



Heuristic Decisions in Static Analysis

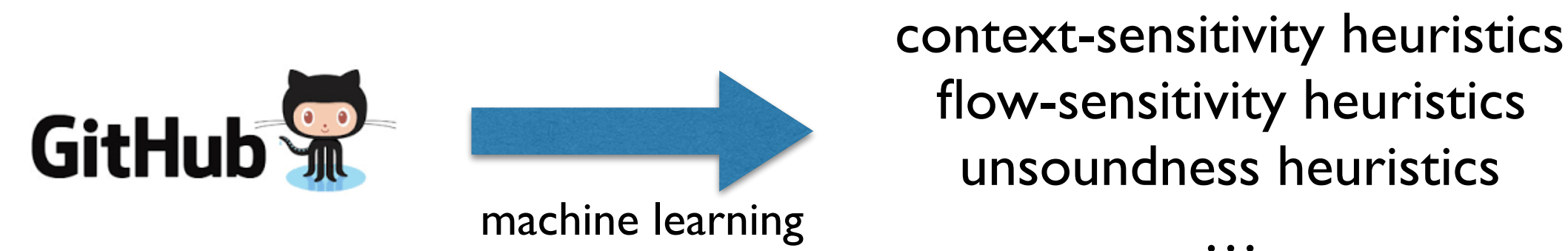


Astrée DOOP TAJS SAFE

- Practical static analyzers involve many heuristics
 - Which procedures should be analyzed context-sensitively?
 - Which relationships between variables should be tracked?
 - Which program parts to analyze unsoundly or soundly?, etc.
- Designing a good heuristic is an art
 - Usually done by trials and error: **nontrivial** and **suboptimal**

Automatically Generating Heuristics from Data

- Automate the process: use data to make heuristic decisions in static analysis
- Automatic:** little reliance on analysis designers



- Powerful:** machine-tuning outperforms hand-tuning
- Stable:** can be generated for arbitrary programs

Selective Context-Sensitivity

```

1 class D{} class E{}
2 class C{
3   void dummy(){}
4   Object id1(Object v){return id2(v);} //4
5   Object id2(Object v){return v;}
6 class B{
7   void m(){
8     C c = new C();
9     D d = (D)c.id1(new D()); //query1 //9
10    E e = (E)c.id1(new E()); //query2 //10
11    c.dummy(); //11
12 public class A{
13   public static void main(String[] args){
14     B b = new B();
15     b.m(); //15
16     b.m(); //16
  
```

- Context-insensitivity fails to prove the queries
- 2-call-site-sensitivity succeeds but not scale

solution

Apply 2-call-sens: {C.id2}
 Apply 1-call-sens: {C.id1}
 Apply insens: {B.m, C.dummy}

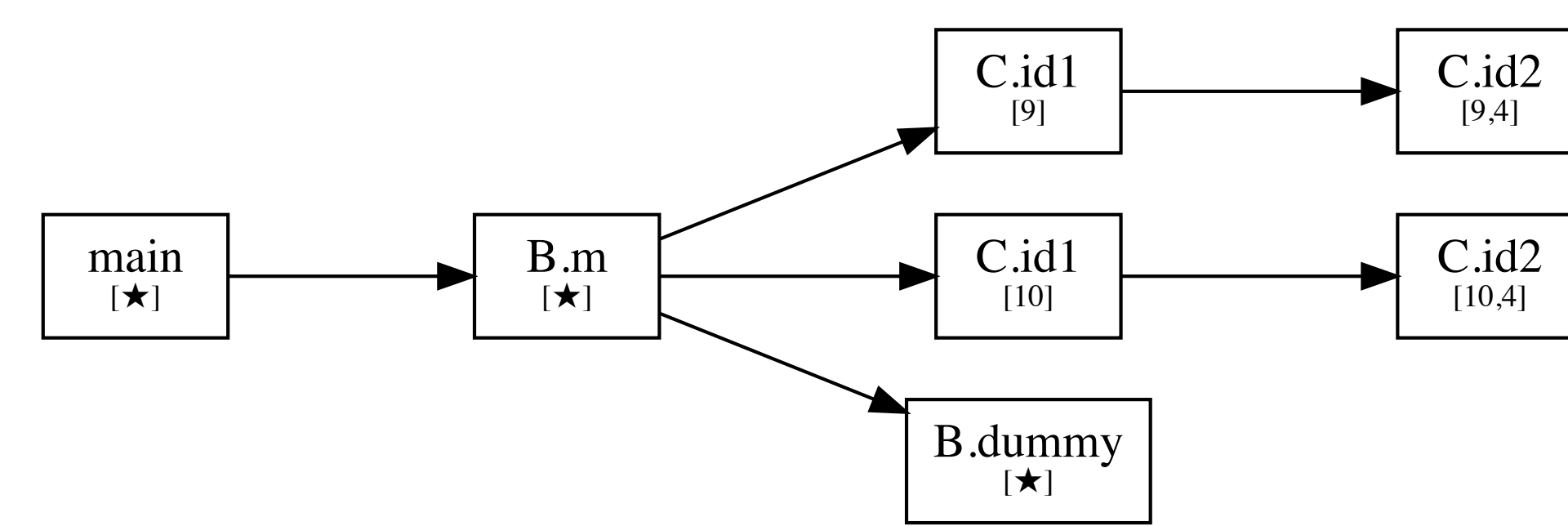


Figure 1: call graph of the solution

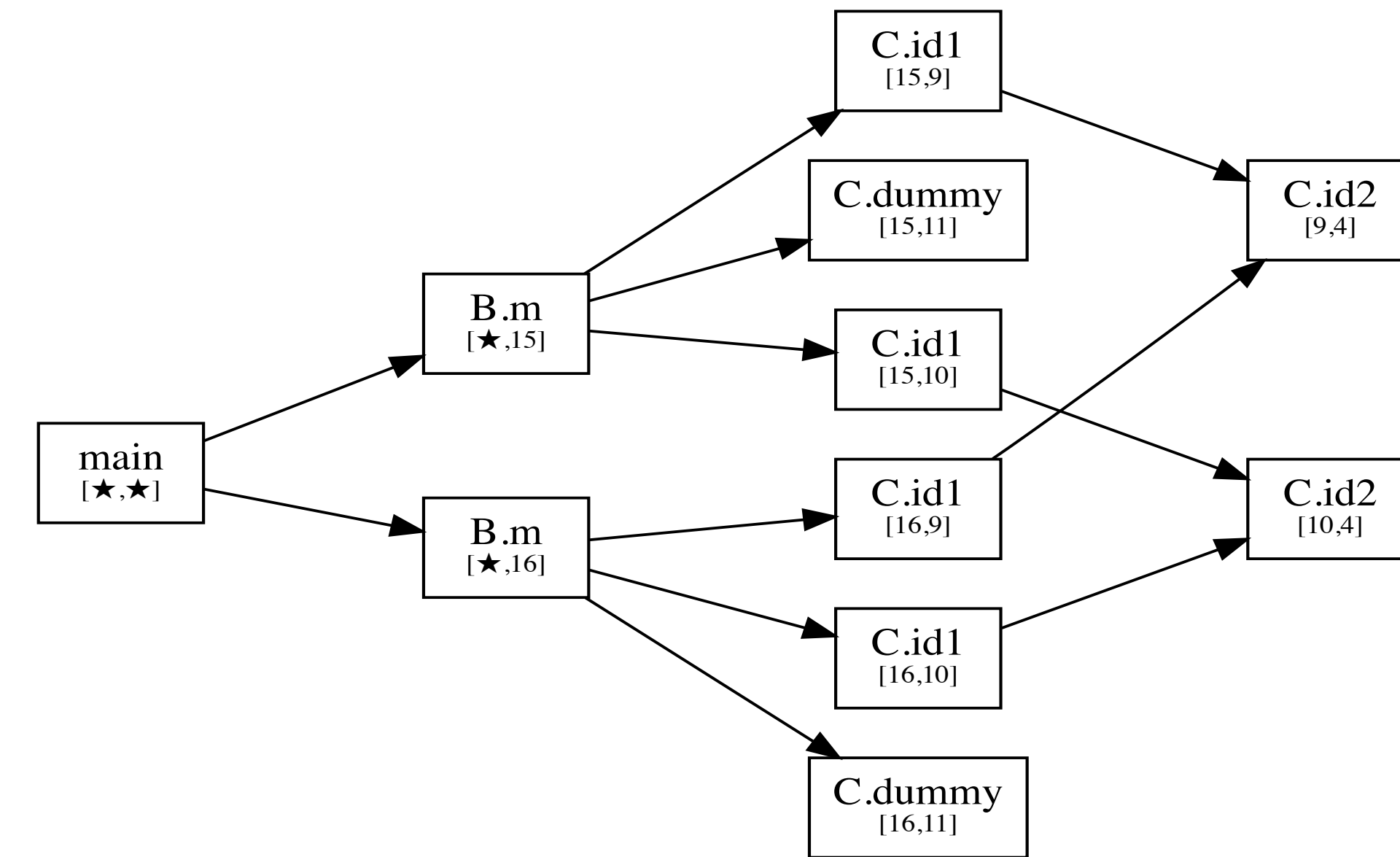
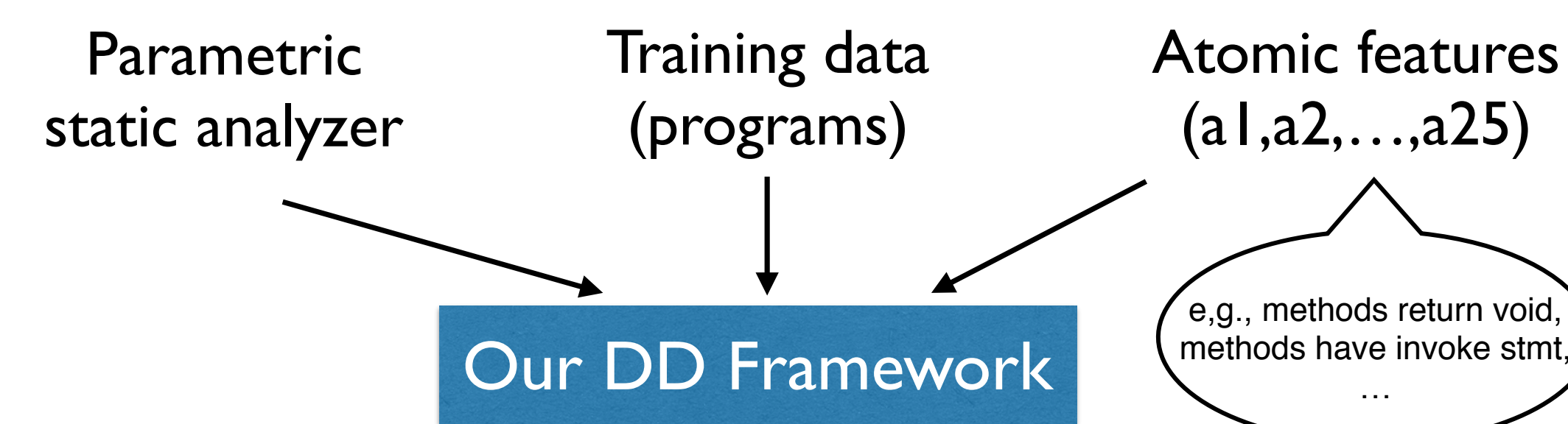


Figure 2: call graph of 2-call-site sensitive

Challenge: How to decide?

⇒ Data-Driven approach

Data-Driven Ctx-Sensitivity



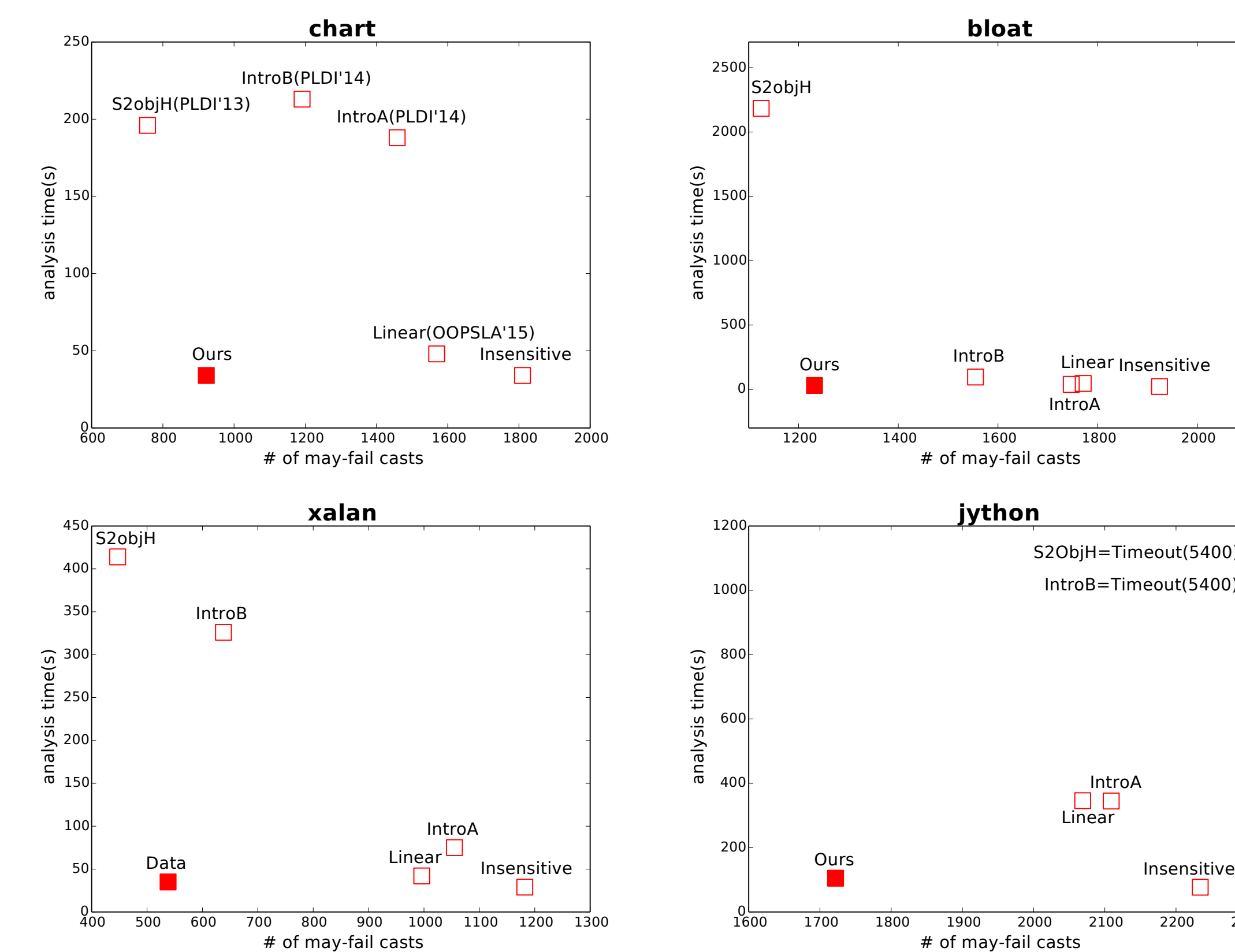
Heuristic for applying (hybrid) object-sensitivity:

```

f2: Methods that require 2-object-sensitivity
1 A ~3 A ~6 A 8 A ~9 A ~16 A ~17 A ~18 A ~19 A ~20 A ~21 A ~22 A ~23 A ~24 A ~25
f1: Methods that require 1-object-sensitivity
(1 A ~3 A ~4 A ~7 A ~8 A 6 A ~9 A ~15 A ~16 A ~17 A ~18 A ~19 A ~20 A ~21 A ~22 A ~23 A ~24 A ~25)V
(~3 A ~4 A ~7 A ~8 A ~9 A 10 A 11 A 12 A 13 A ~16 A ~17 A ~18 A ~19 A ~20 A ~21 A ~22 A ~23 A ~24 A ~25)V
(~3 A ~9 A 13 A 14 A 15 A ~16 A ~17 A ~18 A ~19 A ~20 A ~21 A ~22 A ~23 A ~24 A ~25)V
(1 A 2 A ~3 A 4 A ~5 A ~6 A ~7 A ~8 A ~9 A ~10 A ~13 A ~15 A ~16 A ~17 A ~18 A ~19 A ~20 A ~21 A ~22
A ~23 A ~24 A ~25)
  
```

Performance

- Training with 4 small programs from DaCapo, and applied to 6 large programs
- Machine-tuning outperforms hand-tuning



Key Contributions

We achieve the improvement with two key ideas.

- A new expressive model (Disjunctive Model)
- Learning algorithm for new model

Disjunctive Model

Disjunctive model expresses set with DNF form.

Method : Features Goal = {M₁, M₄}
 M₁ : a₁a₂ Disjunctive Model(Possible):
 M₂ : a₁ (a₁ ∧ a₂) ∨ (¬a₁ ∧ ¬a₂)
 M₃ : a₂ Linear Model(Impossible):
 M₄ : C₁ * a₁ + C₂ * a₂

Figure 3: Disjunctive vs Linear

With {a₁, a₂}, Disjunctive model can express the target methods, but Linear model cannot.

Learning Algorithm

Let $\Pi = \{f_1, \dots, f_k\}$ be parameters. Each f_i expresses methods to be assigned with depth i . We assign deeper depth if a method is in both f_i and $f_j (i \neq j)$. We learn Π by solving the following problem.

Optimization problem

Find parameter $\Pi = \langle f_1, \dots, f_k \rangle$ that minimizes the cost of analysis while satisfying precision constraint over training set.

Challenge

Assuming that $|S|$ is the space of possible boolean formulas over which we learn, search space of original problem is $|S|^k$. We reduce the search space into $k * |S|$ by decomposing the original problem into k subproblems ($\Psi_1 \sim \Psi_k$). Each f_i is obtained from Ψ_i and we solve them from Ψ_k to Ψ_1 .

Decomposed problem Ψ_i

Let $\Pi = \langle true, true, \dots, true, f_i, f_{i+1}, \dots, f_k \rangle$. Find formula f_i that makes Π minimize the cost while satisfying precision constraint over training set.

Learning Algorithm for Ψ_i

To solve Ψ_i , we made a greedy algorithm. Let $\{a_1, \dots, a_n\}$ be atomic features. Our algorithm proceeds in the following steps:

- f_i starts from disjunctions of $2n$ clauses :
 $f_i = a_1 \vee \neg a_1 \vee \dots \vee a_n \vee \neg a_n$
- Choose the most expensive clause c_i to refine.
- Strengthen the clause c_i by conjoining an decent atom a_k with c_i : $f'_i = c_1 \vee \dots \vee (c_i \wedge a_k) \vee \dots \vee c_j$.
- Check if f'_i satisfies precision constraint. If it is, $f_i = f'_i$.
- Repeat 2~4 until f_i cannot be refined.

* co-first authors