

Lecture 10

Pointer Analysis

1. Datalog
2. Context-insensitive, flow-insensitive pointer analysis
3. Context sensitivity

Readings: Chapter 12

Pointer Analysis to Improve Security

- Top web application security vulnerabilities
 - SQL injection, cross-site scripting
- User input accessing databases
- Information flow analysis (taint analysis)
- Sound analysis that found errors in 8 out of 9 apps

PQL

```
 $p_1 = req.getParameter ( );$   
 $stmt.executeQuery (p_2);$ 
```

p_1 and p_2 point to same object?

Pointer alias analysis

Automatic Analysis Generation



Programmer:
Security analysis
in 10 lines

Compiler writer:
Ptr analysis
in 10 lines

1000s of lines
1 year tuning

PQL

Datalog

BDD operations

bddbddb
(**BDD**-based
deductive database)
with
Active Machine Learning

BDD: 10,000s-lines library

Goals of the Lecture

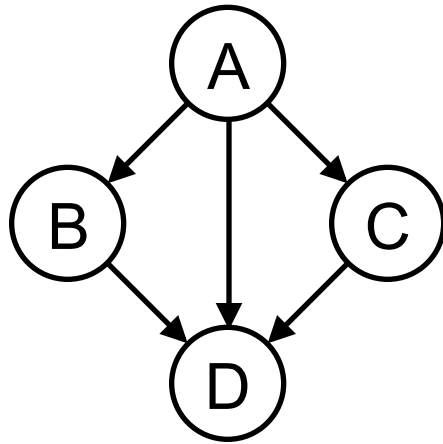
- Pointer analysis
 - Interprocedural, context-sensitive, flow-insensitive
(Dataflow: intraprocedural, flow-sensitive)
- Power of languages and abstractions
- Elegant abstractions
 - Logic programming
 - BDDs: Binary decision diagrams
(Most-cited CS paper a few years ago)

1. Datalog Basics

- $p(X_1, X_2, \dots, X_n)$
 - p is a predicate
 - X_1, X_2, \dots, X_n are terms such as variables or constants
- A predicate can be viewed as a relation

Example: Call graph edges

Predicate vs. Relation



calls(A,B)
calls(A,C)
calls(A,D)
calls(B,D)
calls(C,D)

Predicates

- Calls (x,y): x calls y is true
- Ground atoms: predicates with constant arguments

Relations

- Calls (x,y) : x, y is in a “calls” relationship
- Extensional database: tuples representing facts

Datalog Programs:

Set of Rules (Intensional DB)

- $H :- B_1 \ \& \ B_2 \ \dots \ \& \ B_n$
- LHS is true if RHS is true
 - Rules define the intensional database
- Example: Datalog program to compute call*
 - transitive closure of calls relation
 - $\text{calls}^*(x, y)$ if x calls y directly or indirectly
 - $\text{calls}^*(x, y) :- \text{calls}(x, y)$
 - $\text{calls}^*(x, z) :- \text{calls}^*(x, y) \ \& \ \text{calls}^*(y, z)$
- Result:
 - set of ground atoms inferred by applying the rules until no new inferences can be made

Datalog vs. SQL

- SQL
 - Imperative programming:
 - join, union, projection, selection
 - Explicit iteration
- Datalog: logical database language
 - Declarative programming
 - Recursive definition: fixpoint computation
 - Negation can lead to oscillation
 - Stratified: only negate one “stratum” at a time

2. Flow-insensitive Points-to Analysis

- Alias analysis:
 - Can two pointers point to the same location?
 - `*a, *(a+8)`
- Points-to analysis:
 - What objects does each pointer points to?
 - Two pointers cannot be aliased if they must point to different objects

How to Name Objects?

- Objects are dynamically allocated
- Use finite names to refer to unbounded # objects
- 1 scheme: Name an object by its allocation site

```
main () {                f () {  
    p = f();              A: a = new O ();  
    q = f();              B: b = new O ();  
}                          return a;  
                          }
```

Points-To Analysis for Java

- Variables ($v \in V$): local variables in the program
- Heap-allocated objects ($h \in H$)
 - has a set of fields ($f \in F$)
 - named by allocation site

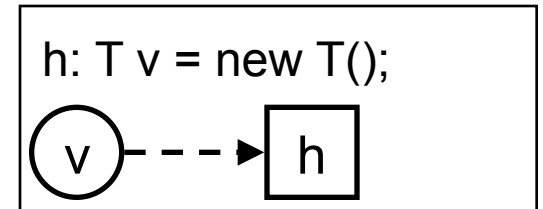
Program Abstraction

- Allocations $h: v = \text{new } c$
- Store $v_1.f = v_2$
- Loads $v_2 = v_1.f$
- Moves, arguments: $v_1 = v_2$
- Assume: a (conservative) call graph is known a priori
 - Call: formal = actual
 - Return: actual = return value

Pointer Analysis Rules

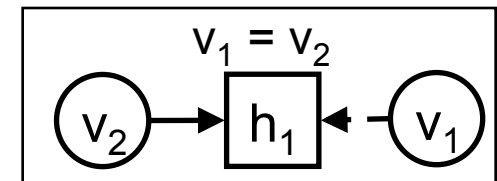
Object creation

$\text{pts}(v, h) \text{ :- "h: T v = new T()".}$



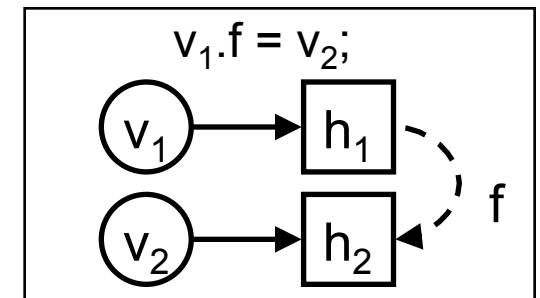
Assignment

$\text{pts}(v_1, h_1) \text{ :- "v}_1 = v_2" \ \& \ \text{pts}(v_2, h_1).$



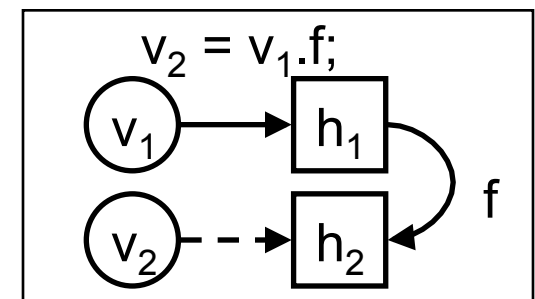
Store

$\text{hpts}(h_1, f, h_2) \text{ :- "v}_1.f = v_2" \ \& \ \text{pts}(v_1, h_1) \ \& \ \text{pts}(v_2, h_2).$



Load

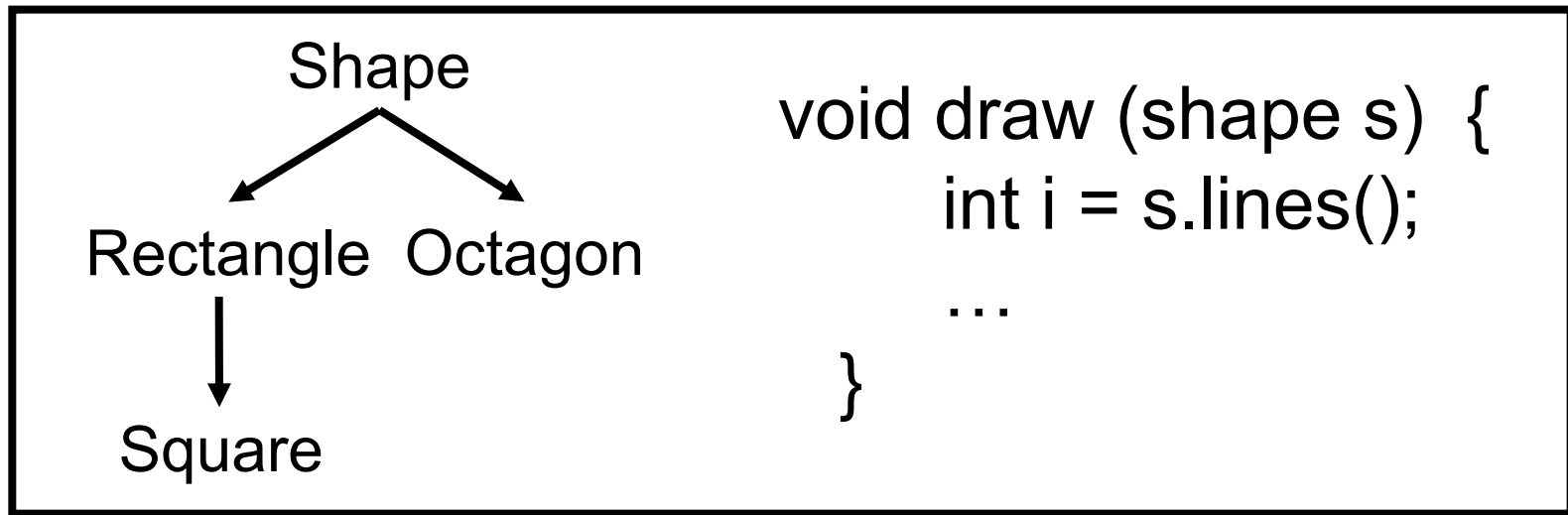
$\text{pts}(v_2, h_2) \text{ :- "v}_2 = v_1.f" \ \& \ \text{pts}(v_1, h_1) \ \& \ \text{hpts}(h_1, f, h_2).$



Pointer Alias Analysis

- Specified by a few Datalog rules
 - Creation sites
 - Assignments
 - Stores
 - Loads
- Apply rules until they converge

Virtual Method Invocation



- Class hierarchy analysis $\text{cha}(t, n, m)$
 - Given an invocation $v.n(\dots)$,
if v points to object of type t ,
then m is the method invoked
 - t 's first superclass that defines n

Pointer Analysis

Can Improve Call Graphs

Discover points-to results and methods invoked on the fly

$\text{hType}(h, t)$: h has type t

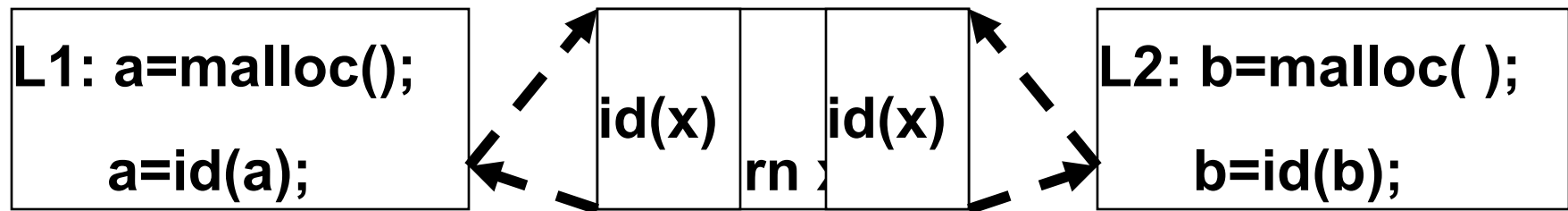
$\text{actual}(s, i, v)$: v is the i th actual parameter in call site s .

$\text{formal}(m, i, v)$: v is the i th formal parameter declared in method m .

$\text{invokes}(s, m) \text{ :- "s: v.n (...)" \& pts(v, h) \&}$
 $\text{hType}(h, t) \& \text{cha}(t, n, m)$

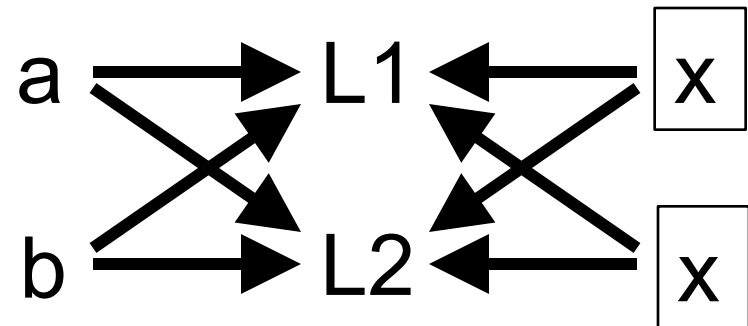
$\text{pts}(v, h) \text{ :- invokes}(s, m) \&$
 $\text{formal}(m, i, v) \& \text{actual}(s, i, w) \&$
 $\text{pts}(w, h)$

3. Context-Sensitive Pointer Analysis

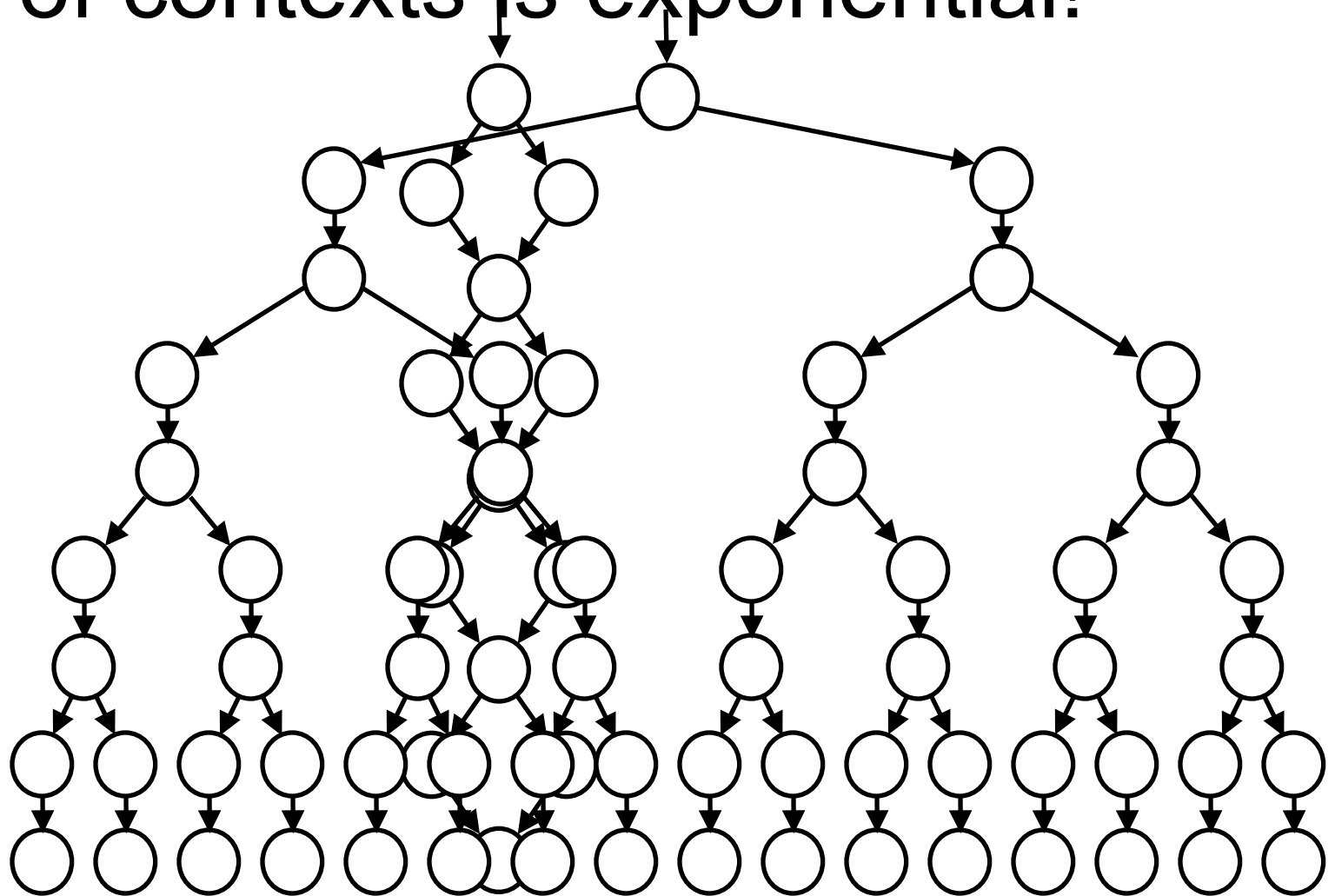


context-sensitive

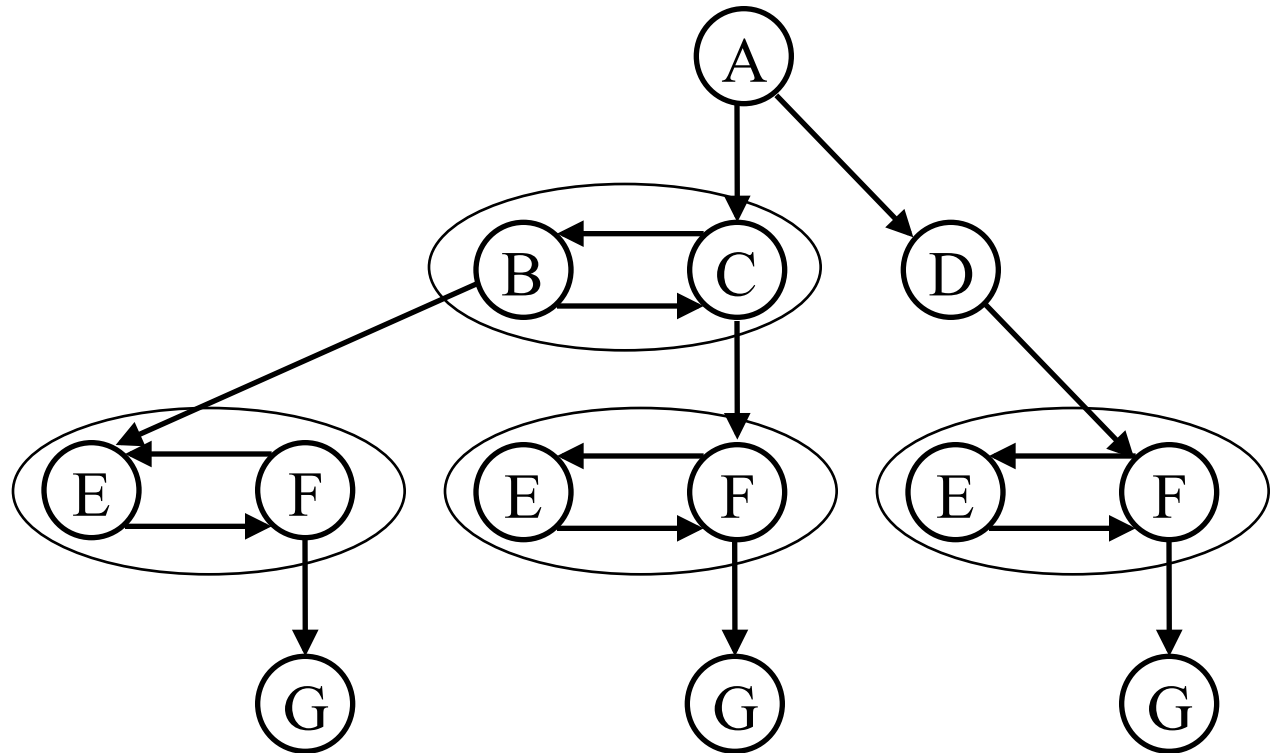
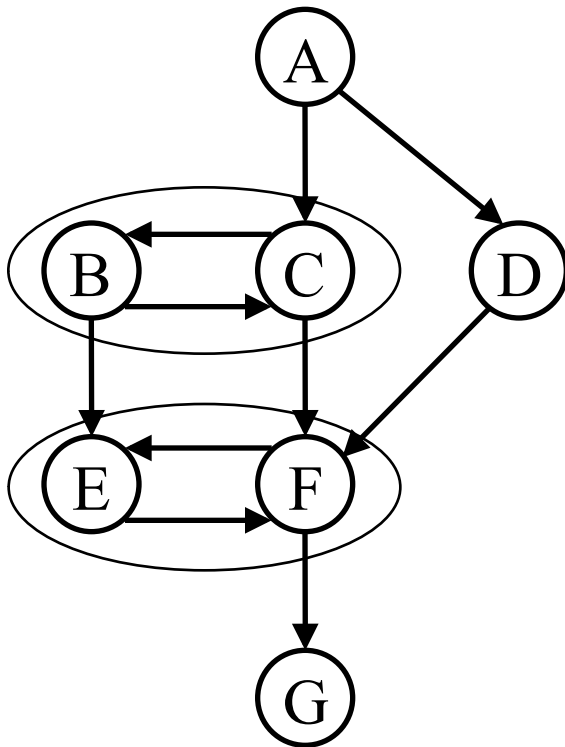
context-insensitive



Even without recursion,
of contexts is exponential!

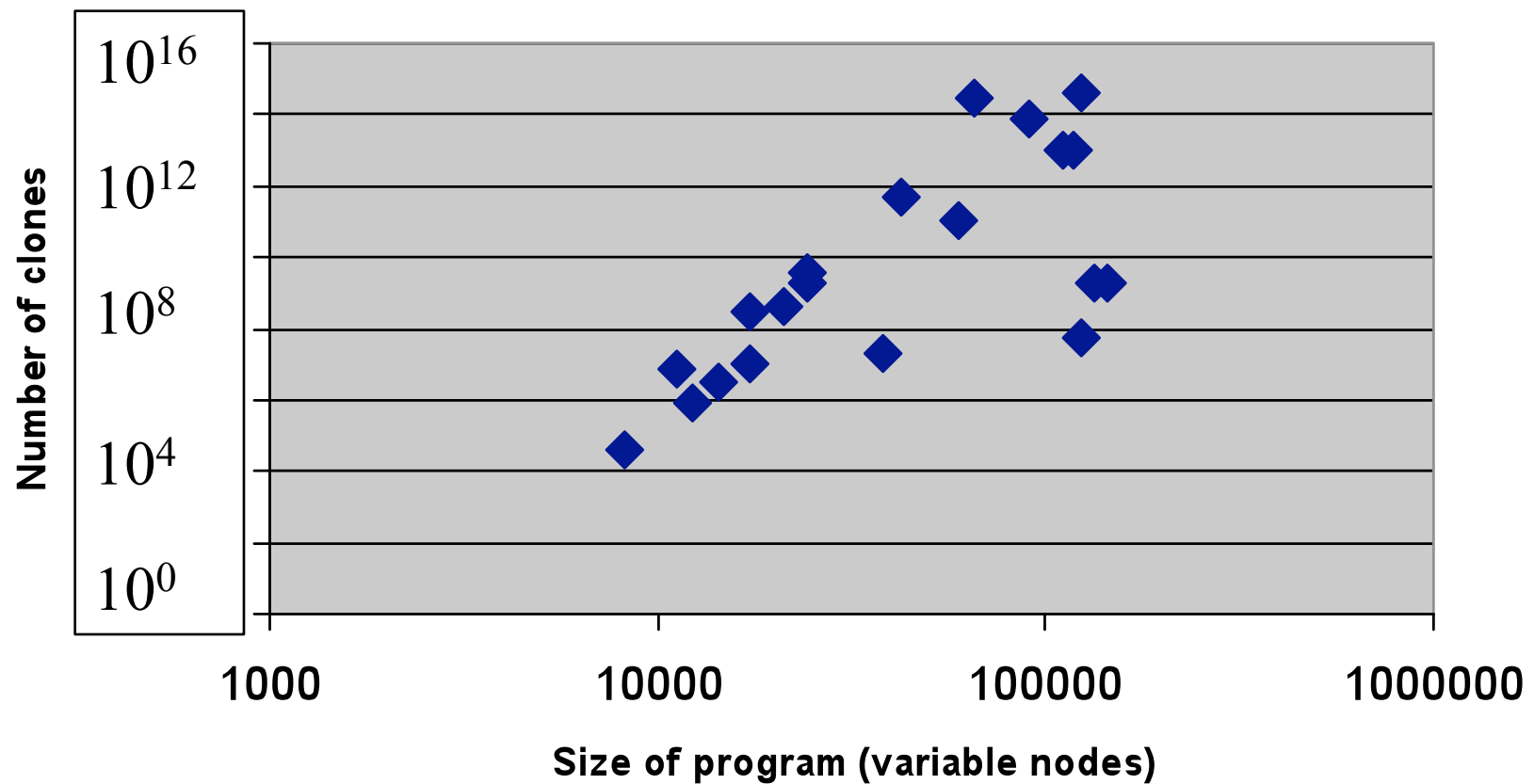


Recursion



Top 20 Sourceforge Java Apps

Number of Clones



Cloning-Based Algorithm

- Whaley&Lam, PLDI 2004 (best paper award)
- Apply the context-insensitive algorithm to the program to discover the call graph
- Find strongly connected components
- Create a “clone” for every context
- Apply the context-insensitive algorithm to cloned call graph
- Lots of redundancy in result
- Exploit redundancy by clever use of BDDs (binary decision diagrams)

Automatic Analysis Generation



Programmer:
Security analysis
in 10 lines

Compiler Writer:
Ptr analysis in 10 lines

1000s of lines
1 year tuning

PQL

Datalog

BDD operations

bddbddb
(**BDD**-based
deductive database)
with
Active Machine Learning

BDD: 10,000s-lines library