

Bilateral Algorithms for Symbolic Abstraction^{*}

Aditya Thakur¹, Matt Elder¹, and Thomas Reps^{1,2**}

¹ University of Wisconsin; Madison, WI, USA

² GrammaTech, Inc.; Ithaca, NY, USA

Abstract. Given a concrete domain \mathcal{C} , a concrete operation $\tau : \mathcal{C} \rightarrow \mathcal{C}$, and an abstract domain \mathcal{A} , a fundamental problem in abstract interpretation is to find the *best abstract transformer* $\tau^\# : \mathcal{A} \rightarrow \mathcal{A}$ that over-approximates τ . This problem, as well as several other operations needed by an abstract interpreter, can be reduced to the problem of *symbolic abstraction*: the symbolic abstraction of a formula φ in logic \mathcal{L} , denoted by $\hat{\alpha}(\varphi)$, is the best value in \mathcal{A} that over-approximates the meaning of φ . When the concrete semantics of τ is defined in \mathcal{L} using a formula φ_τ that specifies the relation between input and output states, the best abstract transformer $\tau^\#$ can be computed as $\hat{\alpha}(\varphi_\tau)$.

In this paper, we present a new framework for performing symbolic abstraction, discuss its properties, and present several instantiations for various logics and abstract domains. The key innovation is to use a *bilateral* successive-approximation algorithm, which maintains both an over-approximation and an under-approximation of the desired answer.

1 Introduction

For several years, we have been investigating connections between abstract interpretation and logic—in particular, how to harness decision procedures to obtain algorithms for several fundamental primitives used in abstract interpretation. Automation ensures correctness and precision of these primitives [3, §1.1], and drastically reduces the time taken to implement the primitives [19, §2.5] This paper presents new results on this topic.

Like several previous papers [25, 15, 11, 34], this paper concentrates on the problem of developing an algorithm for *symbolic abstraction*: the symbolic abstraction of a formula φ in logic \mathcal{L} , denoted by $\hat{\alpha}(\varphi)$, is the best value in a given abstract domain \mathcal{A} that over-approximates the meaning of φ [25]. To be more precise, given a formula $\varphi \in \mathcal{L}$, let $\llbracket \varphi \rrbracket$ denote the meaning of φ —i.e., the set of

^{*} Supported, in part, by NSF under grants CCF-{0810053, 0904371}, by ONR under grants N00014-{09-1-0510, 10-M-0251, 11-C-0447}, by ARL under grant W911NF-09-1-0413, by AFRL under grants FA9550-09-1-0279 and FA8650-10-C-7088; and by DARPA under cooperative agreement HR0011-12-2-0012. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the authors, and do not necessarily reflect the views of the sponsoring agencies.

^{**} T. Reps has an ownership interest in GrammaTech, Inc., which has licensed elements of the technology discussed in this publication.

concrete states that satisfy φ . Then $\widehat{\alpha}(\varphi)$ is the unique value $a \in \mathcal{A}$ such that (i) $\llbracket \varphi \rrbracket \subseteq \gamma(a)$, and (ii) for all $a' \in \mathcal{A}$ for which $\llbracket \varphi \rrbracket \subseteq \gamma(a')$, $a \sqsubseteq a'$. In this paper, we present a new framework for performing symbolic abstraction, discuss its properties, and present several instantiations for various logics and abstract domains.

Several key operations needed by an abstract interpreter can be reduced to symbolic abstraction. For instance, one use of symbolic abstraction is to bridge the gap between concrete semantics and an abstract domain. Cousot and Cousot [5] gave a *specification* of the most-precise abstract interpretation of a concrete operation τ that is possible in a given abstract domain:

Given a Galois connection $\mathcal{C} \xleftrightarrow[\alpha]{\gamma} \mathcal{A}$, the *best abstract transformer*, $\tau^\# : \mathcal{A} \rightarrow \mathcal{A}$, is the most precise abstract operator possible that over-approximates τ . $\tau^\#$ can be expressed as follows: $\tau^\# = \alpha \circ \tau \circ \gamma$.

The latter equation defines the limit of precision obtainable using abstraction \mathcal{A} . However, the definition is non-constructive; it does not provide an *algorithm*, either for applying $\tau^\#$ or for finding a representation of the function $\tau^\#$. In particular, in many cases, the explicit application of γ to an abstract value would yield an intermediate result—a set of concrete states—that is either infinite or too large to fit in computer memory.

In contrast, it is often convenient to use a logic \mathcal{L} to state the concrete semantics of transformer τ as a formula $\varphi_\tau \in \mathcal{L}$ that specifies the relation between input and output states. Then, using an *algorithm for symbolic abstraction*, a representation of $\tau^\#$ can be computed as $\widehat{\alpha}(\varphi_\tau)$.

To see how symbolic abstraction can yield better results than conventional approaches to the creation of abstract transformers, consider an example from machine-code analysis: the x86 instruction “`add bh, al`” adds `al`, the low-order byte of 32-bit register `eax`, to `bh`, the second-to-lowest byte of 32-bit register `ebx`. The semantics of this instruction can be expressed in quantifier-free bit-vector (QFBV) logic as

$$\varphi_I \stackrel{\text{def}}{=} \mathbf{ebx}' = \left(\begin{array}{l} (\mathbf{ebx} \ \& \ 0\text{xFFF00FF}) \\ | \ ((\mathbf{ebx} + 256 * (\mathbf{eax} \ \& \ 0\text{xFF})) \ \& \ 0\text{xFF00}) \end{array} \right) \wedge \mathbf{eax}' = \mathbf{eax}, \quad (1)$$

where “ $\&$ ” and “ $|$ ” denote bitwise-and and bitwise-or. Eqn. (1) shows that the semantics of the instruction involves non-linear bit-masking operations.

Now suppose that abstract domain \mathcal{A} is the domain of affine relations over integers mod 2^{32} [11]. For this domain, $\widehat{\alpha}(\varphi_I)$ is $(2^{16}\mathbf{ebx}' = 2^{16}\mathbf{ebx} + 2^{24}\mathbf{eax}) \wedge (\mathbf{eax}' = \mathbf{eax})$, which captures the relationship between the low-order two bytes of `ebx` and the low-order byte of `eax`. It is the best over-approximation to Eqn. (1) that can be expressed as an affine relation. In contrast, a more conventional approach to creating an abstract transformer for φ_I is to use operator-by-operator reinterpretation of Eqn. (1). The resulting abstract transformer would be $(\mathbf{eax}' = \mathbf{eax})$, which loses all information about `ebx`. Such loss in precision is exacerbated when considering larger loop-free blocks of instructions.

Motivation. Reps, Sagiv, and Yorsh (RSY) [25] presented a framework for computing $\hat{\alpha}$ that applies to any logic and abstract domain that satisfies certain conditions. King and Søndergaard [15] gave a specific $\hat{\alpha}$ algorithm for an abstract domain of Boolean affine relations. Elder et al. [11] extended their algorithm to affine relations in arithmetic modulo 2^w —i.e., for some bit-width w of bounded integers. (When the generalized algorithm is applied to φ_I from Eqn. (1), it finds the $\hat{\alpha}(\varphi_I)$ formula indicated above.) Because the generalized algorithm is similar to the Boolean one, we refer to it as KS. We use RSY[AR] to denote the RSY framework instantiated for the abstract domain of affine relations modulo 2^w .

The RSY[AR] and KS algorithms resemble one another in that they both find $\hat{\alpha}(\varphi)$ via successive approximation from “below”. However, *the two algorithms are not the same*. As discussed in §2, although both the RSY[AR] and KS algorithms issue queries to a decision procedure, compared to the RSY[AR] algorithm, the KS algorithm issues *comparatively inexpensive* decision-procedure queries. Moreover, the differences in the two algorithms cause an order-of-magnitude difference in performance: in our experiments, *KS is approximately ten times faster* than RSY[AR].

These issues motivated us to (i) investigate the fundamental principles underlying the difference between the RSY[AR] and KS algorithms, and (ii) seek a *framework* into which the KS algorithm could be placed, so that its advantages could be transferred to other domains. A third motivating issue was that neither the RSY framework nor the KS algorithm are resilient to timeouts. Because the algorithms maintain only under-approximations of the desired answer, if the successive-approximation process takes too much time and needs to be stopped, they must return \top to be sound. We desired an algorithm that could return a nontrivial (non- \top) value in case of a timeout.

- The outcome of our work is a new *framework* for symbolic abstraction that
- is applicable to any abstract domain that satisfies certain conditions (similar to the RSY algorithm)
 - uses a successive-approximation algorithm that is *parsimonious* in its use of the decision procedure (similar to the KS algorithm)
 - is *bilateral*; that is, it maintains both an under-approximation and a (non-trivial) over-approximation of the desired answer, and hence is resilient to timeouts: the procedure can return the over-approximation if it is stopped at any point (unlike the RSY and KS algorithms).

The key concept used in generalizing the KS algorithm is an operation that we call **AbstractConsequence** (Defn. 1, §3). We show that many abstract domains have an **AbstractConsequence** operation that enables the kind of inexpensive decision-procedure queries that we see in the KS algorithm (Thm. 2, §3).

Our experiments show that the bilateral algorithm for the AR domain improves precision at up to 15% of a program’s control points (i.e., the beginning of a basic block that ends with a branch), and on average is more precise for 3.1% of the control points (computed as the arithmetic mean).

Algorithm 1: $\hat{\alpha}_{\text{RSY}}^\uparrow(\mathcal{L}, \mathcal{A})(\varphi)$	Algorithm 2: $\hat{\alpha}_{\text{KS}}^\uparrow(\varphi)$
<pre> 1 lower ← ⊥ 2 3 while true do 4 5 S ← Model($\varphi \wedge \neg\hat{\gamma}(\text{lower})$) 6 if S is TimeOut then 7 return ⊥ 8 else if S is None then 9 break // $\varphi \Rightarrow \hat{\gamma}(\text{lower})$ 10 else // $S \not\models \hat{\gamma}(\text{lower})$ 11 lower ← lower \sqcup $\beta(S)$ 12 ans ← lower 13 return ans </pre>	<pre> 1 lower ← ⊥ 2 i ← 1 3 while i ≤ rows(lower) do 4 p ← Row(lower, -i) // $p \sqsupseteq \text{lower}$ 5 S ← Model($\varphi \wedge \neg\hat{\gamma}(p)$) 6 if S is TimeOut then 7 return ⊥ 8 else if S is None then 9 i ← i + 1 // $\varphi \Rightarrow \hat{\gamma}(p)$ 10 else // $S \not\models \hat{\gamma}(p)$ 11 lower ← lower \sqcup $\beta(S)$ 12 ans ← lower 13 return ans </pre>

Contributions. The contributions of the paper can be summarized as follows:

- We show how the KS algorithm can be modified into a *bilateral algorithm* that maintains sound under- and over-approximations of the answer (§2).
- We present a framework for symbolic abstraction based on a bilateral algorithm for computing $\hat{\alpha}$ (§3).
- We give several instantiations of the framework (§3 and §4).
- We compare the performance of various algorithms (§2 and §5).

§6 discusses related work. A longer version is available as a technical report [32].

2 Towards a Bilateral Algorithm

Alg. 1 shows the general RSY algorithm ($\hat{\alpha}_{\text{RSY}}^\uparrow(\mathcal{L}, \mathcal{A})$) [25], which is parameterized on logic \mathcal{L} and abstract domain \mathcal{A} . Alg. 2 shows the KS algorithm ($\hat{\alpha}_{\text{KS}}^\uparrow$) [15, 11], which is specific to the QFBV logic and the affine-relations (AR) domain. The following notation is used in the algorithms:

- The operation of *symbolic concretization* (line 5 of Algs. 1 and 2), denoted by $\hat{\gamma}$, maps an abstract value $a \in \mathcal{A}$ to a formula $\hat{\gamma}(a) \in \mathcal{L}$ such that a and $\hat{\gamma}(a)$ represent the same set of concrete states (i.e., $\gamma(a) = \llbracket \hat{\gamma}(a) \rrbracket$).
- Given a formula $\psi \in \mathcal{L}$, $\text{Model}(\psi)$ returns (i) a satisfying model S if a decision procedure was able to determine that ψ is satisfiable in a given time limit, (ii) **None** if a decision procedure was able to determine that ψ is unsatisfiable in a given time limit, and (iii) **TimeOut** otherwise.
- The *representation function* β (line 11 of Algs. 1 and 2) maps a singleton concrete state $S \in \mathcal{C}$ to the least value in \mathcal{A} that over-approximates $\{S\}$.

An abstract value in the AR domain is a conjunction of affine equalities, which can be represented in a normal form as a matrix in which each row expresses a non-redundant affine equality [11]. (Rows are 0-indexed.) Given a matrix m , $\text{rows}(m)$ returns the number of rows of m (line 3 in Alg. 2), and $\text{Row}(m, -i)$, for $1 \leq i \leq \text{rows}(m)$, returns row $(\text{rows}(m) - i)$ of m (line 4 in Alg. 2).

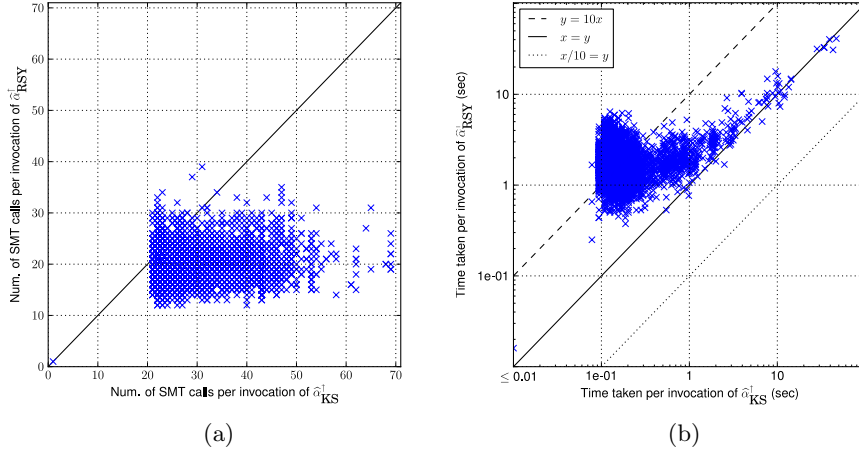


Fig. 1. (a) Scatter plot showing of the number of decision-procedure queries during each pair of invocations of $\hat{\alpha}_{RSY}^\uparrow$ and $\hat{\alpha}_{KS}^\uparrow$, when neither invocation had a decision-procedure timeout. (b) Log-log scatter plot showing the times taken by each pair of invocations of $\hat{\alpha}_{RSY}^\uparrow$ and $\hat{\alpha}_{KS}^\uparrow$, when neither invocation had a decision-procedure timeout.

Both algorithms have a similar overall structure. Both are successive approximation algorithms: they compute a sequence of successively “larger” approximations to the set of states described by φ . Both maintain an under-approximation of the final answer in the variable *lower*, which is initialized to \perp on line 1. Both call a decision procedure (line 5), and having found a model S that satisfies the query, the under-approximation is updated by performing a join (line 11).

The differences between Algs. 1 and 2 are highlighted in gray. The key difference is the nature of the decision-procedure query on line 5. $\hat{\alpha}_{RSY}^\uparrow$ uses *all* of *lower* to construct the query, while $\hat{\alpha}_{KS}^\uparrow$ uses only a single row from *lower* (line 4)—i.e., just a *single affine equality*, which has two consequences. First, $\hat{\alpha}_{KS}^\uparrow$ should issue a larger number of queries, compared with $\hat{\alpha}_{RSY}^\uparrow$. Suppose that the value of *lower* has converged to the final answer via a sequence of joins performed by the algorithm. To discover that convergence has occurred, $\hat{\alpha}_{RSY}^\uparrow$ has to issue just a single decision-procedure query, whereas $\hat{\alpha}_{KS}^\uparrow$ has to confirm it by issuing $\mathbf{rows}(\mathit{lower}) - i$ number of queries, proceeding row-by-row. Second, each individual query issued by $\hat{\alpha}_{KS}^\uparrow$ is simpler than the ones issued by $\hat{\alpha}_{RSY}^\uparrow$. Thus, *a priori*, it is not clear which algorithm will perform better in practice.

We compared the time for $\hat{\alpha}_{RSY}^\uparrow$ (instantiated for QFBV and the AR domain) and $\hat{\alpha}_{KS}^\uparrow$ to compute basic-block transformers for a set of x86 executables. There was no overall timeout imposed on the invocation of the procedures, but each invocation of the decision procedure (line 5 in Algs. 1 and 2) had a timeout of 3 seconds. (Details of the experimental setup are described in §5.) Fig. 1(a) shows a

Algorithm 3: $\hat{\alpha}_{\text{KS}}^\uparrow(\varphi)$	Algorithm 4: $\tilde{\alpha}_{\text{KS}^+}^\uparrow(\varphi)$
1	1 $upper \leftarrow \top$
2 $lower \leftarrow \perp$	2 $lower \leftarrow \perp$
3 $i \leftarrow 1$	3 $i \leftarrow 1$
4 while $i \leq \text{rows}(lower)$ do	4 while $i \leq \text{rows}(lower)$ do
5 $p \leftarrow \text{Row}(lower, -i)$	5 $p \leftarrow \text{Row}(lower, -i)$
// $p \sqsupseteq lower$	// $p \sqsupseteq lower, p \not\sqsupseteq upper$
6 $S \leftarrow \text{Model}(\varphi \wedge \neg\hat{\gamma}(p))$	6 $S \leftarrow \text{Model}(\varphi \wedge \neg\hat{\gamma}(p))$
7 if S is TimeOut then	7 if S is TimeOut then
8 return \top	8 return $upper$
9 else if S is None then	9 else if S is None then
// $\varphi \Rightarrow \hat{\gamma}(p)$	10 $upper \leftarrow upper \sqcap p$ // $\varphi \Rightarrow \hat{\gamma}(p)$
10 $i \leftarrow i + 1$	$i \leftarrow i + 1$
11 else // $S \not\models \hat{\gamma}(p)$	11 else // $S \not\models \hat{\gamma}(p)$
12 $lower \leftarrow lower \sqcup \beta(S)$	12 $lower \leftarrow lower \sqcup \beta(S)$
13 $ans \leftarrow lower$	13 $ans \leftarrow lower$
14 return ans	14 return ans

scatter-plot of the *number of decision-procedure calls* in each invocation of $\hat{\alpha}_{\text{RSY}}^\uparrow$ versus the corresponding invocation of $\hat{\alpha}_{\text{KS}}^\uparrow$, when neither of the procedures had a decision-procedure timeout. $\hat{\alpha}_{\text{RSY}}^\uparrow$ issues fewer decision-procedure queries: on average (computed as an arithmetic mean), $\hat{\alpha}_{\text{KS}}^\uparrow$ invokes 42% more calls to the decision procedure. Fig. 1(b) shows a log-log scatter-plot of the *total time* taken by each invocation of $\hat{\alpha}_{\text{RSY}}^\uparrow$ versus the time taken by $\hat{\alpha}_{\text{KS}}^\uparrow$. $\hat{\alpha}_{\text{KS}}^\uparrow$ is much faster than $\hat{\alpha}_{\text{RSY}}^\uparrow$: overall, computed as the geometric mean of the speedups on each of the x86 executables, $\hat{\alpha}_{\text{KS}}^\uparrow$ is about ten times faster than $\hat{\alpha}_{\text{RSY}}^\uparrow$.

The order-of-magnitude speedup can be attributed to the fact that each of the $\hat{\alpha}_{\text{KS}}^\uparrow$ decision-procedure queries is less expensive than the ones issued by $\hat{\alpha}_{\text{RSY}}^\uparrow$. At line 4 in $\hat{\alpha}_{\text{KS}}^\uparrow$, p is a single constraint; consequently, the decision-procedure query contains the *single* conjunct $\neg\hat{\gamma}(p)$ (line 5). In contrast, at line 5 in $\hat{\alpha}_{\text{RSY}}^\uparrow$, $lower$ is a *conjunction* of constraints, and consequently the decision-procedure query contains $\neg\hat{\gamma}(lower)$, which is a *disjunction* of constraints.

Neither $\hat{\alpha}_{\text{RSY}}^\uparrow$ nor $\hat{\alpha}_{\text{KS}}^\uparrow$ is resilient to timeouts. A decision-procedure query—or the cumulative time for $\hat{\alpha}^\uparrow$ —might take too long, in which case the only safe answer that can be returned is \top (line 6 in Algs. 1 and 2). To remedy this situation, we show how $\hat{\alpha}_{\text{KS}}^\uparrow$ can be modified to maintain a non-trivial over-approximation of the desired answer. Alg. 4 is such a *bilateral* algorithm that maintains both an under-approximation and over-approximation of $\hat{\alpha}(\varphi)$. The original $\hat{\alpha}_{\text{KS}}^\uparrow$ is shown in Alg. 3 for comparison; the differences in the algorithms are highlighted in gray. (Note that line numbers are different in Algs. 2 and 3.)

The $\tilde{\alpha}_{\text{KS}^+}^\uparrow$ algorithm (Alg. 4) initializes the over-approximation ($upper$) to \top on line 1. At any stage in the algorithm $\varphi \Rightarrow \hat{\gamma}(upper)$. On line 10, it is sound to update $upper$ by performing a meet with p because $\varphi \Rightarrow \hat{\gamma}(p)$. Progress is

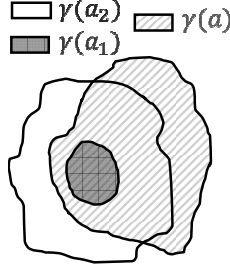


Fig. 2. Abstract Consequence:
 For all $a_1, a_2 \in \mathcal{A}$
 where $\gamma(a_1) \subseteq \gamma(a_2)$, if $a =$
 $\text{AbstractConsequence}(a_1, a_2)$,
 then $\gamma(a_1) \subseteq \gamma(a)$ and
 $\gamma(a) \not\subseteq \gamma(a_2)$.

Algorithm 5: $\tilde{\alpha}^\dagger\langle \mathcal{L}, \mathcal{A} \rangle(\varphi)$

```

1  upper ← ⊤
2  lower ← ⊥
3  while lower ≠ upper ∧ ResourcesLeft do
  // lower ⊈ upper
4  p ← AbstractConsequence(lower, upper)
  // p ⊇ lower, p ⊈ upper
5  S ← Model(φ ∧ ¬γ̂(p))
6  if S is TimeOut then
7    return upper
8  else if S is None then // φ ⇒ γ̂(p)
9    upper ← upper ⊓ p
11 else // S ⊈ γ̂(p)
12   lower ← lower ⊔ β(S)
13 return ans

```

guaranteed because $p \not\subseteq \text{upper}$. In case of a decision-procedure timeout (line 7), Alg. 4 returns upper as the answer (line 8). We use “ \sim ” to emphasize the fact that $\tilde{\alpha}_{\text{KS}^+}^\dagger(\varphi)$ can return an over-approximation of $\hat{\alpha}(\varphi)$ in case of a timeout. However, if the loop exits without a timeout, then $\tilde{\alpha}_{\text{KS}^+}^\dagger(\varphi)$ returns $\hat{\alpha}(\varphi)$.

3 A Parametric Bilateral Algorithm

Like the original KS algorithm, $\tilde{\alpha}_{\text{KS}^+}^\dagger$ applies only to the AR domain. The results presented in §2 provide motivation to generalize $\tilde{\alpha}_{\text{KS}^+}^\dagger$ so that we can take advantage of its benefits with domains other than AR. In this section, we present the bilateral framework we developed. Proofs for all theorems are found in [32].

We first introduce the *abstract-consequence* operation, which is the key operation in our generalized algorithm:

Definition 1. An operation $\text{AbstractConsequence}(\cdot, \cdot)$ is an **acceptable abstract-consequence operation** iff for all $a_1, a_2 \in \mathcal{A}$ such that $a_1 \not\subseteq a_2$, $a = \text{AbstractConsequence}(a_1, a_2)$ implies $a_1 \subseteq a$ and $a \not\subseteq a_2$. \square

Fig. 2 illustrates Defn. 1 graphically, using the concretizations of a_1 , a_2 , and a .

Alg. 5 presents the parametric bilateral algorithm $\tilde{\alpha}^\dagger\langle \mathcal{L}, \mathcal{A} \rangle(\varphi)$, which performs symbolic abstraction of $\varphi \in \mathcal{L}$ for abstract domain \mathcal{A} . The differences between Alg. 5 and Alg. 4 are highlighted in gray.

The assumptions placed on the logic and the abstract domain are as follows:

1. There is a Galois connection $\mathcal{C} \xrightleftharpoons[\alpha]{\gamma} \mathcal{A}$ between \mathcal{A} and concrete domain \mathcal{C} .
2. Given $a_1, a_2 \in \mathcal{A}$, there are algorithms to evaluate $a_1 \sqcup a_2$ and $a_1 \sqcap a_2$, and to check $a_1 = a_2$.

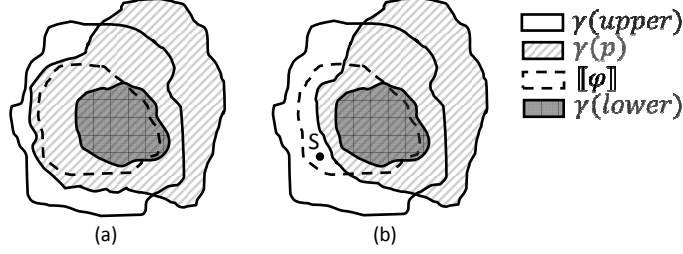


Fig. 3. The two cases arising in Alg. 5: $\varphi \wedge \neg \widehat{\gamma}(p)$ is either (a) unsatisfiable, or (b) satisfiable with $S \models \varphi$ and $S \not\models \widehat{\gamma}(p)$. (Note that although $lower \sqsubseteq \widehat{\alpha}(\varphi) \sqsubseteq upper$ and $\llbracket \varphi \rrbracket \subseteq \gamma(upper)$ are invariants of Alg. 5, $\gamma(lower) \subseteq \llbracket \varphi \rrbracket$ does not necessarily hold, as depicted above.)

3. There is a symbolic-concretization operation $\widehat{\gamma}$ that maps an abstract value $a \in \mathcal{A}$ to a formula $\widehat{\gamma}(a)$ in \mathcal{L} .
4. There is a decision procedure for the logic \mathcal{L} that is also capable of returning a model satisfying a formula in \mathcal{L} .
5. The logic \mathcal{L} is closed under conjunction and negation.
6. There is an acceptable abstract-consequence operation for \mathcal{A} (Defn. 1).

The abstract value p returned by **AbstractConsequence** (line 4 of Alg. 5) is used to generate the decision-procedure query (line 5); Fig. 3 illustrates the two cases arising based on whether $\varphi \wedge \neg \widehat{\gamma}(p)$ is satisfiable or unsatisfiable. The overall resources, such as time, used by Alg. 5 can be controlled via the **ResourcesLeft** flag (line 3).

Theorem 1. [Correctness of Alg. 5] *Suppose that \mathcal{L} and \mathcal{A} satisfy requirements 1–6, and $\varphi \in \mathcal{L}$. Let $a \in \mathcal{A}$ be the value returned by $\widetilde{\alpha}^\dagger(\mathcal{L}, \mathcal{A})(\varphi)$. Then*

1. *a over-approximates $\widehat{\alpha}(\varphi)$; i.e., $\widehat{\alpha}(\varphi) \sqsubseteq a$.*
2. *If \mathcal{A} has neither infinite ascending nor infinite descending chains and $\widetilde{\alpha}^\dagger(\mathcal{L}, \mathcal{A})(\varphi)$ returns with no timeout, then $a = \widehat{\alpha}(\varphi)$. \square*

Defn. 1 allows **AbstractConsequence**(a_1, a_2) to return any $a \in \mathcal{A}$ as long as a satisfies $a_1 \sqsubseteq a$ and $a \not\sqsupseteq a_2$. Thus, for a given abstract domain \mathcal{A} there could be multiple implementations of the **AbstractConsequence** operation. In particular, **AbstractConsequence**(a_1, a_2) can return a_1 , because $a_1 \sqsubseteq a_1$ and $a_1 \not\sqsupseteq a_2$. If this particular implementation of **AbstractConsequence** is used, then Alg. 5 reduces to the RSY algorithm (Alg. 1). However, as illustrated in §2, the decision-procedure queries issued by the RSY algorithm can be very expensive.

Conjunctive domains. We now define a class of *conjunctive domains*, for which **AbstractConsequence** can be implemented by the method presented as Alg. 6. The benefit of Alg. 6 is that it causes Alg. 5 to issue the kind of inexpensive queries that we see in $\widehat{\alpha}_{KS}^\uparrow$. Let Φ be a given set of formulas expressed in \mathcal{L} . A *conjunctive domain* over Φ is an abstract domain \mathcal{A} such that:

- For any $a \in \mathcal{A}$, there exists a finite subset $\Psi \subseteq \Phi$ such that $\widehat{\gamma}(a) = \bigwedge \Psi$.

Algorithm 6: AbstractConsequence(a_1, a_2) for conjunctive domains

```

1 if  $a_1 = \perp$  then return  $\perp$ 
2 Let  $\Psi \subseteq \Phi$  be the set of formulas such that  $\hat{\gamma}(a_1) = \bigwedge \Psi$ 
3 foreach  $\psi \in \Psi$  do
4    $a \leftarrow \mu\hat{\alpha}(\psi)$ 
5   if  $a \not\sqsupseteq a_2$  then return  $a$ 
    
```

- For any finite $\Psi \subseteq \Phi$, there exists an $a \in \mathcal{A}$ such that $\gamma(a) = \llbracket \bigwedge \Psi \rrbracket$.
- There is an algorithm $\mu\hat{\alpha}(\varphi)$ (“micro- $\hat{\alpha}$ ”) that, for each singleton formula $\varphi \in \Phi$, returns $a_\varphi \in \mathcal{A}$ such that $\hat{\alpha}(\varphi) = a_\varphi$.
- There is an algorithm that, for all $a_1, a_2 \in \mathcal{A}$, checks $a_1 \sqsubseteq a_2$.

Many common domains are conjunctive domains. For example, using v, v_i for program variables and c, c_i for constants:

Domain	Φ
Interval domain	inequalities of the form $c_1 \leq v$ and $v \leq c_2$
Octagon domain [20]	inequalities of the form $\pm v_1 \pm v_2 \leq c$
Polyhedral domain [7]	linear inequalities over reals or rationals
KS domain [15, 11]	linear equalities over integers mod 2^w

Theorem 2. When \mathcal{A} is a conjunctive domain over Φ , Alg. 6 is an acceptable abstract-consequence operation. \square

Discussion. We can weaken part 2 of Thm. 1 to allow \mathcal{A} to have infinite descending chains by modifying Alg. 5 slightly. The modified algorithm has to ensure that it does not get trapped updating *upper* along an infinite descending chain, and that it exits when *lower* has converged to $\hat{\alpha}(\varphi)$. We can accomplish these goals by forcing the algorithm to perform the basic RSY iteration step at least once every N iterations, for some fixed N . A version of Alg. 5 that implements this strategy is presented in [32].

As presented, Alg. 5 exits and returns the value of *upper* the first time the decision procedure times out. We can improve the precision of Alg. 5 by not exiting after the first timeout, and instead trying other abstract consequences. The algorithm will exit and return *upper* only if it cannot find an abstract consequence for which the decision-procedure terminates within the time bound. For conjunctive domains, Alg. 5 can be modified to enumerate all conjuncts of *lower* that are abstract consequences; to implement this strategy, lines 4–7 of Alg. 5 are replaced with

```

progress ← false // Initialize progress
foreach  $p$  such that  $p = \text{AbstractConsequence}(\text{lower}, \text{upper})$  do
   $S \leftarrow \text{Model}(\varphi \wedge \neg \hat{\gamma}(p))$ 
  if  $S$  is not TimeOut then
    progress ← true // Can make progress
    break
if  $\neg \text{progress}$  then return upper // Could not make progress
    
```

Henceforth, when we refer to $\tilde{\alpha}^\dagger$, we mean Alg. 5 with the above two changes.

Relationship of AbstractConsequence to interpolation. To avoid the potential for confusion, we now discuss how the notion of abstract consequence differs from the well-known concept of *interpolation* [8]:

A logic \mathcal{L} *supports interpolation* if for all $\varphi_1, \varphi_2 \in \mathcal{L}$ such that $\varphi_1 \Rightarrow \varphi_2$, there exists a formula I such that (i) $\varphi_1 \Rightarrow I$, (ii) $I \Rightarrow \varphi_2$, and (iii) I uses only symbols in the shared vocabulary of φ_1 and φ_2 .

Although condition (i) is part of Defn. 1, the restrictions imposed by conditions (ii) and (iii) are not part of Defn. 1. From an operational standpoint, condition (iii) in the definition of interpolation serves as a heuristic that generally allows interpolants to be expressed as small formulas. In the context of $\tilde{\alpha}^\dagger$, we are interested in obtaining small formulas to use in the decision-procedure query (line 5 of Alg. 5). Thus, given $a_1, a_2 \in \mathcal{A}$, it might appear plausible to use an interpolant I of $\hat{\gamma}(a_1)$ and $\hat{\gamma}(a_2)$ in $\tilde{\alpha}^\dagger$ instead of the abstract consequence of a_1 and a_2 . However, there are a few problems with such an approach:

- There is no guarantee that I will indeed be simple; for instance, if the vocabulary of $\hat{\gamma}(a_1)$ is a subset of the vocabulary of $\hat{\gamma}(a_2)$, then I could be $\hat{\gamma}(a_1)$ itself, in which case Alg. 5 performs the more expensive RSY iteration step.
- Converting the formula I into an abstract value $p \in \mathcal{A}$ for use in line 9 of Alg. 5 itself requires performing $\hat{\alpha}$ on I .

As discussed above, many domains are conjunctive domains, and for conjunctive domains is it always possible to find a *single conjunct* that is an abstract consequence (see Thm. 2). Moreover, such a conjunct is not necessarily an interpolant.

4 Instantiations

4.1 Herbrand-Equalities Domain

Herbrand equalities are used in analyses for partial redundancy elimination, loop-invariant code motion [30], and strength reduction [31]. In these analyses, arithmetic operations (e.g., $+$ and $*$) are treated as term constructors. Two program variables are known to hold equal values if the analyzer determines that the variables hold equal terms. Herbrand equalities can also be used to analyze programs whose types are user-defined algebraic data-types.

Basic definitions. Let \mathcal{F} be a set of function symbols. The function *arity*: $\mathcal{F} \rightarrow \mathbb{N}$ yields the number of parameters of each function symbol. *Terms* over \mathcal{F} are defined in the usual way; each function symbol f always requires $\text{arity}(f)$ parameters. Let $\mathcal{T}(\mathcal{F}, X)$ denote the set of finite terms generated by \mathcal{F} and variable set X . The *Herbrand universe* of \mathcal{F} is $\mathcal{T}(\mathcal{F}, \emptyset)$, the set of *ground terms* over \mathcal{F} .

A *Herbrand state* is a mapping from program variables \mathcal{V} to ground terms (i.e., a function in $\mathcal{V} \rightarrow \mathcal{T}(\mathcal{F}, \emptyset)$). The concrete domain consists of all sets of Herbrand states: $\mathcal{C} \stackrel{\text{def}}{=} \mathcal{P}(\mathcal{V} \rightarrow \mathcal{T}(\mathcal{F}, \emptyset))$. We can apply a Herbrand state σ to a term $t \in \mathcal{T}(\mathcal{F}, \mathcal{V})$ as follows:

$$\sigma[t] \stackrel{\text{def}}{=} \begin{cases} \sigma(t) & \text{if } t \in \mathcal{V} \\ f(\sigma[t_1], \dots, \sigma[t_k]) & \text{if } t = f(t_1, \dots, t_k) \end{cases}$$

The Herbrand-equalities domain. Sets of Herbrand states can be abstracted in several ways. One way is to use conjunctions of equations among terms (whence the name “Herbrand-equalities domain”). Such systems of equations can be represented using Equivalence DAGs [30]. A different, but equivalent, approach is to use a representation based on *idempotent substitutions*: $\mathcal{A} = (\mathcal{V} \rightarrow \mathcal{T}(\mathcal{F}, \mathcal{V}))_{\perp}$. Idempotence means that for each $\sigma \neq \perp$ and $v \in \mathcal{V}$, $\sigma[\sigma(v)] = \sigma(v)$. The meaning of an idempotent substitution $\sigma \in \mathcal{A}$ is given by its concretization, $\gamma: \mathcal{A} \rightarrow \mathcal{C}$, where $\gamma(\perp) = \emptyset$, and otherwise

$$\gamma(\sigma) = \{\rho: \mathcal{V} \rightarrow \mathcal{T}(\mathcal{F}, \emptyset) \mid \forall v \in \mathcal{V}: \rho(v) = \rho[\sigma(v)]\}. \quad (2)$$

We now show that the Herbrand-equalities domain satisfies the requirements of the bilateral framework. We will assume that the logical language \mathcal{L} has all the function symbols and constant symbols from \mathcal{F} , equality, and a constant symbol for each element from \mathcal{V} . (In a minor abuse of notation, the set of such constant symbols will also be denoted by \mathcal{V} .) The logic’s universe is the Herbrand universe of \mathcal{F} (i.e., $\mathcal{T}(\mathcal{F}, \emptyset)$). An interpretation maps the constants in \mathcal{V} to terms in $\mathcal{T}(\mathcal{F}, \emptyset)$. To be able to express $\hat{\gamma}(p)$ and $\neg\hat{\gamma}(p)$ (see item 5 below), we assume that \mathcal{L} contains at least the following productions:

$$\begin{aligned} F &::= F \wedge F \mid \neg F \mid v = T \text{ for } v \in \mathcal{V} \mid \text{false} \\ T &::= v \in \mathcal{V} \mid f(T_1, \dots, T_k) \text{ when } \textit{arity}(f) = k \end{aligned} \quad (3)$$

1. There is a Galois connection $\mathcal{C} \xrightleftharpoons[\alpha]{\gamma} \mathcal{A}$:
 - The ordering on \mathcal{C} is the subset relation on sets of Herbrand states.
 - $\gamma(\sigma)$ is given in Eqn. (2).
 - $\alpha(S) = \prod \{a \mid \gamma(a) \supseteq S\}$.
 - For $a, b \in \mathcal{A}$, $a \sqsubseteq b$ iff $\gamma(a) \subseteq \gamma(b)$.
2. Meet is most-general unification of substitutions, computed by standard unification techniques [18, Thm. 3.1].
3. Join is most-specific generalization, computed by “dual unification” or “anti-unification” [23, 26], [18, Thm. 5.8].
4. Equality checking is described by Lassez et al. [18, Prop. 4.10].
5. $\hat{\gamma}: \hat{\gamma}(\perp) = \text{false}$; otherwise, $\hat{\gamma}(\sigma)$ is $\bigwedge_{v \in \mathcal{V}} v = \sigma(v)$.
6. One can obtain a decision procedure for \mathcal{L} formulas using the built-in datatype mechanism of, e.g., Z3 [9] or Yices [10], and obtain the necessary decision procedure using an existing SMT solver.
7. \mathcal{L} is closed under conjunction and negation.
8. **AbstractConsequence:** The domain is a conjunctive domain, as can be seen from the definition of $\hat{\gamma}$.

Thm. 1 ensures that Alg. 5 returns $\hat{\alpha}(\varphi)$ when abstract domain \mathcal{A} has neither infinite ascending nor infinite descending chains. The Herbrand-equalities domain has no infinite ascending chains [18, Lem. 3.15]. The domain described here also has no infinite descending chains, essentially because every right-hand

term in every Herbrand state has no variables but those in \mathcal{V} . (Worked examples of $\tilde{\alpha}^\dagger$ (Alg. 5) for the Herbrand-equalities domain are given in [32].)

4.2 Polyhedral Domain

An element of the polyhedral domain [7] is a convex polyhedron, bounded by hyperplanes. It may be unbounded in some directions. The symbolic concretization of a polyhedron is a conjunction of linear inequalities. The polyhedral domain is a conjunctive domain:

- Each polyhedron can be expressed as some conjunction of linear inequalities (“half-spaces”) from the set $\mathcal{F} = \{\sum_{v \in \mathcal{V}} c_v v \geq c \mid c, c_v \text{ are constants}\}$.
- Every finite conjunction of facts from \mathcal{F} can be represented as a polyhedron.
- $\mu\hat{\alpha}$: Each formula in \mathcal{F} corresponds to a simple, one-constraint polyhedron.
- There is an algorithm for comparing two polyhedra [7].

In addition, there are algorithms for join, meet, and checking equality.

The logic QF-LRA (quantifier-free linear real arithmetic) supported by SMT solvers provides a decision procedure for the fragment of logic that is required to express negation, conjunction, and $\hat{\gamma}$ of a polyhedron. Consequently, the polyhedral domain satisfies the bilateral framework, and therefore supports the $\tilde{\alpha}^\dagger$ algorithm. The polyhedral domain has both infinite ascending chains and infinite descending chains, and hence Alg. 5 is only guaranteed to compute an over-approximation of $\hat{\alpha}(\varphi)$.

Because the polyhedral domain is a conjunctive domain, if $a_1 \sqsubseteq a_2$, then some single constraint a of a_1 satisfies $a \sqsupseteq a_2$. For instance, for the polyhedra a_1 and a_2 in Fig. 4, the region a above the dotted line is an acceptable abstract consequence.

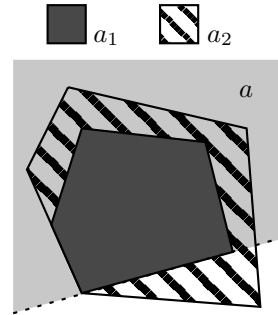


Fig. 4. Abs. conseq. for polyhedra.
 $a = \text{AbstractConsequence}(a_1, a_2)$

5 Experiments

In this section, we compare two algorithms for performing symbolic abstraction for the affine-relations (AR) domain [15, 11]:

- the $\tilde{\alpha}_{\text{KS}}^\dagger$ procedure of Alg. 2 [11].
- the $\tilde{\alpha}^\dagger(\text{AR})$ procedure that is the instantiation of Alg. 5 for the affine-relations (AR) domain and QFBV logic.

Although the bilateral algorithm $\tilde{\alpha}^\dagger(\text{AR})$ benefits from being resilient to time-outs, it maintains *both* an over-approximation and an under-approximation. Thus, the experiments were designed to understand the trade-off between performance and precision. In particular, the experiments were designed to answer the following questions:

Prog. name	Measures of size				Performance (x86)			Better
					$\hat{\alpha}_{KS}^\dagger$		$\tilde{\alpha}^\dagger\langle AR \rangle$	$\tilde{\alpha}^\dagger\langle AR \rangle$
	instrs	procs	BBs	brs	WPDS	t/o	WPDS	precision
finger	532	18	298	48	104.0	4	138.9	6.3%
subst	1093	16	609	74	196.7	4	214.6	0%
label	1167	16	573	103	146.1	2	171.6	0%
chkdsk	1468	18	787	119	377.2	16	417.9	0%
convert	1927	38	1013	161	287.1	10	310.5	0%
route	1982	40	931	243	618.4	14	589.9	2.5%
logoff	2470	46	1145	306	611.2	16	644.6	15.0%
setup	4751	67	1862	589	1499	60	1576	1.0%

Fig. 5. WPDS experiments. The columns show the number of instructions (instrs); the number of procedures (procs); the number of basic blocks (BBs); the number of branch instructions (brs); the times, in seconds, for $\hat{\alpha}_{KS}^\dagger$ and $\tilde{\alpha}^\dagger\langle AR \rangle$ WPDS construction; the number of invocations of $\hat{\alpha}_{KS}^\dagger$ that had a decision procedure timeout (t/o); and the degree of improvement gained by using $\tilde{\alpha}^\dagger\langle AR \rangle$ -generated ARA weights rather than $\hat{\alpha}_{KS}^\dagger$ weights (measured as the percentage of control points whose inferred one-vocabulary affine relation was strictly more precise under $\tilde{\alpha}^\dagger\langle AR \rangle$ -based analysis).

1. How does the speed of $\tilde{\alpha}^\dagger\langle AR \rangle$ compare with that of $\hat{\alpha}_{KS}^\dagger$?
2. How does the precision of $\tilde{\alpha}^\dagger\langle AR \rangle$ compare with that of $\hat{\alpha}_{KS}^\dagger$?

To address these questions, we performed affine-relations analysis (ARA) on x86 machine code, computing affine relations over the x86 registers. Our experiments were run on a single core of a quad-core 3.0 GHz Xeon computer running 64-bit Windows XP (SP2), configured so that a user process has 4GB of memory. We analyzed a corpus of Windows utilities using the WALi [14] system for weighted pushdown systems (WPDSs). For the $\hat{\alpha}_{KS}^\dagger$ -based ($\tilde{\alpha}^\dagger\langle AR \rangle$ -based) analysis we used a weight domain of $\hat{\alpha}^\dagger$ -generated ($\tilde{\alpha}^\dagger\langle AR \rangle$ -generated) ARA transformers. The weight on each WPDS rule encodes the ARA transformer for a basic block B of the program, including a jump or branch to a successor block. A formula φ_B is created that captures the concrete semantics of B , and then the ARA weight for B is obtained by performing $\hat{\alpha}(\varphi_B)$. We used EWPDS merge functions [17] to preserve caller-save and callee-save registers across call sites. The post* query used the FWPDS algorithm [16].

Fig. 5 lists several size parameters of the examples (number of instructions, procedures, basic blocks, and branches).³ Prior research [11] shows that the calls to $\hat{\alpha}$ during WPDS construction dominate the total time for ARA. Although the overall time taken by $\hat{\alpha}$ is not limited by a timeout, we use a 3-second timeout for each invocation of the decision procedure (as in Elder et al. [11]). Column 7 of Fig. 5 lists the number invocations of $\hat{\alpha}_{KS}^\dagger$ that had a decision-procedure timeout, and hence returned \top . (Note that, in general, $\hat{\alpha}_{KS}^\dagger$ implements an over-approximating $\tilde{\alpha}$ operation.)

³ Due to the high cost of the ARA-based 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)`”.

Columns 6 and 8 of Fig. 5 list the time taken, in seconds, for $\hat{\alpha}_{\text{KS}}^\uparrow$ and $\tilde{\alpha}^\uparrow\langle\text{AR}\rangle$ WPDS construction. We observe that on average $\tilde{\alpha}^\uparrow\langle\text{AR}\rangle$ is about 10% slower than $\hat{\alpha}_{\text{KS}}^\uparrow$ (computed as the geometric mean), which answers question 1.

To answer question 2 we compared the precision of the WPDS analysis when using $\hat{\alpha}_{\text{KS}}^\uparrow$ with the precision obtained using $\tilde{\alpha}^\uparrow\langle\text{AR}\rangle$. In particular, we compare the affine-relation invariants computed by the $\hat{\alpha}_{\text{KS}}^\uparrow$ -based and $\tilde{\alpha}^\uparrow\langle\text{AR}\rangle$ -based analyses for each *control point*—i.e., the beginning of a basic block that ends with a branch. The last column of Fig. 5 shows the percentage of control points for which the $\tilde{\alpha}^\uparrow\langle\text{AR}\rangle$ -based analysis computed a strictly more precise affine relation. We see that the $\tilde{\alpha}^\uparrow\langle\text{AR}\rangle$ -based analysis improves precision at up to 15% of control points, and, on average, the $\tilde{\alpha}^\uparrow\langle\text{AR}\rangle$ -based analysis is more precise for 3.1% of the control points (computed as the arithmetic mean), which answers question 2.

6 Related Work

6.1 Related Work on Symbolic Abstraction

Previous work on symbolic abstraction falls into three categories:

1. algorithms for specific domains [24, 3, 2, 15, 11]
2. algorithms for parameterized abstract domains [12, 35, 28, 22]
3. abstract-domain frameworks [25, 34].

What distinguishes category 3 from category 2 is that each of the results cited in category 2 applies to a specific *family* of abstract domains, defined by a *parameterized Galois connection* (e.g., with an abstraction function equipped with a readily identifiable parameter for controlling the abstraction). In contrast, the results in category 3 are defined by an *interface*; for any abstract domain that satisfies the requirements of the interface, one has a method for symbolic abstraction. The approach presented in this paper falls into category 3.

Algorithms for specific domains. Regehr and Reid [24] present a method that constructs abstract transformers for machine instructions, for interval and bitwise abstract domains. Their method does not call a SAT solver, but instead uses the physical processor (or a simulator of a processor) as a black box.

Brauer and King [3] developed a method that works from below to derive abstract transformers for the interval domain. Their method is based on an approach due to Monniaux [22] (see below), but they changed two aspects:

1. They express the concrete semantics with a Boolean formula (via “bit-blasting”), which allows a formula equivalent to $\forall x.\varphi$ to be obtained from φ (in CNF) by removing the x and $\neg x$ literals from all of the clauses of φ .
2. Whereas Monniaux’s method performs abstraction and then quantifier elimination, Brauer and King’s method performs quantifier elimination on the concrete specification, and then performs abstraction.

Barrett and King [2] describe a method for generating range and set abstractions for bit-vectors that are constrained by Boolean formulas. For range analysis, the algorithm separately computes the minimum and maximum value

of the range for an n -bit bit-vector using $2n$ calls to a SAT solver, with each SAT query determining a single bit of the output. The result is the best over-approximation of the value that an integer variable can take on (i.e., $\hat{\alpha}$).

Algorithms for parameterized abstract domains. Graf and Saïdi [12] showed that decision procedures can be used to generate best abstract transformers for predicate-abstraction domains. Other work has investigated more efficient methods to generate approximate transformers that are not best transformers, but approach the precision of best transformers [1, 4].

Yorsh et al. [35] developed a method that works from above to perform $\tilde{\alpha}(\varphi)$ for the kind of abstract domains used in shape analysis (i.e., “canonical abstraction” of logical structures [27]).

Template Constraint Matrices (TCMs) are a parametrized family of linear-inequality domains for expressing invariants in linear real arithmetic. Sankaranarayanan et al. [28] gave a parametrized meet, join, and set of abstract transformers for all TCM domains. Monniaux [22] gave an algorithm that finds the best transformer in a TCM domain across a straight-line block (assuming that concrete operations consist of piecewise linear functions), and good transformers across more complicated control flow. However, the algorithm uses quantifier elimination, and no polynomial-time elimination algorithm is known for piecewise-linear systems.

Abstract-domain frameworks. Thakur and Reps [34] recently discovered a new framework for performing symbolic abstraction from “above”: $\tilde{\alpha}^\downarrow$. The $\tilde{\alpha}^\downarrow$ framework builds upon the insight that Stålmarck’s algorithm for propositional validity checking [29] can be explained using abstract-interpretation terminology [33]. The $\tilde{\alpha}^\downarrow$ framework adapts the same algorithmic components of this generalization to perform symbolic abstraction. Because $\tilde{\alpha}^\downarrow$ maintains an over-approximation of $\hat{\alpha}$, it is resilient to timeouts.

The $\tilde{\alpha}^\downarrow$ framework is based on much different principles from the RSY and bilateral frameworks. The latter frameworks use an *inductive-learning approach* to learn from examples, while the $\tilde{\alpha}^\downarrow$ framework uses a *deductive approach* by using inference rules to deduce the answer. Thus, they represent two different classes of frameworks, with different requirements for the abstract domain.

6.2 Other Related Work

Cover algorithms. Gulwani and Musuvathi [13] defined what they termed the “cover problem”, which addresses *approximate existential quantifier elimination*: Given a formula φ in logic \mathcal{L} , and a set of variables V , find the strongest quantifier-free formula $\bar{\varphi}$ in \mathcal{L} such that $\llbracket \exists V : \varphi \rrbracket \subseteq \llbracket \bar{\varphi} \rrbracket$. They presented cover algorithms for the theories of uninterpreted functions and linear arithmetic, and showed that covers exist in some theories that do not support quantifier elimination.

The notion of a cover has similarities to the notion of symbolic abstraction, but the two notions are distinct. Our technical report [32] discusses the differences in detail, describing symbolic abstraction as over-approximating a formula

φ using an impoverished logic fragment (e.g., approximating an arbitrary QFBV formula, such as Eqn. (1), using conjunctions of modular-arithmetic affine equalities) while a cover algorithm only removes variables V from the vocabulary of φ . The two approaches yield different over-approximations of φ , and the over-approximation obtained by a cover algorithm does not, in general, yield suitable abstract values and abstract transformers.

Logical abstract domains. Cousot et al. [6] define a method of abstract interpretation based on using particular sets of logical formulas as abstract-domain elements (so-called *logical abstract domains*). They face the problems of (i) performing abstraction from unrestricted formulas to the elements of a logical abstract domain [6, §7.1], and (ii) creating abstract transformers that transform input elements of a logical abstract domain to output elements of the domain [6, §7.2]. Their problems are particular cases of $\widehat{\alpha}(\varphi)$. They present heuristic methods for creating over-approximations of $\widehat{\alpha}(\varphi)$.

Connections to machine-learning algorithms. In [25], a connection was made between symbolic abstraction (in abstract interpretation) and the problem of *concept learning* (in machine learning). In machine-learning terms, an abstract domain \mathcal{A} is a *hypothesis space*; each domain element corresponds to a *concept*. Given a formula φ , the symbolic-abstraction problem is to find the most specific concept that explains the meaning of φ .

$\widehat{\alpha}_{\text{RSY}}^{\uparrow}$ (Alg. 1) is related to the Find-S algorithm [21, §2.4] for concept learning. Both algorithms start with the most-specific hypothesis (i.e., \perp) and work bottom-up to find the most-specific hypothesis that is consistent with positive examples of the concept. Both algorithms generalize their current hypothesis each time they process a (positive) training example that is not explained by the current hypothesis. A major difference is that Find-S receives a sequence of positive and negative examples of the concept (e.g., from nature). It discards negative examples, and its generalization steps are based solely on the positive examples. In contrast, $\widehat{\alpha}_{\text{RSY}}^{\uparrow}$ repeatedly calls a decision procedure to generate the next positive example; $\widehat{\alpha}_{\text{RSY}}^{\uparrow}$ never sees a negative example.

A similar connection exists between $\widetilde{\alpha}^{\dagger}$ (Alg. 5) and a different concept-learning algorithm, called the Candidate-Elimination algorithm [21, §2.5]. Both algorithms maintain two approximations of the concept, one that is an over-approximation and one that is an under-approximation.

References

1. T. Ball, A. Podelski, and S. Rajamani. Boolean and Cartesian abstraction for model checking C programs. In *TACAS*, pages 268–283, 2001.
2. E. Barrett and A. King. Range and set abstraction using SAT. *ENTCS*, 267(1), 2010.
3. J. Brauer and A. King. Automatic abstraction for intervals using Boolean formulae. In *SAS*, 2010.
4. E. Clarke, D. Kroening, N. Sharygina, and K. Yorav. Predicate abstraction of ANSI-C programs using SAT. *FMSD*, 25(2–3), 2004.

5. P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In *POPL*, pages 269–282, 1979.
6. P. Cousot, R. Cousot, and L. Mauborgne. Logical abstract domains and interpretations. In *The Future of Software Engineering*, 2011.
7. P. Cousot and N. Halbwachs. Automatic discovery of linear constraints among variables of a program. In *POPL*, 1978.
8. W. Craig. Three uses of the Herbrand-Gentzen theorem in relating model theory and proof theory. *J. Sym. Logic*, 22(3), Sept. 1957.
9. L. de Moura and N. Bjørner. Z3: An efficient SMT solver. In *TACAS*, 2008.
10. B. Dutertre and L. de Moura. Yices: An SMT solver, 2006. yices.csl.sri.com/.
11. M. Elder, J. Lim, T. Sharma, T. Andersen, and T. Reps. Abstract domains of affine relations. In *SAS*, 2011.
12. S. Graf and H. Saïdi. Construction of abstract state graphs with PVS. In *CAV*, 1997.
13. S. Gulwani and M. Musuvathi. Cover algorithms and their combination. In *ESOP*, 2008.
14. N. Kidd, A. Lal, and T. Reps. WALi: The Weighted Automaton Library, 2007. www.cs.wisc.edu/wpis/wpds/download.php.
15. A. King and H. Søndergaard. Automatic abstraction for congruences. In *VMCAI*, 2010.
16. A. Lal and T. Reps. Improving pushdown system model checking. In *CAV*, 2006.
17. A. Lal, T. Reps, and G. Balakrishnan. Extended weighted pushdown systems. In *CAV*, 2005.
18. J. Lassez, M. Maher, and K. Marriott. Unification revisited. In *Foundations of Logic and Functional Programming*, volume 306, pages 67–113. Springer, 1988.
19. J. Lim and T. Reps. A system for generating static analyzers for machine instructions. In *CC*, 2008.
20. A. Miné. The octagon abstract domain. In *WCRE*, pages 310–322, 2001.
21. T. Mitchell. *Machine Learning*. WCB/McGraw-Hill, Boston, MA, 1997.
22. D. Monniaux. Automatic modular abstractions for template numerical constraints. *Logical Methods in Comp. Sci.*, 6(3), 2010.
23. G. Plotkin. A note on inductive generalization. In *Machine Intelligence*, volume 5, pages 153–165. Edinburgh Univ. Press, 1970.
24. J. Regehr and A. Reid. HOIST: A system for automatically deriving static analyzers for embedded systems. In *ASPLOS*, 2004.
25. T. Reps, M. Sagiv, and G. Yorsh. Symbolic implementation of the best transformer. In *VMCAI*, pages 252–266, 2004.
26. J. Reynolds. Transformational systems and the algebraic structure of atomic formulas. *Machine Intelligence*, 5(1):135–151, 1970.
27. M. Sagiv, T. Reps, and R. Wilhelm. Parametric shape analysis via 3-valued logic. *TOPLAS*, 24(3):217–298, 2002.
28. S. Sankaranarayanan, H. Sipma, and Z. Manna. Scalable analysis of linear systems using mathematical programming. In *VMCAI*, 2005.
29. M. Sheeran and G. Stålmarck. A tutorial on Stålmarck’s proof procedure for propositional logic. *FMSD*, 16(1):23–58, 2000.
30. B. Steffen, J. Knoop, and O. Rüthing. The value flow graph: A program representation for optimal program transformations. In *ESOP*, 1990.
31. B. Steffen, J. Knoop, and O. Rüthing. Efficient code motion and an adaption to strength reduction. In *TAPSOFT*, 1991.

32. A. Thakur, M. Elder, and T. Reps. Bilateral algorithms for symbolic abstraction. TR 1713, CS Dept., Univ. of Wisconsin, Madison, WI, Mar. 2012. www.cs.wisc.edu/wpis/papers/tr1713.pdf.
33. A. Thakur and T. Reps. A Generalization of Stålmarck's Method. In *SAS*, 2012.
34. A. Thakur and T. Reps. A method for symbolic computation of precise abstract operations. In *CAV*, 2012.
35. G. Yorsh, T. Reps, and M. Sagiv. Symbolically computing most-precise abstract operations for shape analysis. In *TACAS*, 2004.