cast}_{\mathbb{Q}} \left(\left[\begin{array}{c|c} 1 & \vec{d} \\ \hline 0 & C \end{array} \right] \right) = \sum_{i=1}^l \lambda_i \text{cast}_{\mathbb{Q}}(t_i), \text{ for some } \lambda_1, \lambda_2, \dots, \lambda_l \in \mathbb{Q} \text{ such that}$$

(E.14)

$$\bigwedge_{i=1}^l (0 \leq \lambda_i \leq 1) \wedge (\sum_{i=1}^l \lambda_i = 1) \wedge \bigwedge_{i=1}^l \left(t_i = \left[\begin{array}{c|c} 1 & \vec{d}_i \\ \hline 0 & C_i \end{array} \right] \right), \text{ and}$$

$$\text{cast}_{\mathbb{Q}} \left(\left[\begin{array}{c|c} 1 & \vec{d}' \\ \hline 0 & C' \end{array} \right] \right) = \sum_{i=1}^{l'} \lambda'_i \text{cast}_{\mathbb{Q}}(t'_i), \text{ for some } \lambda'_1, \lambda'_2, \dots, \lambda'_{l'} \in \mathbb{Q} \text{ such that}$$

(E.15)

$$\bigwedge_{i=1}^{l'} (0 \leq \lambda'_i \leq 1) \wedge (\sum_{i=1}^{l'} \lambda'_i = 1) \wedge \bigwedge_{i=1}^{l'} \left(t'_i = \left[\begin{array}{c|c} 1 & \vec{d}'_i \\ \hline 0 & C'_i \end{array} \right] \right).$$

To show that $m \in \gamma(\mathbf{a}' \circ_{\text{WC}} \mathbf{a})$, we consider two cases.

Overflows in matrix multiplication. If there is an overflow encountered in any matrix multiplication $t_i \times t'_j$, then $\text{base}(a'') = \top_{\mathcal{B}}$ and consequently, $m \in \gamma(a' \circ_{\text{WC}} a)$ is true trivially.

No overflows in matrix multiplication. If there is no overflow encountered in any of the matrix multiplications $t_i \times t'_j$, then it suffices to prove that

$$(C'' : \vec{d}'') \in \bigsqcup_{i=1}^l \bigsqcup_{j=1}^{l'} \{t_i \times t'_j\}. \quad (\text{E.16})$$

Eqn. (E.16) translates to proving that for some $\{\lambda''_1, \lambda''_2, \dots, \lambda''_{l''}\}$:

$$\left(\text{cast}_{\mathbb{Q}} \left(\left[\begin{array}{c|c} 1 & \vec{d}'' \\ \hline 0 & C'' \end{array} \right] \right) = \sum_{i=1}^{l''} \lambda''_i \text{cast}_{\mathbb{Q}}(t''_i) \right), \text{ for some } \lambda''_1, \lambda''_2, \dots, \lambda''_{l''} \in \mathbb{Q} \text{ such that} \quad (\text{E.17})$$

$$\bigwedge_{i=1}^{l''} (0 \leq \lambda''_i \leq 1) \wedge \left(\sum_{i=1}^{l''} \lambda''_i = 1 \right) \wedge \bigwedge_{i=1}^{l''} \left(t''_i = \left[\begin{array}{c|c} 1 & \vec{d}''_i \\ \hline 0 & C''_i \end{array} \right] \right).$$

$$\begin{aligned} & \text{cast}_{\mathbb{Q}} \left(\left[\begin{array}{c|c} 1 & \vec{d}'' \\ \hline 0 & C'' \end{array} \right] \right) \\ &= \text{cast}_{\mathbb{Q}} \left(\left[\begin{array}{c|c} 1 & \vec{d} \\ \hline 0 & C \end{array} \right] \times \left[\begin{array}{c|c} 1 & \vec{d}' \\ \hline 0 & C' \end{array} \right] \right) \text{ (by Eqn. (E.13))} \\ &= \text{cast}_{\mathbb{Q}} \left(\left[\begin{array}{c|c} 1 & \vec{d} \\ \hline 0 & C \end{array} \right] \right) \times \text{cast}_{\mathbb{Q}} \left(\left[\begin{array}{c|c} 1 & \vec{d}' \\ \hline 0 & C' \end{array} \right] \right) \\ & \quad \text{(by Lem. E.4, because } \left[\begin{array}{c|c} 1 & \vec{d} \\ \hline 0 & C \end{array} \right] \times \left[\begin{array}{c|c} 1 & \vec{d}' \\ \hline 0 & C' \end{array} \right] \text{ does not overflow)} \\ &= \sum_{i=1}^l \lambda_i \text{cast}_{\mathbb{Q}}(t_i) \times \sum_{j=1}^{l'} \lambda'_j \text{cast}_{\mathbb{Q}}(t'_j) \\ & \quad \text{(by Eqn. (E.14) and Eqn. (E.15).)} \\ &= \sum_{i=1}^l \sum_{j=1}^{l'} \lambda_i \lambda'_j \text{cast}_{\mathbb{Q}}(t_i) \times \text{cast}_{\mathbb{Q}}(t'_j) \\ & \quad \text{(by distributivity of matrix multiplication over addition)} \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^l \sum_{j=1}^{l'} \lambda_i \lambda'_j \text{cast}_{\mathbb{Q}}(\mathbf{t}_i \times \mathbf{t}'_j) \\
&\quad \text{(by Lem. E.2)} \\
&= \sum_{m=1}^{l''} \lambda''_m \text{cast}_{\mathbb{Q}}(\mathbf{t}''_m)
\end{aligned}$$

where $l'' = l \cdot l'$, $\bigwedge_{m=1}^{l''} (1 \leq (\lambda''_m = \lambda_{\lfloor m/l' \rfloor} \cdot \lambda'_{(m-1)\%l'+1}) \leq 1)$, and

$$\sum_{m=1}^{l''} \lambda''_m = \sum_{i=1}^{l'} \lambda_i \cdot \sum_{j=1}^{l'} \lambda'_j = 1 \cdot 1 = 1.$$

□

F SOUNDNESS OF THE MERGE OPERATION

In this section, we show that the merge operation defined in §4.3.3 is sound. Recall that the merge function is defined as:

$$\begin{aligned}
 \text{Merge}(a, a') &= a'', \text{ where} & (F.1) \\
 \text{base}(a'') &= (b_g \sqcap \text{havoc}(\text{base}(\text{Id}), \text{gsyms})) \circ a \\
 b_g &= \text{havoc}(\text{base}(a'), \text{lsyms}), \\
 \text{lsyms} &= \text{symbols}(C_{g1}) \cup \text{symbols}(C_{11}) \cup \\
 &\quad \text{symbols}(d_1) \cup \text{symbols}(C_{1g}) \\
 \text{gsyms} &= \text{symbols}(C_{gg}) \cup \text{symbols}(d_g)
 \end{aligned}$$

As mentioned in Eqn. (4.6), the exact merge-function semantics are specified as follows:

$$\text{MERGE}(\gamma(a), \gamma(a')) = \text{REVERTLOCALS}(\gamma(a')) \circ \gamma(a) \quad (F.2)$$

Theorem F.1.

$$\text{MERGE}(\gamma(a), \gamma(a')) \subseteq \gamma(\text{Merge}(a, a')). \quad (F.3)$$

Proof. We will prove Thm. F.1 by contradiction. Consider a model $m = (\vec{g}_m, \vec{l}_m; \vec{g}_m', \vec{l}_m')$, such that $m \in \text{MERGE}(\gamma(a), \gamma(a'))$ and $m \notin \gamma(\text{Merge}(a, a'))$. Let $a_{\text{REvLocs}} \in \text{ATA}[\mathcal{B}]$ be an abstract domain value such that

$$\text{base}(a_{\text{REvLocs}}) = \text{havoc}(\text{base}(a'), \text{lsyms}) \sqcap \text{havoc}(\text{base}(\text{Id}), \text{gsyms}) \quad (F.4)$$

By the soundness of abstract composition, existence of m implies existence of $n = (\vec{g}_n, \vec{l}_n; \vec{g}_n', \vec{l}_n')$, such that $n \in \text{REVERTLOCALS}(\gamma(a'))$ and $n \notin \gamma(a_{\text{REvLocs}})$. We will show that n cannot exist. Consequently, m

cannot exist, and thus merge is sound.

$$\begin{aligned}
& n \in \text{REVERTLOCALS}(\gamma(\alpha')) \\
\Leftrightarrow & (\vec{g}_n, \vec{l}_n; \vec{g}_n', \vec{l}_n') \in \{(\vec{g}, \vec{q}, \vec{g}', \vec{q}') \mid (\vec{g}, \vec{l}, \vec{g}', \vec{l}') \in \gamma(\alpha')\} \text{ (by Eqn. (4.7))} \\
\Leftrightarrow & \exists \left(\begin{bmatrix} C_{gg} & C_{gl} \\ C_{lg} & C_{ll} \end{bmatrix} : (\vec{d}_g, \vec{d}_l) \right) \in \gamma(\text{base}(\alpha')), \vec{l}, \vec{l}', \vec{q} : \\
& (\vec{g}_n' = \vec{g}_n \cdot C_{gg} + \vec{l} \cdot C_{lg} + \vec{d}_g) \wedge (\vec{l}' = \vec{g}_n \cdot C_{gl} + \vec{l} \cdot C_{ll} + \vec{d}_l), \\
& \wedge (\vec{l}_n = \vec{q}) \wedge (\vec{l}_n' = \vec{q}') \wedge (\vec{l} = \vec{0}) \\
& (\vec{l} \text{ are initialized to 0 in each function}) \\
\Leftrightarrow & \exists \left(\begin{bmatrix} C_{gg} & C_{gl} \\ C_{lg} & C_{ll} \end{bmatrix} : (\vec{d}_g, \vec{d}_l) \right) \in \gamma(\text{base}(\alpha')) : \\
& (\vec{g}_n' = \vec{g}_n \cdot C_{gg} + \vec{d}_g) \wedge (\vec{l}_n = \vec{l}_n') \\
& (\text{by removing the existential variables } \vec{l}, \vec{l}' \text{ and } \vec{q}) \\
\Leftrightarrow & \exists \left(\begin{bmatrix} C_{gg} & C_{gl} \\ C_{lg} & C_{ll} \end{bmatrix} : (\vec{d}_g, \vec{d}_l) \right) \in \gamma(\text{base}(\alpha')) : \\
& \text{havoc} \left(\left[1 \mid \vec{g}_n \quad \vec{l}_n \right] \left[\begin{array}{c|cc} 1 & \vec{d}_g & \vec{d}_l \\ \hline 0 & C_{gg} & C_{gl} \\ 0 & 0 & C_{ll} \end{array} \right] = \left[1 \mid \vec{g}_n' \quad \vec{l}_n' \right], \text{lsyms} \right) \wedge (\vec{l}_n' = \vec{l}_n) \\
& (\text{because } (\vec{g}_n' = \vec{g}_n \cdot C_{gg} + \vec{d}_g) \text{ is the result of havoc on } \text{lsyms} \text{ for} \\
& (\vec{g}_n' = \vec{g}_n \cdot C_{gg} + \vec{d}_g) \wedge (\vec{l}_n' = \vec{g}_n \cdot C_{gl} + \vec{l}_n \cdot C_{ll} + \vec{d}_l)) \\
\Leftrightarrow & \exists \left(\begin{bmatrix} C_{gg} & C_{gl} \\ C_{lg} & C_{ll} \end{bmatrix} : (\vec{d}_g, \vec{d}_l) \right) \in \gamma(\text{base}(\alpha')) : \\
& \text{havoc} \left(\left[1 \mid \vec{g}_n \quad \vec{l}_n \right] \left[\begin{array}{c|cc} 1 & \vec{d}_g & \vec{d}_l \\ \hline 0 & C_{gg} & C_{gl} \\ 0 & C_{gl} & C_{ll} \end{array} \right] = \left[1 \mid \vec{g}_n' \quad \vec{l}_n' \right], \text{lsyms} \right) \wedge (\vec{l}_n' = \vec{l}_n) \\
& (\text{because } \vec{l}_n \text{ are initialized to zero})
\end{aligned}$$

$$\Leftrightarrow \exists \left(\left[\begin{array}{cc} C_{gg} & C_{gl} \\ C_{lg} & C_{ll} \end{array} \right] : (\vec{d}_g, \vec{d}_l) \right) \in \gamma(\text{base}(a')) :$$

$$\text{havoc} \left(\left[1 \mid \vec{g}_n \quad \vec{l}_n \right] \left[\begin{array}{c|cc} 1 & \vec{d}_g & \vec{d}_l \\ \hline 0 & C_{gg} & C_{gl} \\ 0 & C_{gl} & C_{ll} \end{array} \right] = \left[1 \mid \vec{g}_n' \quad \vec{l}_n' \right], \text{lsyms} \right) \wedge$$

$$\text{havoc} \left(\left[1 \mid \vec{g}_n \quad \vec{l}_n \right] \left[\begin{array}{c|cc} 1 & 0 & 0 \\ \hline 0 & I & 0 \\ 0 & 0 & I \end{array} \right] = \left[1 \mid \vec{g}_n' \quad \vec{l}_n' \right], \text{gsyms} \right)$$

(because havoc of *gsyms* on the identity transformation yields $\vec{l}_n' = \vec{l}_n$)

$$\Leftrightarrow \left(\exists \left(\left[\begin{array}{cc} C_{gg} & C_{gl} \\ C_{lg} & C_{ll} \end{array} \right] : (\vec{d}_g, \vec{d}_l) \right) \in \gamma(\text{havoc}(\text{base}(a'), \text{lsyms})) : \right.$$

$$\left. \left[1 \mid \vec{g}_n \quad \vec{l}_n \right] \left[\begin{array}{c|cc} 1 & \vec{d}_g & \vec{d}_l \\ \hline 0 & C_{gg} & C_{gl} \\ 0 & C_{gl} & C_{ll} \end{array} \right] = \left[1 \mid \vec{g}_n' \quad \vec{l}_n' \right] \right) \wedge$$

$$\left(\exists \left(\left[\begin{array}{cc} C_{gg} & C_{gl} \\ C_{lg} & C_{ll} \end{array} \right] : (\vec{d}_g, \vec{d}_l) \right) \in \gamma(\text{havoc}(\text{base}(\text{Id}), \text{gsyms})) : \right.$$

$$\left. \left[1 \mid \vec{g}_n \quad \vec{l}_n \right] \left[\begin{array}{c|cc} 1 & \vec{d}_g & \vec{d}_l \\ \hline 0 & C_{gg} & C_{gl} \\ 0 & C_{gl} & C_{ll} \end{array} \right] = \left[1 \mid \vec{g}_n' \quad \vec{l}_n' \right] \right)$$

$$\Leftrightarrow (\vec{g}_n, \vec{l}_n, \vec{g}_n', \vec{l}_n') \in \gamma(\mathbf{a}_{\text{RevLocs}})$$

$$\Leftrightarrow \text{Model}(\vec{g}_n, \vec{l}_n, \vec{g}_n', \vec{l}_n') \text{ does not exist.}$$

$$\Rightarrow \text{Model}(\vec{g}_m, \vec{l}_m, \vec{g}_m', \vec{l}_m') \text{ does not exist. (Contradiction.)}$$

(by soundness of abstract composition)

If the abstract composition operation is exact, then the implication in the last step of the proof becomes biconditional. Thus, if abstract composition is exact then the merge operation is exact. \square

REFERENCES

- [1] Bach, E. 1992. Linear algebra modulo n . Unpublished manuscript.
- [2] Bagnara, R., K. Dobson, P.M. Hill, M. Mundell, and E. Zaffanella. 2006. Grids: A domain for analyzing the distribution of numerical values. In *Int. Workshop on Logic Based Prog. Dev. and Transformation*.
- [3] Bagnara, R., P. M. Hill, and E. Zaffanella. 2008. The Parma Polyhedra Library: Toward a complete set of numerical abstractions for the analysis and verification of hardware and software systems. *Science of Computer Programming* 72(1–2):3–21.
- [4] Bagnara, R., P.M. Hill, and E. Zaffanella. 2006. Widening operators for powerset domains. *Int. Journal on Software Tools for Technology Transfer* 8(4/5):449–466.
- [5] Balakrishnan, G., and T. Reps. 2010. WYSINWYX: What You See Is Not What You eXecute. *Trans. on Prog. Lang. and Syst.*
- [6] BBC. 2014. Gangnam style music video 'broke' youtube view limit. <http://www.bbc.com/news/world-asia-30288542>.
- [7] Beyer, D. 2015. Software verification and verifiable witnesses - (report on sv-comp 2015). In *Int. Conf. on Tools and Algorithms for the Construction and Analysis of Systems*.
- [8] Bjørner, N., A. Phan, and L. Fleckenstein. 2015. νz -an optimizing smt solver. In *Int. Conf. on Tools and Algorithms for the Construction and Analysis of Systems*.
- [9] Bloch, Joshua. Extra, extra - read all about it: Nearly all binary searches and mergesorts are broken. "googleresearch.blogspot.com/2006/06/extra-extra-read-all-about-it-nearly.html".

- [10] Bouajjani, A., J. Esparza, and T. Touili. 2003. A generic approach to the static analysis of concurrent programs with procedures. In *Symposium on Principles of Programming Languages*.
- [11] Brauer, J., and A. King. 2010. Automatic abstraction for intervals using Boolean formulae. In *Static Analysis Symposium*.
- [12] ———. 2011. Transfer function synthesis without quantifier elimination. In *European Symp. on Programming*.
- [13] ———. 2012. Transfer function synthesis without quantifier elimination. *Logical Methods in Comp. Sci.* 8(3).
- [14] Burch, Jerry, and David Dill. 1994. Automatic verification of pipelined microprocessor control. In *Int. Conf. on Computer Aided Verification*.
- [15] Bygde, S., B. Lisper, and N. Holsti. 2011. Fully bounded polyhedral analysis of integers with wrapping. In *Int. Workshop on Numerical and Symbolic Abstract Domains*.
- [16] Chang, B.-Y.E., and K.R.M. Leino. 2005. Abstract interpretation with alien expressions and heap structures. In *Int. Conf. on Verification, Model Checking, and Abstract Interpretation*.
- [17] Chen, L., A. Miné, and P. Cousot. 2008. A sound floating-point polyhedra abstract domain. In *Asian Symp. on Prog. Lang. and Systems*.
- [18] Chen, L., A. Miné, J. Wang, and P. Cousot. 2009. Interval polyhedra: An abstract domain to infer interval linear relationships. In *Static Analysis Symposium*.
- [19] Cousot, P., and R. Cousot. 1976. Static determination of dynamic properties of programs. In *Int. Symp on Programming*.

- [20] ———. 1977. Abstract interpretation: A unified lattice model for static analysis of programs by construction of approximation of fixed points. In *Symposium on Principles of Programming Languages*, 238–252.
- [21] ———. 1979. Systematic design of program analysis frameworks. In *Symposium on Principles of Programming Languages*, 269–282.
- [22] Cousot, P., R. Cousot, and L. Mauborgne. 2011. The reduced product of abstract domains and the combination of decision procedures. In *Int. Conf. on Foundations of Software Science and Computation Structures*.
- [23] Cousot, P., and N. Halbwachs. 1978. Automatic discovery of linear constraints among variables of a program. In *Symposium on Principles of Programming Languages*, 84–96.
- [24] Dietz, W., P. Li, J. Regehr, and V. Adve. 2012. Understanding integer overflow in C/C++. In *Int. Conf. on Software Engineering*.
- [25] dSPACE. dSPACE TargetLink. www.dspace.com/en/pub/home/products/sw/pcgs/targetli.cfm.
- [26] Dutertre, B., and L. de Moura. 2006. Yices: An SMT solver. <http://yices.csl.sri.com>.
- [27] Dutertre, Bruno, and Leonardo de Moura. 2006. A fast linear-arithmetic solver for dpll(t). In *Int. Conf. on Computer Aided Verification*.
- [28] EDN Network. 2013. Toyota’s killer firmware: Bad design and its consequences. <http://www.edn.com/design/automotive/4423428/Toyota-s-killer-firmware--Bad-design-and-its-consequences>.

- [29] Elder, M., J. Lim, T. Sharma, T. Andersen, and T. Reps. 2011. Abstract domains of affine relations. In *Static Analysis Symposium*.
- [30] ———. 2014. Abstract domains of affine relations. *Trans. on Prog. Lang. and Syst.*
- [31] Engadget. 2015. To keep a boeing dreamliner flying, reboot once every 248 days. <https://www.engadget.com/2015/05/01/boeing-787-dreamliner-software-bug>.
- [32] FDA. 2014. Fda statement on radiation overexposures in panama. <https://www.fda.gov/radiation-emittingproducts/radiationsafety/alertsandnotices/ucm116533.htm>.
- [33] Ganesh, V., and D.L. Dill. 2007. A decision procedure for bit-vectors and arrays. In *Int. Conf. on Computer Aided Verification*.
- [34] Gange, G., J. Navas, P. Schachte, H. Søndergaard, and P. Stuckey. 2013. Abstract interpretation over non-lattice abstract domains. In *Static Analysis Symposium*.
- [35] Garner, H. L. 1978. Theory of computer addition and overflow. *IEEE Trans. on Computers* C-27(4).
- [36] GCN. 1998. Software glitches leave navy smart ship dead in the water. <https://gcn.com/Articles/1998/07/13/Software-glitches-leave-Navy-Smart-Ship-dead-in-the-water.aspx>.
- [37] Ghorbal, K., F. Ivančić, G. Balakrishnan, N. Maeda, and A. Gupta. 2012. Donut domains: Efficient non-convex domains for abstract interpretation. In *Int. Conf. on Verification, Model Checking, and Abstract Interpretation*.
- [38] Gleick, J. 1996. A bug and a crash. <https://around.com/ariane.html>.

- [39] Granger, P. 1989. Static analysis of arithmetic congruences. *Int. J. of Comp. Math.*
- [40] ———. 1991. Analyses semantiques de congruence. Ph.D. thesis, Ecole Polytechnique.
- [41] Gulwani, S., and G.C. Necula. 2003. Discovering affine equalities using random interpretation. In *Symposium on Principles of Programming Languages*.
- [42] ———. 2005. Precise interprocedural analysis using random interpretation. In *Symposium on Principles of Programming Languages*.
- [43] Hafner, J., and K. McCurley. 1991. Asymptotically fast triangularization of matrices over rings. *SIAM J. Comput.* 20(6).
- [44] Howell, J.A. 1986. Spans in the module $(\mathbb{Z}_m)^s$. *Linear and Multilinear Algebra* 19.
- [45] Jeannet, B. New Polka. <http://pop-art.inrialpes.fr/~bjeannet/newpolka/index.html>.
- [46] Jeannet, B., D. Gopan, and T. Reps. 2005. A relational abstraction for functions. In *Static Analysis Symposium*.
- [47] Jones, N.D., and A. Mycroft. 1986. Data flow analysis of applicative programs using minimal function graphs. In *Symposium on Principles of Programming Languages*, 296–306.
- [48] Jourdan, Jacques-Henri, Vincent Laporte, Sandrine Blazy, Xavier Leroy, and David Pichardie. 2015. A formally-verified c static analyzer. In *Symposium on Principles of Programming Languages*.
- [49] Karr, M. 1976. Affine relationship among variables of a program. *Acta Inf.* 6:133–151.

- [50] Kidd, N., A. Lal, and T. Reps. 2007. WALi: The Weighted Automaton Library. www.cs.wisc.edu/wpis/wpds/download.php.
- [51] Kildall, G.A. 1973. A unified approach to global program optimization. In *Symposium on Principles of Programming Languages*, 194–206.
- [52] King, A., and H. Søndergaard. 2008. Inferring congruence equations with SAT. In *Int. Conf. on Computer Aided Verification*.
- [53] ———. 2010. Automatic abstraction for congruences. In *Int. Conf. on Verification, Model Checking, and Abstract Interpretation*.
- [54] Knoop, J., and B. Steffen. 1992. The interprocedural coincidence theorem. In *Int. Conf. on Compiler Construction*.
- [55] Lal, A., and T. Reps. 2006. Improving pushdown system model checking. In *Int. Conf. on Computer Aided Verification*.
- [56] Lal, A., T. Reps, and G. Balakrishnan. 2005. Extended weighted pushdown systems. In *Int. Conf. on Computer Aided Verification*.
- [57] Lattner, C., and V.S. Adve. 2004. LLVM: A compilation framework for lifelong program analysis & transformation. In *Int. Symp. on Code Generation and Optimization*.
- [58] Laviro, V., and F. Logozzo. 2009. Subpolyhedra: A (more) scalable approach to infer linear inequalities. In *Int. Conf. on Verification, Model Checking, and Abstract Interpretation*.
- [59] Li, Y., A. Albarghouthi, Z. Kincaid, A. Gurfinkel, and M. Chechik. 2014. Symbolic optimization with smt solvers. In *Symposium on Principles of Programming Languages*.
- [60] Lim, J. 2011. Transformer Specification Language: A system for generating analyzers and its applications. Ph.D. thesis, Comp. Sci. Dept., Univ. of Wisconsin, Madison, WI. Tech. Rep. 1689.

- [61] Lim, J., and T. Reps. 2008. A system for generating static analyzers for machine instructions. In *Int. Conf. on Compiler Construction*.
- [62] ———. 2013. TSL: A system for generating abstract interpreters and its application to machine-code analysis. *Trans. on Prog. Lang. and Syst.* 35(1).
- [63] llvm. LLVM: Low level virtual machine. llvm.org.
- [64] Malmkjær, K. 1993. Abstract interpretation of partial-evaluation algorithms. Ph.D. thesis, Dept. of Comp. and Inf. Sci., Kansas State Univ., Manhattan, Kansas.
- [65] Masdupuy, F. 1992. Array abstractions using semantic analysis of trapezoid congruences. In *Int. Conf. Supercomputing*, 226–235.
- [66] ———. 1992. Array abstractions using semantic analysis of trapezoid congruences. In *Int. Conf. Supercomputing*.
- [67] Meyer, C.D. 2000. *Matrix analysis and applied linear algebra*. Philadelphia, PA: SIAM.
- [68] Miné, A. 2001. The octagon abstract domain. In *Working Conference on Reverse Engineering*, 310–322.
- [69] ———. 2002. A few graph-based relational numerical abstract domains. In *Static Analysis Symposium*.
- [70] ———. 2004. Relational abstract domains for the detection of floating-point run-time errors. In *European Symp. on Programming*.
- [71] ———. 2006. The octagon abstract domain. *J. Higher-Order and Symbolic Computation* 19(1):31–100.
- [72] ———. 2012. Abstract domains for bit-level machine integer and floating-point operations. In *Int. Joint Conf. on Automated Reasoning*.

- [73] Monniaux, D. 2010. Automatic modular abstractions for template numerical constraints. *Logical Methods in Comp. Sci.*
- [74] de Moura, L., and N. Bjørner. 2008. Z3: An efficient SMT solver. In *Int. Conf. on Tools and Algorithms for the Construction and Analysis of Systems*.
- [75] Müller-Olm, M., and H. Seidl. 2004. Precise interprocedural analysis through linear algebra. In *Symposium on Principles of Programming Languages*.
- [76] ———. 2005. Analysis of modular arithmetic. In *European Symp. on Programming*.
- [77] ———. 2005. Personal communication.
- [78] ———. 2007. Analysis of modular arithmetic. *Trans. on Prog. Lang. and Syst.* 29(5).
- [79] Mycroft, A., and N.D. Jones. 1985. A relational framework for abstract interpretation. In *Programs as data objects*.
- [80] Nelson, G., and D. Oppen. 1979. Simplification by cooperating decision procedures. *Trans. on Prog. Lang. and Syst.* 1(2).
- [81] Nielson, F. 1989. Two-level semantics and abstract interpretation. *Theor. Comp. Sci.* 69:117–242.
- [82] PPL. PPL: The Parma polyhedra library. www.cs.unipr.it/ppl/.
- [83] Reps, T., G. Balakrishnan, and J. Lim. 2006. Intermediate-representation recovery from low-level code. In *Part. eval. and semantics-based prog. manip.*
- [84] Reps, T., M. Sagiv, and G. Yorsh. 2004. Symbolic implementation of the best transformer. In *Int. Conf. on Verification, Model Checking, and Abstract Interpretation*.

- [85] Reps, T., S. Schwoon, S. Jha, and D. Melski. 2005. Weighted push-down systems and their application to interprocedural dataflow analysis. *Science of Computer Programming* 58(1–2).
- [86] Reps, T., and A. Thakur. 2016. Automating abstract interpretation. In *Int. Conf. on Verification, Model Checking, and Abstract Interpretation*.
- [87] Sagiv, M., T. Reps, and S. Horwitz. 1996. Precise interprocedural dataflow analysis with applications to constant propagation. *Theor. Comp. Sci.* 167:131–170.
- [88] Sankaranarayanan, S., F. Ivančić, I. Shlyakhter, and A. Gupta. 2006. Static analysis in disjunctive numerical domains. In *Static Analysis Symposium*.
- [89] Sankaranarayanan, S., H.B. Sipma, and Z. Manna. 2005. Scalable analysis of linear systems using mathematical programming. In *Int. Conf. on Verification, Model Checking, and Abstract Interpretation*.
- [90] Schmidt, D.A. 1986. *Denotational semantics*. Boston, MA: Allyn and Bacon, Inc.
- [91] Sen, R., and Y.N. Srikant. 2007. Executable analysis using abstract interpretation with circular linear progressions. In *Memocode*.
- [92] Sharir, M., and A. Pnueli. 1981. Two approaches to interprocedural data flow analysis. In *Program flow analysis: Theory and applications*. Prentice-Hall.
- [93] Sharma, T., and T. Reps. 2017. A new abstraction framework for affine transformers. In *Static Analysis Symposium*.
- [94] ———. 2017. Sound bit-precise numerical domains. In *Int. Conf. on Verification, Model Checking, and Abstract Interpretation*.

- [95] Sharma, T., A. Thakur, and T. Reps. 2013. An abstract domain for bit-vector inequalities. Tech. Rep. TR-1789, Comp. Sci. Dept., Univ. of Wisconsin, Madison, WI.
- [96] Simon, A., and A. King. 2007. Taming the wrapping of integer arithmetic. In *Static Analysis Symposium*.
- [97] Simon, A., A. King, and J.M. Howe. 2002. Two variables per linear inequality as an abstract domain. In *Int. Workshop on Logic Based Prog. Dev. and Transformation*, 71–89.
- [98] Spolsky, Joel. 2007. Explaining the excel bug. <https://www.joelonsoftware.com/2007/09/26/explaining-the-excel-bug>.
- [99] Storjohann, A. 2000. Algorithms for matrix canonical forms. Ph.D. thesis, ETH Zurich, Zurich, Switzerland. Diss. ETH No. 13922.
- [100] Thakur, A., M. Elder, and T. Reps. 2012. Bilateral algorithms for symbolic abstraction. In *Static Analysis Symposium*.
- [101] Thakur, A., and T. Reps. 2012. A method for symbolic computation of abstract operations. In *Int. Conf. on Computer Aided Verification*.
- [102] Warren, H.S., Jr. 2003. *Hacker's delight*. Addison-Wesley.