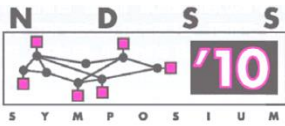# Reverse Engineering of Data Structures and Types

# Automatic Reverse Engineering of Program Data Structures from Binary Execution
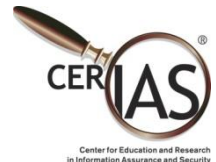
Zhiqiang Lin
Xiangyu Zhang,  Dongyan Xu
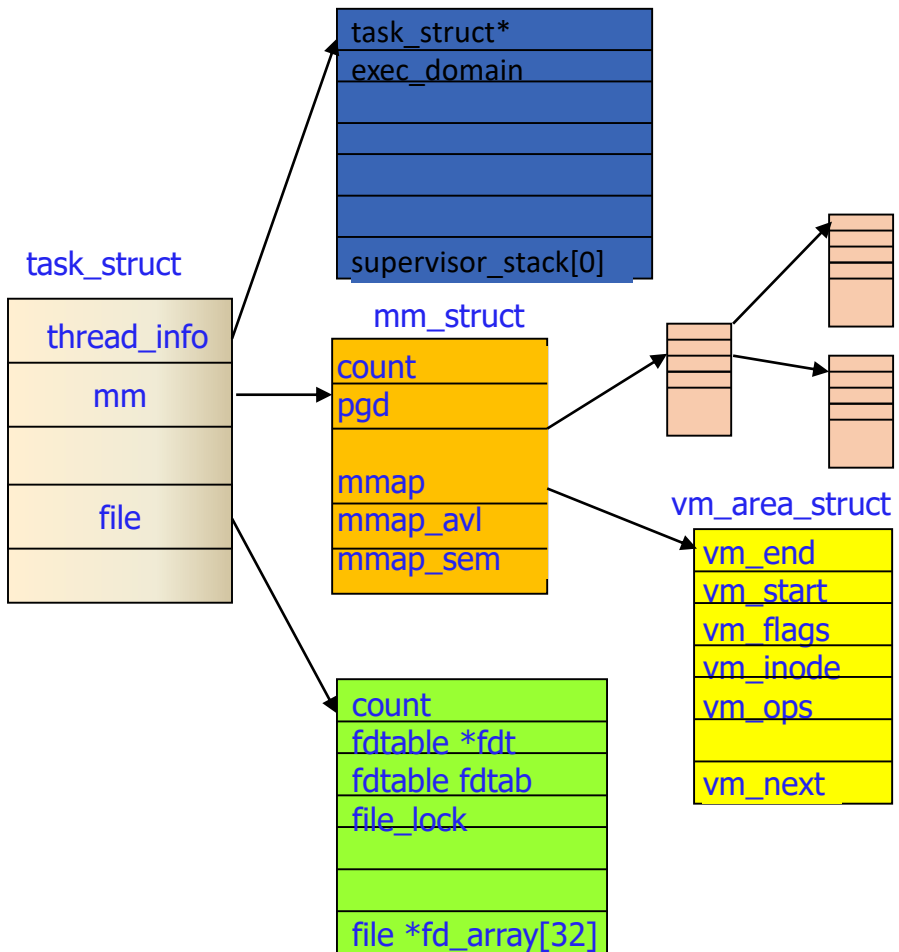
Dept. of Computer Science and CERIAS
Purdue University

March 3rd, 2010
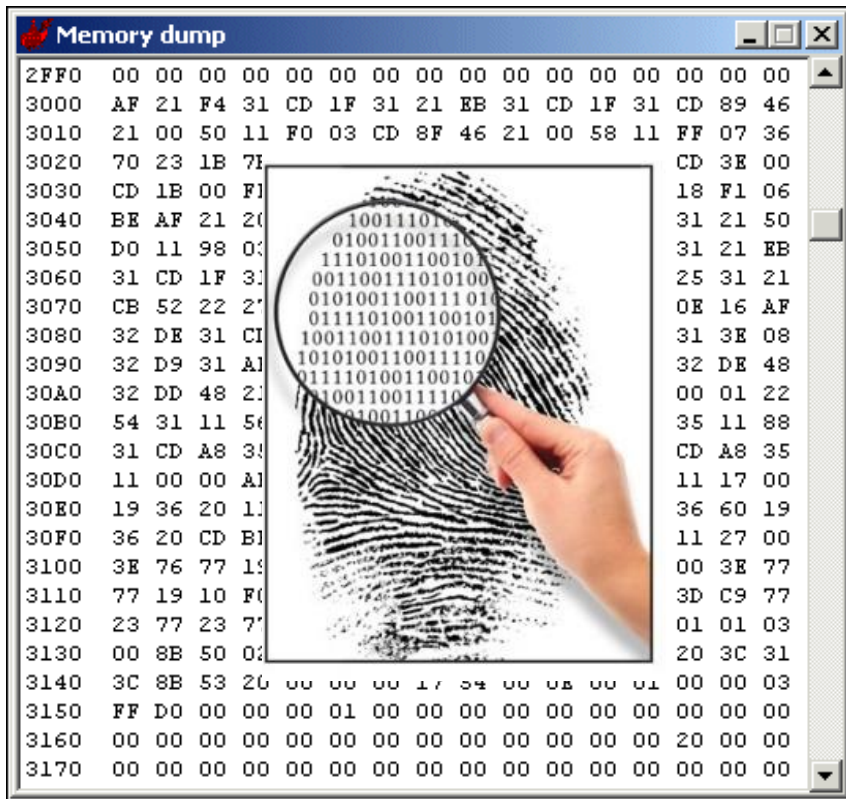
PURDUE UNIVERSITY

CERIAS

# Problem Definition

- Recover data structure specifications from binary
  - Syntactic
    - Layout
    - Offset
    - Size

  - Semantic
    - ip_addr_t, pid_t, …
    - input_t
    - malloc_arg_t, format_string_t…

# Security Applications -- Memory Forensics



task_struct*
exec_domain

supervisor_stack[0]

task_struct

thread_info

mm

file

mm_struct

count
pgd

mmap
mmap_avl
mmap_sem

vm_area_struct

vm_end
vm_start
vm_flags
vm_inode
vm_ops

vm_next

count
fdtable *fdt
fdtable fdtab
file_lock

file *fd_array[32]

# Security Applications -- Vulnerability Discovery

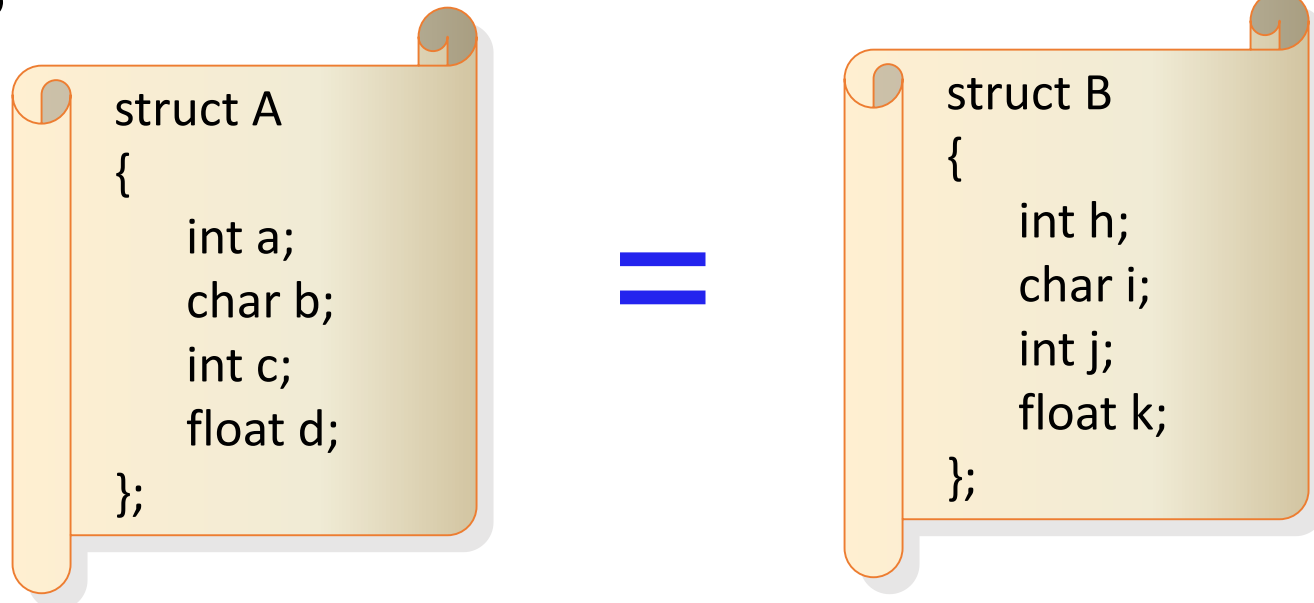```
1  void main(int argc, char* argv[])
2  {
3      char   tempname[1024];
4      strcpy(tempname, argv[1]);
5  }
```

$ ./a.out aaa'\n'

| |
|---|
| char [4] |
| unused[1020] |
| stack_frame_pointer |
| return_address |
| char* |

input_t

$ ./a.out aaaaaaaaaaaaaaaaa...a'\n'

# Security Applications -- Program Signature

```
struct A
{
    int a;
    char b;
    int c;
    float d;
};
```

$=$

```
struct B
{
    int h;
    char i;
    int j;
    float k;
};
```

Laika [Cozzie et al., OSDI'08] Using Data structure as a classifier → Malware Signature

# Overview



Application binary → Type Resolution Points ⇕ Data flow tracking, Type Resolution → Data Structure Specification

**REWARDS**

**R**everse **E**ngineering **W**ork for **A**utomatic **R**evelation of **D**ata **S**tructures

# Key Ideas

```
1  80480a0:  call    0x80480b4<foo>
2  80480a5:  mov    $0x1,%eax
3  80480aa:  mov    $0x0,%ebx
 ...
10 80480c1:  call    0x8048110<getpid>
11 80480c6:  mov    eax, 0x8049124
 ...
36 8048110:  mov    $0x14,%eax
37 8048115:  int    $0x80
38 8048117:  ret
```

```
1 struct {
2     unsigned int pid;
3     char data[16];
4 }test;
5
6 void foo(){
7     char *p="hello world";
8     test.pid=getpid();
9     strcpy(test.data,p);
10 }
...
```

```
fun_0x080480b4{                  foo
-28:     unused[20],
-08:     char *,
-04:     stack_frame_t,
+00:     ret_addr_t
}
```

```
bss_0x08049124{               global
+00:     pid_t,               var test
.04:     char[12],
         unused[4]
}
```

Data Flow

Type Resolution
Point

```
fun_0x08048110{        getpid
+00:   ret_addr_t
}
```

# Type Resolution Point - I

```
<getpid>
36  8048110:   mov    $0x14,%eax
37  8048115:   int    $0x80
38  8048117:   ret
```

- System calls
  - Syscall num →
    - Syscall_enter: Type parameter passing registers (i.e., ebx, ecx, edx, esi, edi, and ebp) if they involved

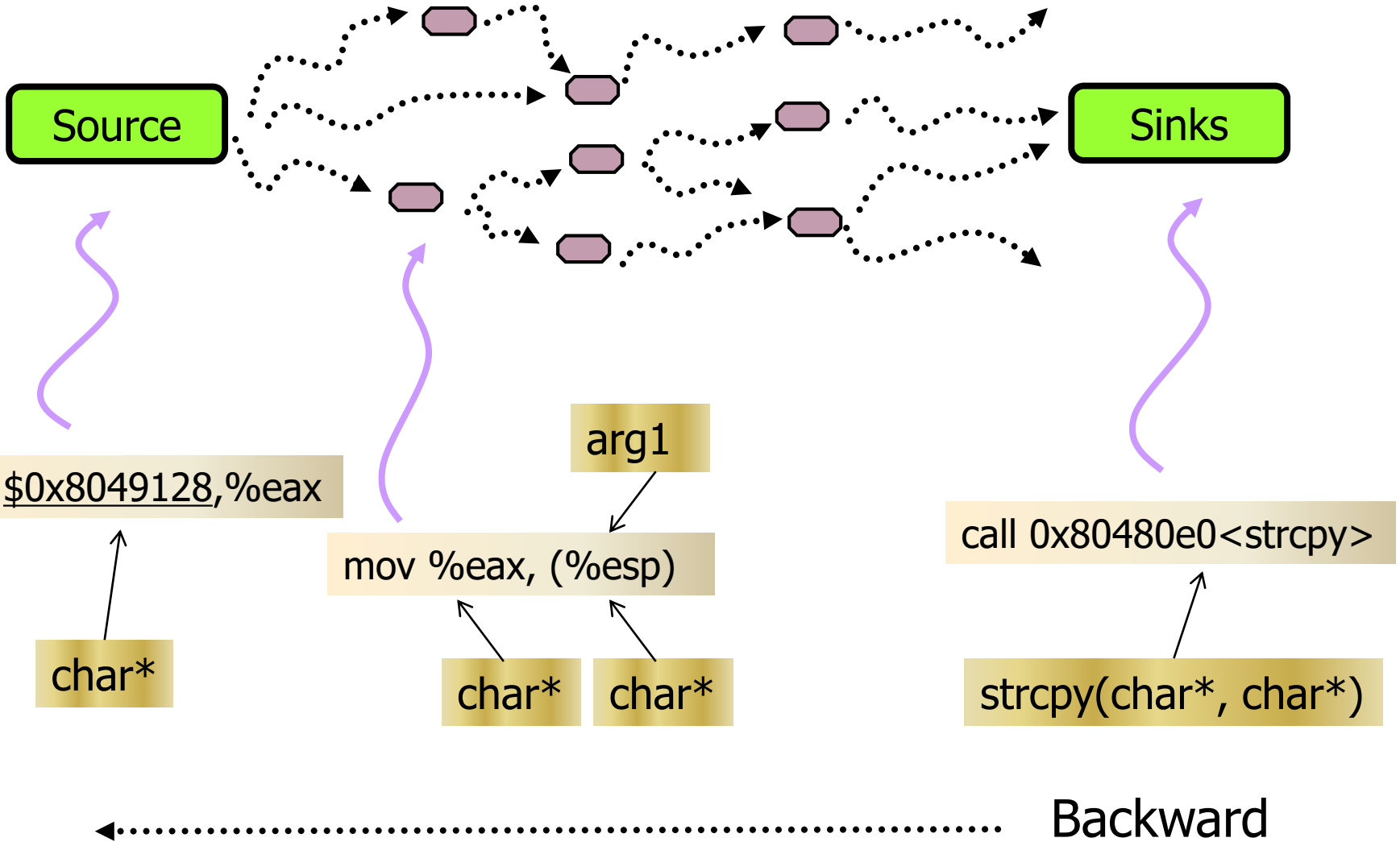    - Syscall_exit: type return value (eax)

# Type Resolution Point - II

```
13  80480ce:      mov %eax,0x4(%esp)       <arg2>
14  80480d2:      movl $0x8049128, (%esp)  <arg1>
15  80480d9:      call 0x80480e0           <strcpy>
```

- Standard Library Call (API)
  - Type corresponding argument and return value
  - More wealth than System call
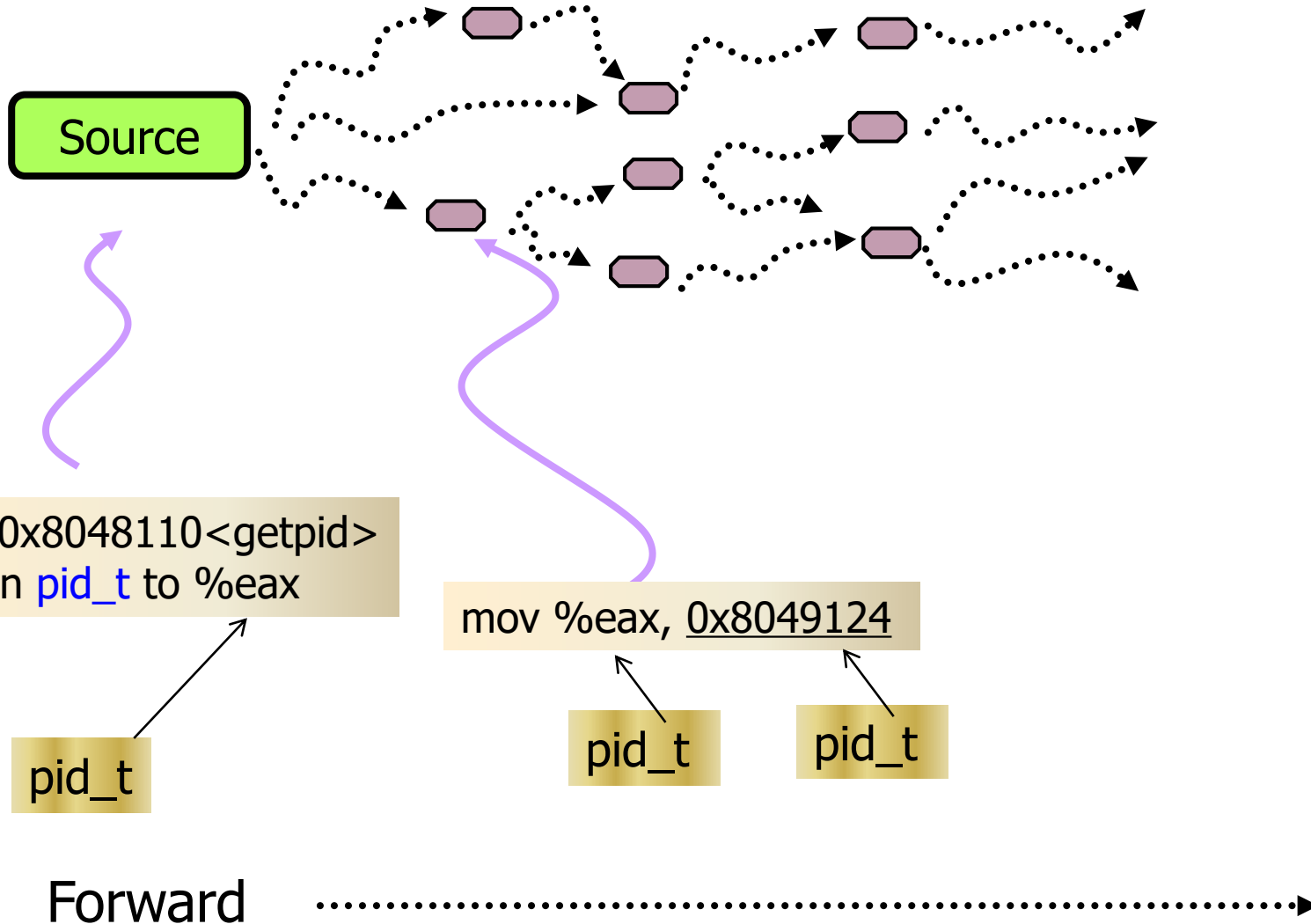    - 2016 APIs in Libc.so.6  vs. 289 sys call (2.6.15)

# Type Resolution Point - III

- Other type revealing instructions in binary code
  - String related (e.g., MOVS/B/D/W, LOADS/STOS/B/D/W)
  - Floating-point related (e.g., FADD, FABS, FST)
  - Pointer-related (e.g., MOV (%edx), %ebx)

# Backward Type Resolution



movl $0x8049128,%eax

char*

arg1

mov %eax, (%esp)

char*  char*

call 0x80480e0<strcpy>

strcpy(char*, char*)

Backward

# Forward Type Resolution



Source

call  0x8048110<getpid>
return pid_t to %eax

mov %eax, 0x8049124

pid_t

pid_t

pid_t

Forward

# Data Flow Propagation

- Similar to taint analysis, using shadow memory to keep the variable attributes and track the propagation

- New challenges
  - Two-way type resolution
  - Dynamic life time of stack and heap variables
    - Memory locations will be reused
  - Multiple instances of the same type

# Implementations

- Dynamic Binary Instrumentation
  - A pin-tool on top of Pin-2.6
    - http://www.pintool.org/



- Data structures of interest:
  - File information
  - Network communication
  - Process information
  - Input-influenced (for vulnerability discovery)

# Evaluations

- Experiment set up
    - 10 utility binaries (e.g., ls, ps, ping, netstat...)
    - Linux 2.6.18, gcc-3.4, 2G mem

- Ground truth → debug information

# Experimental Results

- False negative: the data structure we missed (compared with the ground truth)
  - Global: 70%
  - Heap: 55%
  - Stack: 60%

- Why false negative
  - Dynamic analysis

# Experimental Results

- False positive (we get the wrong data type)
  - Global: 3%
  - Heap: 0%
  - Stack: 15%

# Sources of False Positives

- Loss of data structure hierarchy
  - There is no corresponding type resolution point with the same hierarchical type
  - Heuristics exist (e.g., AutoFormat [Lin et al, NDSS'08])
- Path-sensitive memory reuse
  - Compiler might assign different local variables declared in different program paths to the same memory address
  - Path-sensitive analysis

- Type cast
  - int a=(float)b;

# Application I: Memory Forensics

```
...
08052170   b0 5b fe b7   b0 5b fe b7   05 00 00 00   02 00 92 7e
08052180   0a 00 00 0b 00 00 00 00   00 00 00 00   c7 b0 af 4a
08052190   c7 b0 af 4a 00 00 00 00   58 2a 05 08   00 00 00 00
080521a0   00 00 00 00 00 00 00 00   00 00 00 00 00 00 00 00
...
```

struct 0x804dd4f {
| 00: | pthread_t; | b0 5b fe b7 |
| 04: | int; | b0 5b fe b7 |
| 08: | socket; | 05 00 00 00 |
| 12: | struct sockaddr_in; | 02 00 92 7e 0a 00 00 0b .. 00 00 |
| 28: | time_t; | c7 b0 af 4a |
| 32: | time_t; | c7 b0 af 4a |
| 36: | unused[4]; | 00 00 00 00 |
| 40: | struct 0x804ddfb*; | 58 2a 05 08 |
};

# Application I: Memory Forensics

```
...
08052170   b0 5b fe b7 b0 5b fe b7   05 00 00 00  02 00  92 7e
08052180   0a 00 00 0b  00 00 00 00   00 00 00 00  c7 b0 af 4a
08052190   c7 b0 af 4a 00 00 00 00   58 2a 05 08 00 00 00 00
080521a0   00 00 00 00 00 00 00 00   00 00 00 00 00 00 00 00
...
```

struct 0x804dd4f {      {
  00:      sin_family;           02 00
  04:      sin_port;             92 7e
  08:      sin_addr;             0a 00 00 0b
  08:      struct sockaddr_in;   02 00 92 7e 0a 00 00 0b .. 00 00
                                 00 00 00 00 00 00 00 00 00 00
  28:      time_t;
  32:      time_t;
  36:      unused[4];
  40:      struct 0x804ddfb*;    ip_addr_t: 10.0.0.11
};

# Application II: Vulnerability Discovery

| Program | buffer overflow #total | #real | integer overflow #total | #real | Format string #total | #real |
|---|---|---|---|---|---|---|
| ncompress-4.2.4 | 1 | 1 | 0 | 0 | 0 | 0 |
| bftpd-1.0.11 | 3 | 1 | 0 | 0 | 0 | 0 |
| gzip-1.2.4 | 3 | 1 | 0 | 0 | 0 | 0 |
| nullhttpd-0.5.0 | 5 | 1 | 2 | 1 | 0 | 0 |
| xzgv-5.8 | 3 | 0 | 8 | 1 | 0 | 0 |
| gnuPG-1.4.3 | 0 | 1 | 3 | 1 | 0 | 0 |
| ipgrab-0.9.9 | 0 | 1 | 5 | 1 | 0 | 0 |
| cfingerd-1.4.3 | 4 | 0 | 0 | 0 | 1 | 1 |
| ngircd-0.8.2 | 12 | 0 | 0 | 0 | 1 | 1 |

Suspects    True

# Discussion

- Dynamic analysis
  - Full coverage of data structures
- Only user-level data structure
  - OS kernel
- Obfuscated binary can cheat REWARDS
  - No library call
  - Memory cast
- Lack of other application-specific type resolution points
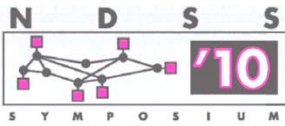  - Other APIs
    - E.g., browser-related

# Summary

- REWARDS
  - Binary only
  - Dynamic analysis
  - Data flow tracking

- Key insight
  - Using system call/API/Type revealing instruction as type resolution point
  - Two-way type resolution
  - Unique benefits to memory forensics and vulnerability discovery
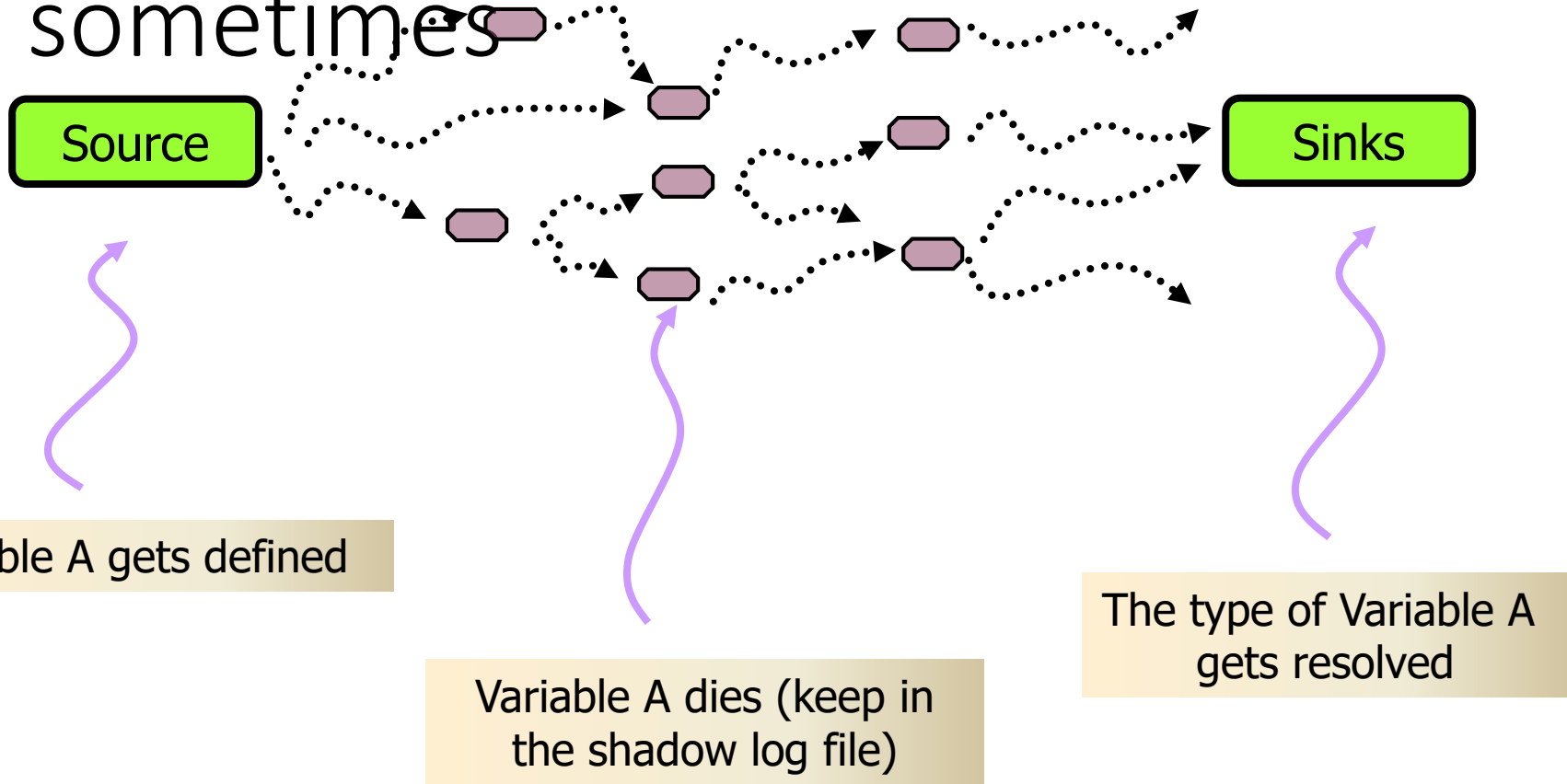
# Backup: Path-sensitive memory reuse

```
919          case 'e': {
920              if (strlen(params[0]) > 1) switch (tolower(params[0][1])) {
921                  case 's': {
922                      struct escan_rec rc;
923                      if (num_params <= 2) Log("escan <dest>\n");
924                      else {
925                          memset((void*)&rc,0,sizeof(struct escan_rec));


1061     if (udp.len != 0) if (!audp_recv(&udp,&tmpclient,buf,3000)) {
1062         struct header *rc;
1063         struct llheader *llrc;
1064         char k=0;
1065         buf+=sizeof(struct llheader);
1066         udp.len-=sizeof(struct llheader);
1067         llrc=(struct llheader *)(buf-sizeof(struct llheader));
1068         rc=(struct header *)buf;
1069         if (llrc->type == 0) {
1070             struct llheader ll;
1071             memset((void*)&ll,0,sizeof(struct llheader));
```

(pudclient.c)

# Backup: offline is needed sometimes



Source

Sinks

Variable A gets defined

Variable A dies (keep in the shadow log file)

The type of Variable A gets resolved

Backward

# Backup: Vulnerability Discovery

- Buffer overflow

```
fun_0x08048e76{
- 1052:    char[13],                      input_t   argv
- 1039:    unused[1023],...
- 0004:    stack_frame_t,
+0000:    ret_addr_t,
+0004:    char*}
```

- Format string

```
fun 0x0805f9a5 { ...,
- 0284:    format_string_t[76],           input_t   recv
- 0208:    unused[204],
- 0004:    stack_frame_t,
+0000:    ret_addr_t,...}
```

- Integer overflow

```
bss 0x0809ac80 { ...
+91952:    int (malloc_arg_t)             input_t  fread
+91956:    int (malloc_arg_t)
...}
```