

Analyzing Memory Accesses in x86 Executables

Gogul Balakrishnan

Thomas Reps

University of Wisconsin

Motivation

- Basic infrastructure for language-based security
 - buffer-overflow detection
 - information-flow vulnerabilities
 - ...
- What if we do not have source code?
 - viruses, worms, mobile code, etc.
 - legacy code (w/o source)
- Limitations of existing tools
 - over-conservative treatment of memory accesses
 - ⇒ Many false positives
 - unsafe treatment of pointer arithmetic
 - ⇒ Many false negatives

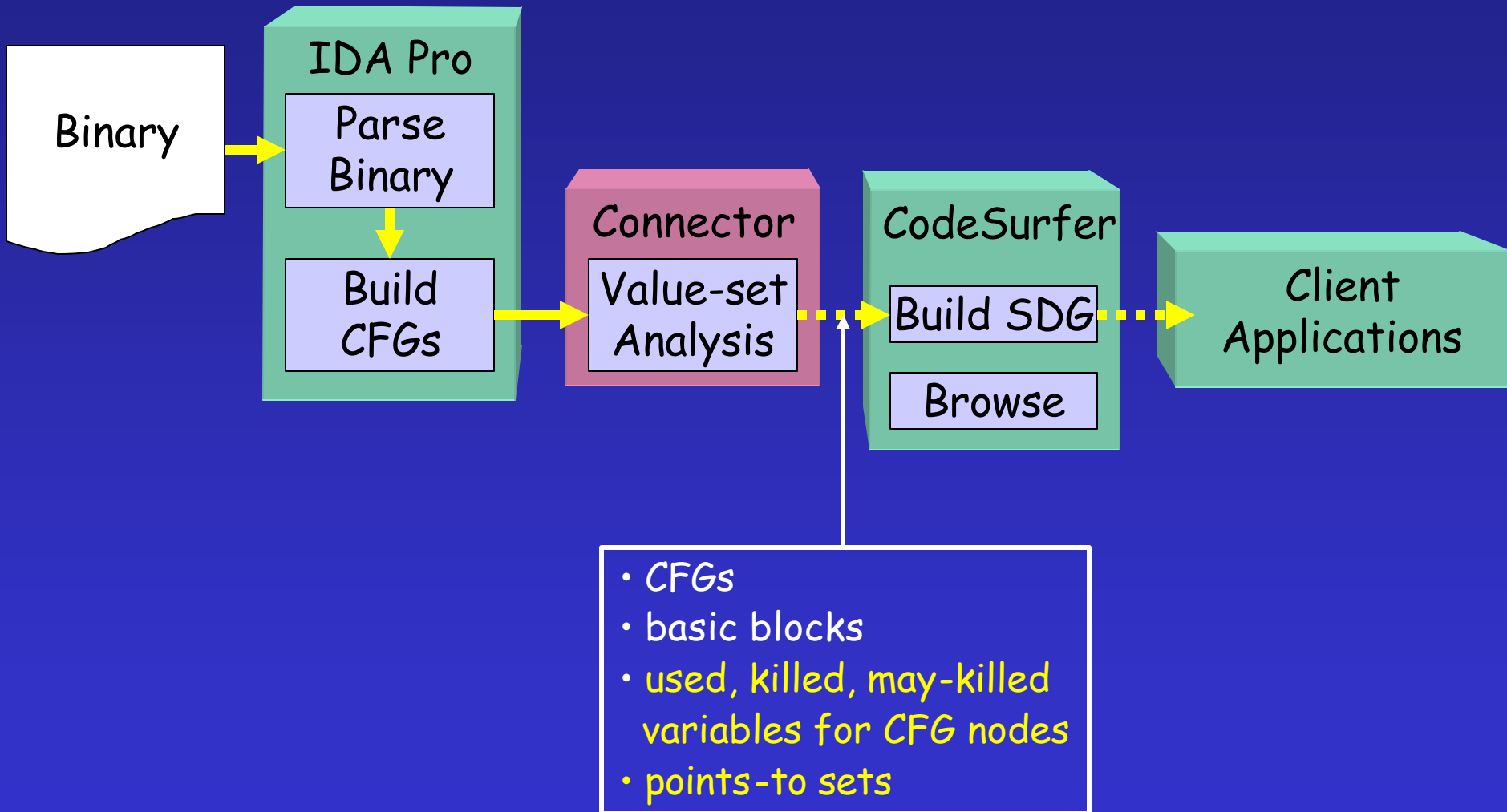
Goal (1)

- Create an **intermediate representation (IR)** that is similar to the IR used in a compiler
 - CFGs
 - used, killed, may-killed variables for CFG nodes
 - points-to sets
 - call-graph
- Why?
 - a tool for a security analyst
 - a general infrastructure for binary analysis

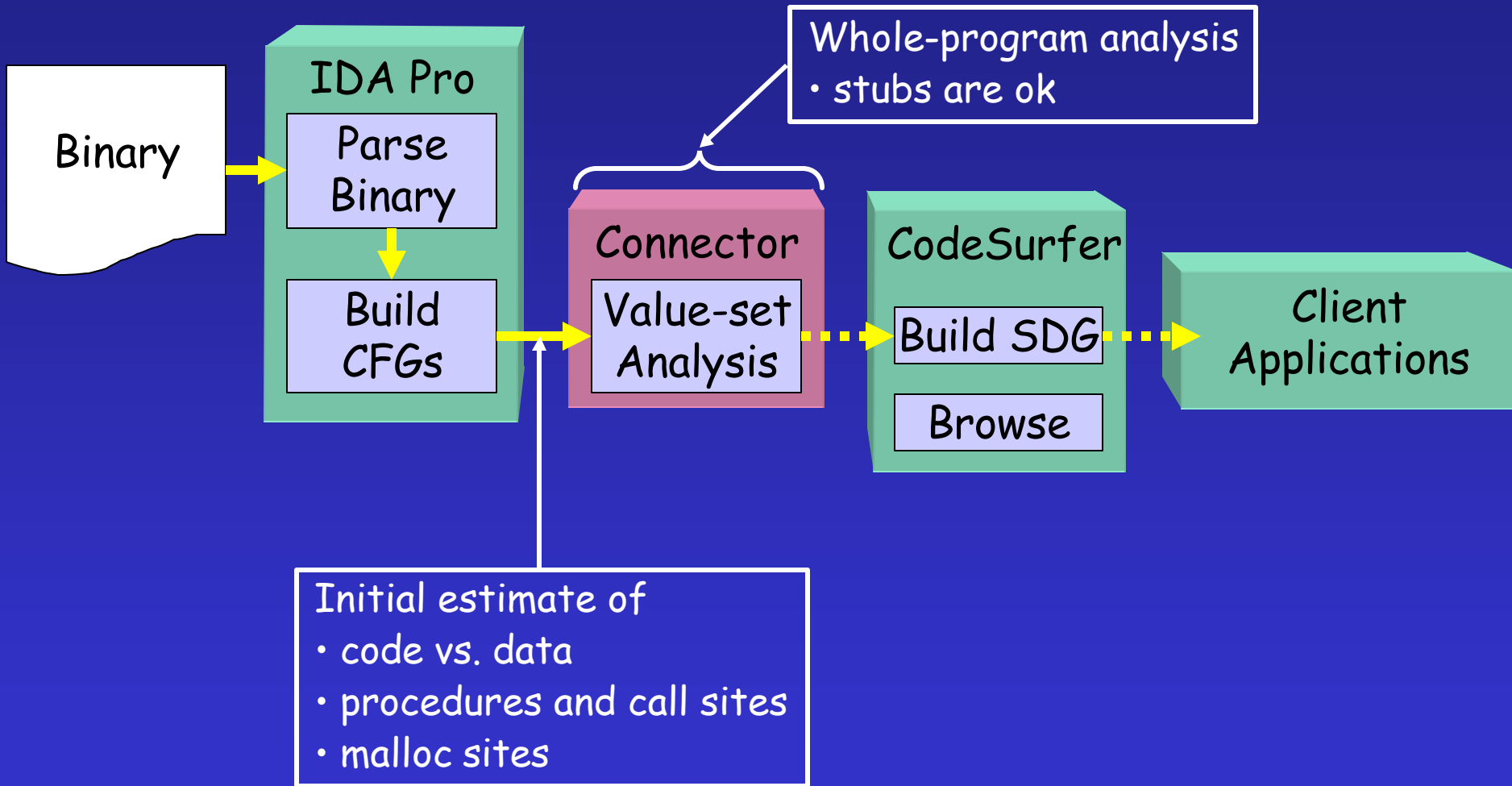
Goal (2)

- Scope: programs that conform to a "standard compilation model"
 - data layout determined by compiler
 - some variables held in registers
 - global variables \rightarrow absolute addresses
 - local variables \rightarrow offsets in `esp`-based stack frame
- Report violations
 - violations of stack protocol
 - return address modified within procedure

Codesurfer/x86 Architecture



Codesurfer/x86 Architecture



Outline

- Example
- Challenges
- Value-set analysis
- Performance
- [Future work]

Running Example

```
int arrVal=0, *pArray2;

int main() {
    int i, a[10], *p;
    /* Initialize pointers */
    pArray2=&a[2];
    p=&a[0];
    /* Initialize Array*/
    for(i=0; i<10; ++i) {
        *p=arrVal;
        p++;
    }
    /* Return a[2] */
    return *pArray2;
}
```

```
; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

sub     esp, 40           ;adjust stack
lea     edx, [esp+8]     ;
mov     [4], edx         ;pArray2=&a[2]
lea     ecx, [esp]       ;p=&a[0]
mov     edx, [0]         ;

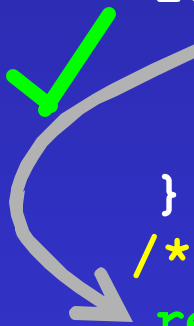
loc_9:
mov     [ecx], edx       ;*p=arrVal
add     ecx, 4           ;p++
inc     ebx              ;i++
cmp     ebx, 10          ;i<10?
j1     short loc_9      ;

mov     edi, [4]         ;
mov     eax, [edi]       ;return *pArray2
add     esp, 40
retn
```


Running Example

```
int arrVal=0, *pArray2;

int main() {
    int i, a[10], *p;
    /* Initialize pointers */
    pArray2=&a[2];
    p=&a[0];
    /* Initialize Array*/
    for(i=0; i<10; ++i) {
        *p=arrVal;
        p++;
    }
    /* Return a[2] */
    return *pArray2;
}
```

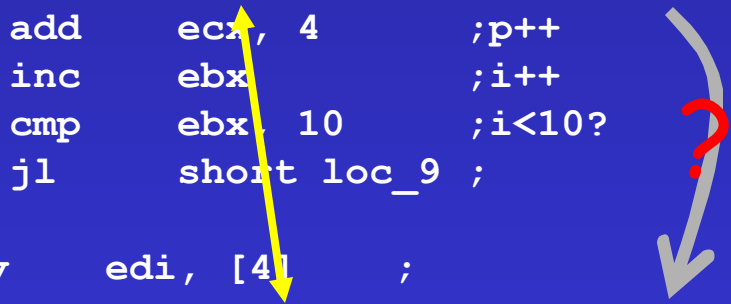


```
; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

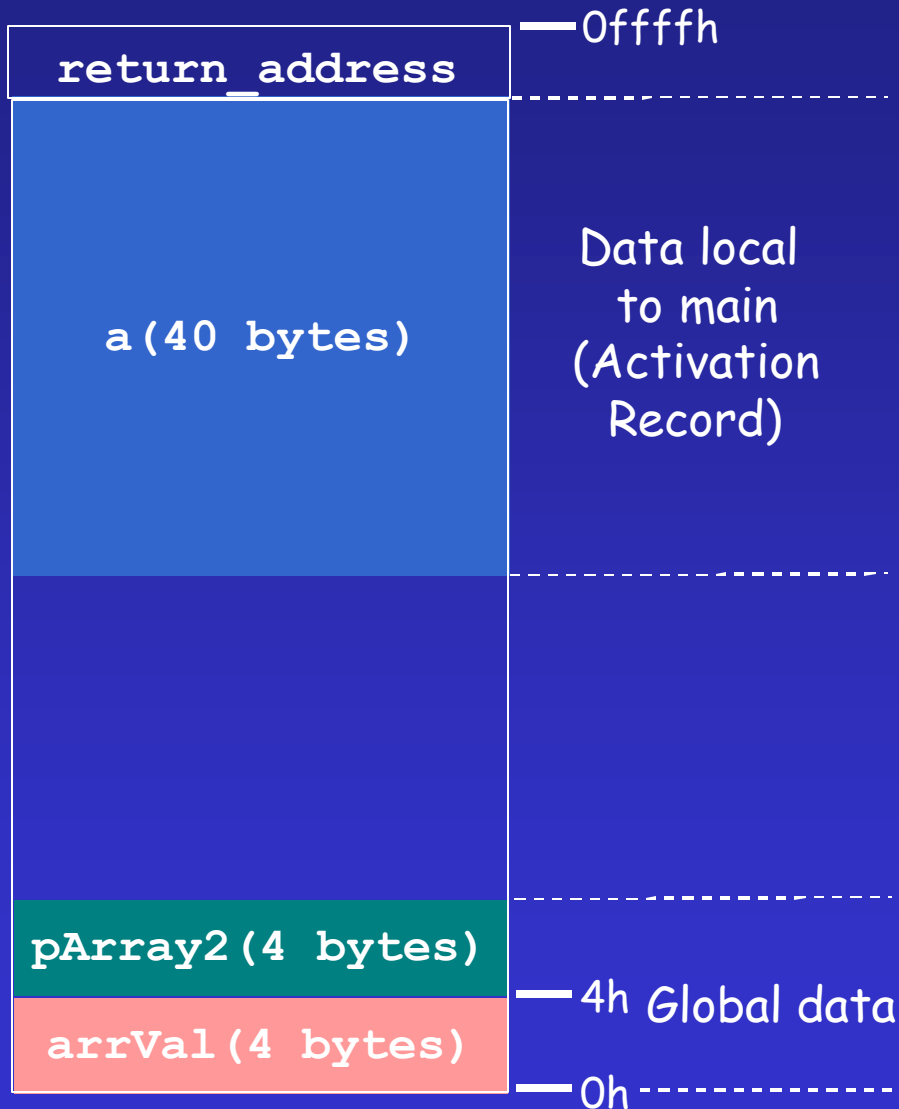
sub     esp, 40           ;adjust stack
lea     edx, [esp+8]     ;
mov     [4], edx         ;pArray2=&a[2]
lea     ecx, [esp]       ;p=&a[0]
mov     edx, [0]         ;

loc_9:
mov     [ecx], edx       ;*p=arrVal
add     ecx, 4           ;p++
inc     ebx              ;i++
cmp     ebx, 10         ;i<10?
j1     short loc_9      ;

mov     edi, [4]         ;
mov     eax, [edi] ;return *pArray2
add     esp, 40
retn
```



Running Example - Address Space



```
; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

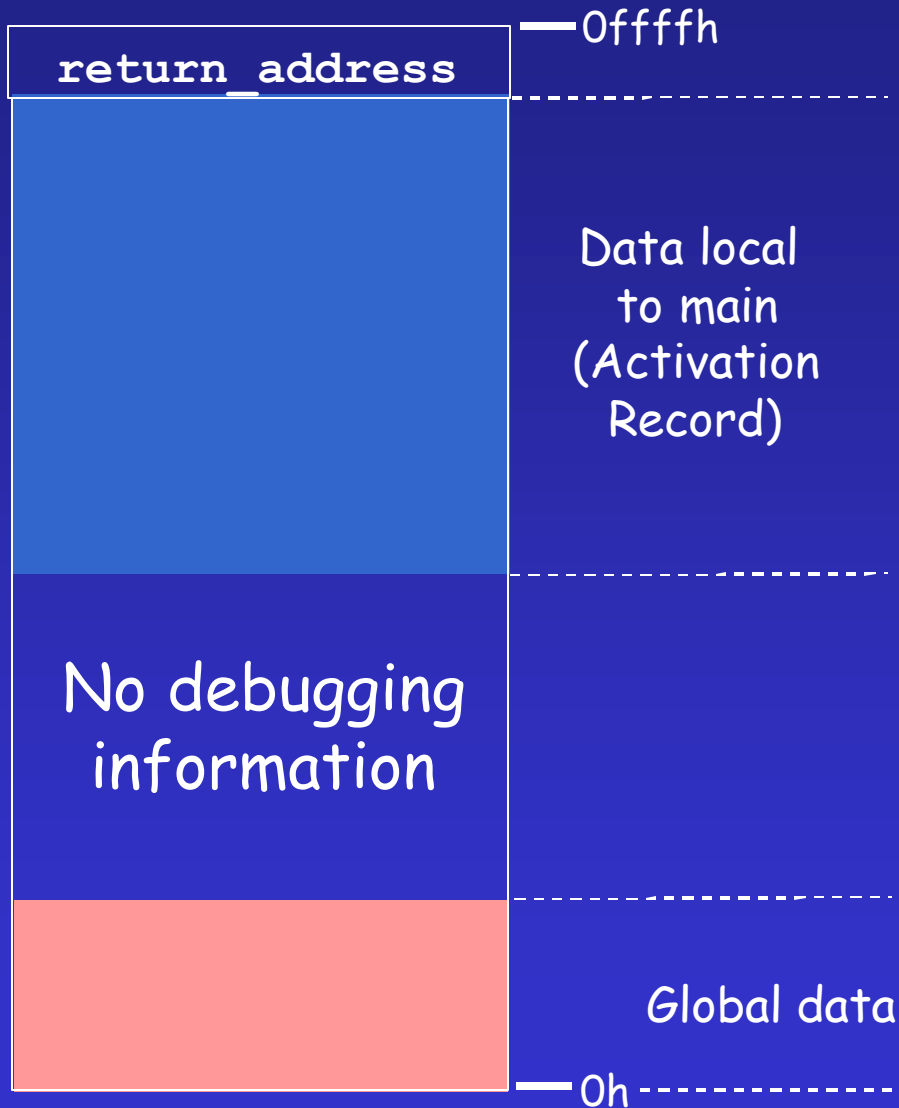
sub     esp, 40      ;adjust stack
lea     edx, [esp+8] ;
mov     [4], edx    ;pArray2=&a[2]
lea     ecx, [esp]  ;p=&a[0]
mov     edx, [0]    ;

loc_9:
mov     [ecx], edx  ;*p=arrVal
add     ecx, 4      ;p++
inc     ebx         ;i++
cmp     ebx, 10    ;i<10?
j1     short loc_9 ;

mov     edi, [4]    ;
mov     eax, [edi] ;return *pArray2
add     esp, 40
retn
```



Running Example - Address Space



```
; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

sub     esp, 40      ;adjust stack
lea     edx, [esp+8] ;
mov     [4], edx    ;pArray2=&a[2]
lea     ecx, [esp]  ;p=&a[0]
mov     edx, [0]    ;

loc_9:
mov     [ecx], edx  ;*p=arrVal
add     ecx, 4      ;p++
inc     ebx         ;i++
cmp     ebx, 10     ;i<10?
j1     short loc_9 ;

mov     edi, [4]    ;
mov     eax, [edi] ;return *pArray2
add     esp, 40
retn
```

A red question mark is placed to the right of the `loc_9` block, with a grey arrow pointing downwards from it towards the `mov eax, [edi]` instruction.

Challenges (1)

- No debugging/symbol-table information
- Explicit memory addresses
 - need something similar to C variables
 - **a-locs**
- Only have an initial estimate of
 - code, data, procedures, call sites, malloc sites
 - **extend IR on-the-fly**
 - disassemble data, add to CFG, ...
 - similar to elaboration of CFG/call-graph in a compiler because of calls via function pointers

Challenges (2)

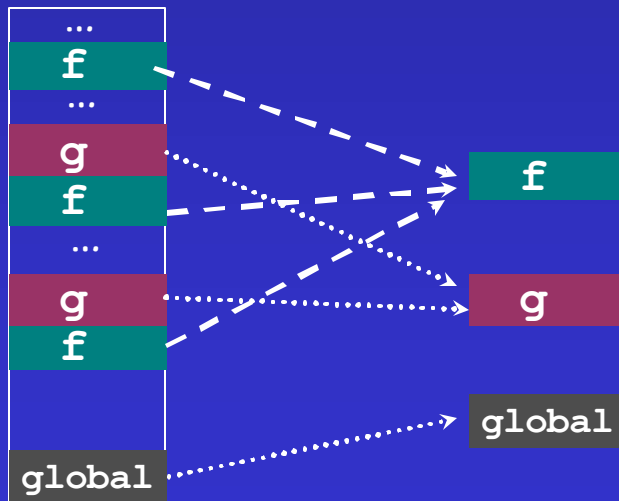
- Indirect-addressing mode
 - need "pointer analysis"
 - value-set analysis
- Pointer arithmetic
 - need numeric analysis (e.g., range analysis)
 - value-set analysis
- Checking for non-aligned accesses
 - pointer forging?
 - keep stride information in value-sets

Not Everything is Bad News !

- Multiple source languages OK
- Some optimizations make our task easier
 - optimizers try to use registers, not memory
 - **deciphering memory operations is the hard part**

Memory-regions

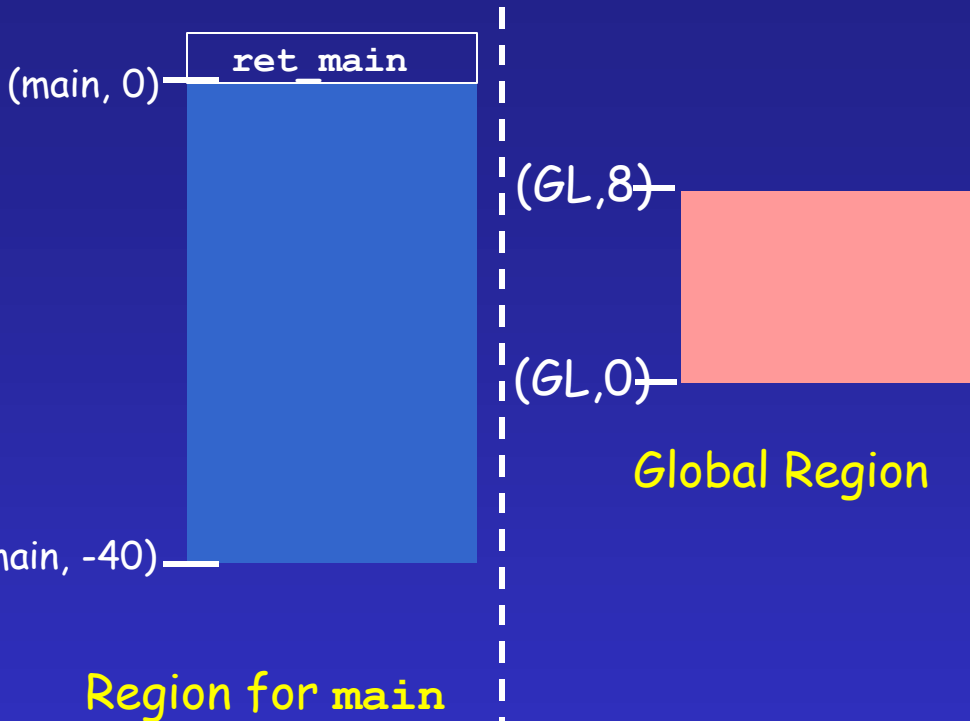
- An abstraction of the address space
- **Idea: group similar runtime addresses**
 - collapse the runtime ARs for each procedure



Memory-regions

- An abstraction of the address space
- **Idea: group similar runtime addresses**
 - collapse the runtime ARs for each procedure
- Similarly,
 - one region for all global data
 - one region for each malloc site

Example - Memory-regions



```
; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

sub    esp, 40      ;adjust stack
lea    edx, [esp+8] ;
mov    [4], edx     ;pArray2=&a[2]
lea    ecx, [esp]   ;p=&a[0]
mov    edx, [0]     ;

loc_9:
    mov    [ecx], edx ;*p=arrVal
    add    ecx, 4     ;p++
    inc    ebx        ;i++
    cmp    ebx, 10    ;i<10?
    jl    short loc_9 ;

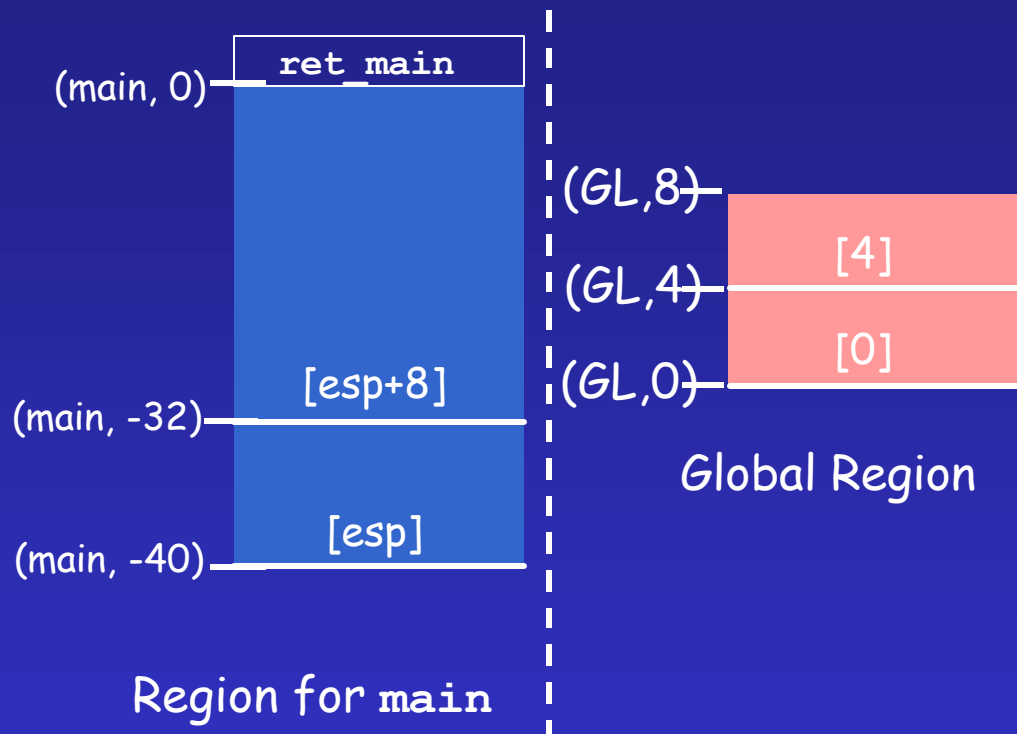
    mov    edi, [4]   ;
    mov    eax, [edi] ;return *pArray2
    add    esp, 40
    retn
```



"Need Something Similar to C Variables"

- Standard compilation model
 - some variables held in registers
 - global variables → absolute addresses
 - local variables → offsets in stack frame
- A-locs
 - locations between consecutive addresses
 - locations between consecutive offsets
 - registers
- Use a-locs instead of variables in static analysis
 - e.g., killed a-loc ≈ killed variable

Example - A-locs



```

; ebx ̂ i
; ecx ̂ variable p

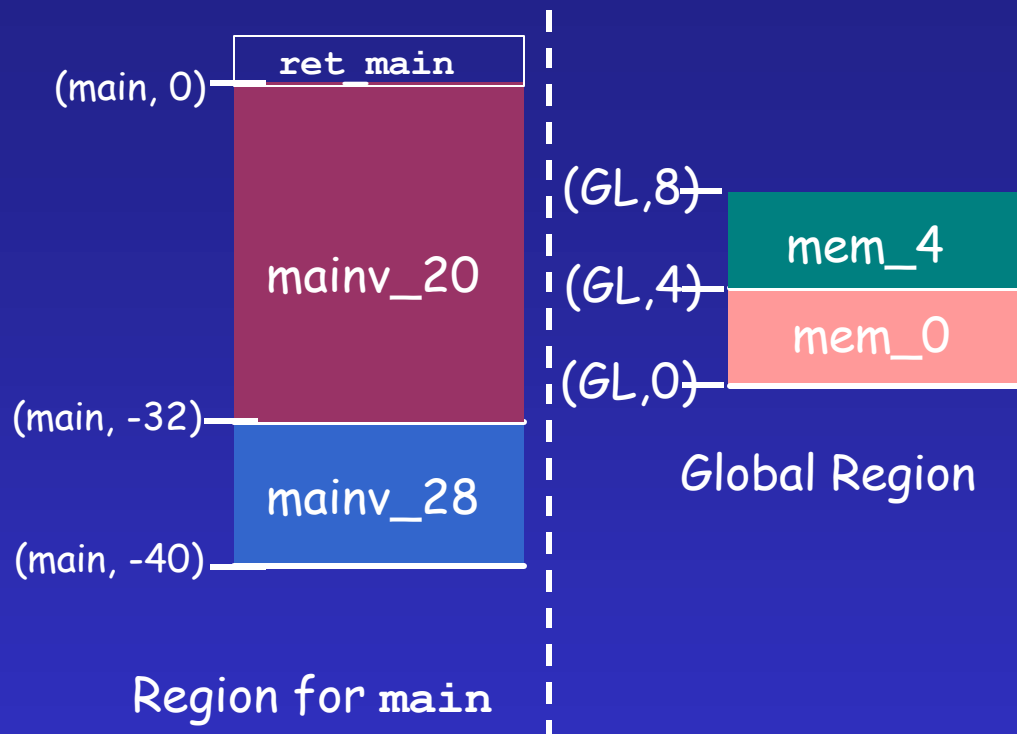
sub esp, 40 ;adjust stack
lea edx, [esp+8] ;
mov [4], edx ;pArray2=&a[2]
lea ecx, [esp] ;p=&a[0]
mov edx, [0] ;

loc_9:
mov [ecx], edx ;*p=arrVal
add ecx, 4 ;p++
inc ebx ;i++
cmp ebx, 10 ;i<10?
j1 short loc_9 ;

mov edi, [4] ;
mov eax, [edi] ;return *pArray2
add esp, 40
retn
    
```



Example - A-locs



```

; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

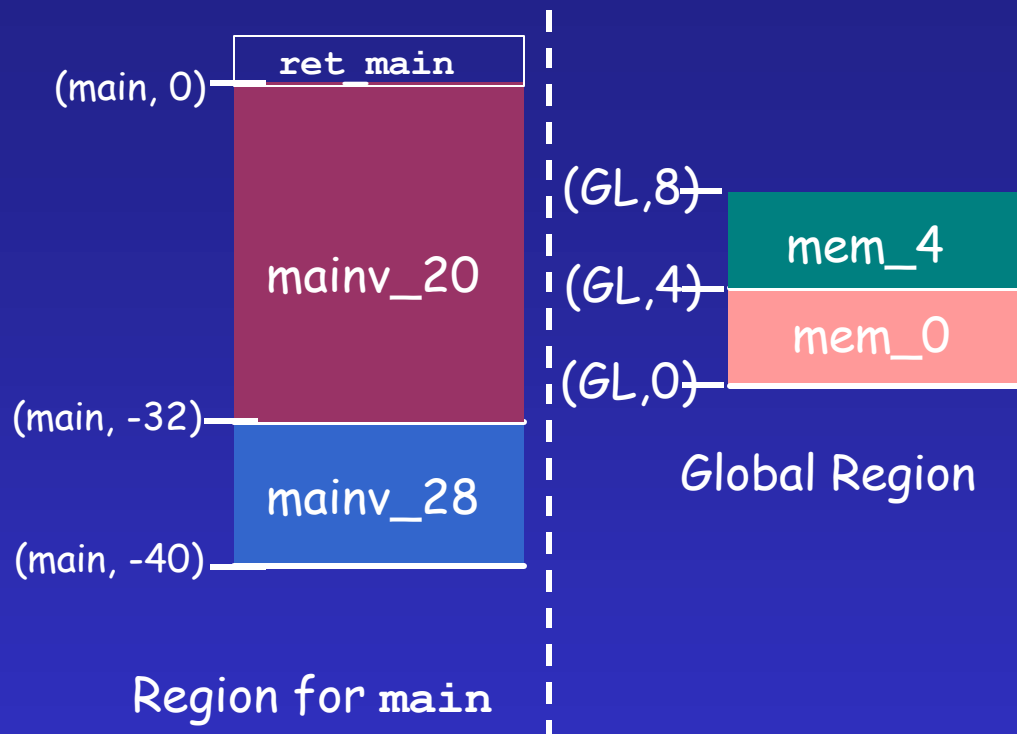
sub esp, 40 ;adjust stack
lea edx, [esp+8] ;
mov [4], edx ;pArray2=&a[2]
lea ecx, [esp] ;p=&a[0]
mov edx, [0] ;

loc_9:
mov [ecx], edx ;*p=arrVal
add ecx, 4 ;p++
inc ebx ;i++
cmp ebx, 10 ;i<10?
j1 short loc_9 ;

mov edi, [4] ;
mov eax, [edi] ;return *pArray2
add esp, 40
retn
    
```



Example - A-locs



```

; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

sub esp, 40 ;adjust stack
lea edx, &mainv_2 ;
mov mem_4, edx ;pArray2=&a[2]
lea ecx, &mainv_2 ;p=&a[0]
mov edx, mem_0 ;

loc_9:
mov [ecx], edx ;*p=arrVal
add ecx, 4 ;p++
inc ebx ;i++
cmp ebx, 10 ;i<10?
jl short loc_9 ;

mov edi, mem_4 ;
mov eax, [edi] ;return *pArray2
add esp, 40
retn
    
```



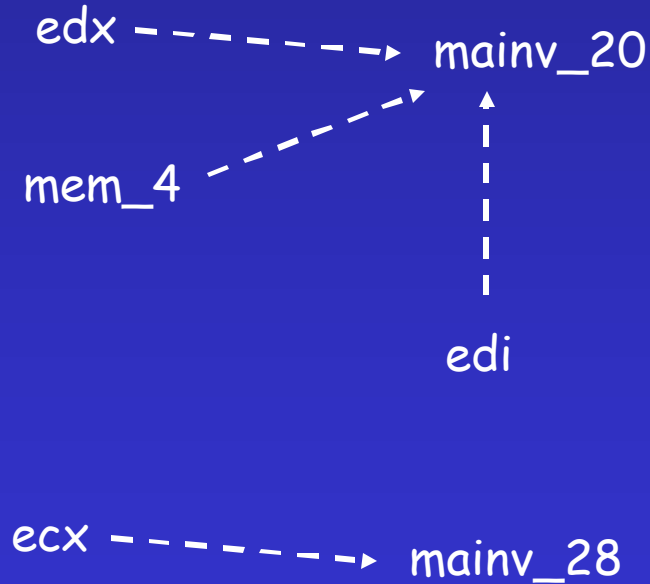
Example - A-locs

locals: mainv_28, mainv_20

{a[0], a[2]}

globals: mem_0, mem_4

{arrVal, pArray2}



```
; ebx  $\hat{U}$  i  
; ecx  $\hat{U}$  variable p
```

```
sub    esp, 40      ;adjust stack  
lea    edx, &mainv_20;  
mov    mem_4, edx   ;pArray2=&a[2]  
lea    ecx, &mainv_28;p=&a[0]  
mov    edx, mem_0   ;
```

```
loc_9:  
mov    [ecx], edx   ;*p=arrVal  
add    ecx, 4       ;p++  
inc    ebx          ;i++  
cmp    ebx, 10      ;i<10?  
j1     short loc_9 ;
```

```
mov    edi, mem_4   ;  
mov    eax, [edi] ;return *pArray2  
add    esp, 40  
retn
```



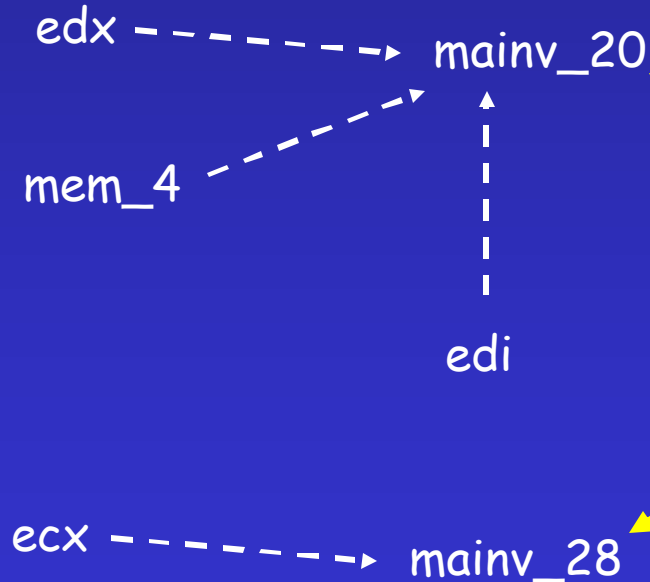
Example - A-locs

locals: mainv_28, mainv_20

{a[0], a[2]}

globals: mem_0, mem_4

{arrVal, pArray2}



```
; ebx  $\hat{U}$  i
```

```
; ecx  $\hat{U}$  variable p
```

```
sub    esp, 40      ;adjust stack
lea    edx, &mainv_20;
mov    mem_4, edx   ;pArray2=&a[2]
lea    ecx, &mainv_28;p=&a[0]
mov    edx, mem_0   ;
```

```
loc_9:
```

```
mov    [ecx], edx  ;*p=arrVal
```

```
add    ecx, 4      ;p++
```

```
inc    ebx         ;i++
```

```
cmp    ebx, 10    ;i<10?
```

```
jnl   short loc_9 ;
```

```
mov    edi, mem_4 ;
```

```
mov    eax, [edi] ;return *pArray2
```

```
add    esp, 40
```

```
retn
```

Value-Set Analysis

- **Resembles a pointer-analysis algorithm**
 - interprets pointer-manipulation operations
 - pointer arithmetic, too
- **Resembles a numeric-analysis algorithm**
 - over-approximate the set of values/addresses held by an a-loc
 - range information
 - stride information
 - interprets arithmetic operations on sets of values/addresses

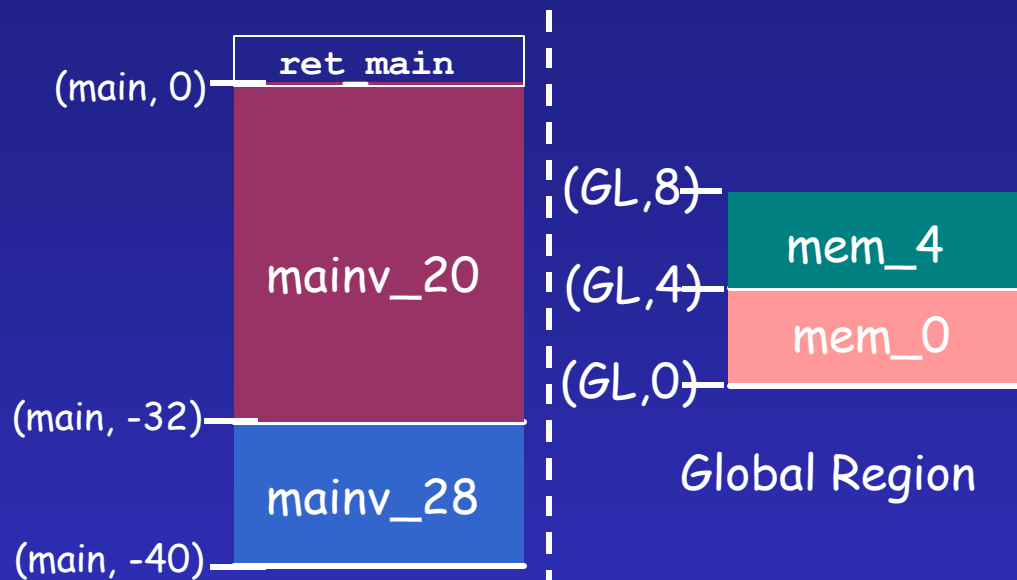
Value-set

- An a-loc \approx a variable
 - the address of an a-loc
(memory-region, offset within the region)
- An a-loc \approx an aggregate variable
 - addresses of elements of an a-loc
($\text{rgn}, \{o_1, o_2, \dots, o_n\}$)
- Value-set = a set of such addresses
 $\{(\text{rgn}_1, \{o_1, o_2, \dots, o_n\}), \dots, (\text{rgn}_r, \{o_1, o_2, \dots, o_m\})\}$
"r" - number of regions in the program

Value-set

- Set of addresses: $\{(rgn_1, \{o_1, \dots, o_n\}), \dots, (rgn_r, \{o_1, \dots, o_m\})\}$
- Idea: approximate $\{o_1, \dots, o_k\}$ with a numeric domain
 - $\{1, 3, 5, 9\}$ represented as $2[0,4]+1$
 - Reduced Interval Congruence (RIC)
 - common stride
 - lower and upper bounds
 - displacement
- Set of addresses is an r-tuple: (ric_1, \dots, ric_r)
 - ric_1 : offsets in global region
 - a set of numbers: $(ric_1, \wedge, \dots, \wedge)$

Example - Value-set analysis



```

; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

sub    esp, 40      ;adjust stack
lea    edx, [esp+8] ;
mov    [4], edx     ;pArray2=&a[2]
lea    ecx, [esp]   ;p=&a[0]
mov    edx, [0]     ;

```

```

loc_9:
mov    [ecx], edx   ;*p=arrVal 1
add    ecx, 4       ;p++
inc    ebx          ;i++
cmp    ebx, 10      ;i<10?
j1     short loc_9 ;

```

```

mov    edi, [4]     ;
mov    eax, [edi]   ;return *pArr 2
add    esp, 40
retn

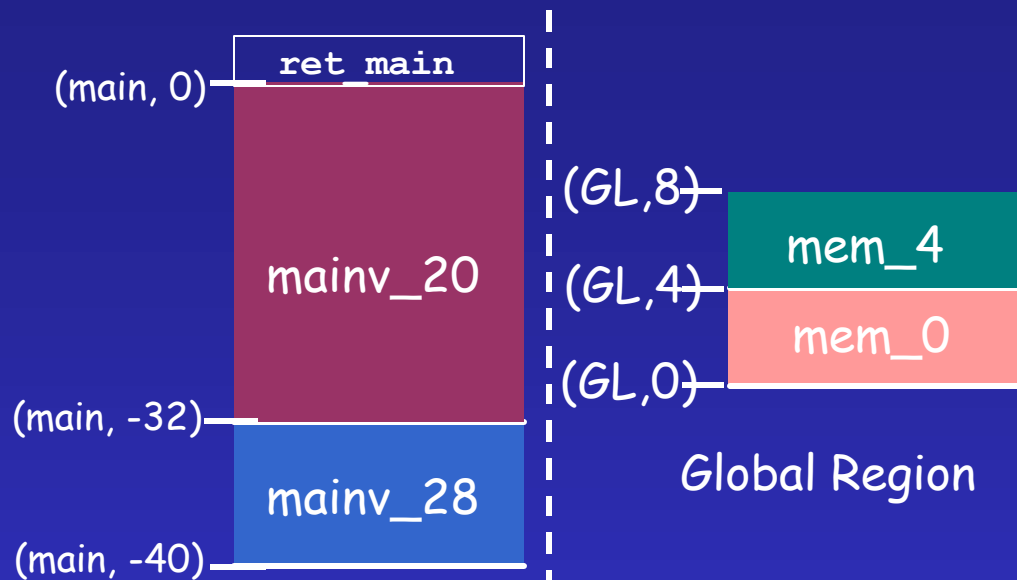
```

Region for main

- 1 \rightarrow $ecx \rightarrow (\hat{\quad}, 4[0,8]-40)$
- $ebx \rightarrow (1[0,9], \hat{\quad})$
- $esp \rightarrow (\hat{\quad}, -40)$
- 2 \rightarrow $edi \rightarrow (\hat{\quad}, -32)$
- $esp \rightarrow (\hat{\quad}, -40)$



Example - Value-set analysis



```

; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

sub     esp, 40      ;adjust stack
lea     edx, [esp+8] ;
mov     [4], edx    ;pArray2=&a[2]
lea     ecx, [esp]  ;p=&a[0]
mov     edx, [0]    ;
    
```

```

loc_9:
mov     [ecx], edx  ;*p=arrVal 1
add     ecx, 4      ;p++
inc     ebx         ;i++
cmp     ebx, 10    ;i<10?
j1     short loc_9 ;
    
```

```

mov     edi, [4]    ;
mov     eax, [edi] ;return *pArr 2
add     esp, 40
retn
    
```

Region for main
 1 $ecx \rightarrow (\wedge, 4[0,8]-40)$

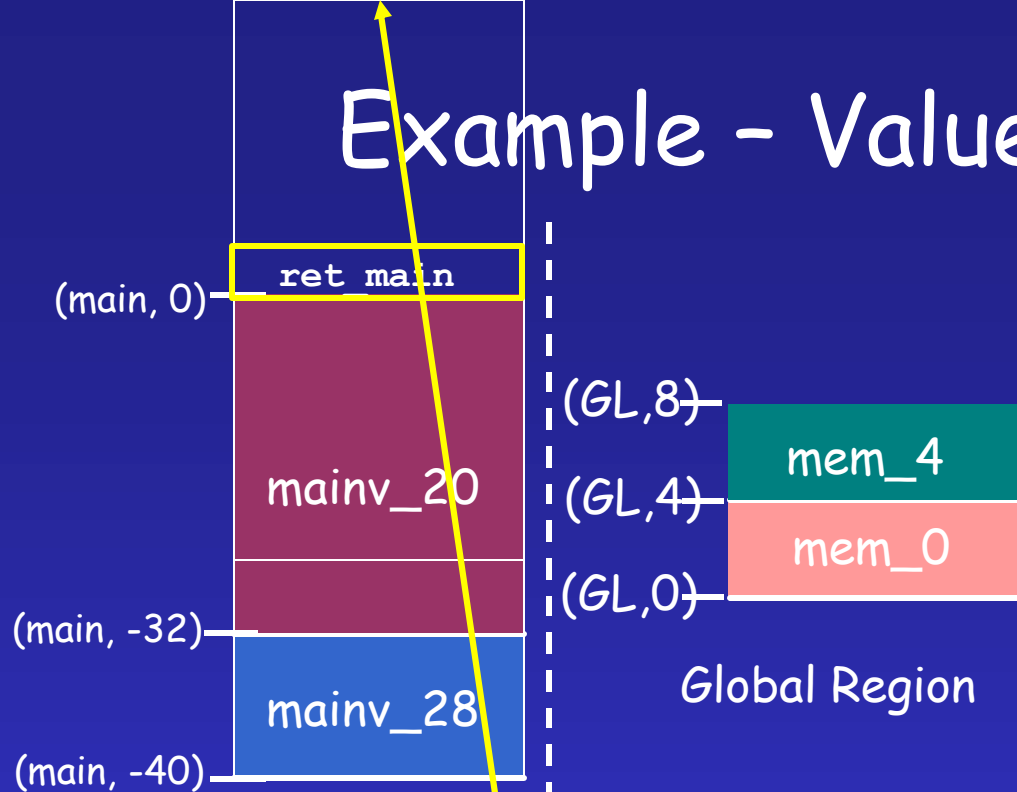


$$= (\wedge, -32)^1 \mathcal{A}$$

2 $edi \rightarrow (\wedge, -32)$



Example - Value-set analysis



```

; ebx  $\hat{U}$  i
; ecx  $\hat{U}$  variable p

sub     esp, 40           ;adjust stack
lea    edx, [esp+8]      ;
mov     [4], edx         ;pArray2=&a[2]
lea    ecx, [esp]        ;p=&a[0]
mov     edx, [0]         ;

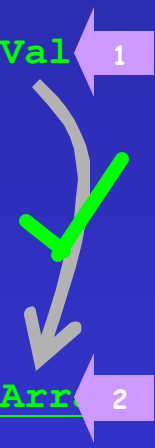
loc_9:
mov     [ecx], edx       ;*p=arrVal 1
add     ecx, 4           ;p++
inc     ebx              ;i++
cmp     ebx, 10          ;i<10?
j1     short loc_9      ;

mov     edi, [4]         ;
mov     eax, [edi]       ;return *pArr 2
add     esp, 40
retn
    
```

1 $ecx \rightarrow (^, 4[0]8)-40)$

2 $edi \rightarrow (^, -32)$

A stack-smashing attack?



Affine-Relation Analysis

- Value-set domain is **non-relational**
 - cannot capture relationships among a-locs
- Imprecise results
 - e.g. no upper bound for `ecx` at `loc_9`
 - `ecx` \rightarrow (\wedge , 4[0,8]-40)

```
loc_9:
    mov     [ecx], edx    ;*p=arrVal
    add     ecx, 4        ;p++
    inc     ebx          ;i++
    cmp     ebx, 10      ;i<10?
    jl     short loc_9 ;
    . . .
```

Affine-Relation Analysis

- Obtain affine relations via static analysis
- Use affine relations to improve precision
 - e.g., at `loc_9`

$ecx = esp + (4 \cdot ebx)$, $ebx = ([0, 9], ^)$, $esp = (^, -40)$

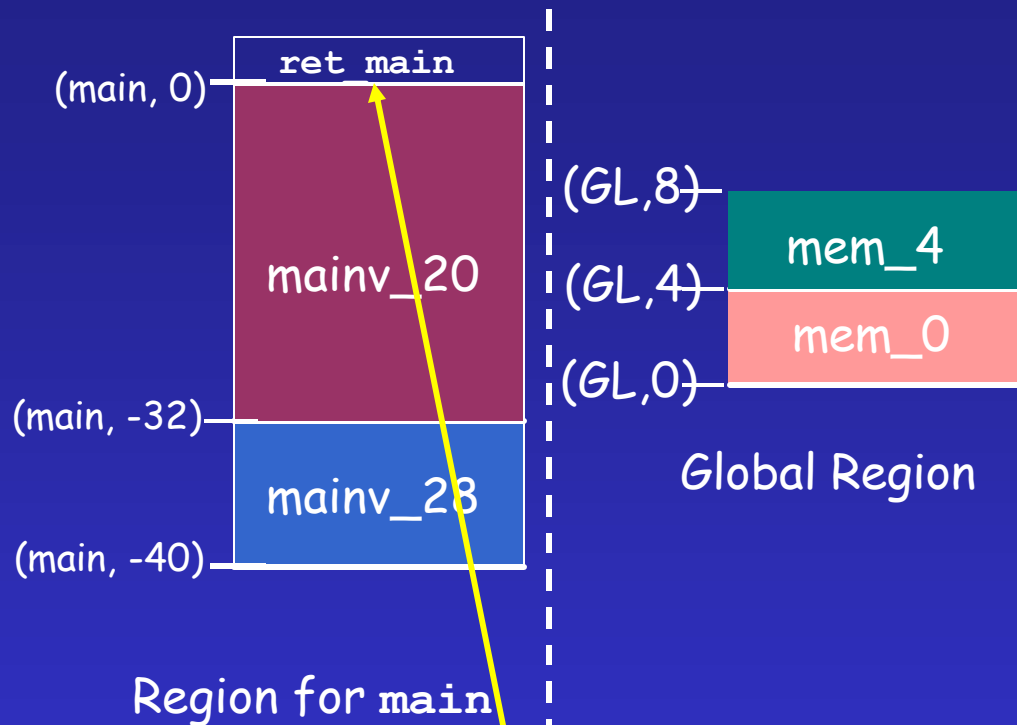
$\mathbb{P} \quad ecx = (^, -40) + 4([0, 9])$

$\mathbb{P} \quad ecx = (^, 4[0, 9] - 40)$

\mathbb{P} upper bound for `ecx` at `loc_9`

```
...  
loc_9:  
mov     [ecx], edx    ;*p=arrVal  
add     ecx, 4        ;p++  
inc     ebx           ;i++  
cmp     ebx, 10       ;i<10?  
jl     short loc_9 ;  
...
```

Example - Value-set analysis



1 **ecx** → (⊥, 4[0,9]-40)

No stack-smashing attack reported

```

; ebx ̂ i
; ecx ̂ variable p

sub     esp, 40           ;adjust stack
lea     edx, [esp+8]     ;
mov     [4], edx         ;pArray2=&a[2]
lea     ecx, [esp]       ;p=&a[0]
mov     edx, [0]         ;

loc_9:
mov     [ecx], edx       ;*p=arrVal 1
add     ecx, 4           ;p++
inc     ebx              ;i++
cmp     ebx, 10          ;i<10?
jnl    short loc_9      ;

mov     edi, [4]         ;
mov     eax, [edi]       ;return *pArr 2
add     esp, 40
retn
    
```



Affine-Relation Analysis

- Affine relation
 - x_1, x_2, \dots, x_n - a-locs
 - a_0, a_1, \dots, a_n - integer constants
 - $a_0 + \sum_{i=1..n}(a_i x_i) = 0$
- **Idea: determine affine relations on registers**
 - use such relations to improve precision
- Implemented using WPDS++

Performance

Program	nProc	nInsts	Value-set analysis (seconds)	Affine-relations (seconds)
javac	36	3,555	42	36
cat(2.0.14)	123	3,892	51	32
cut(2.0.14)	129	4,329	28	50
grep(2.4.2)	245	16,808	85	78
flex(2.5.4)	239	23,373	200	376
tar(1.13.19)	587	50,347		210
awk(3.1.0)	595	69,927		1,507
winhlp32 (5.00.2195.2014)	1,018	108,380		2,002



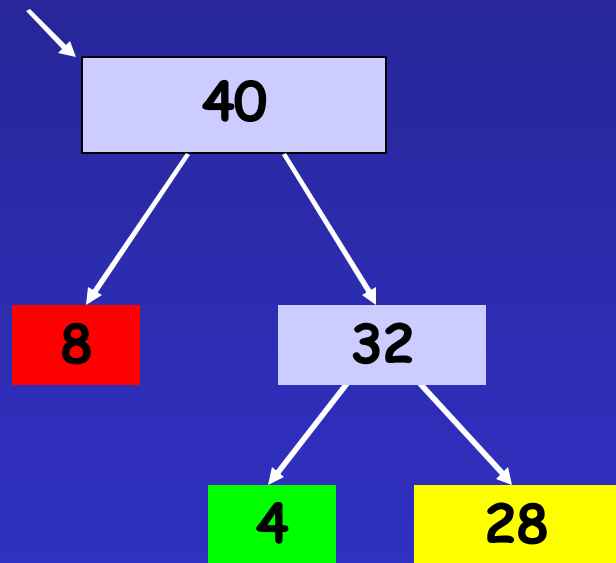
Future Work

- Aggregate Structure Identification
 - Ramalingam et al. [POPL 99]
 - Ignore declarative information
 - Identify fields from the access patterns
 - Useful for
 - improving the a-loc abstraction
 - discovering type information



Future Work

AR[-40:-1]



```
; ebx  $\hat{U}$  i  
; ecx  $\hat{U}$  variable p
```

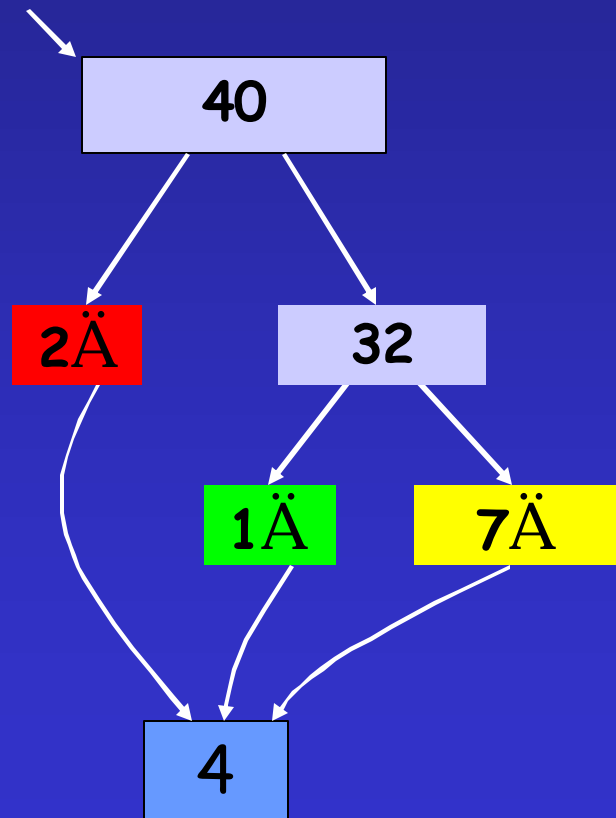
```
sub    esp, 40        ;adjust stack  
lea    edx, [esp+8]  ;  
mov    [4], edx      ;pArray2=&a[2]  
lea    ecx, [esp]    ;p=&a[0]  
mov    edx, [0]      ;
```

```
loc_9:  
mov    [ecx], edx    ;*p=arrVal  
add    ecx, 4        ;p++  
inc    ebx           ;i++  
cmp    ebx, 10       ;i<10?  
jl     short loc_9  ;
```

```
mov    edi, [4]      ;  
➔ mov    eax, [edi]  ;return *pArray2  
add    esp, 40  
retn
```

Future Work

AR[-40:-1]

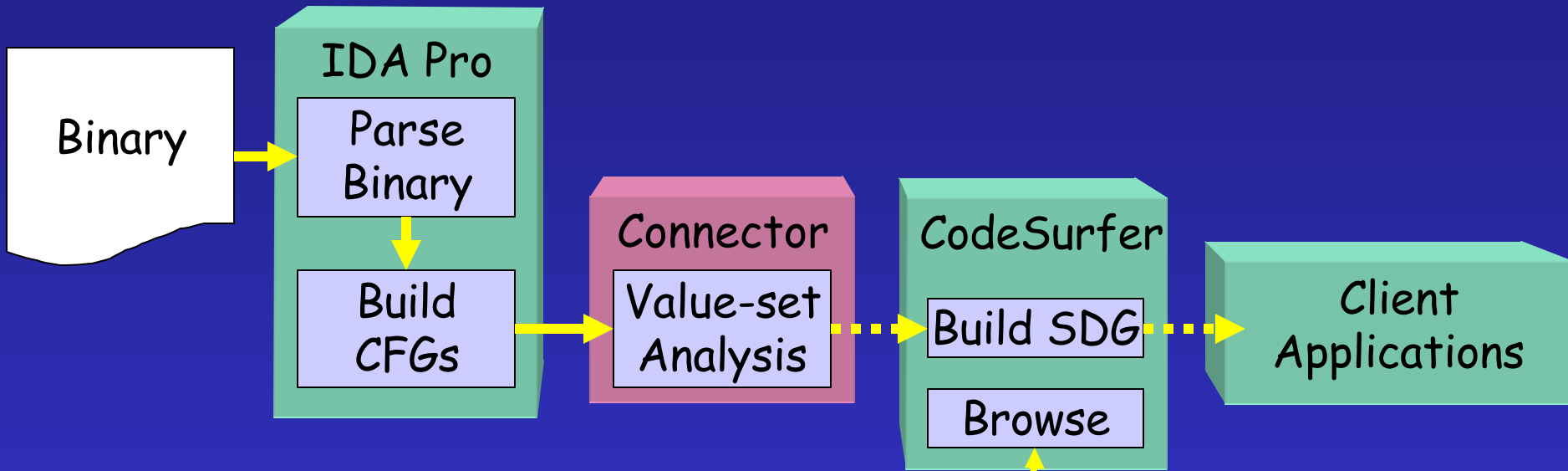


```
; ebx  $\hat{U}$  i  
; ecx  $\hat{U}$  variable p
```

```
sub    esp, 40        ;adjust stack  
lea    edx, [esp+8]   ;  
mov    [4], edx       ;pArray2=&a[2]  
lea    ecx, [esp]     ;p=&a[0]  
mov    edx, [0]       ;
```

```
loc_9:  
mov    [ecx], edx     ;*p=arrVal  
add    ecx, 4         ;p++  
inc    ebx            ;i++  
cmp    ebx, 10        ;i<10?  
jl     short loc_9 ;  
  
mov    edi, [4]       ;  
mov    eax, [edi]     ;return *pArray2  
add    esp, 40  
retn
```

Codesurfer/x86 Architecture



For more details

- [Gogul Balakrishnan's demo](#)
- Gogul Balakrishnan's poster
- Consult UW-TR 1486
[<http://www.cs.wisc.edu/~reps/#tr1486>]

Analyzing Memory Accesses in x86 Executables

Gogul Balakrishnan

Thomas Reps

University of Wisconsin