



Data-Driven Context-Sensitivity for Points-to Analysis

Sehun Jeong, Minseok Jeon,
Sungdeok Cha, and Hakjoo Oh

Korea University

OOPSLA 2017

Heuristics in Static Analysis



Astrée **DOOP** **TAJS** **SAFE**

- Modern static analyzers use many heuristics
 - Partial context-sensitivity
 - Partial octagon analysis
 - Partial flow-sensitivity
 - Partial path-sensitivity
 - Partial unsoundness
- Designing such heuristics is an art

Automatic Heuristic Generation



Context-sensitivity heuristics
Flow-sensitivity heuristics
Path-sensitivity heuristics
...

- **Automatic:** Learning algorithms infer heuristics from data.
- **Powerful:** Machine-tuned heuristics outperform hand-tuned ones.
- **Stable:** Learned heuristics perform well across training data.

Context-sensitivity

```
1 class D {} class E {}
2 class C {
3     void dummy(){}
4     Object id1(Object v){ return id2(v); }
5     Object id2(Object v){ return v; }
6 }
7 class B {
8     void m (){
9         C c = new C();
10        D d = (D)c.id1(new D()); //Query 1
11        E e = (E)c.id1(new E()); //Query 2
12        c.dummy();
13    }
14 }
15 public class A {
16     public static void main(String[] args){
17         B b = new B();
18         b.m();
19         b.m();
20     }
21 }
```

- Context-insensitivity fails to prove queries.
 - It merges `id1` invocations.
- 2-call-site sensitivity succeeds but not scale.
 - `m` and `dummy` are not worthy.

Selective 2-callsite-Sensitive Analysis

```
1 class D {} class E {}
2 class C {
3     void dummy(){}
4     Object id1(Object v){ return id2(v); }
5     Object id2(Object v){ return v; }
6 }
7 class B {
8     void m (){
9         C c = new C();
10        D d = (D)c.id1(new D()); //Query 1
11        E e = (E)c.id1(new E()); //Query 2
12        c.dummy();
13    }
14 }
15 public class A {
16     public static void main(String[] args){
17         B b = new B();
18         b.m();
19         b.m();
20     }
21 }
```

- Apply 2-callsite-sens: C.id2
- Apply 1-callsite-sens: C.id1
- Apply context-insens: main, B.m, C.dummy
- Proves all queries.

Selective 2-callsite-Sensitive Analysis

```
1 class D {} class E {}
2 class C {
3     void dummy(){
4         Object id1(Object v){ return id2(v); }
5         Object id2(Object v){ return v; }
6     }
7 class B {
8     void m (){
9         C c = new C();
10        D d = (D)c.id1(new D()); //Query 1
11        E e = (E)c.id1(new E()); //Query 2
12        c.dummy();
13    }
14 }
15 public class
16     public stat
17     B b = new
18     b.m();
19     b.m();
20 }
21 }
```

- Apply 2-callsite-sens: C.id2
- Apply 1-callsite-sens: C.id1
- Apply context-insens: main, B.m, C.dummy

Who decide?

Data-driven approach

This talk

Data-Driven Context Sensitivity

Parametric
Static Analyzer

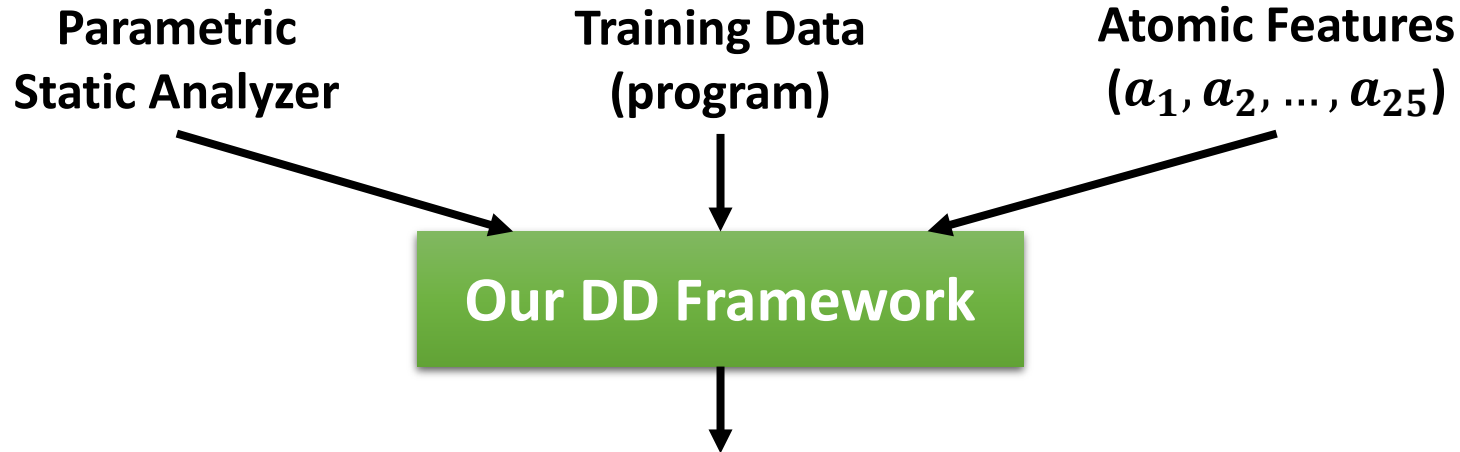
Training Data
(program)

Atomic Features
(a_1, a_2, \dots, a_{25})

Our DD Framework

Predicate over methods
e.g.) has return stmt?,
has voidreturn stmt?

Data-Driven Context Sensitivity

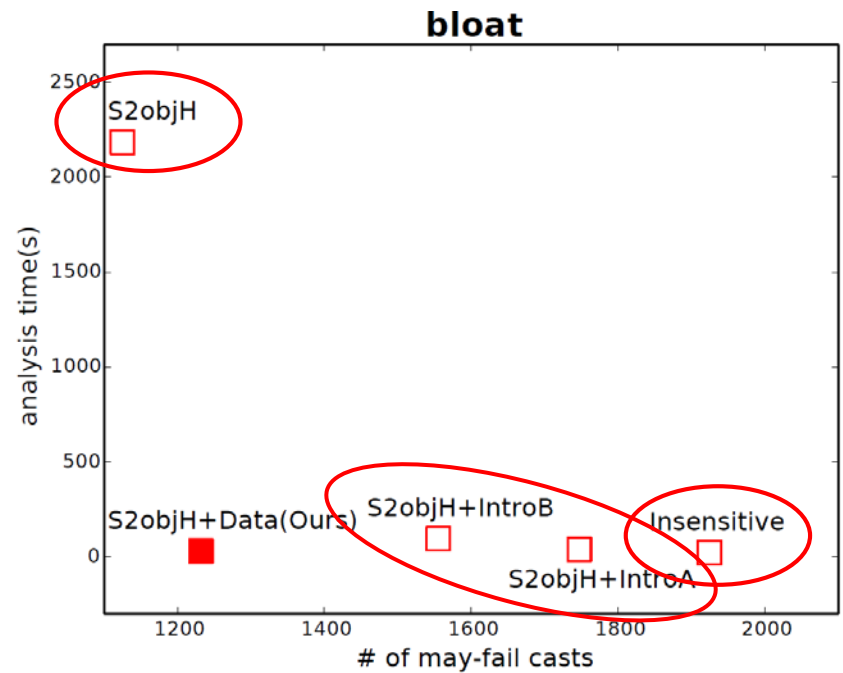
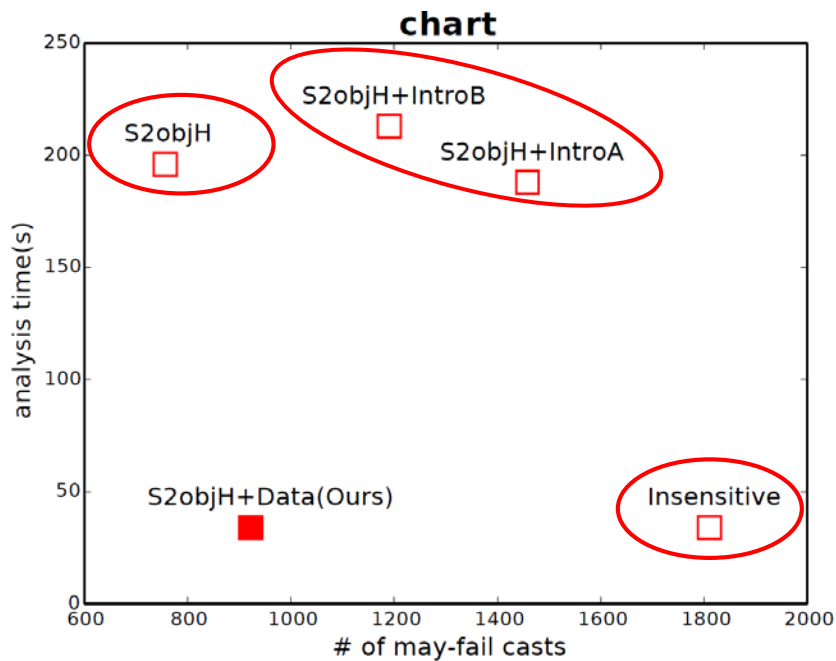


Heuristic for applying 2-hybrid-object-sensitivity

- **Methods that require 2-hybrid-object-sensitivity**
 $1 \wedge \neg 3 \wedge \neg 6 \wedge 8 \wedge \neg 9 \wedge \neg 16 \wedge \neg 17 \wedge \neg 18 \wedge \neg 19 \wedge \neg 20 \wedge \dots \wedge \neg 24 \wedge \neg 25$
- **Methods that require 1-hybrid-object-sensitivity**
 $(1 \wedge \neg 3 \wedge \neg 4 \wedge \neg 7 \wedge \neg 8 \wedge 6 \wedge \neg 9 \wedge \neg 15 \wedge \neg 16 \wedge \neg 17 \wedge \dots \wedge \neg 24 \wedge \neg 25) \vee$
 $(\neg 3 \wedge \neg 4 \wedge \neg 7 \wedge \neg 8 \wedge \neg 9 \wedge 10 \wedge 11 \wedge 12 \wedge 13 \wedge \neg 16 \wedge \dots \wedge \neg 24 \wedge \neg 25) \vee$
 $(\neg 3 \wedge \neg 9 \wedge 13 \wedge 14 \wedge 15 \wedge \neg 16 \wedge \neg 17 \wedge \neg 18 \wedge \neg 19 \wedge \dots \wedge \neg 24 \wedge \neg 25) \vee$
 $(1 \wedge 2 \wedge \neg 3 \wedge 4 \wedge \neg 5 \wedge \neg 6 \wedge \neg 7 \wedge \neg 8 \wedge \neg 9 \wedge \neg 10 \wedge \neg 13 \wedge \dots \wedge \neg 24 \wedge \neg 25)$

Performance Highlight

- Well-trained disjunctive heuristics tops hand-tuning.



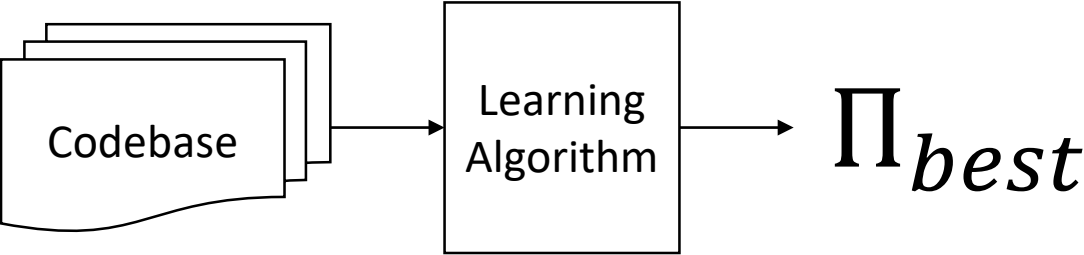
Details

Data-Driven Context Sensitivity

- 1. Parameterized disjunctive heuristics.

$$H_{\Pi}: M_P \rightarrow \{0, 1, \dots, k\}$$

- 2. Learning disjunctive heuristics for context-sensitivity.



Parametric Static Analyzer

$$F_P: A_P \rightarrow 2^Q \times Cost$$

- Context-sensitivity abstraction
- $a \in A_P: M_P \rightarrow \{0, 1, \dots, k\}$
 - k is the maximum context depth
 - M_P is a set of methods in program P

Parametric Static Analyzer

$$F_P(M_P \rightarrow \{0\})$$

- Context-insensitive analysis

Parametric Static Analyzer

$$F_P (M_P \rightarrow \{2\})$$

- Conventional 2-context-sensitive analysis
 - 2-object-sensitivity
 - 2-call-site-sensitivity
 - hybrid-2-object-sensitivity
 - 2-type sensitivity
 - ...

Parametric Static Analyzer

$$F_P(H_\Pi)$$

- Ours uses machine-tuned heuristic H .

Parametric Disjunctive Heuristics

$$H_{\Pi}: M_P \rightarrow \{0, 1, \dots, k\}$$

- Takes two steps:
 1. **Feature extraction**: Represent M_P as feature sets.
 2. **Decision making**: Assign context-depth to each methods.

Atomic Features

- Predicates over method definitions

$$\mathbb{A} = \{a_1, a_2, \dots, a_{25}\}$$

$$a_i: M_P \rightarrow \mathbb{B}$$

Atomic Features

“Does the method has a specific word in its signature string?”

Signature features

#1	“java”	#3	“sun”	#5	“void”	#7	“int”	#9	“String”
#2	“lang”	#4	“()”	#6	“security”	#8	“util”	#10	“init”

Statement features

#11	AssignStmt	#16	BreakpointStmt	#21	LookupStmt
#12	IdentityStmt	#17	EnterMonitorStmt	#22	NopStmt
#13	InvokeStmt	#18	ExitMonitorStmt	#23	RetStmt
#14	ReturnStmt	#19	GotoStmt	#24	ReturnVoidStmt
#15	ThrowStmt	#20	IfStmt	#25	TableSwitchStmt

“Does the method has a specific type of statement in its body?”

(1) Feature Extraction

$$a(m) = \{a_i \in \mathbb{A} \mid a_i(m) = \text{true}\}$$

$$a(\text{M1}) = \{a_1\}$$

$$a(\text{M2}) = \{a_2\}$$

$$a(\text{M3}) = \{a_3, a_4\}$$

- Represent methods into feature sets.

(2) Decision Making

Extracted Feature

M1 : { a_1 }
M2 : { a_2 }
M3 : { a_3, a_4 }

Learned Parameter

$\Pi = \langle f_2 = (a_1 \wedge \neg a_2), f_1 = a_1 \vee a_2 \rangle$



Initial Decision

$f_2: \{M1\}$
 $f_1: \{M2, M1\}$

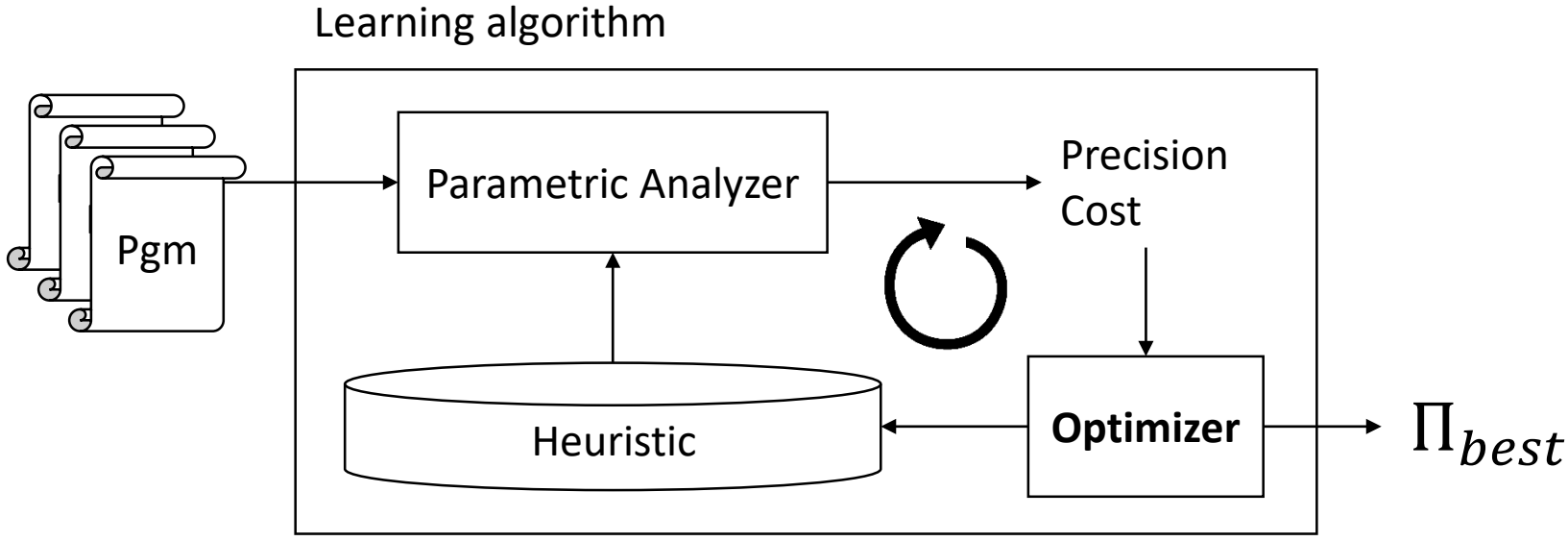
Context-sensitivity Abstraction

$H_\Pi: \{\{M1\} \mapsto 2, \{M2\} \mapsto 1, \{M3\} \mapsto 0\}$



Learning Disjunctive Heuristics

- Our learning process searches good heuristic in iterative fashion.



Optimization Problem

- $\Pi = \langle f_1, f_2, \dots, f_k \rangle$
- Find Π that minimizes

$$\sum_{P \in P_{gm}} \text{cost}(F_P(H_\Pi(P)))$$

- while satisfying

$$\frac{\sum_{P \in P_{gm}} |\text{proved}(F_P(H_\Pi(P)))|}{\sum_{P \in P_{gm}} |\text{proved}(F_P(k))|} \geq \gamma$$

Problem Decomposition Reduces Search Space

Goal

Find Π that minimizes cost
while satisfying target precision

Learning Π as a whole

$$\Pi = \langle f_1, f_2, \dots, f_{k-1}, f_k \rangle \rightarrow |S|^k$$

Set of all possible
Boolean formulas

Our approach: Learning one f at a time

$$\Pi = \langle true, \dots, true, f_k \rangle$$

Problem Decomposition Reduces Search Space

Goal

Find Π that minimizes cost
while satisfying target precision

Learning Π as a whole

Set of all possible
Boolean formulas

$$\Pi = \langle f_1, f_2, \dots, f_{k-1}, f_k \rangle \rightarrow |S|^k$$

Our approach: Learning one f at a time

$$\Pi = \langle true, \dots, true, f_{k-1}, f_k \rangle$$

Problem Decomposition Reduces Search Space

Goal

Find Π that minimizes cost
while satisfying target precision

Learning Π as a whole

Set of all possible
Boolean formulas

$$\Pi = \langle f_1, f_2, \dots, f_{k-1}, f_k \rangle \rightarrow |S|^k$$

Our approach: Learning one f at a time

$$\Pi = \langle true, f_2, \dots, f_{k-1}, f_k \rangle$$

Problem Decomposition Reduces Search Space

Goal

Find Π that minimizes cost
while satisfying target precision

Learning Π as a whole

Set of all possible
Boolean formulas

$$\Pi = \langle f_1, f_2, \dots, f_{k-1}, f_k \rangle \rightarrow |S|^k$$

Our approach: Learning one f at a time

$$\Pi = \langle f_1, f_2, \dots, f_{k-1}, f_k \rangle$$

Problem Decomposition Reduces Search Space

Goal

Find Π that minimizes cost while satisfying target precision

Learning Π as a whole

Set of all possible Boolean formulas

$$\Pi = \langle f_1, f_2, \dots, f_{k-1}, f_k \rangle \rightarrow |S|^k$$

Our approach: Learning one f at a time

$$\Pi = \langle f_1, f_2, \dots, f_{k-1}, f_k \rangle \rightarrow k \cdot |S|$$

Has Same Power? Yes.

- We have a theorem for it.
- Please consult with our paper.

Illustration of Algorithm for Searching Boolean Formula f

Initial formula

$$(a_1) \vee (\neg a_1) \vee (a_2) \vee (\neg a_2)$$

The most general formula

GOAL: Find f that proves all queries

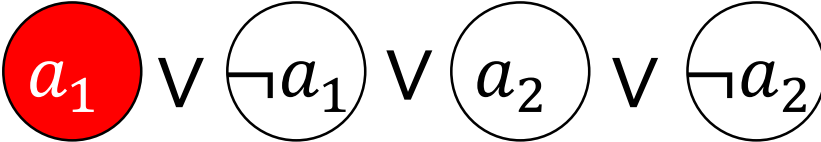
	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Performance Table

$$W = \{a_1, \neg a_1, a_2, \neg a_2\}$$

Refinement Targets

Iteration 1

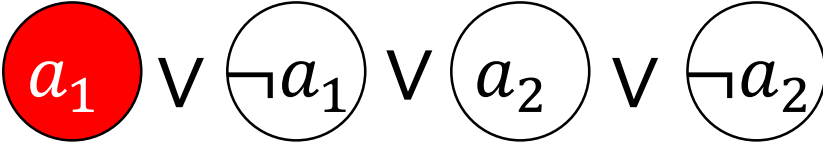


	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Pick the most expensive clause

$$W = \{a_1, \neg a_1, a_2, \neg a_2\}$$

Iteration 1

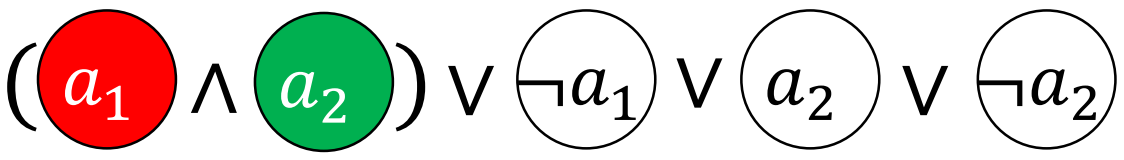


	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Remove the clause from the workset

$$W = \{\cancel{a_1}, \neg a_1, a_2, \neg a_2\}$$

Iteration 1

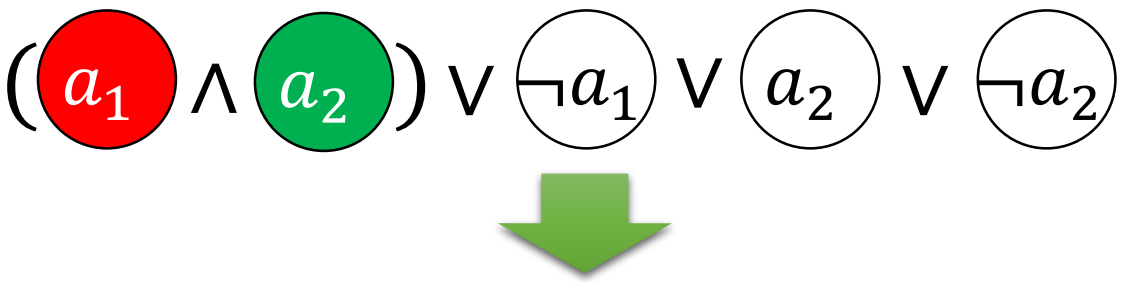


	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Refine the clause conservatively

$$W = \{\neg a_1, a_2, \neg a_2\}$$

Iteration 1



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

$$\frac{\sum_{P \in P_{gm}} |\text{proved}(F_P(H_{\Pi}(P)))|}{\sum_{P \in P_{gm}} |\text{proved}(F_P(k))|} = 1$$

Check whether the formula satisfies the precision goal

$$W = \{\neg a_1, a_2, \neg a_2\}$$

Iteration 1

$$(\textcircled{a_1} \wedge \textcircled{a_2}) \vee \textcircled{\neg a_1} \vee \textcircled{a_2} \vee \textcircled{\neg a_2}$$



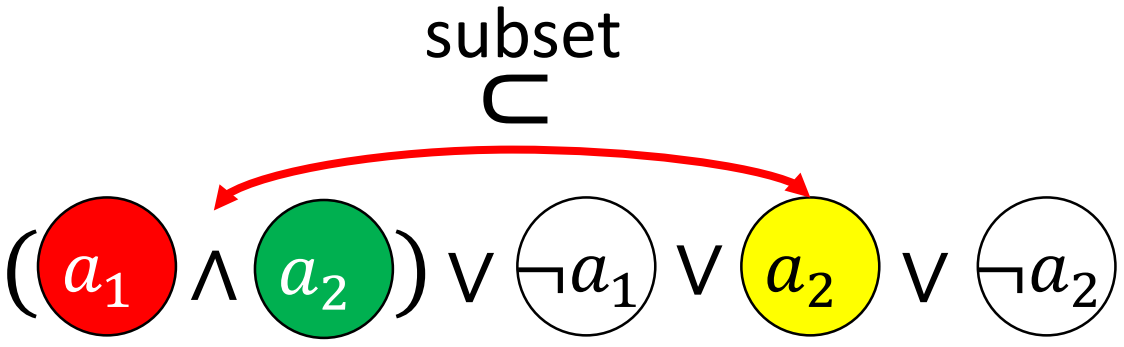
$$F_P(H_{\Pi}(P)) : \{Q1, Q2, Q3, Q4\}$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Current formula proves all queries

$$W = \{\neg a_1, a_2, \neg a_2\}$$

Iteration 1

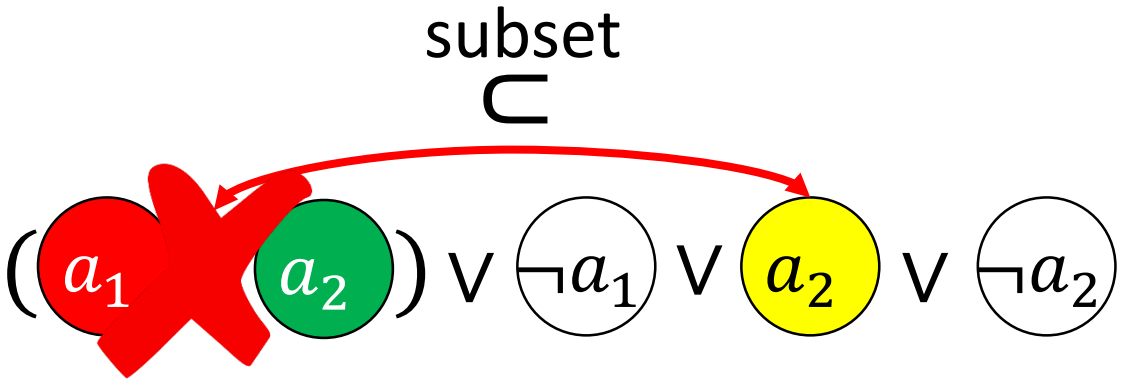


	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

However, the refinement didn't refine the formula

$$W = \{\neg a_1, a_2, \neg a_2\}$$

Iteration 1



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Drop the refined clause from the formula

$$W = \{\neg a_1, a_2, \neg a_2\}$$

Iteration 2



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

$$W = \{a_1, \neg a_1, a_2\}$$

Iteration 2



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Again, pick the most expensive clause

$$W = \{a_1, \neg a_1, a_2\}$$

Iteration 2



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Remove the clause from the workset

$$W = \{a_1, \neg a_1, \mathbf{a_2}\}$$

Iteration 2

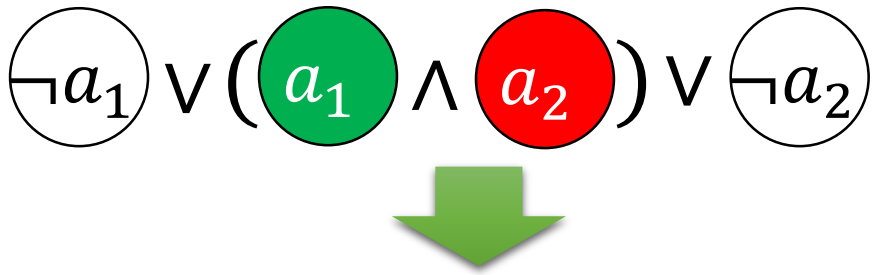
$$\neg a_1 \vee (a_1 \wedge a_2) \vee \neg a_2$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Refine the clause conservatively

$$W = \{a_1, \neg a_1\}$$

Iteration 2



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

$$\frac{\sum_{P \in P_{gm}} |\text{proved}(F_P(H_{\Pi}(P)))|}{\sum_{P \in P_{gm}} |\text{proved}(F_P(k))|} = 1$$

Check whether the formula satisfies the precision goal

$$W = \{a_1, \neg a_1\}$$

Iteration 2

$$\neg a_1 \vee (a_1 \wedge a_2) \vee \neg a_2$$



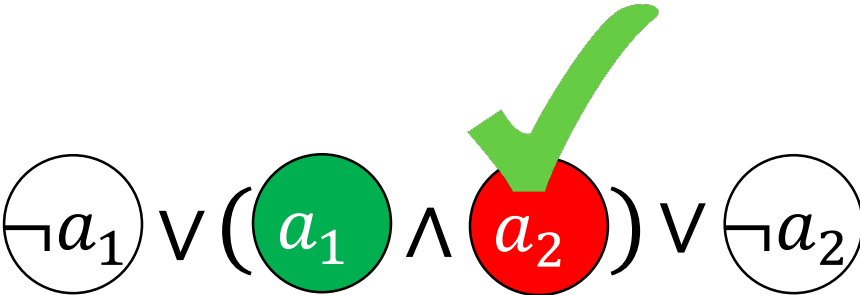
$$F_P(H_{\Pi}(P)) : \{Q1, Q2, Q3, Q4\}$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Still, the formula proves all queries

$$W = \{a_1, \neg a_1\}$$

Iteration 2

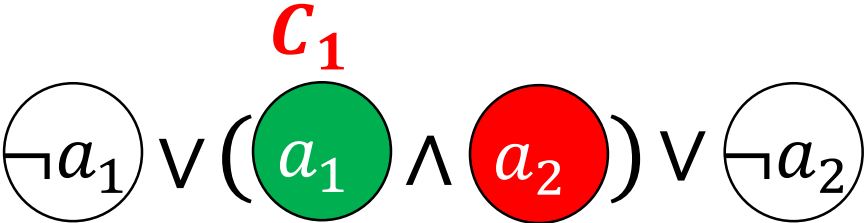


	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10

Refined clause passes subset checking

$$W = \{a_1, \neg a_1\}$$

Iteration 2

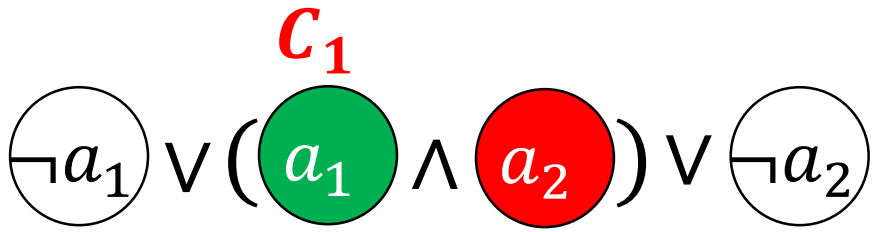


	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13

Record the refined clause C_1

$$W = \{a_1, \neg a_1\}$$

Iteration 2



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
<i>C</i> ₁	Q1, Q2, Q4	13

...and add *C*₁ to the workset for further refinement

$$W = \{a_1, \neg a_1, (a_1 \wedge a_2)\}$$

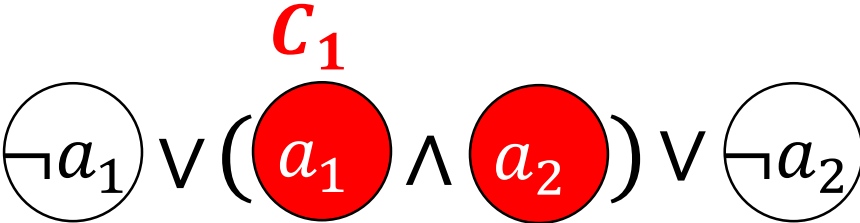
Iteration 3

$$\neg a_1 \vee \left(\overset{C_1}{a_1} \wedge a_2 \right) \vee \neg a_2$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13

$$W = \{a_1, \neg a_1, (a_1 \wedge a_2)\}$$

Iteration 3



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13

Pick the most expensive clause

$$W = \{a_1, \neg a_1, (a_1 \wedge a_2)\}$$

Iteration 3

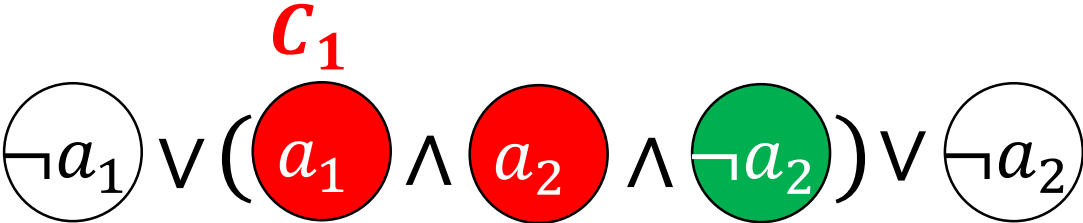
$$\neg a_1 \vee (\overset{C_1}{a_1} \wedge a_2) \vee \neg a_2$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13

Remove the clause from the workset

$$W = \{a_1, \neg a_1, \cancel{(a_1 \wedge a_2)}\}$$

Iteration 3

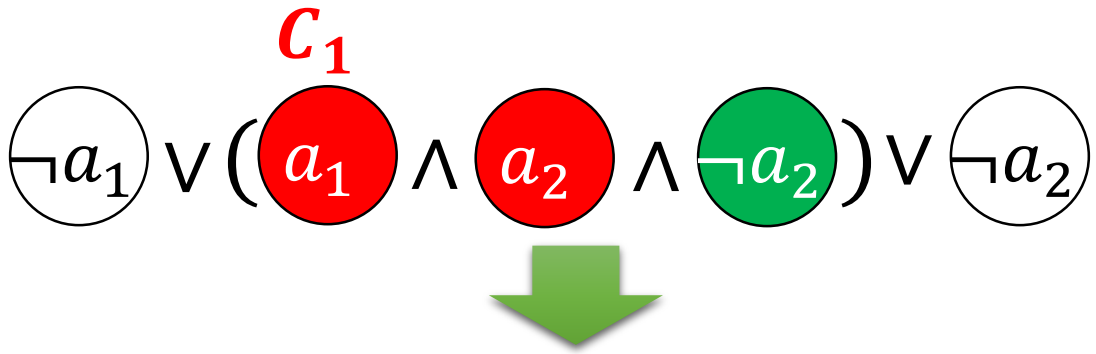


	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13

Refine the clause conservatively

$$W = \{a_1, \neg a_1\}$$

Iteration 3



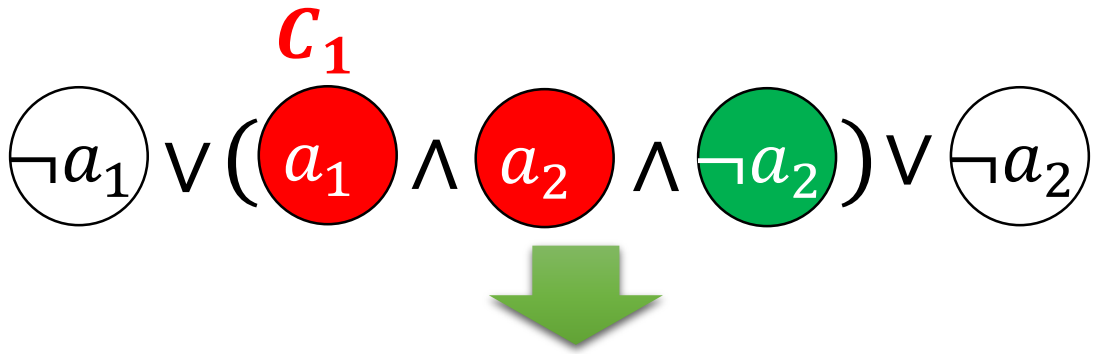
	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13

$$\frac{\sum_{P \in P_{gm}} |\text{proved}(F_P(H_{\Pi}(P)))|}{\sum_{P \in P_{gm}} |\text{proved}(F_P(k))|} = 1$$

Check whether the formula satisfies the precision goal

$$W = \{a_1, \neg a_1\}$$

Iteration 3



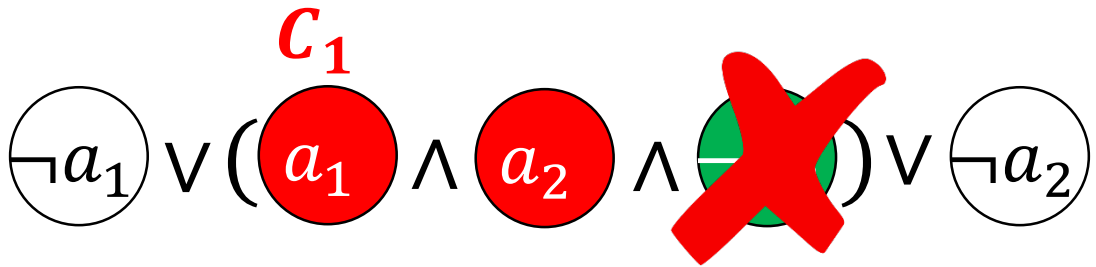
$F_P(H_{\Pi}(P)) : \{Q1, Q2, \cancel{Q3}, Q4\}$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13

This refinement failed to prove all queries

$W = \{a_1, \neg a_1\}$

Iteration 3



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13

Revert to the last state

$$W = \{a_1, \neg a_1\}$$

After iteration 4 and 5...

$$C_1 \left((a_1) \wedge (a_2) \right) \vee \left((\neg a_2) \wedge (\neg a_1) \right) C_2$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13
C_2	Q1, Q2, Q3	

$$W = \{(\neg a_2 \wedge \neg a_1)\}$$

Iteration 6

$$C_1 \quad C_2$$

$$\left((a_1) \wedge (a_2) \right) \vee \left(\neg a_2 \wedge \neg a_1 \right)$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13
C_2	Q1, Q2, Q3	7

Pick the most expensive clause in the workset

$$W = \{(\neg a_2 \wedge \neg a_1)\}$$

Iteration 6

$$C_1 \quad C_2$$

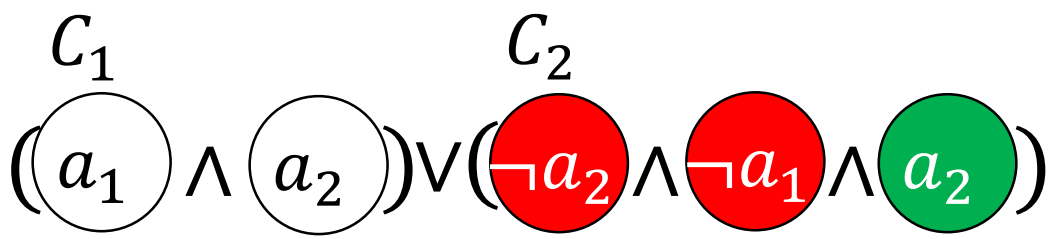
$$\left((a_1) \wedge (a_2) \right) \vee \left(\neg a_2 \wedge \neg a_1 \right)$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13
C_2	Q1, Q2, Q3	7

Remove the clause from the workset

$$W = \{ \langle \neg a_2 \wedge \neg a_1 \rangle \}$$

Iteration 6

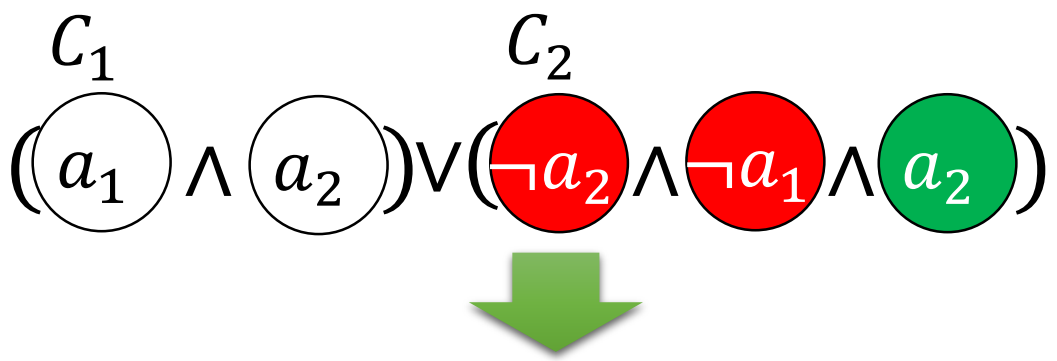


	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13
C_2	Q1, Q2, Q3	7

Refine the clause conservatively

$$W = \emptyset$$

Iteration 6



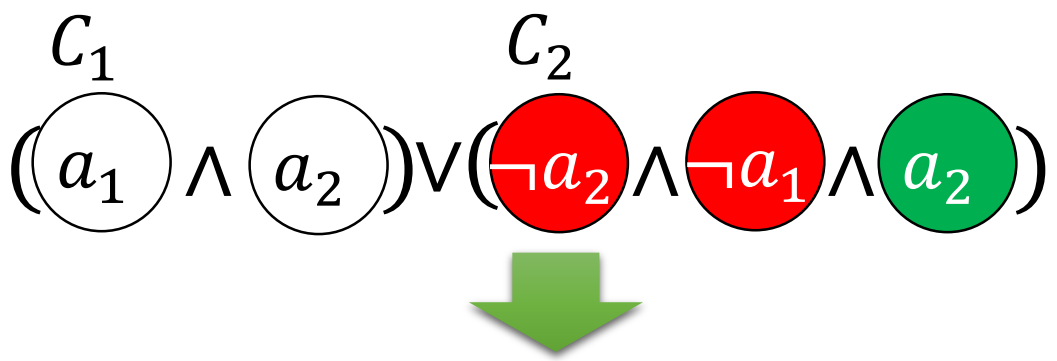
$$\frac{\sum_{P \in P_{gm}} |\text{proved}(F_P(H_{\Pi}(P)))|}{\sum_{P \in P_{gm}} |\text{proved}(F_P(k))|} = 1$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13
C_2	Q1, Q2, Q3	7

Check whether the formula satisfies the precision goal

$$W = \emptyset$$

Iteration 6



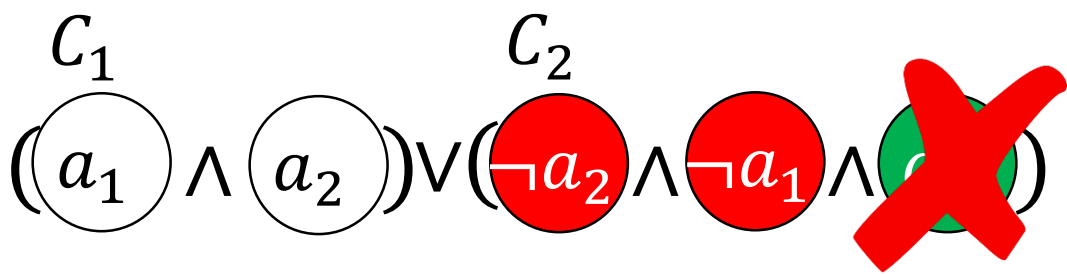
$F_P(H_{\Pi}(P)) : \{Q1, Q2, \del{Q3}, Q4\}$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13
C_2	Q1, Q2, Q3	7

This refinement failed to prove all queries

$W = \emptyset$

Iteration 6



	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13
C_2	Q1, Q2, Q3	7

Revert to the last state

$W = \emptyset$

END

$$C_1 \quad C_2$$

$$\left(a_1 \wedge a_2 \right) \vee \left(\neg a_2 \wedge \neg a_1 \right)$$

	Proven Qs	Cost
a_1	Q1, Q3	20
$\neg a_1$	Q2	7
a_2	Q1, Q2, Q3	15
$\neg a_2$	Q1, Q2, Q4	10
C_1	Q1, Q2, Q4	13
C_2	Q1, Q2, Q3	7

Algorithm ends when there is no refineable claus.

$$W = \emptyset$$

Experiments

Settings

- Data-Driven Doop
 - Our version of Doop points-to analysis framework
- Four context sensitivities, One human-tuned heuristic
 - Selective hybrid object, object, type, call-site
 - Introspective analysis
- DaCapo Benchmark suite
 - Four small programs for training
 - One large programs for validation
 - Five large programs for testing

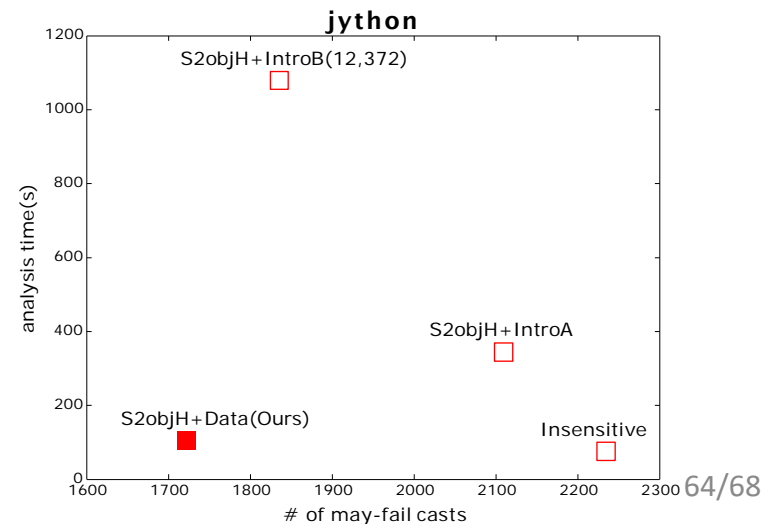
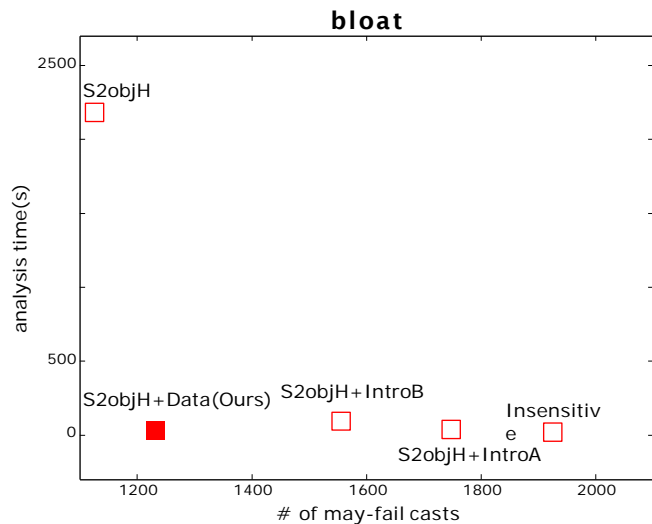
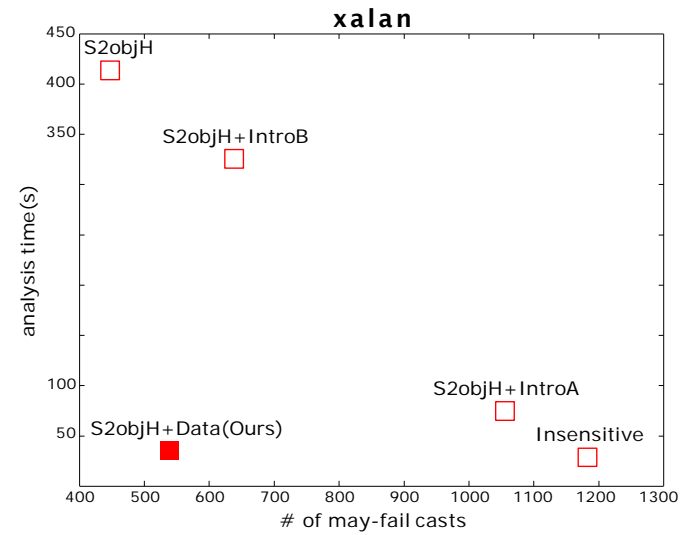
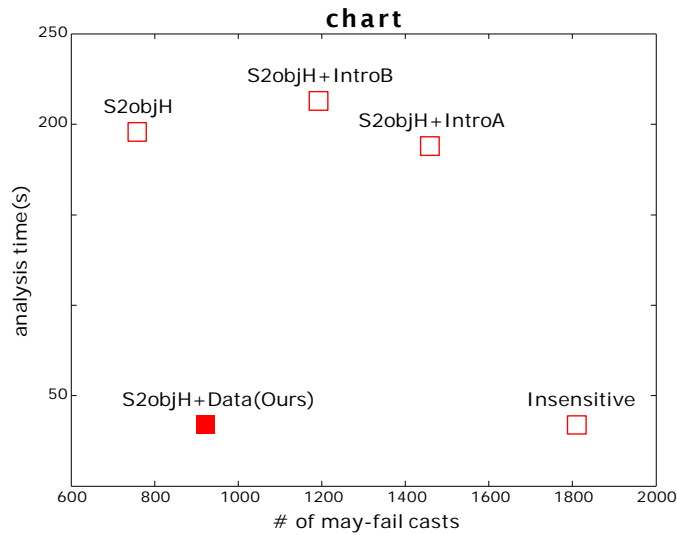
Research Questions

RQ1: Effectiveness

RQ2: Model adequacy

RQ3: Learned Heuristics

RQ1: Performance



RQ2: vs. Linear Model

- On same time budget for training, disjunctive model outperforms non-disjunctive model

Benchmarks	Non-disjunctive		Disjunctive(Ours)	
	may-fail casts	time(s)	may-fail casts	time(s)
eclipse	946	25	596	21
chart	1,569	48	937	33
bloat	1,771	46	1,232	27
xalan	996	42	539	33
jython	2069	346	1,738	104
Total	7,352	346	5,042	218

RQ3: Learned Features

- Object vs. Type
- Object sensitivity strictly more precise than Type.
- Our learning algorithm revealed the relationship.

$$\begin{array}{l} f_1 \text{ for } 2objH+Data : (1 \wedge 2 \wedge \neg 3 \wedge \neg 6 \wedge \neg 7 \wedge \neg 8 \wedge \neg 9 \wedge \neg 16 \wedge \dots \wedge \neg 22 \wedge \neg 23 \wedge \neg 24 \wedge \neg 25) \vee \\ \quad (\neg 1 \wedge \neg 2 \wedge 8 \wedge 5 \wedge \neg 9 \wedge 11 \wedge 12 \wedge \dots \wedge \neg 21 \wedge \neg 22 \wedge \neg 23 \wedge \neg 24 \wedge \neg 25) \vee \\ \quad (\neg 3 \wedge \neg 4 \wedge \neg 7 \wedge \neg 8 \wedge \neg 9 \wedge 10 \wedge 11 \wedge \dots \wedge \neg 21 \wedge \neg 22 \wedge \neg 23 \wedge \neg 24 \wedge \neg 25) \\ \hline f_1 \text{ for } 2typeH+Data : 1 \wedge 2 \wedge \neg 3 \wedge \neg 6 \wedge \neg 7 \wedge \neg 8 \wedge \neg 9 \wedge \neg 15 \wedge \neg 16 \wedge \dots \wedge \neg 22 \wedge \neg 23 \wedge \neg 24 \wedge \neg 25 \end{array}$$

RQ3: Learned Features

- Our algorithm revealed commonality among object-related context sensitivities.

$$f_2 \text{ for } S2objH+Data : 1 \wedge \neg 3 \wedge \neg 6 \wedge 8 \wedge \neg 9 \wedge \neg 16 \wedge \neg 17 \wedge \neg 18 \wedge \dots \wedge \neg 25$$

$$f_2 \text{ for } 2objH+Data : 1 \wedge \neg 3 \wedge \neg 6 \wedge 8 \wedge \neg 9 \wedge \neg 16 \wedge \neg 17 \wedge \neg 18 \wedge \dots \wedge \neg 25$$

$$f_2 \text{ for } 2typeH+Data : 1 \wedge \neg 3 \wedge \neg 6 \wedge 8 \wedge \neg 9 \wedge \neg 16 \wedge \neg 17 \wedge \neg 18 \wedge \dots \wedge \neg 25$$

$$f_2 \text{ for call-site-sensitivity} : 1 \wedge \neg 6 \wedge \neg 7 \wedge 11 \wedge 12 \wedge 13 \wedge \neg 16 \wedge \neg 17 \wedge \neg 18 \wedge \dots \wedge \neg 25$$

Conclusion

- Our approach uses a heuristic rule for context-sensitivity parameterized by Boolean formulas.
 - can express disjunctive properties of methods
- Good parameters for the heuristic can be learned from codebase by our learning algorithm
- The learned heuristics make context-sensitive points-to analysis more practical.