# Automatic Mining of Functionally Equivalent Code Fragments via Random Testing
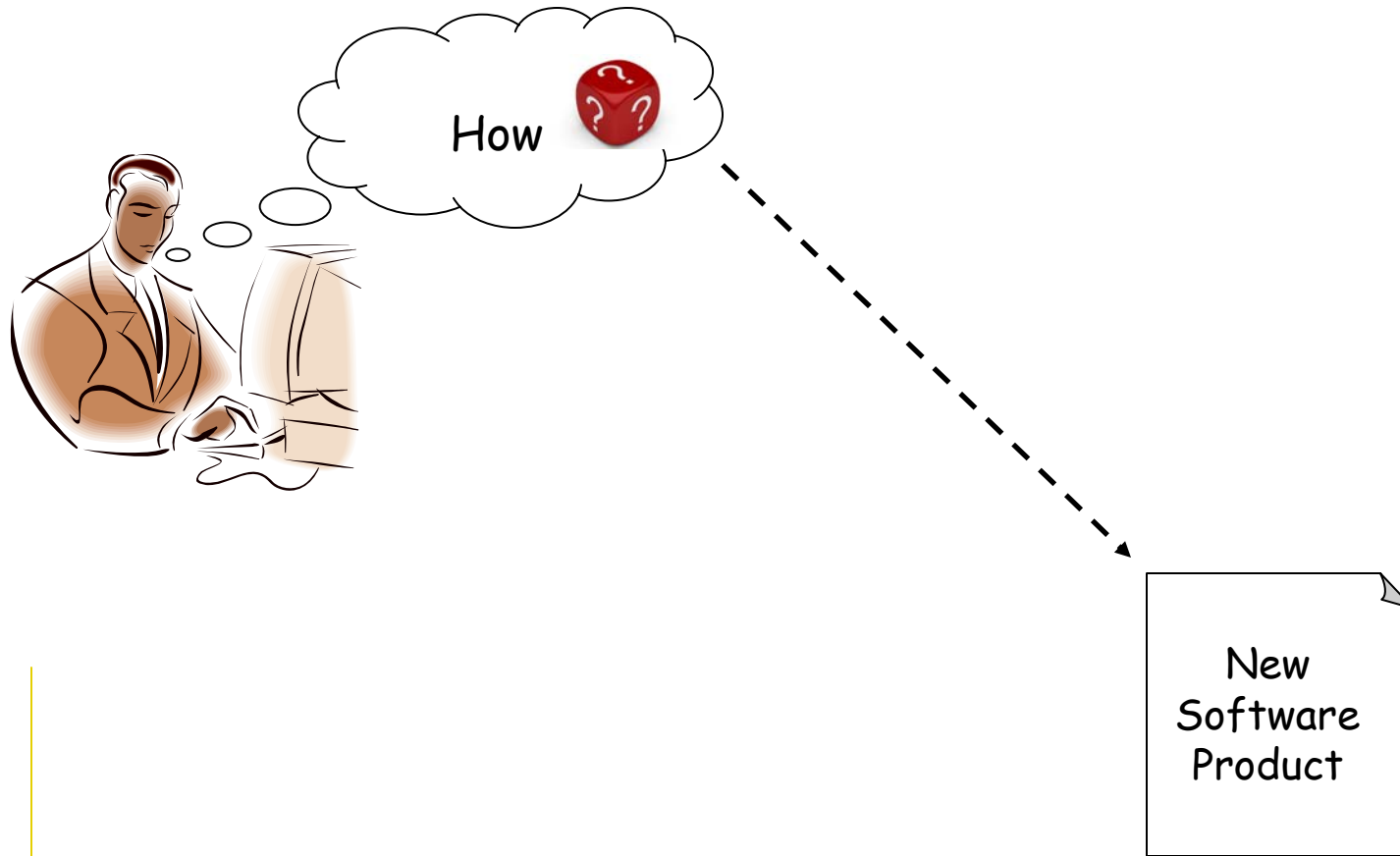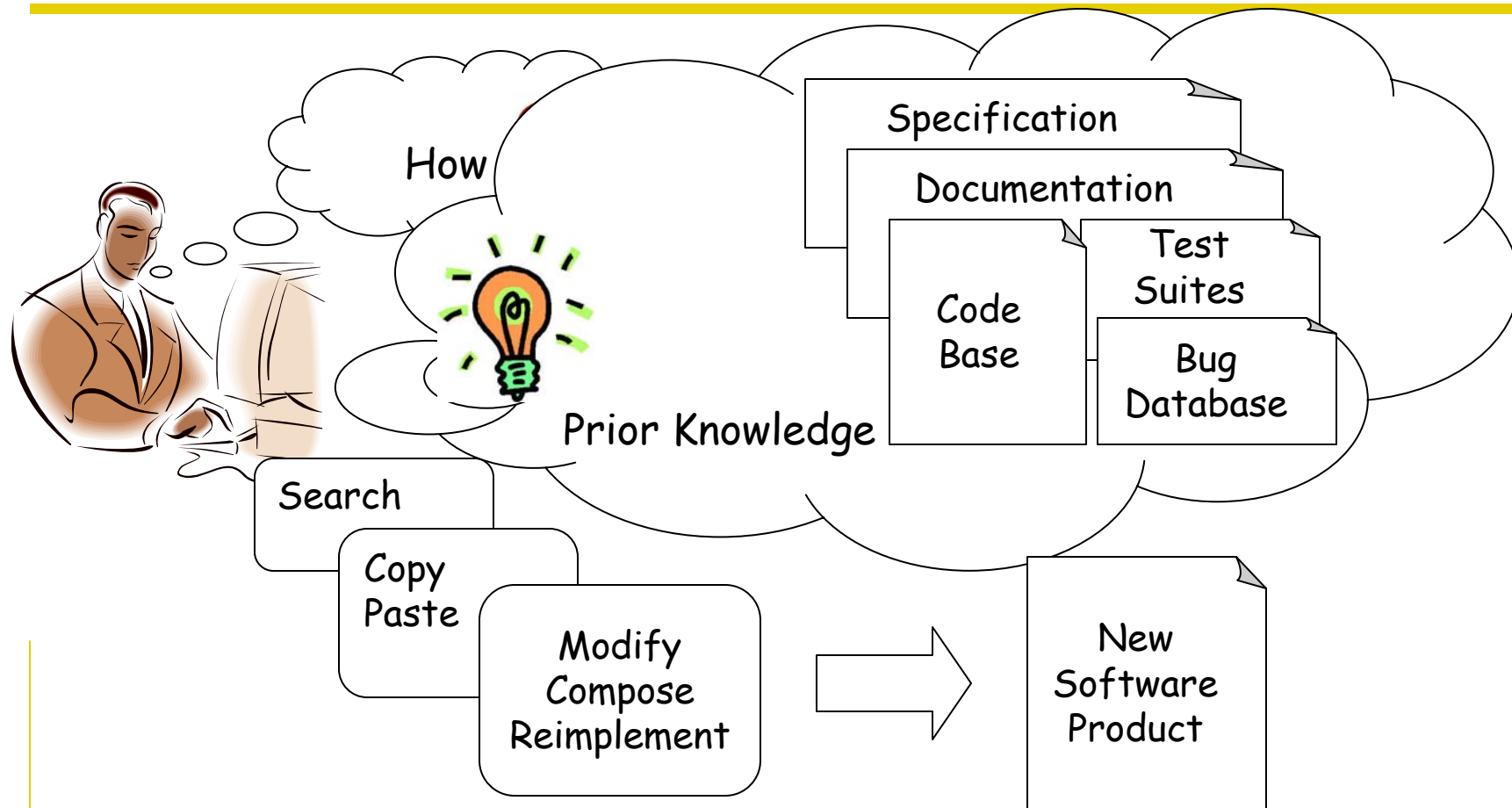
Lingxiao Jiang  and  Zhendong Su

# Cloning in Software Development

How ❓

New
Software
Product

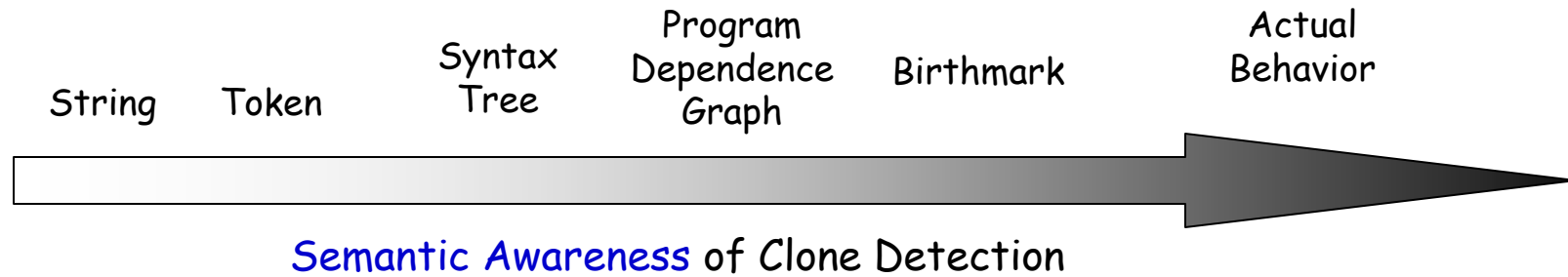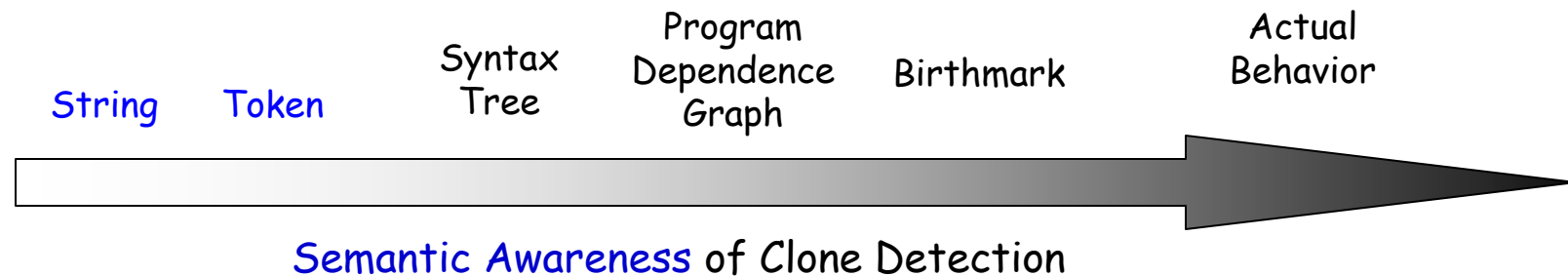# Cloning in Software Development

# Applications of Clone Detection

- Refactoring

- Pattern mining

- Reuse

- Debugging

- Evolution study

- Plagiarism detection

# A Spectrum of Clone Detection

String     Token     Syntax Tree     Program Dependence Graph     Birthmark     Actual Behavior

**Semantic Awareness** of Clone Detection

# A Spectrum of Clone Detection

String          Token          Syntax Tree          Program Dependence Graph          Birthmark          Actual Behavior
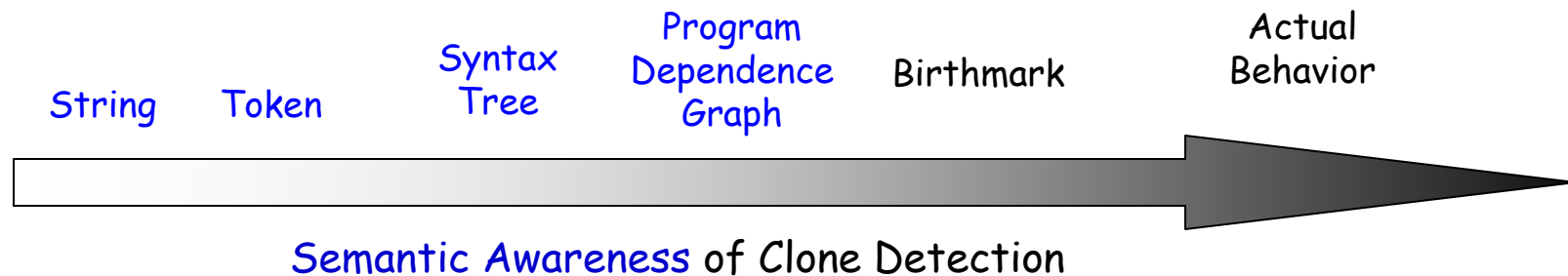
**Semantic Awareness** of Clone Detection

- 1992: Baker, parameterized string algorithm
- 2002: Kamiya et al., CCFinder
- 2004: Li et al., CP-Miner
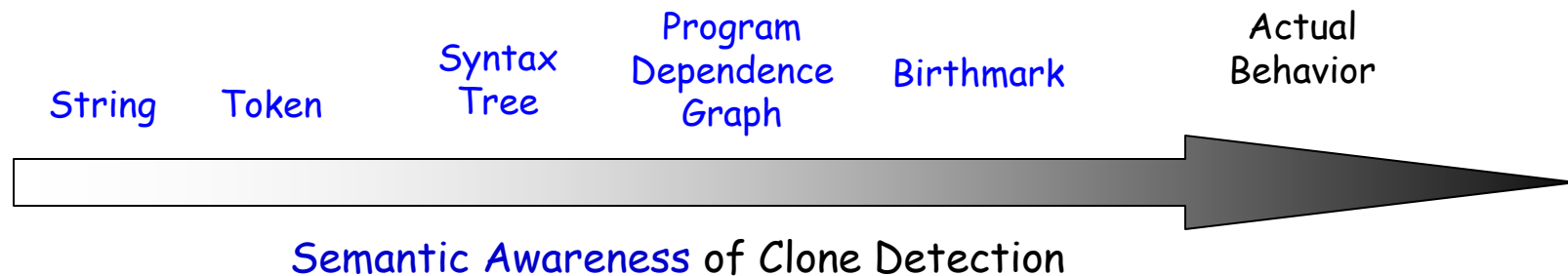- 2007: Basit et al., Repeated Tokens Finder

# A Spectrum of Clone Detection

String    Token    Syntax Tree    Program Dependence Graph    Birthmark    Actual Behavior

**Semantic Awareness** of Clone Detection

- 1998: Baxter et al., CloneDR
- 2004: Wahler et al., XML-based
- 2007: Jiang et al., Deckard
- 2000, 2001: Komondoor et al.
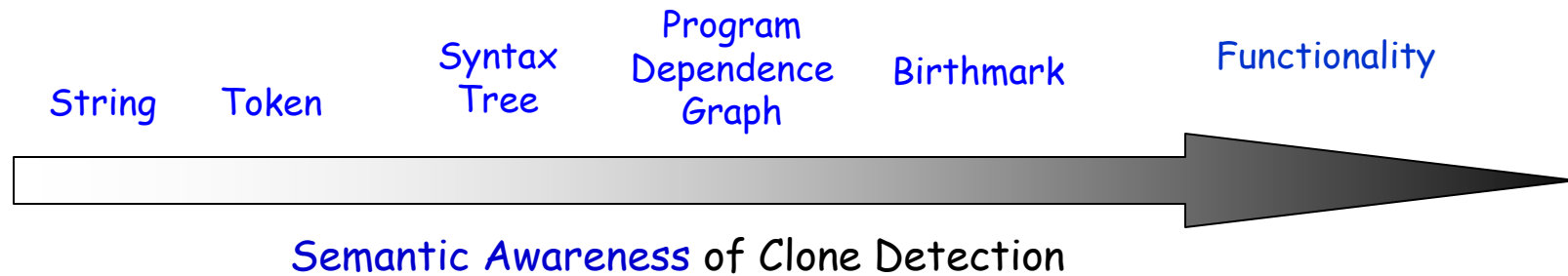- 2006: Liu et al., GPLAG
- 2008: Gabel et al.

**CSSR**
Center for Software and Systems Research

# A Spectrum of Clone Detection

String    Token    Syntax Tree    Program Dependence Graph    Birthmark    Actual Behavior

**Semantic Awareness** of Clone Detection

- 1999: Collberg et al., Software watermarking
- 2007: Schuler et al., Dynamic birthmarking
- 2008: Lim et al., Static birthmarking
- 2008: Zhou et al., Combined approach

**CSSR**
Center for Software and Systems Research

# A Spectrum of Clone Detection

String          Token          Syntax Tree          Program Dependence Graph          Birthmark          Functionality

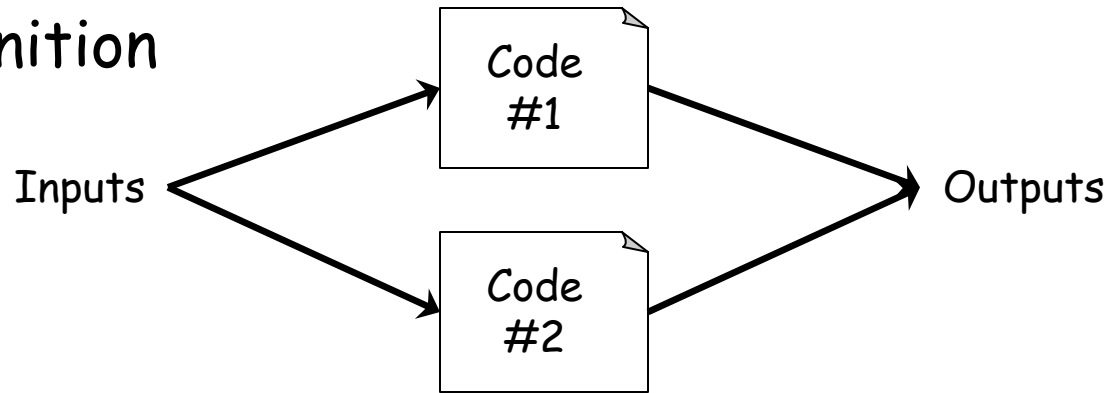Semantic Awareness of Clone Detection

- ## Functional equivalence
  - ### How extensive is its existence

# Functional Equivalence

- ## Definition



- ## Applicability: arbitrary piece of code
  - – Source and binary
  - – From whole program to whole function to code fragments

- ## Example: sorting algorithms
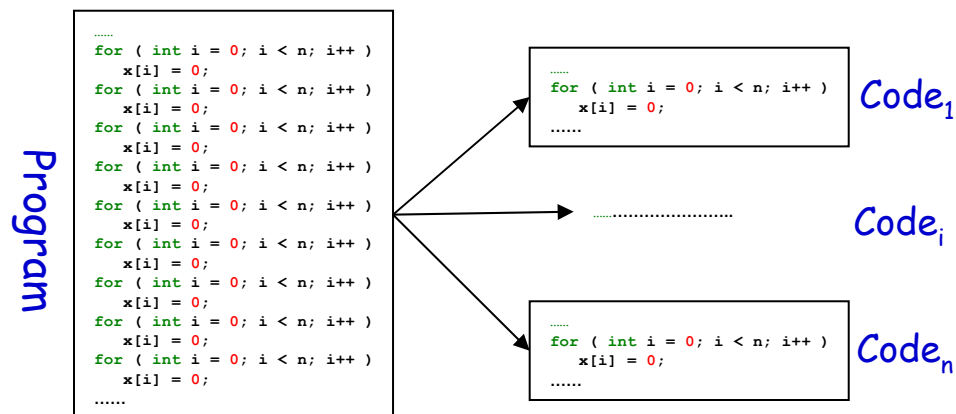  - – Bubble, selection, merge, quick, heap

# Previous Work on Program Equivalence

- [Cousineau 1979; Raoult 1980; Zakharov 1987; Crole 1995; Pitts 2002; Bertran 2005; Matsumoto 2006; Siegel 2008; …]

- Many based on formal semantics
- Consider whole programs or functions only
  - Not arbitrary code fragments
- Check equivalence among given pieces of code
  - Not scalable detection

CSSR
Center for Software and Systems Research

# Our Objectives

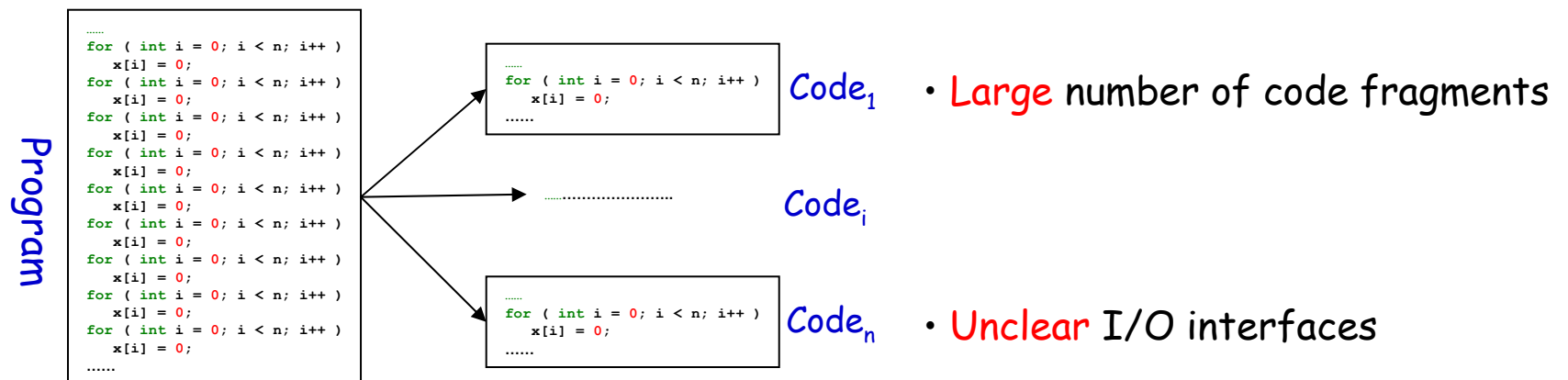- ## Detect functionally equivalent code fragments



- ## Compare I/O behaviors directly
  - Run each piece of code with random inputs

# Our Objectives — Challenges

- ## Detect functionally equivalent code fragments



- **Large** number of code fragments

- **Unclear** I/O interfaces

- ## Compare I/O behaviors directly
  - Run each piece of code with random inputs

- **Huge** number of code executions

# Key 1: Semantic-Aware I/O Identification

- Identify input and output variables based on data flows in the code:
  - Variables used before defined are inputs
  - Variables defined but may not used are outputs

```
1   min = i;
2   j = i+1;
3   while (1) {
4   if(j >= LENGTH)
5   break;
6   if(data[j] < data[min])
7   min = j;
8   j++; }
9   if(min > i) {
10  tmp = data[min];
11  data[min] = data[i];
12  data[i] = tmp; }
```

Input variables: i and data

Output variables: data

CSSR
Center for Software and Systems Research

# Key 2: Limit Number of Inputs

- Schwartz-Zippel lemma: polynomial identities can be tested with few random values
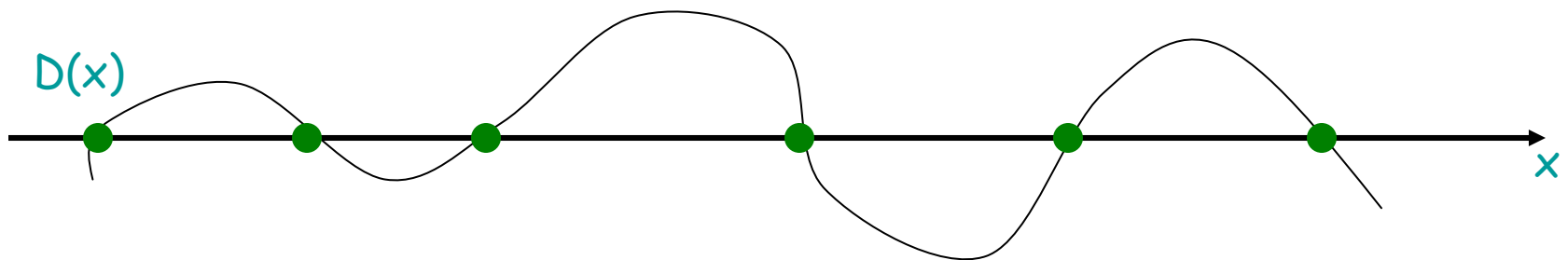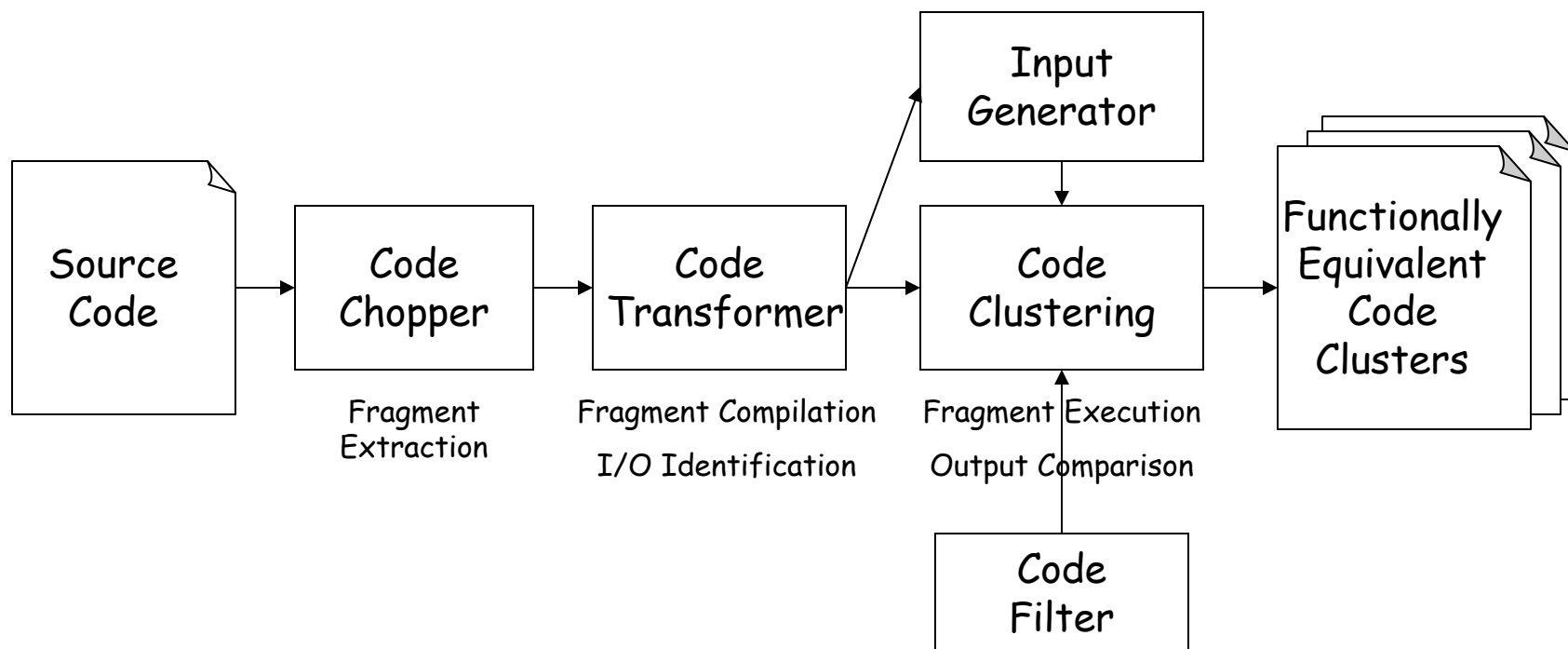  - Let $D(x)$ be $p_1(x) - p_2(x)$
  - If $p_1(x) = p_2(x)$,

$D(x)$



  - If $p_1(x) \neq p_2(x)$,
    - $D(x) = 0$ has at most finite number $d$ of roots
    - Prob ( $D(v) = 0$ ) is bounded by $d$, for any random value $v$ from the domain of $x$.

$D(x)$

# EqMiner

# Code Chopper

- Sliding windows of various sizes on serialized statements

```
1   min = i;
2
3   for(j=i+1; j<LENGTH; j++)
4   {
5
6     if(data[j] < data[min])
7         min = j;
8   }
9   if (min > i) {
10      int tmp = data[min];
11      data[min] = data[i];
12      data[i] = tmp; }
```

```
1   min = i;
2   j = i+1;
3   while (1) {
4   if(j >= LENGTH)
5   break;
6   if(data[j] < data[min])
7   min = j;
8   j++; }
9   if(min > i) {
10  tmp = data[min];
11  data[min] = data[i];
12  data[i] = tmp; }
```

CSSR
Center for Software and Systems Research

# Code Transformer

- Declare undeclared variables, labels
- Define all used types
- Remove assembly code
- Replace goto, return statements
- Replace function calls
  - Replace each call with a random input variable
  - Ignore side effects, only consider return values
- Read inputs
- Dump outputs

# Input Generation

- In order to share concrete input values among input variables for different code fragments, separate the generation into two phases:

  1. Construct bounded memory pools filled with random primary values and pointers. E.g.,

| Primary value pool (bytes): | 100 | -78 | ...... |
|---|---|---|---|
| Pointer value pool (0/1): | 1 | 0 | ...... |

  2. Initialize each variable with values from the pools. E.g.,

```
struct { int x, y; } X;
```

Input variables: X* x; int* y;

```
x = malloc(sizeof(X));
x.x = 100; x.y = -78;

y = 0;
```

# Code Clustering

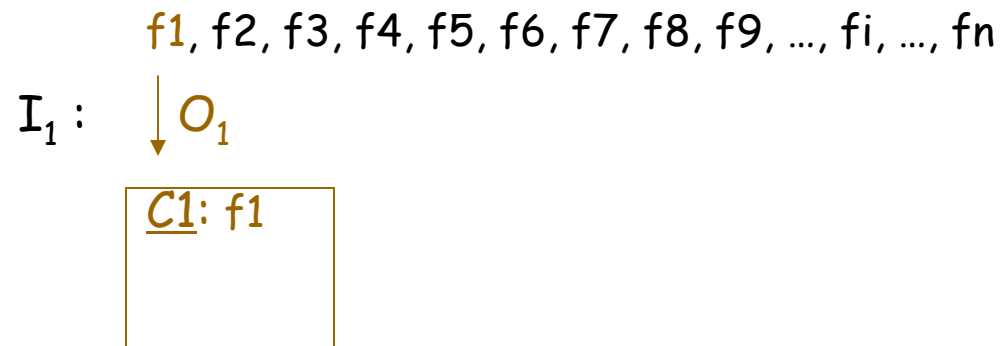- **Eager partitioning** of code fragments for a set of random inputs

$$f1, f2, f3, f4, f5, f6, f7, f8, f9, ..., fi, ..., fn$$

$I_1$ :

# Code Clustering

- Eager partitioning of code fragments for a set of random inputs

$f1$, f2, f3, f4, f5, f6, f7, f8, f9, …, fi, …, fn

$I_1$ :     $\downarrow O_1$

C1: f1

# Code Clustering

- Eager partitioning of code fragments for a set of random inputs

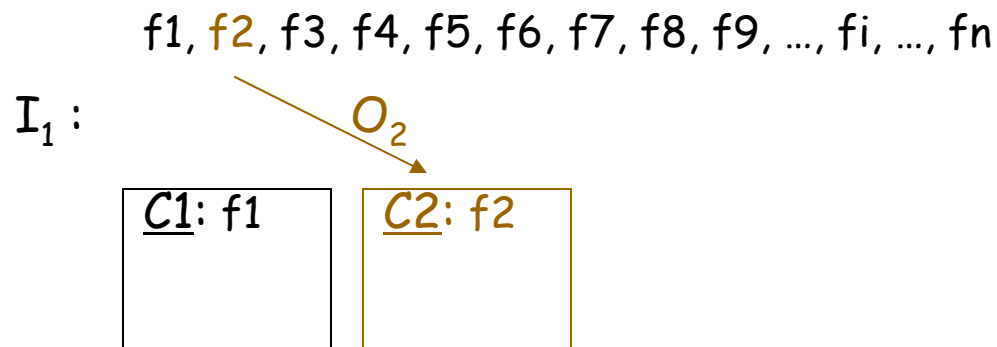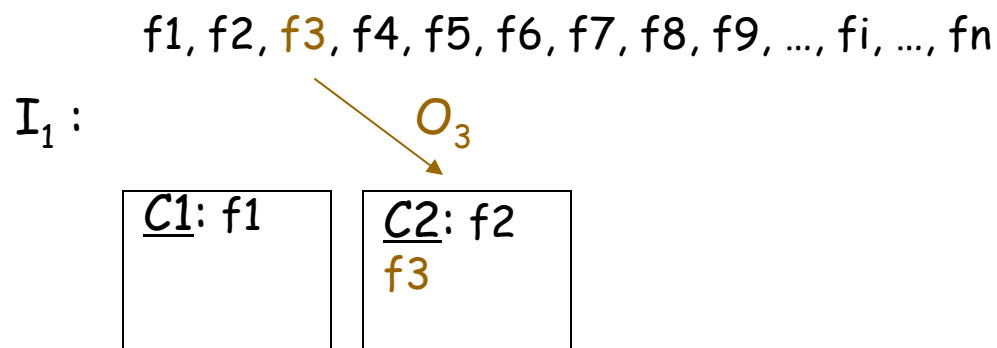f1, f2, f3, f4, f5, f6, f7, f8, f9, …, fi, …, fn

$I_1$ :　　　　　　　　　$O_2$

| $\underline{C1}$: f1 | $\underline{C2}$: f2 |
|---|---|

CSSR

Center for Software and Systems Research

# Code Clustering

- Eager partitioning of code fragments for a set of random inputs

f1, f2, f3, f4, f5, f6, f7, f8, f9, …, fi, …, fn

$I_1$ :                          $O_3$

C1: f1          C2: f2
                f3

# Code Clustering

- Eager partitioning of code fragments for a set of random inputs

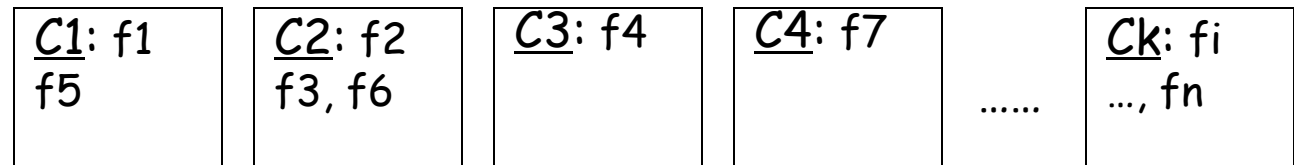f1, f2, f3, f4, f5, f6, f7, f8, f9, …, fi, …, fn

$I_1$ :

$O_4$

| $C1$: f1 | $C2$: f2 f3 | $C3$: f4 |

# Code Clustering

- Eager partitioning of code fragments for a set of random inputs

f1, f2, f3, f4, f5, f6, f7, f8, f9, ..., fi, ..., fn

$I_1$ :

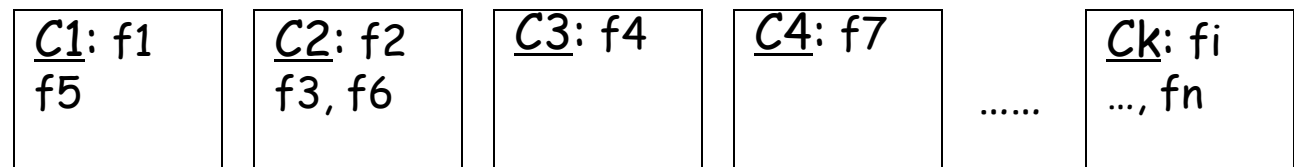| C1: f1 f5 | C2: f2 f3, f6 | C3: f4 | C4: f7 | ...... | Ck: fi ..., fn |
|---|---|---|---|---|---|

# Code Clustering

- Eager partitioning of code fragments for a set of random inputs

f1, f2, f3, f4, f5, f6, f7, f8, f9, ..., fi, ..., fn

$I_1$ :

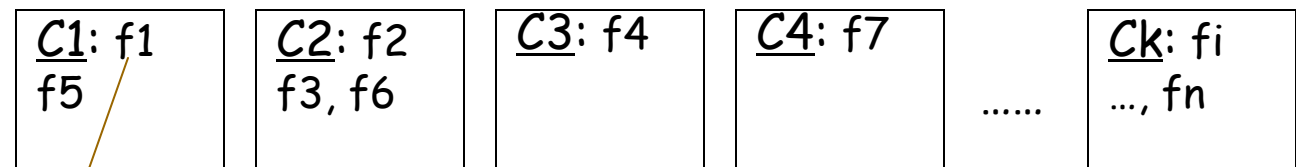| C1: f1 f5 | C2: f2 f3, f6 | C3: f4 | C4: f7 | ...... | Ck: fi ..., fn |

$I_2$ : repeat the same for each intermediate cluster

# Code Clustering

- Eager partitioning of code fragments for a set of random inputs

f1, f2, f3, f4, f5, f6, f7, f8, f9, …, fi, …, fn

$I_1$ :

| C1: f1 f5 | C2: f2 f3, f6 | C3: f4 | C4: f7 | …… | Ck: fi …, fn |

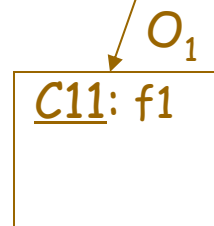$I_2$ : repeat the same for each intermediate cluster

$O_1$

C11: f1

# Code Clustering

- Eager partitioning of code fragments for a set of random inputs

f1, f2, f3, f4, f5, f6, f7, f8, f9, …, fi, …, fn

$I_1$ :

| $\underline{C1}$: f1 f5 | $\underline{C2}$: f2 f3, f6 | $\underline{C3}$: f4 | $\underline{C4}$: f7 | …… | $\underline{Ck}$: fi …, fn |

$I_2$ : repeat the same for each intermediate cluster
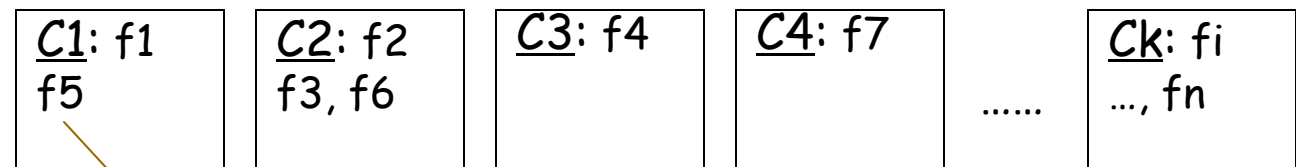
$O_5$

| $\underline{C11}$: f1 | $\underline{C12}$: f5 |

# Code Clustering
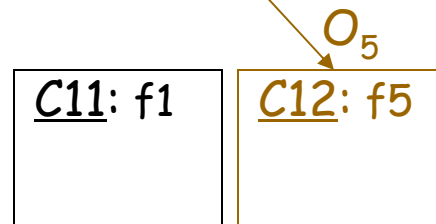
- Eager partitioning of code fragments for a set of random inputs

f1, f2, f3, f4, f5, f6, f7, f8, f9, …, fi, …, fn

$I_1$ :

| C1: f1 f5 | C2: f2 f3, f6 | C3: f4 | C4: f7 | …… | Ck: fi …, fn |

$I_2$ : repeat the same for each intermediate cluster
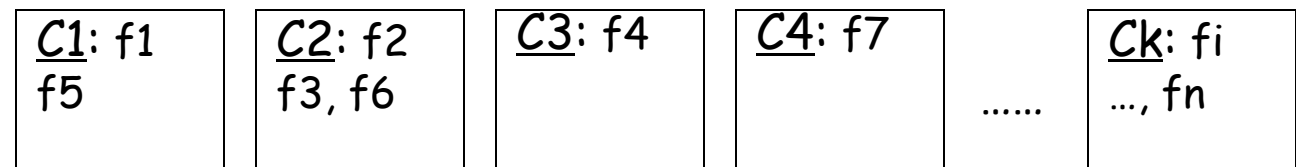
| C11: f1 | C12: f5 | …… | Ck1: fi fj | Ck2: fl …, fp | Ckx: fq …, fn |

# Code Clustering
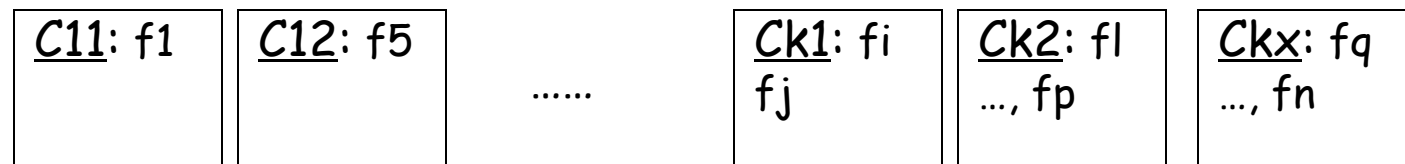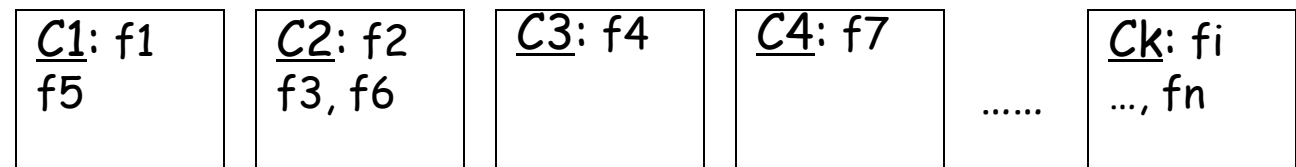
- Eager partitioning of code fragments for a set of random inputs

f1, f2, f3, f4, f5, f6, f7, f8, f9, …, fi, …, fn

$I_1$ :

| C1: f1<br>f5 | C2: f2<br>f3, f6 | C3: f4 | C4: f7<br><br>…… | Ck: fi<br>…, fn |
|---|---|---|---|---|

$I_2$ : repeat the same for each intermediate cluster

| C11: f1 | C12: f5 | …… | Ck1: fi<br>fj | Ck2: fl<br>…, fp | Ckx: fq<br>…, fn |
|---|---|---|---|---|---|

$I_s$ : until only one code fragment is left for each cluster, or until a reasonable number $s$ of inputs are used

CSSR
Center for Software and Systems Research

# EqMiner

# Results on Sorting Algorithms

- 5 sorting algorithms with both recursive and non-recursive versions
  - ~350 LoC
  - ~200 code fragments


- s = 10

  - 69 clone clusters reported
    - Most are portions of the algorithms
    - 4 non-recursive versions are in a same cluster

CSSR
Center for Software and Systems Research

# Results on the Linux Kernel

- s = 10

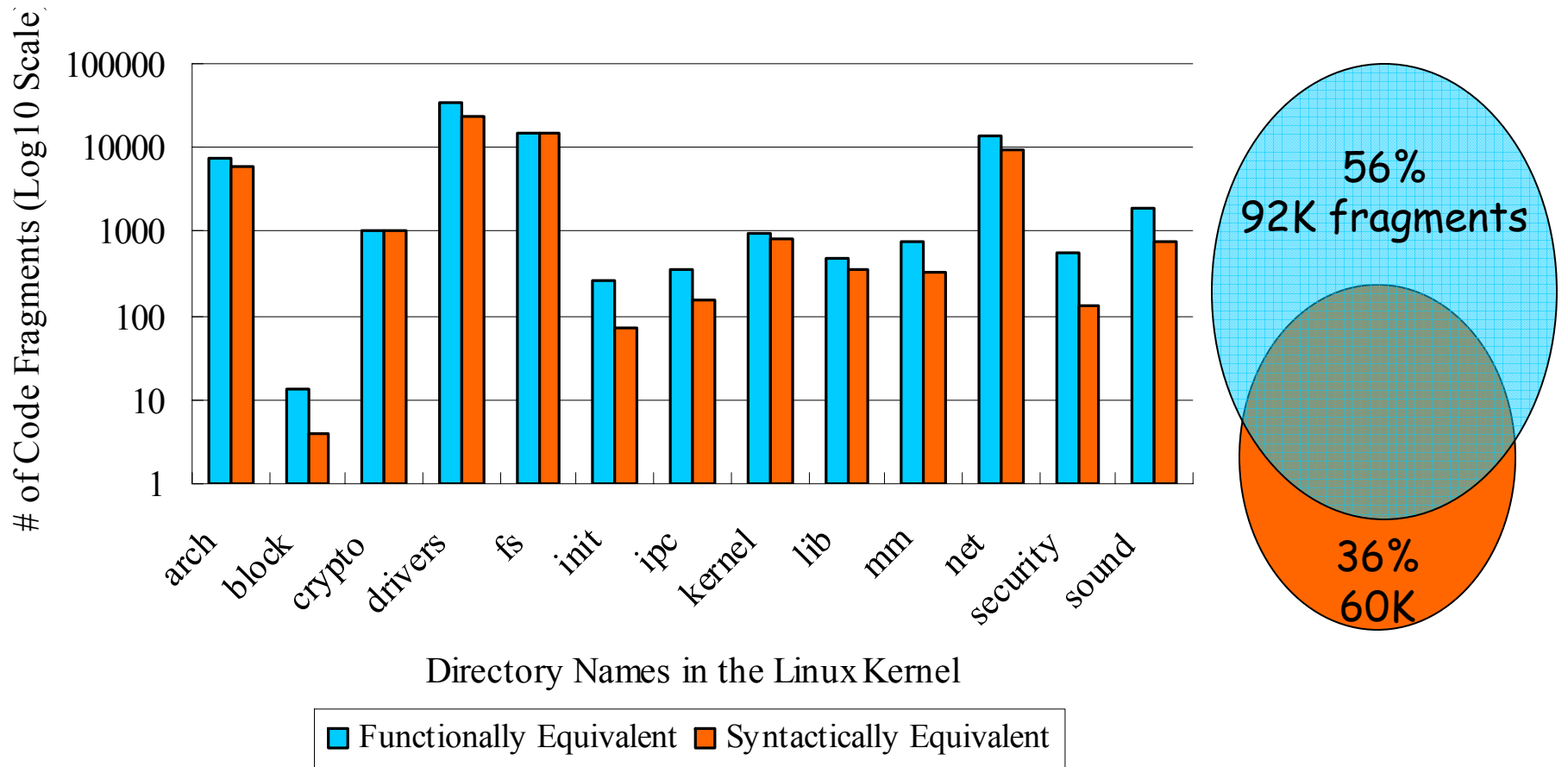  - >800K code fragments were separated into 32K non-trivial clusters

# Results on the Linux Kernel

- s = 10
  - >800K code fragments were separated into 32K non-trivial clusters


- Additional 100 for 128 semi-randomly selected clusters

  - 3% of all of the code fragments became singletons

- 100 more tests
  - 0.5% additional

# Differences from Syntactic Clones



**# of Code Fragments (Log10 Scale)**

y-axis: 100000, 10000, 1000, 100, 10, 1

x-axis: arch, block, crypto, drivers, fs, init, ipc, kernel, lib, mm, net, security, sound

**Directory Names in the Linux Kernel**

■ Functionally Equivalent   ■ Syntactically Equivalent

56%
92K fragments

36%
60K

Center for Software and Systems Research

# Differences from Syntactic Clones

- **False positives**
  - Function calls

- Macro related + few outputs

```
if ( ALWAYS_FALSE ) {
    ……

} else
    output = input;                      output = input;
```

- Lexical differences

```
output = input + 10;                  output = input + 100;

output = 0;                           output = 0;
if ( output < input ) {               if ( output < input ) {
    ...                                   ...
    output = input + 1;                   output = output + 1;
}                                     }
```

# Conclusion & Future Work

- First scalable detection of functionally equivalent code based on random testing

- Confirm the existence of many functional clones which complement syntactic clones
  – Enable further studies on functional clone patterns
  – Explore utilities of functional equivalent code

# Thank you!

Questions?
jiangl@cs.ucdavis.edu

CSSR
Center for Software and Systems Research