# Identifying Variables in x86 Executables

Gogul Balakrishnan

Thomas Reps

University of Wisconsin

# Motivation

- Code-inspection tools for security analysts
  - dependence-based navigation ("code surfing")

- Analyses for identifying
  - security vulnerabilities and bugs
  - malicious code
  - commonalities and differences

- Platform for
  - code obfuscation and de-obfuscation
  - de-compilation
  - installation of protection mechanisms
  - remediation of security vulnerabilities

# Why Executables?

- Reflects actual behaviors that may arise
- Allows platform-specific artifacts to be taken into account
  - memory layout
  - register usage
  - execution order
  - compiler bugs
  - Thompson-style attack
- Source code hides the low-level (actual) behaviors that implement high-level abstractions
- Source-code analyses typically make unsafe assumptions (e.g., that the program is ANSI-C compliant)
  - loss of fidelity can allow vulnerabilities to escape notice

# Puzzle

```
int callee(int a, int b)
  int local;
  if (local == 5) return
  else return 2;
}
```

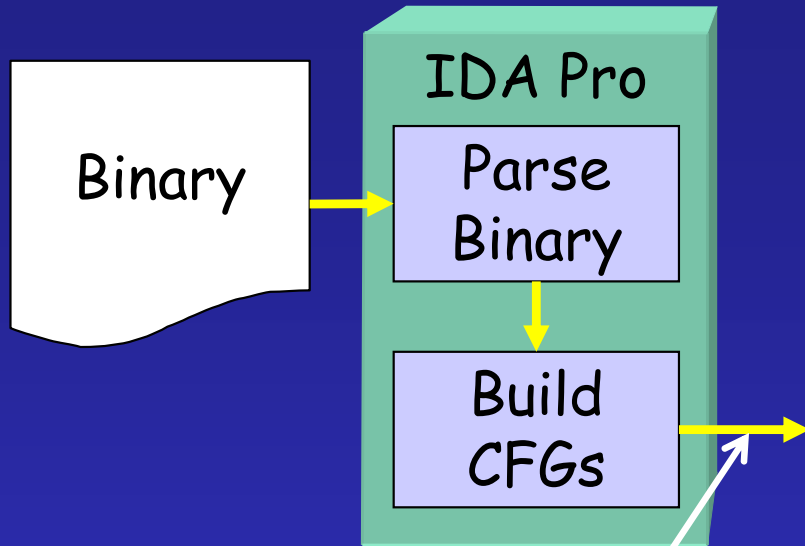| Standard prolog | | Prolog for 1 local | |
|---|---|---|---|
| push | ebp | push | ebp |
| mov | ebp, esp | mov | ebp, esp |
| sub | esp, 4 | push | ecx |

Answer: 1
(for the Microsoft compiler)

```
int main() {
  int c = 5;
  int d = 7;

  int v = callee(c,d);
  // What is the value of v here?
  return 0;
}
```

```
mov    [ebp+var_8], 5
mov    [ebp+var_C], 7
mov    eax, [ebp+var_C]
push   eax
mov    ecx, [ebp+var_8]
push   ecx
call   _callee
. . .
```
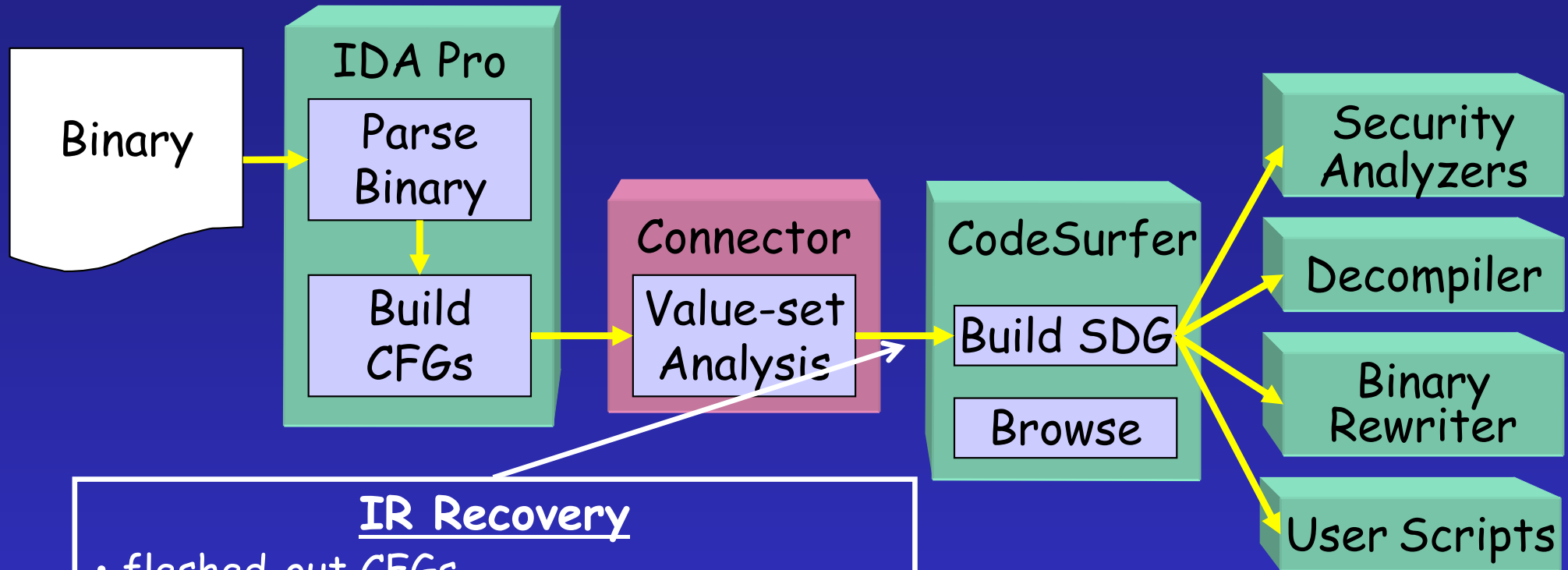
# CodeSurfer/x86 Architecture

Binary

**IDA Pro**

Parse Binary

Build CFGs

Initial estimate of
- code vs. data
- procedures
- call sites
- malloc sites

# CodeSurfer/x86 Architecture

Binary

**IDA Pro**
- Parse Binary
- Build CFGs

**Connector**
- Value-set Analysis

**CodeSurfer**
- Build SDG
- Browse

- Security Analyzers
- Decompiler
- Binary Rewriter
- User Scripts

**IR Recovery**
- fleshed-out CFGs
- fleshed-out call graph
- used, killed, may-killed variables for CFG nodes
- points-to sets
- reports of violations
- [variables]
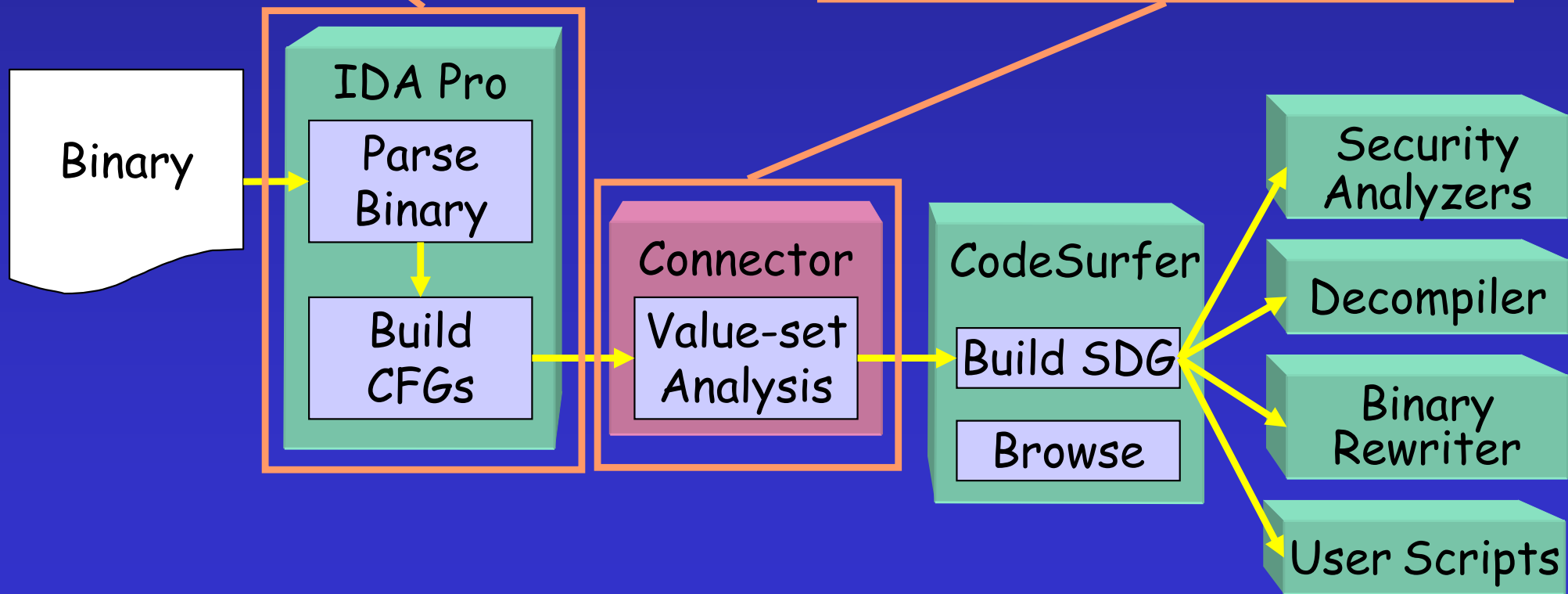- [types: base types, pointer types, structs, and classes]

9

# Scope

- Programs that conform to a "standard compilation model"
  - procedures
  - activation records
  - global data region
  - heap, etc.

- Report violations
  - violations of stack protocol
  - return address modified within procedure

# Technical Challenges

- Distinguishing between code and data
- Identifying variables

- Identifying parameters
- Resolving indirect jumps
- Resolving indirect calls
- Identifying may-aliases

Binary

**IDA Pro**
- Parse Binary
- Build CFGs

**Connector**
- Value-set Analysis

**CodeSurfer**
- Build SDG
- Browse

- Security Analyzers
- Decompiler
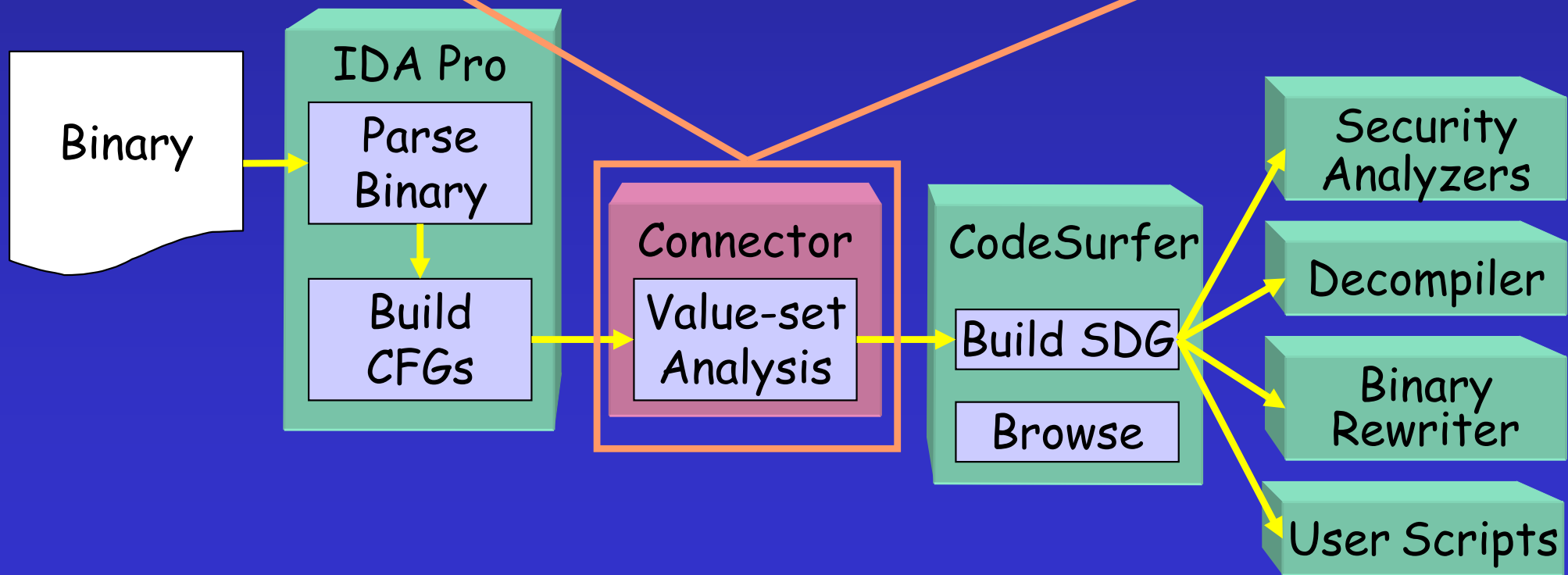- Binary Rewriter
- User Scripts

# Technical Challenges

- Distinguishing between code and data
- Identifying variables

- Identifying parameters
- Resolving indirect jumps
- Resolving indirect calls
- Identifying may-aliases

Binary → IDA Pro [ Parse Binary → Build CFGs ] → Connector [ Value-set Analysis ] → CodeSurfer [ Build SDG, Browse ] → Security Analyzers, Decompiler, Binary Rewriter, User Scripts

13

# Running Example

```c
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;
}
```

```asm
; ebx ⇔ variable i
; ecx ⇔ variable p

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

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, [8]     ;
mov     eax, [edi]   ;return *pArray2
add     esp, 40
retn
```
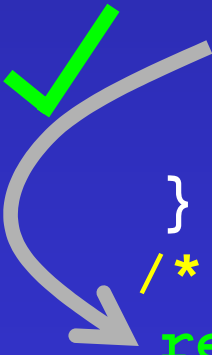
14

# 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;
}
```
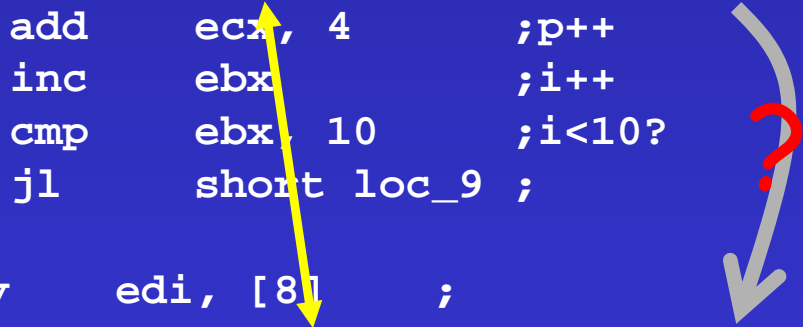
```
; ebx ⇔ variable i
; ecx ⇔ variable p

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

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, [8]       ;
mov     eax, [edi]   ;return *pArray2
add     esp, 40
retn
```
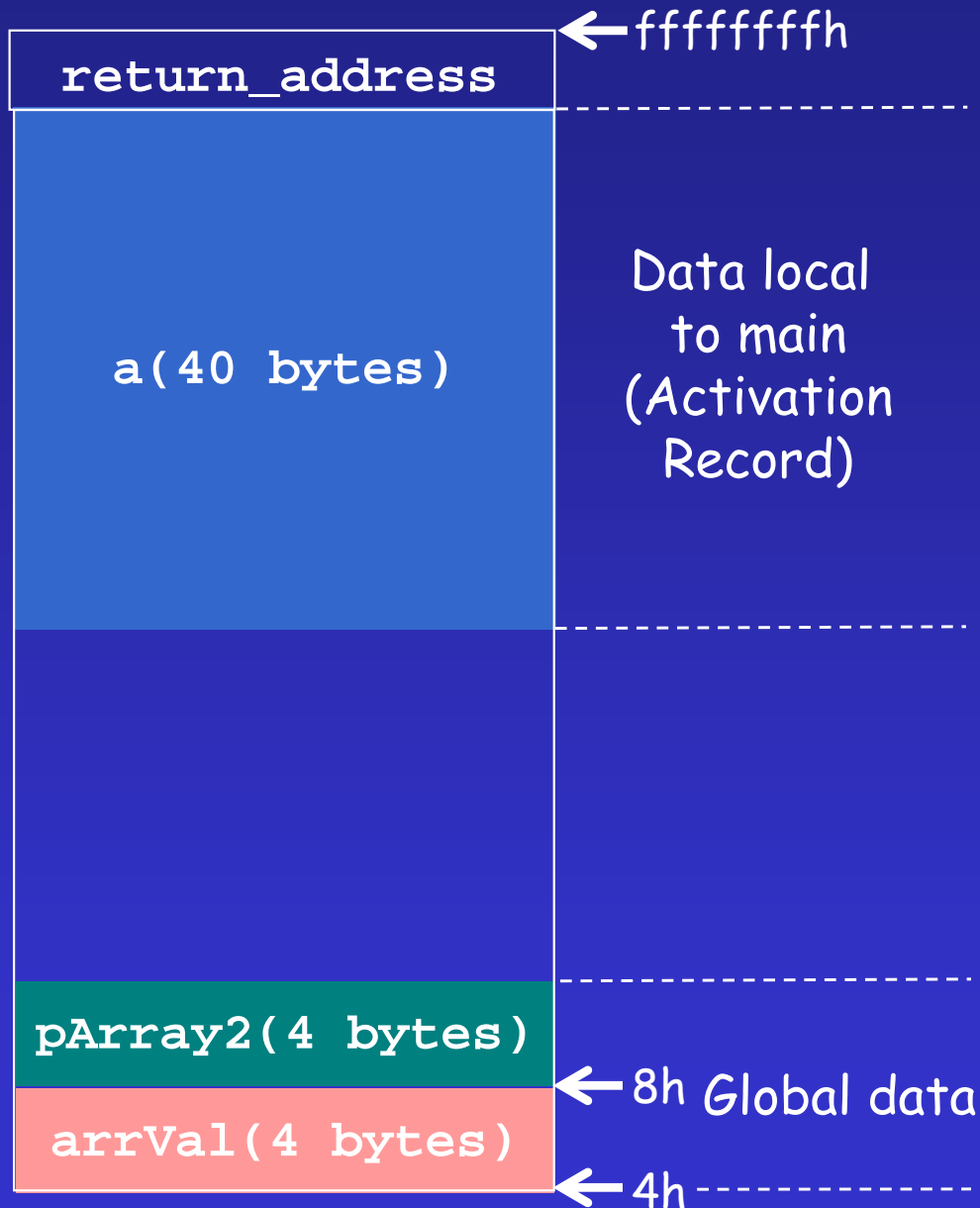
15

# Running Example – Address Space

| |
|---|
| return_address |
| a(40 bytes) |
| pArray2(4 bytes) |
| arrVal(4 bytes) |

←ffffffffh

Data local
to main
(Activation
Record)

←8h Global data

←4h

```
; ebx ⟺ variable i
; ecx ⟺ variable p

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

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, [8]         ;
mov     eax, [edi]   ;return *pArray2
add     esp, 40
retn
```
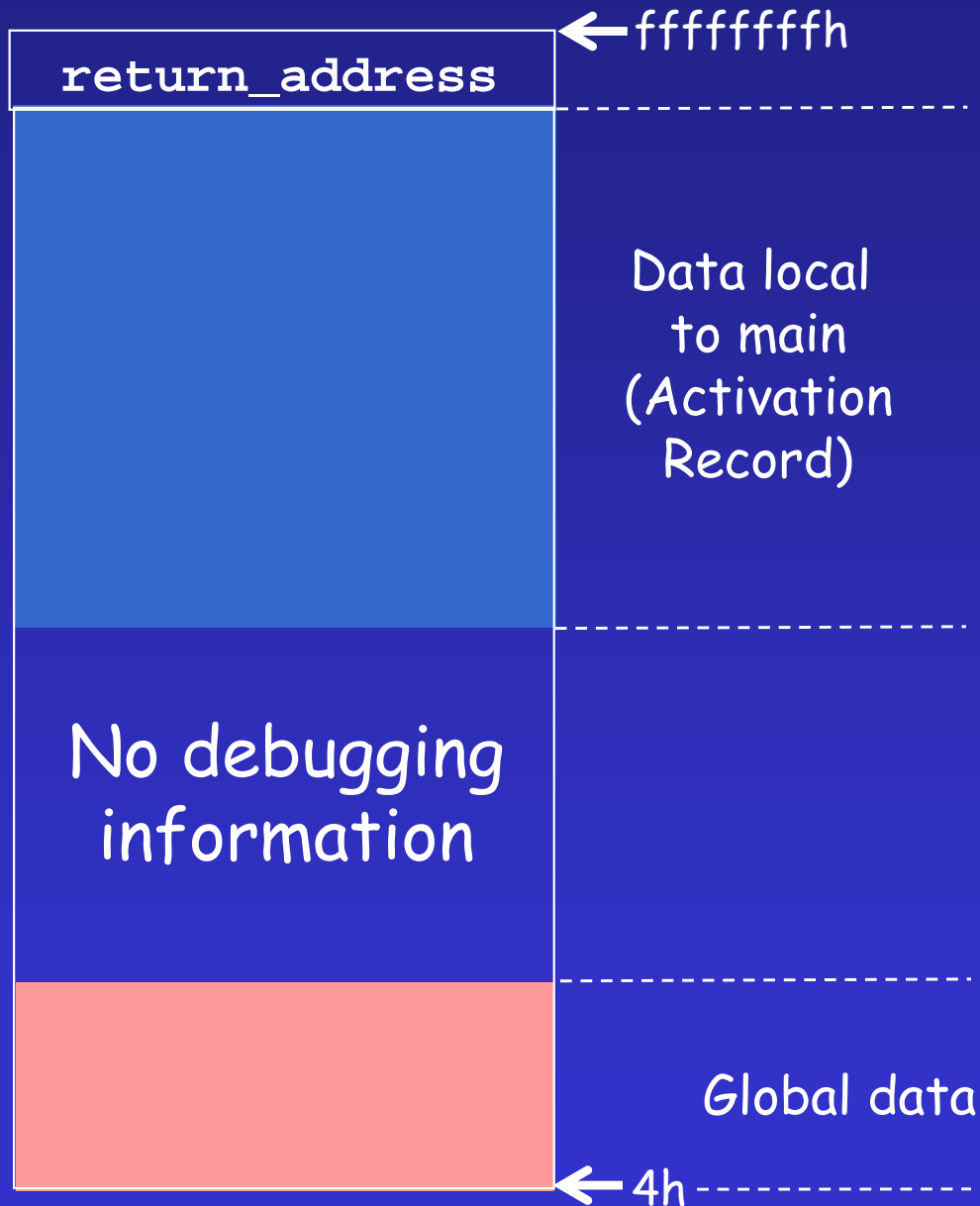
# Running Example – Address Space

| return_address |
| --- |

← fffffffffh

Data local
to main
(Activation
Record)

No debugging
information

Global data

← 4h

```
; ebx  ⇔  variable i
; ecx  ⇔  variable p

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

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, [8]        ;
mov     eax, [edi]      ;return *pArray2
add     esp, 40
retn
```
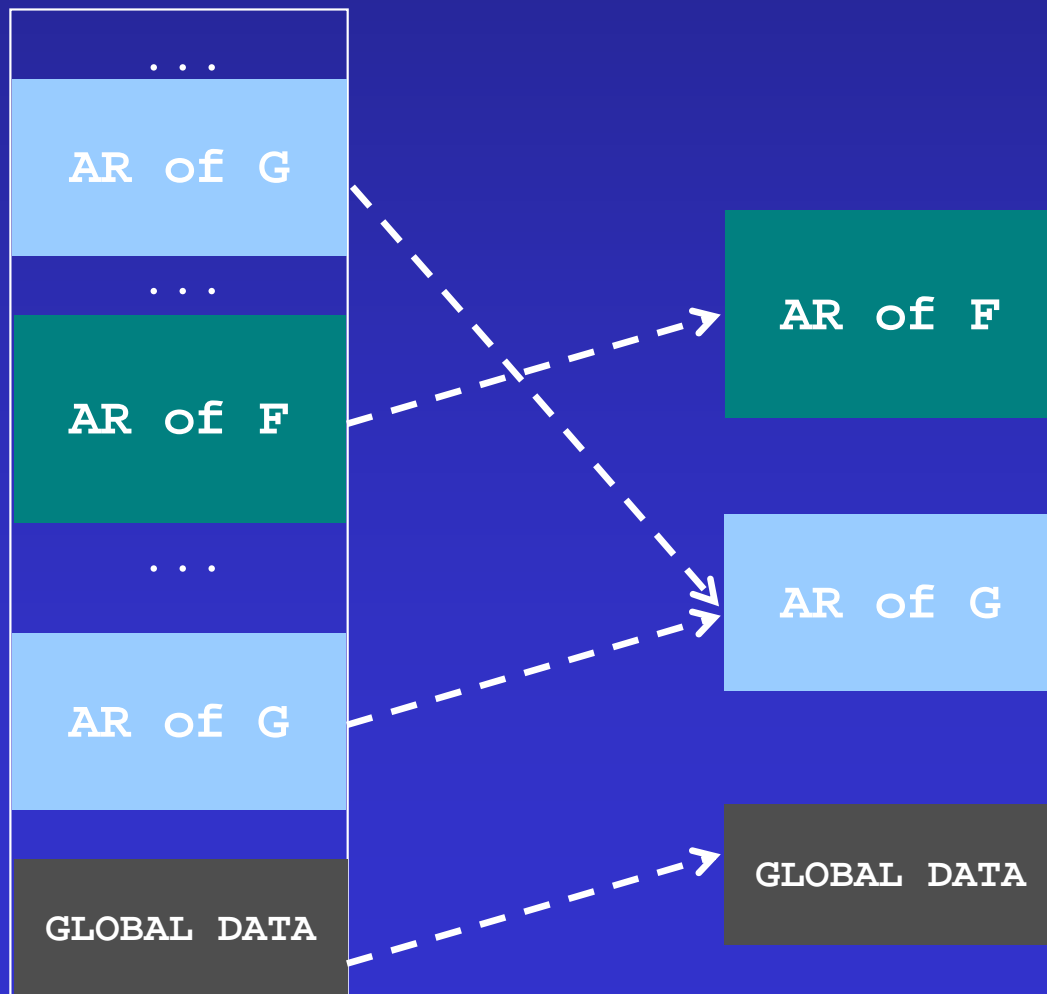
?
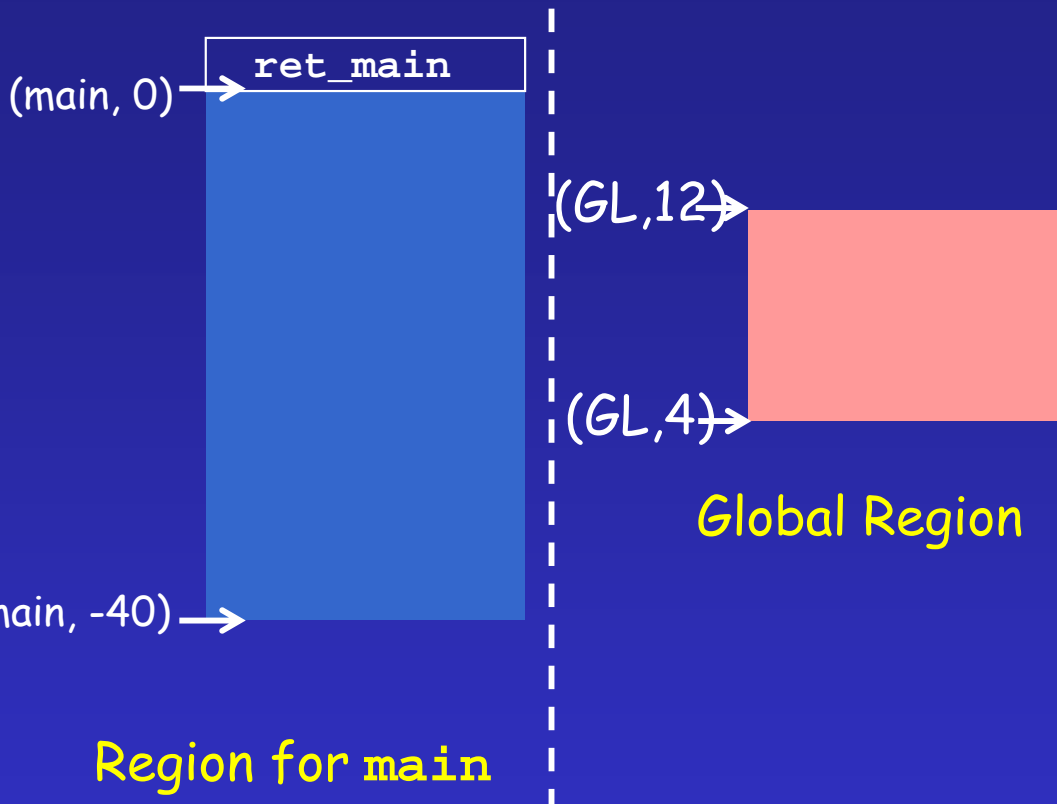
# Identifying Variables

- An abstraction of concrete memory configurations
  - Memory regions

- Infer layout of memory regions
  - A-locs (like variables)

# Memory Regions

- An abstraction of concrete memory configurations
  - Idea: group similar runtime addresses
  - e.g., collapse the runtime ARs for each procedure, malloc-sites, global data

# Example – Memory Regions

ret_main

(main, 0)

(GL,12)

(GL,4)

Global Region

(main, -40)

Region for **main**

```
; ebx ⇔ variable i
; ecx ⇔ variable p

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

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, [8]      ;
mov     eax, [edi]   ;return *pArray2
add     esp, 40
retn
```
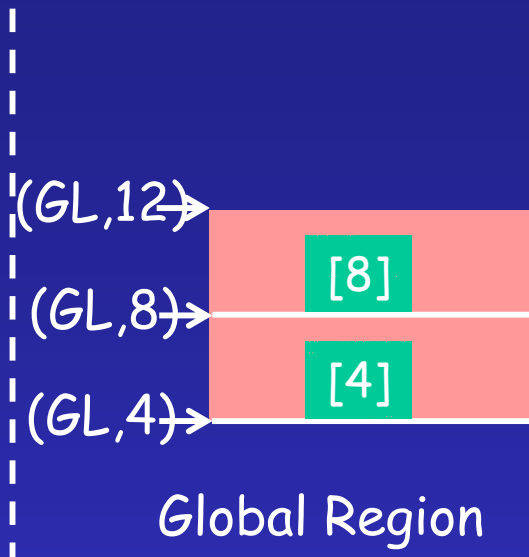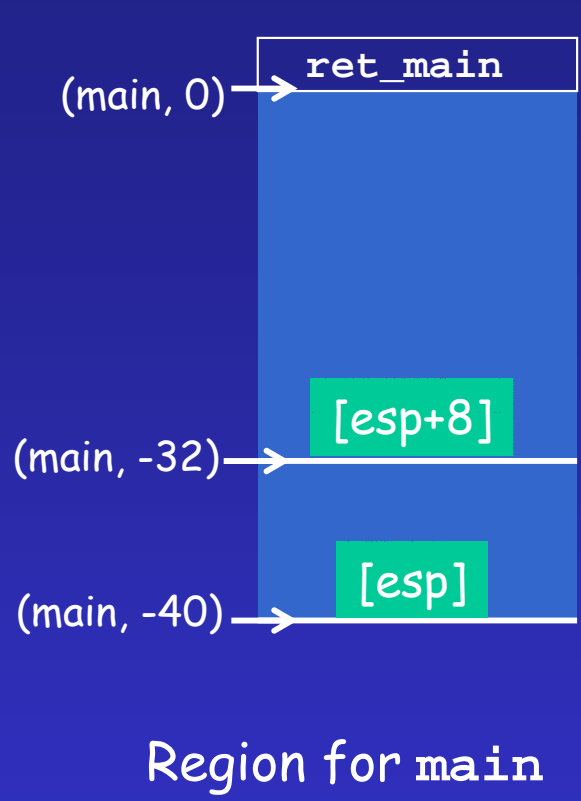
20

# Infer Layout of Memory Regions

* Data-layout known at assembly/compile time
  - 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

# Example – A-locs

(main, 0) → **ret_main**

(GL,12) →

`; ebx ⇔ variable i`
`; ecx ⇔ variable p`

**[8]**

(GL,8) →

**[esp+8]**

(main, -32) →

(GL,4) →

**[4]**

Global Region

(main, -40) →

**[esp]**

Region for **main**

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

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, [8]       ;
mov      eax, [edi]   ;return *pArray2
add      esp, 40
retn
```

?

22

# Example – A-locs



(main, 0) → **ret_main**

**mainv_32**

(main, -32) →

**mainv_40**

(main, -40) →

**Region for `main`**

(GL,12) →
(GL,8) → **mem_8**
(GL,4) → **mem_4**

**Global Region**

```
; ebx ⟺ variable i
; ecx ⟺ variable p

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

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, [8]     ;
mov     eax, [edi]   ;return *pArray2
add     esp, 40
retn
```
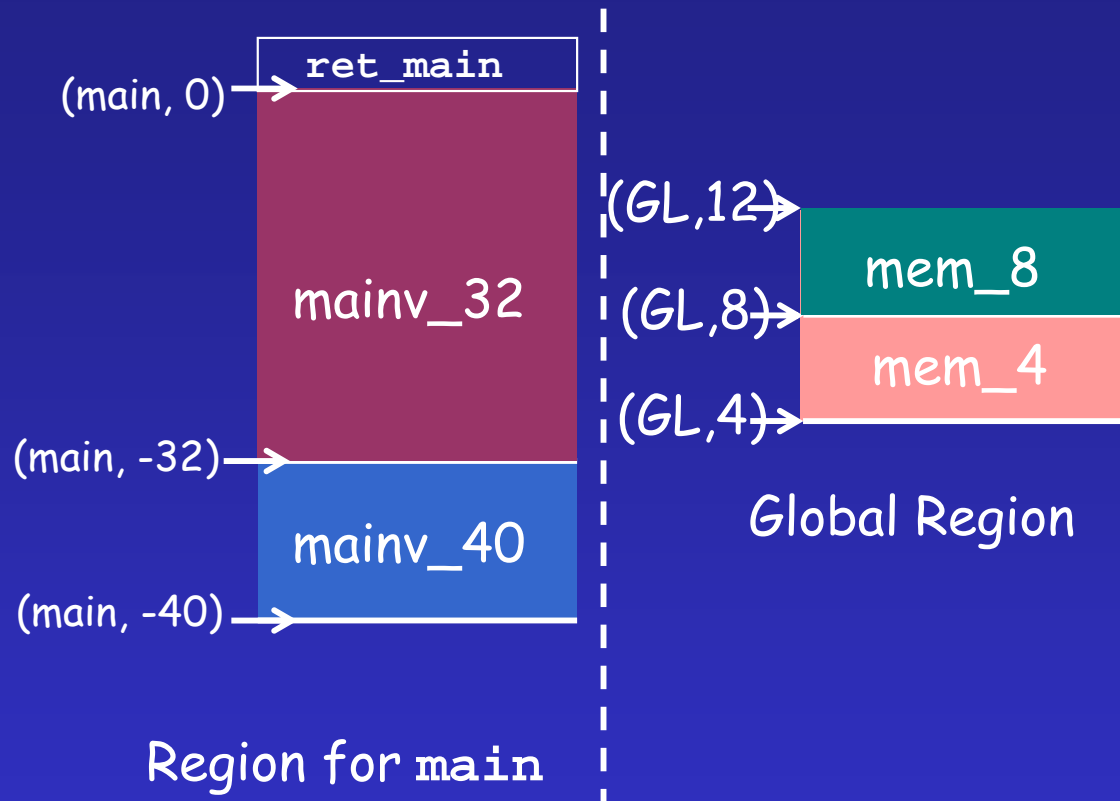
23

# Example – A-locs



```
; ebx ⇔ variable i
; ecx ⇔ variable p

sub      esp, 40          ;adjust stack
lea      edx, &mainv_32;
mov      mem_8, edx       ;pArray2=&a[2]
lea      ecx, &mainv_40;p=&a[0]
mov      edx, mem_4       ;

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_8       ;
mov      eax, [edi]   ;return *pArray2
add      esp, 40
retn
```
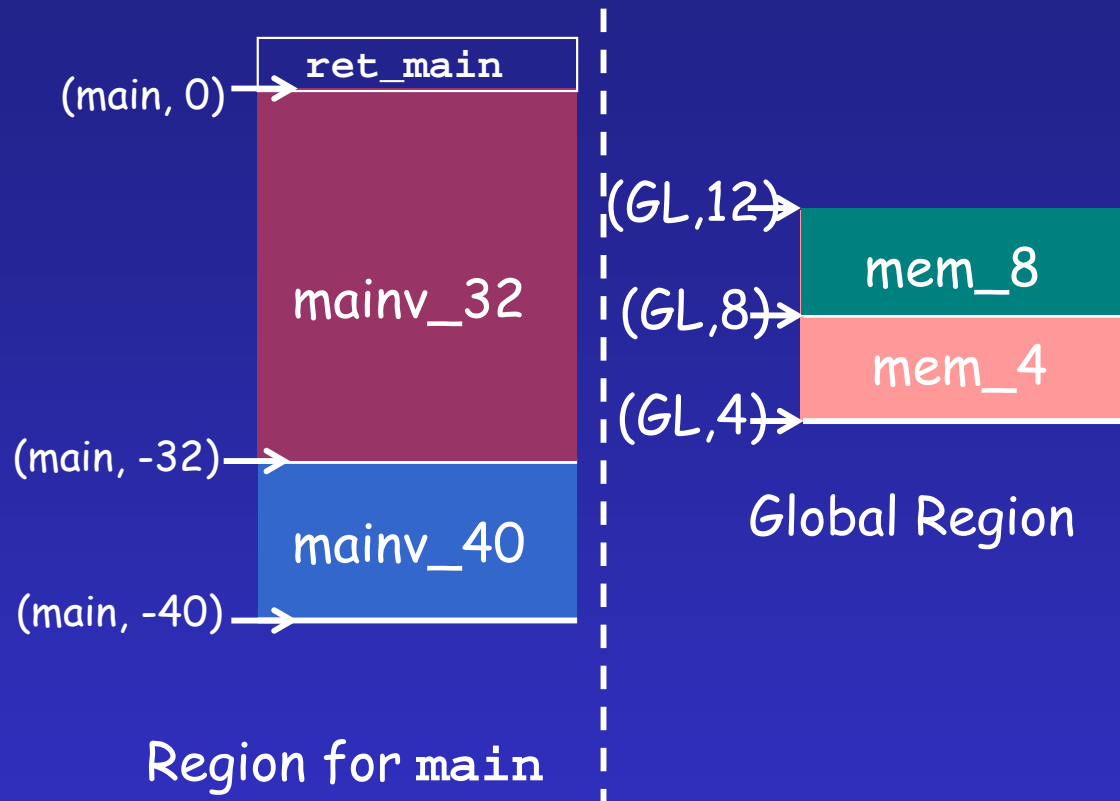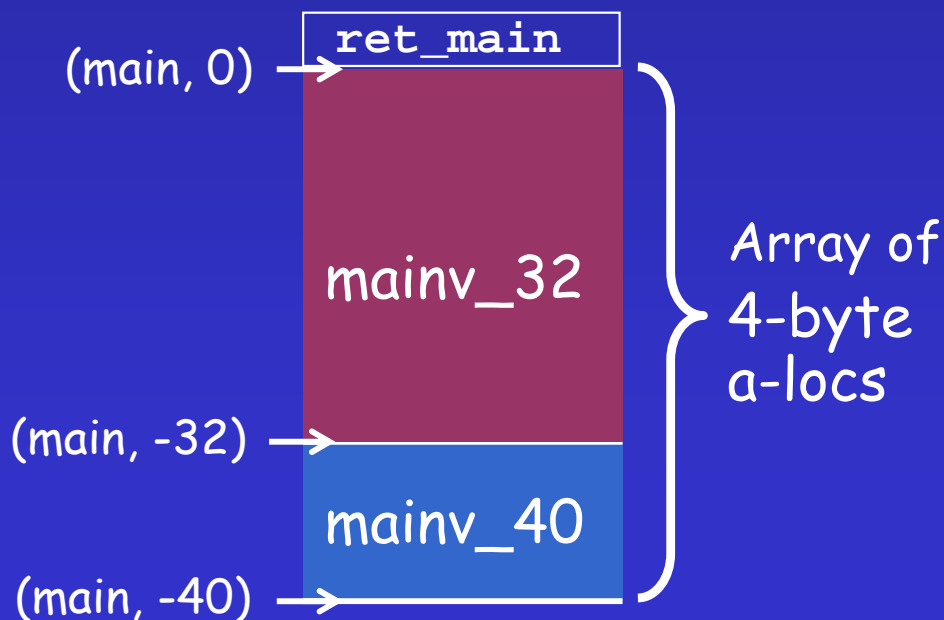
ret_main

(main, 0)

mainv_32

(main, -32)

mainv_40

(main, -40)

Region for main

(GL,12)

(GL,8)

(GL,4)

mem_8

mem_4

Global Region

24

# Better Identification of Variables

- IDAPro A-locs
  - Based on explicitly specified addresses/offsets
- VSA provides access patterns for indirect operands
  - $ecx \rightarrow (\perp, 4[0,9]-40)$

```
                    . . .

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
```
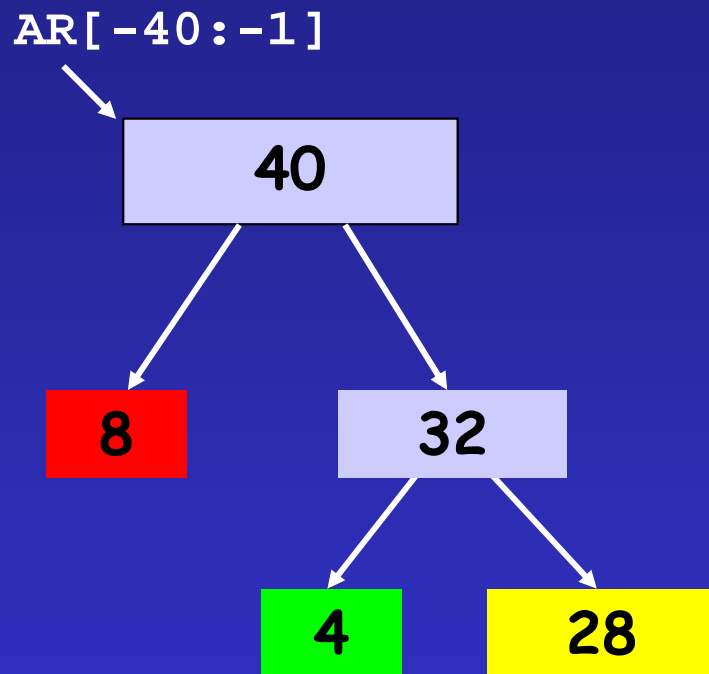
ret_main

(main, 0)

mainv_32

Array of
4-byte
a-locs

(main, -32)

mainv_40

(main, -40)

# Aggregate Structure Identification

- Partition aggregates according to the program's memory-access patterns
  - original motivation: Y2K [Ramalingam et al. POPL 99]
- Uses in our context
  - improved identification of variables
    - identifies a better set of a-locs
      $\Rightarrow$ better IR $\Rightarrow$ fewer false alarms
  - recovery of type information
    - identifies structs and arrays
    - propagates type information from known parameter types (system calls & library functions)
      $\Rightarrow$ better de-compilation

# Aggregate Structure Identification

AR[-40:-1]

```
40
```

```
8        32
```
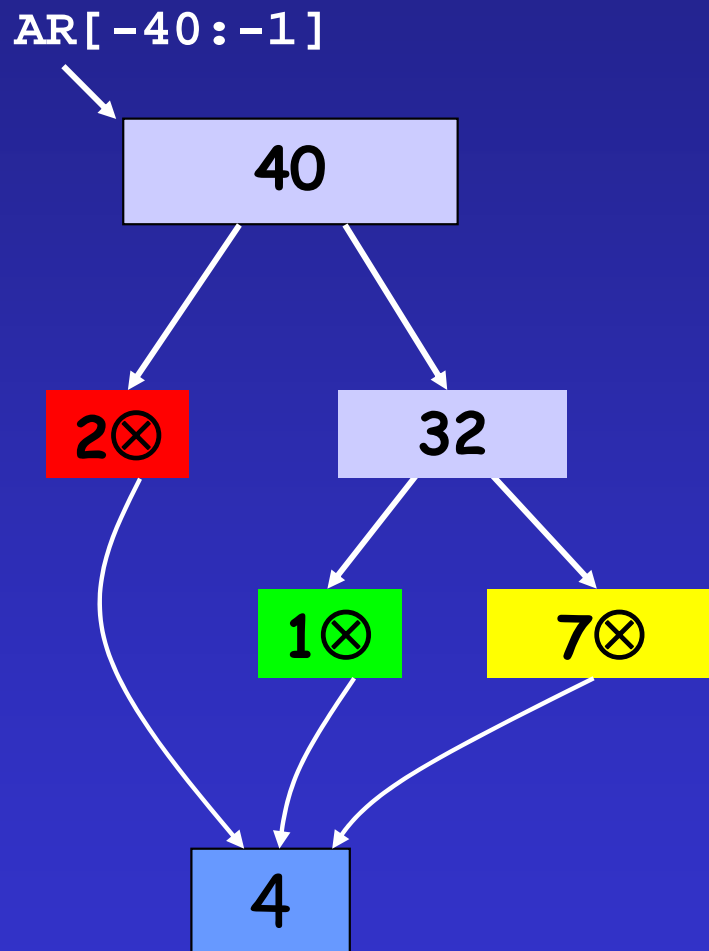
```
4        28
```

```
; ebx ⇔ variable 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?
    jl      short loc_9 ;

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

# Aggregate Structure Identification

**AR[-40:-1]**

```
; ebx ⇔ variable 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?
    jl      short loc_9 ;

mov     edi, [4]       ;
mov     eax, [edi]     ;return *pArray2
add     esp, 40
retn
```
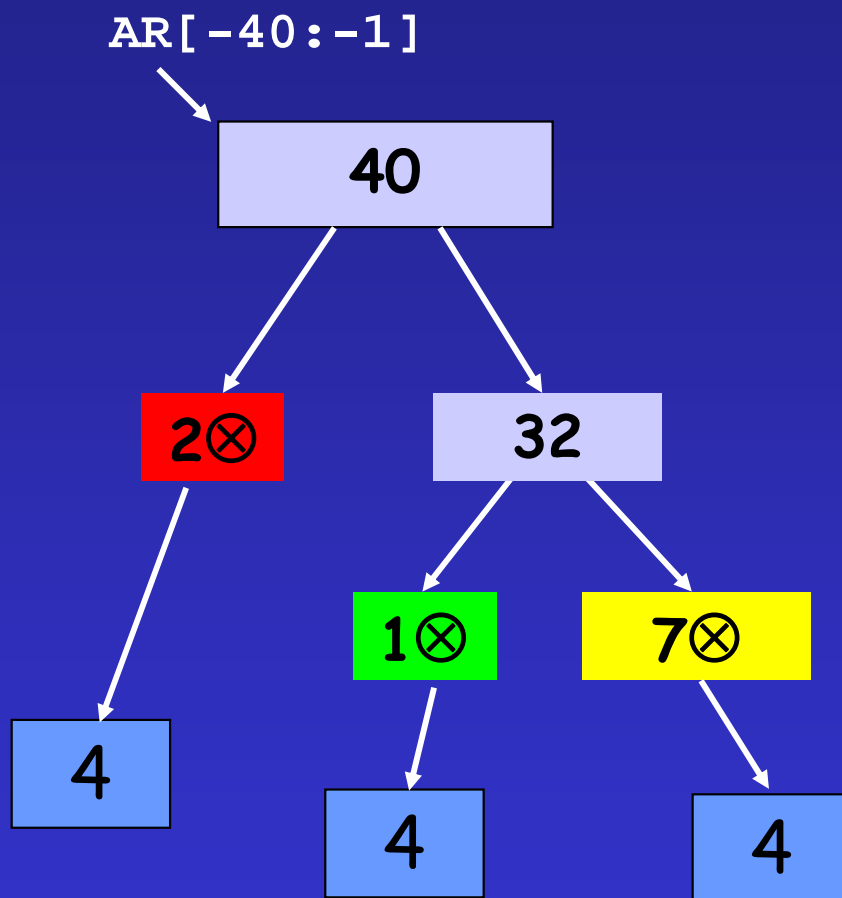
Diagram nodes: 40, 32, 2⊗, 1⊗, 7⊗, 4

# Aggregate Structure Identification
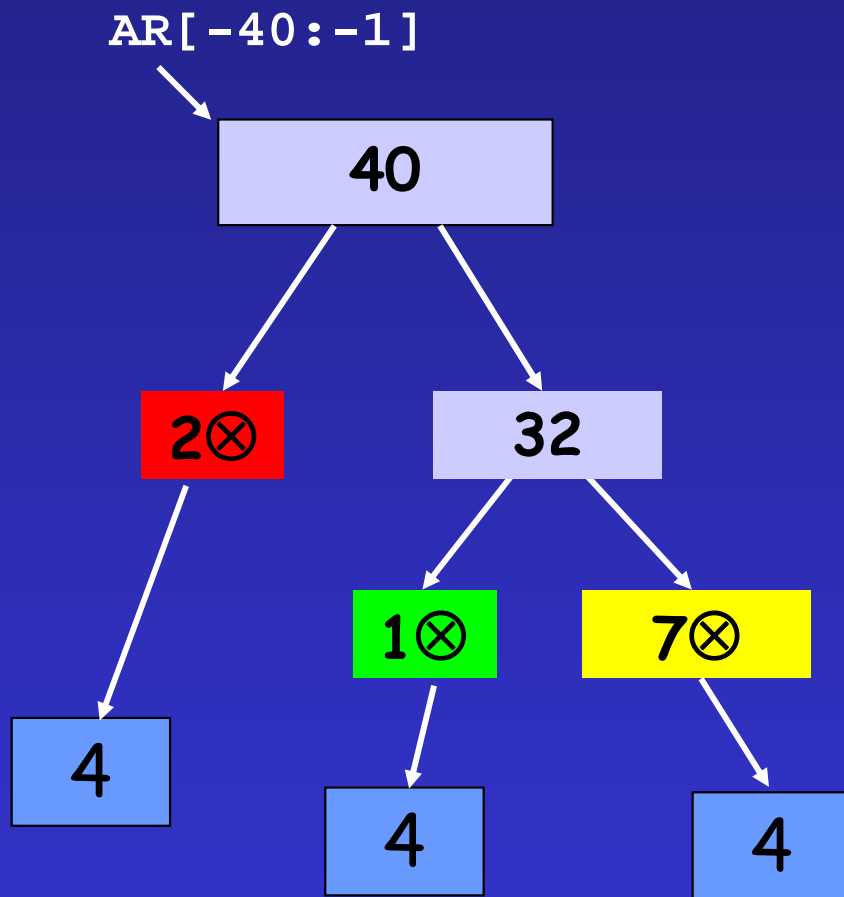
```
AR[-40:-1]
```



```
; ebx ⇔ variable 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?
    jl       short loc_9 ;

mov      edi, [4]        ;
mov      eax, [edi]    ;return *pArray2
add      esp, 40
retn
```
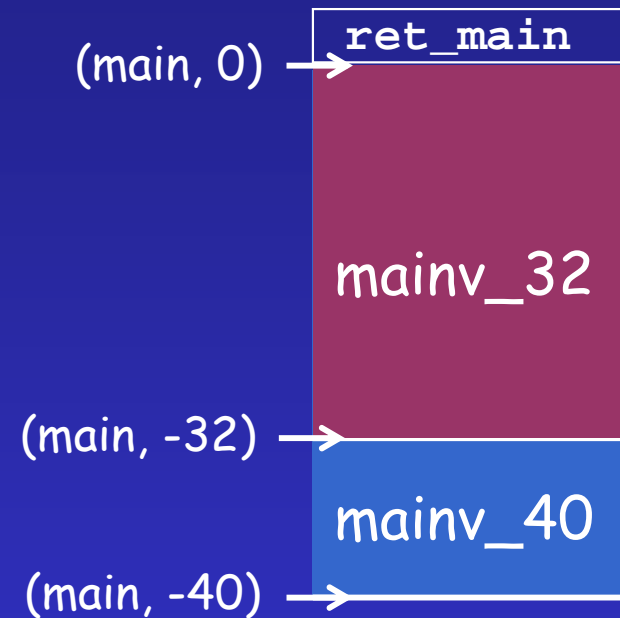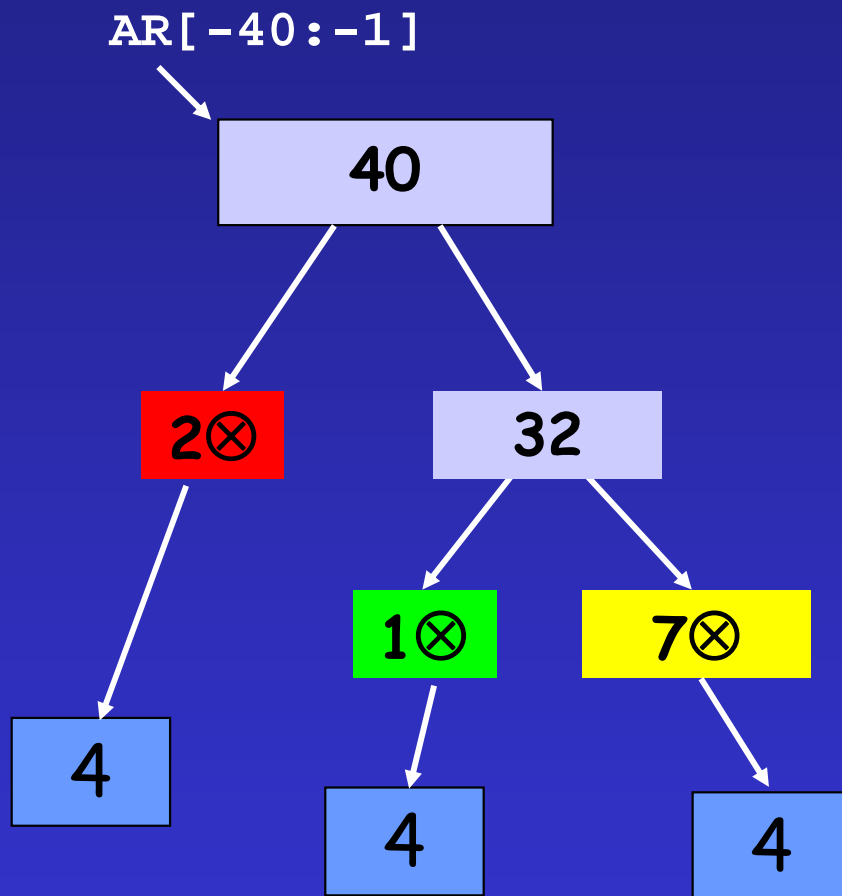
ASI: two arrays; one scalar

29

# Aggregate Structure Identification

AR[-40:-1]

```
        40
       /  \
      /    \
    2⊗      32
   /       /  \
  4      1⊗    7⊗
         |      |
         4      4
```

ASI: two arrays;
one scalar

```
(main, 0) →   | ret_main |
              |          |
              |          |
              | mainv_32 |
              |          |
(main, -32) → |          |
              |          |
              | mainv_40 |
(main, -40) → |          |
```

Region for **main**

IDA Pro
one 8-byte a-loc
one 32-byte a-loc
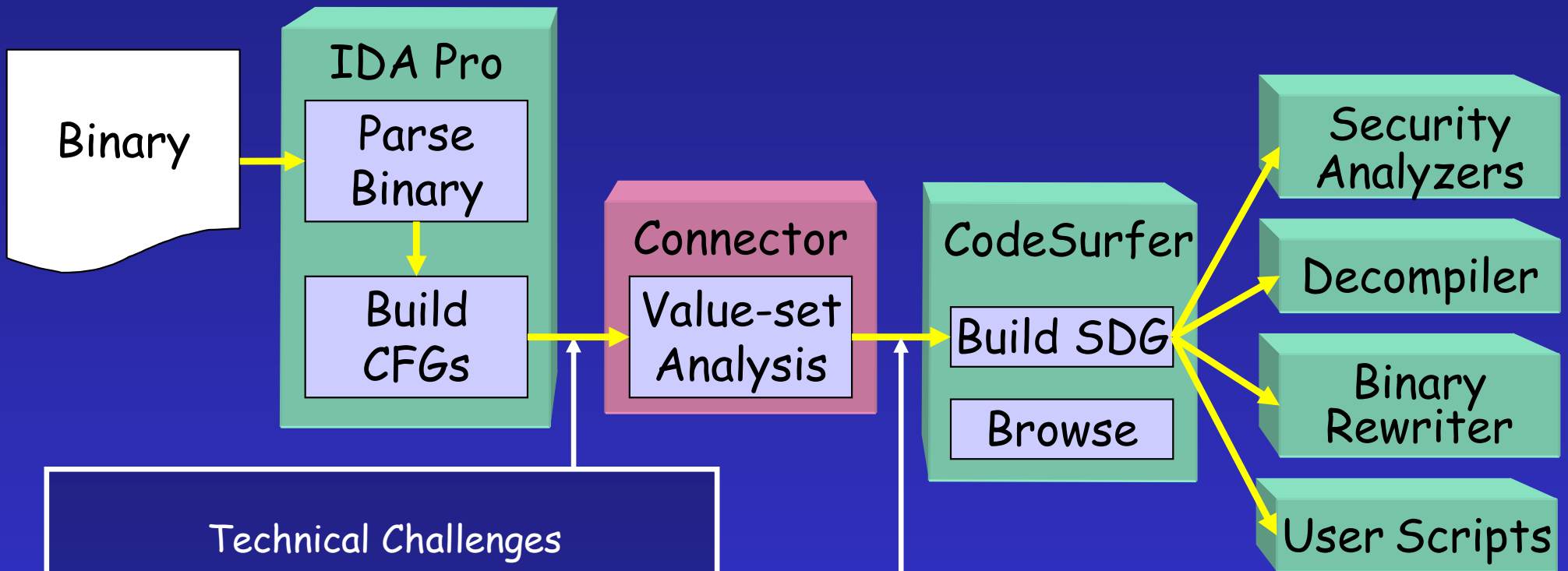
# Aggregate Structure Identification



`AR[-40:-1]`

High level type:

```
struct {
    int a[2];
    int b;
    int c[7];
};
```

ASI: two arrays; one scalar

31

# CodeSurfer/x86 Architecture

Binary → IDA Pro [ Parse Binary → Build CFGs ] → Connector [ Value-set Analysis ] → CodeSurfer [ Build SDG / Browse ] → Security Analyzers / Decompiler / Binary Rewriter / User Scripts
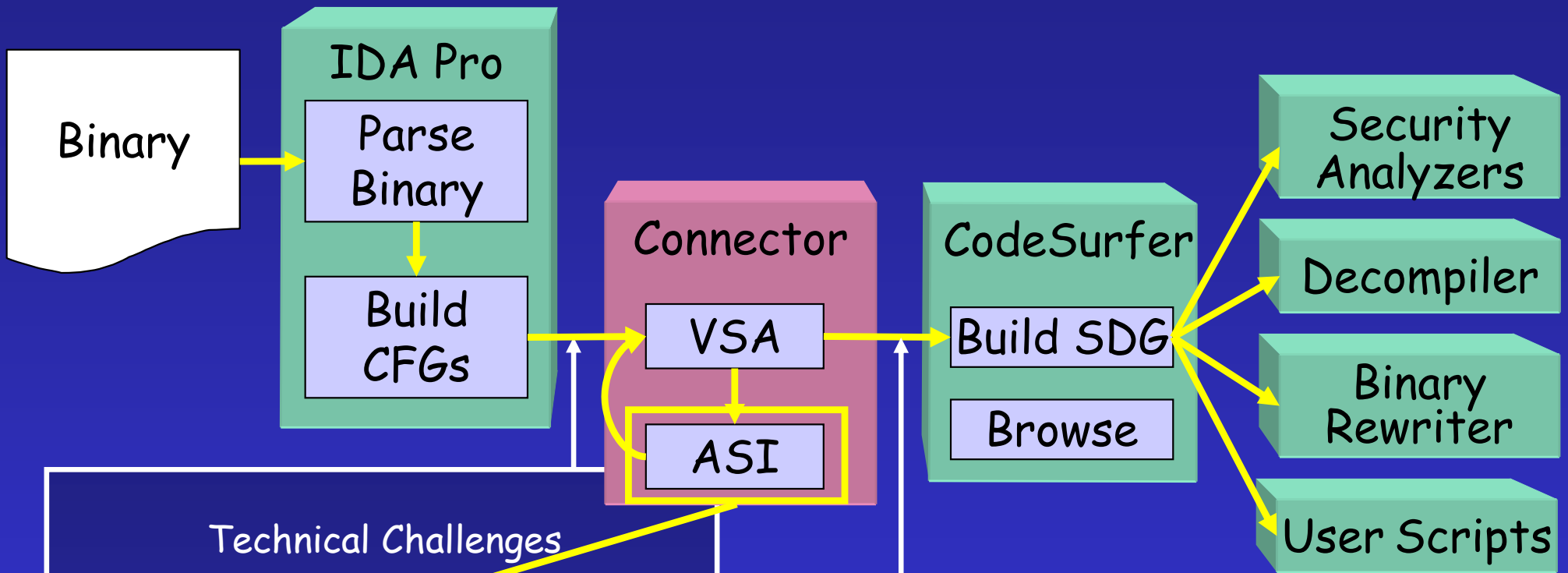
**Technical Challenges**

- Distinguishing between code and data
- Identifying variables
- Identifying parameters
- Resolving indirect jumps
- Resolving indirect calls
- Identifying may-aliases

- fleshed-out CFGs
- fleshed-out call graph
- used, killed, may-killed variables for CFG nodes
- points-to sets
- reports of violations

32

# CodeSurfer/x86 Architecture

Binary → **IDA Pro** [ Parse Binary → Build CFGs ] → **Connector** [ VSA ⇄ ASI ] → **CodeSurfer** [ Build SDG, Browse ] → Security Analyzers / Decompiler / Binary Rewriter / User Scripts
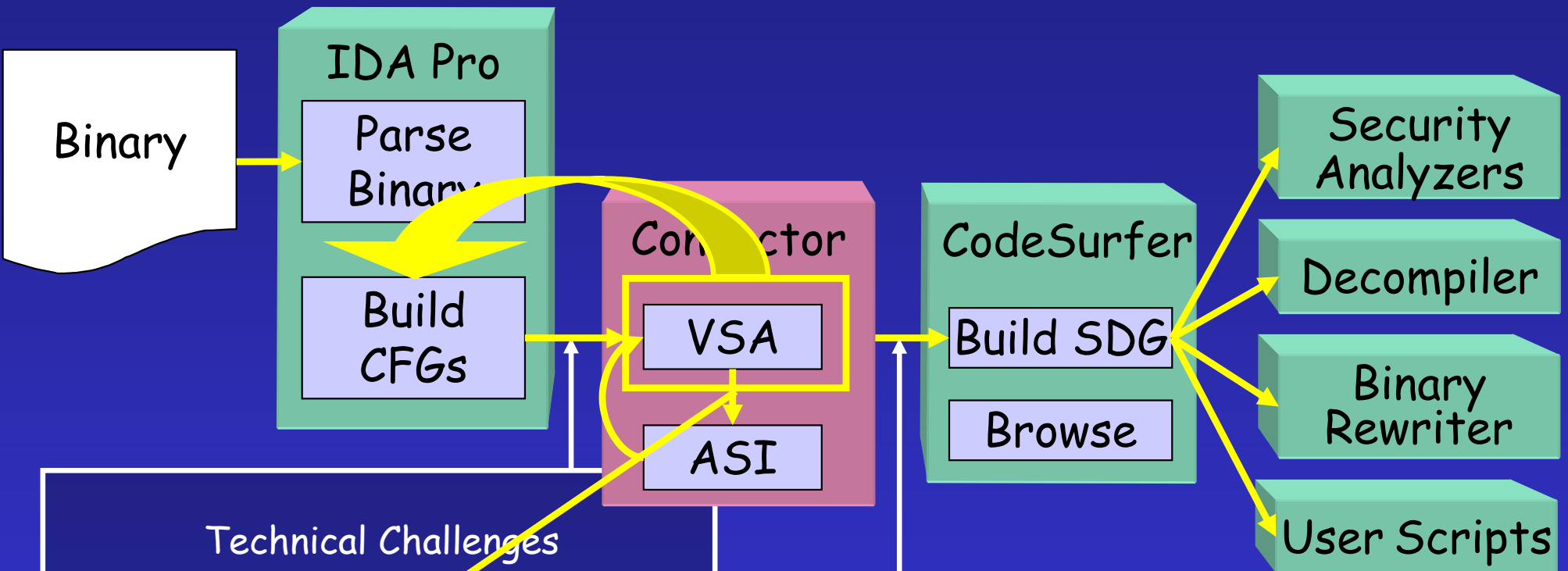
**Technical Challenges**

- Distinguishing between code and data
- Identifying variables
- Identifying parameters
- Resolving indirect jumps
- Resolving indirect calls
- Identifying may-aliases

- fleshed-out CFGs
- fleshed-out call graph
- used, killed, may-killed variables for CFG nodes
- points-to sets
- reports of violations

33

# CodeSurfer/x86 Architecture

Binary

IDA Pro
- Parse Binary
- Build CFGs

Connector
- VSA
- ASI

CodeSurfer
- Build SDG
- Browse

Security Analyzers

Decompiler

Binary Rewriter

User Scripts

## Technical Challenges

- Distinguishing between code and data
- Identifying variables
- Identifying parameters
- Resolving indirect jumps
- Resolving indirect calls
- Identifying may-aliases

- fleshed-out CFGs
- fleshed-out call graph
- used, killed, may-killed variables for CFG nodes
- points-to sets
- reports of violations

34

# Wrap Up

- Code-inspection tools for security analysts
  - dependence-based navigation ("code surfing")

- Analyses for identifying
  - security vulnerabilities and bugs
  - malicious code
  - commonalities and differences

- Platform for
  - code obfuscation and de-obfuscation
  - de-compilation
  - installation of protection mechanisms
  - remediation of security vulnerabilities

35

# Identifying Variables in x86 Executables

## Gogul Balakrishnan
## Thomas Reps

### University of Wisconsin