# Learning a Variable-Clustering Strategy for Octagon from Labeled Data Generated by a Static Analysis

Kihong Heo[1],  Hakjoo Oh[2], Hongseok Yang[3]

Seoul National University[1]
Korea University[2]
University of Oxford[3]
SAS 2016 @Edinburgh

# Long Term Goal

- Self-evolving static analysis by learning big data

  - data : similar codes, old versions, user-feedbacks, bug reports, test results, etc
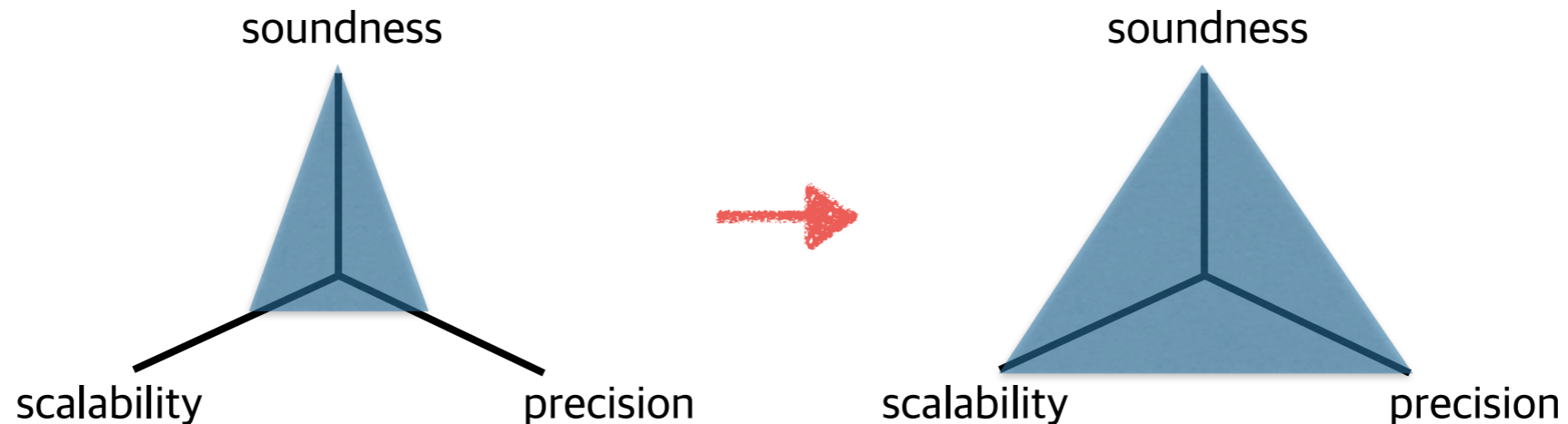
  - mature in other fields :   amazon …



Big Data + Static Analyzer

# Long Term Goal



$$F \in Pgm \times \underline{\Pi} \to \mathcal{A}$$

- Finding a good abstraction for adaptive static analysis

  - **Machine Learning** (learner) + **Static Analysis** (teacher)

  - e.g.) **relation**, context, flow, etc

# Relational Analysis

- Tracking relationships among variables

- e.g.) octagon analysis : $(\pm x) - (\pm y) \leq c$

```
1  int a = b;
2  int c = input();            // User input
3  for (i = 0; i < b; i++) {
4    assert (i < a);           // Query 1
5    assert (i < c);           // Query 2
6  }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | ∞ | ∞ | ∞ |
| b | ∞ | 0 | ∞ | ∞ |
| c | ∞ | ∞ | 0 | ∞ |
| i | ∞ | ∞ | ∞ | 0 |

{a, b, c, i}

*Consider x-y ≤ c only,
for simplicity

4

# Relational Analysis

- Tracking relationships among variables

- e.g.) octagon analysis : $(\pm x) - (\pm y) \leq c$

```
1   int a = b;
2   int c = input();              // User input
3   for (i = 0; i < b; i++) {
4     assert (i < a);             // Query 1
5     assert (i < c);             // Query 2
6   }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | 0 | ∞ | ∞ |
| b | 0 | 0 | ∞ | ∞ |
| c | ∞ | ∞ | 0 | ∞ |
| i | ∞ | ∞ | ∞ | 0 |

b - a ≤ 0

a - b ≤ 0

{a, b, c, i}

5

# Relational Analysis

- Tracking relationships among variables

- e.g.) octagon analysis : $(\pm x) - (\pm y) \leq c$

```
1   int a = b;
2   int c = input();              // User input
3   for (i = 0; i < b; i++) {
4     assert (i < a);             // Query 1
5     assert (i < c);             // Query 2
6   }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | 0 | ∞ | ∞ |  c - a ≤ ∞
| b | 0 | 0 | ∞ | ∞ |  c - b ≤ ∞
| c | ∞ | ∞ | 0 | ∞ |
| i | ∞ | ∞ | ∞ | 0 |

a - c ≤ ∞
b - c ≤ ∞

{a, b, c, i}

6

# Relational Analysis

- Tracking relationships among variables

- e.g.) octagon analysis : $(\pm x) - (\pm y) \leq c$

```
1   int a = b;
2   int c = input();              // User input
3   for (i = 0; i < b; i++) {
4     assert (i < a);             // Query 1
5     assert (i < c);             // Query 2
6   }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | 0 | ∞ | ∞ |
| b | 0 | 0 | ∞ | -1 |
| c | ∞ | ∞ | 0 | ∞ |
| i | ∞ | ∞ | ∞ | 0 |

i - b ≤ -1

{a, b, c, i}

# Relational Analysis

- Tracking relationships among variables

- e.g.) octagon analysis : $(\pm x) - (\pm y) \leq c$

```
1  int a = b;
2  int c = input();              // User input
3  for (i = 0; i < b; i++) {
4    assert (i < a);             // Query 1
5    assert (i < c);             // Query 2
6  }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | 0 | ∞ | -1 |
| b | 0 | 0 | ∞ | -1 |
| c | ∞ | ∞ | 0 | ∞ |
| i | ∞ | ∞ | ∞ | 0 |

i - a ≤ -1

{a, b, c, i}

8

# Relational Analysis

- Tracking relationships among variables

- e.g.) octagon analysis : $(\pm x) - (\pm y) \leq c$

```
1   int a = b;
2   int c = input();              // User input
3   for (i = 0; i < b; i++) {
4      assert (i < a);            // Query 1
5      assert (i < c);            // Query 2
6   }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | 0 | ∞ | -1 |
| b | 0 | 0 | ∞ | -1 |
| c | ∞ | ∞ | 0 | ∞ |
| i | ∞ | ∞ | ∞ | 0 |

i - c ≤ ∞

{a, b, c, i}

# Relational Analysis

- Tracking relationships among variables

- e.g.) octagon analysis : $(\pm x) - (\pm y) \leq c$

```
1  int a = b;
2  int c = input();              // User input
3  for (i = 0; i < b; i++) {
4    assert (i < a);             // Query 1
5    assert (i < c);             // Query 2
6  }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | 0 | ∞ | -1 |
| b | 0 | 0 | ∞ | -1 |
| c | ∞ | ∞ | 0 | ∞ |
| i | ∞ | ∞ | ∞ | 0 |

{a, b, ϲ, i}

Do we need c?

# Selective Relational Analysis

- **Selectively** tracking relationships among variables

  - within the same cluster

```
1   int a = b;
2   int c = input();              // User input
3   for (i = 0; i < b; i++) {
4     assert (i < a);            // Query 1
5     assert (i < c);            // Query 2
6   }
```

|   | a | b | i  |
|---|---|---|----|
| a | 0 | 0 | −1 |
| b | 0 | 0 | −1 |
| i | ∞ | ∞ | 0  |

$+$

$-\infty \leq c \leq +\infty$
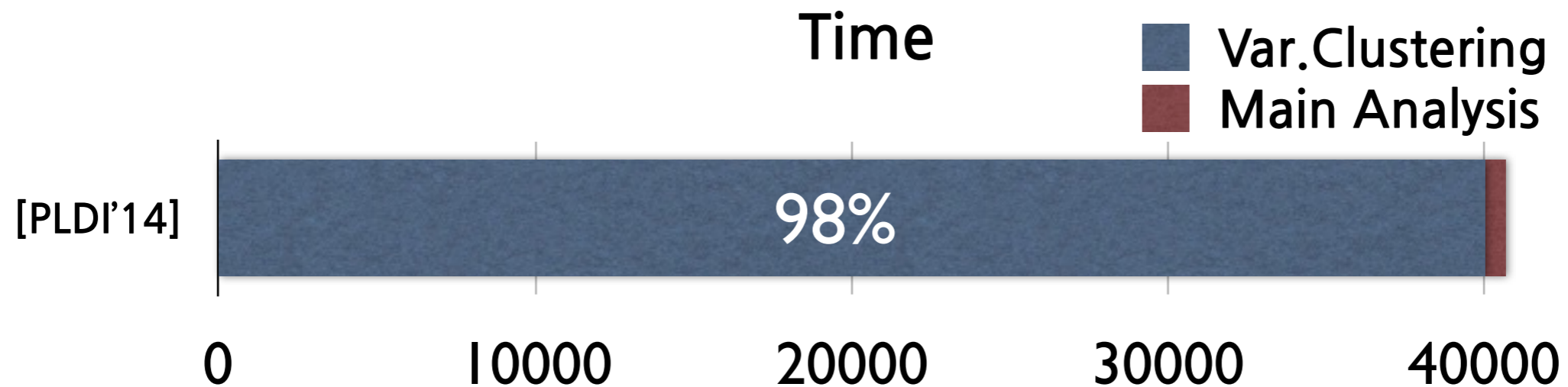
$\{a,b,i\}$        $\{c\}$

# Previous Solution

- Variable clustering by impact pre-analysis

  - estimating the impact of relationships

  - more scalable than the baseline Octagon analysis

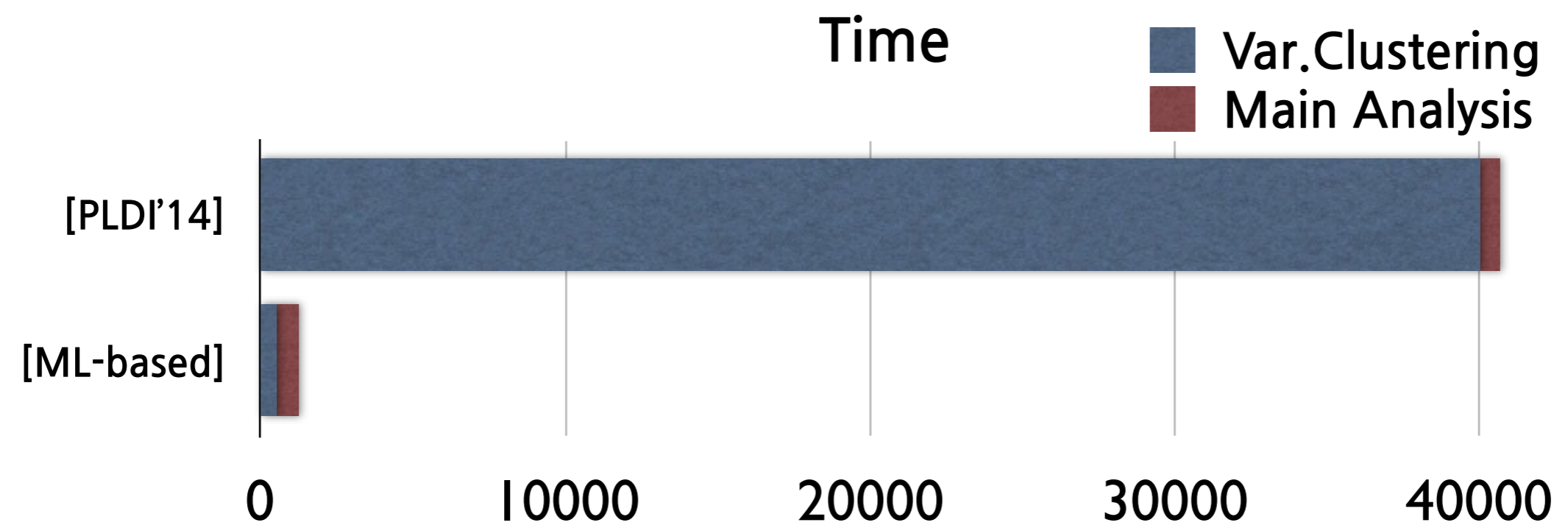  - more scalable & precise than other clustering methods

# Problem

- Variable clustering by impact pre-analysis

  - fully relational pre-analysis as an online estimator

  - e.g.) 17 open source benchmarks (~100KLOC)

Time

■ Var.Clustering
■ Main Analysis

[PLDI'14]  98%

0    10000   20000   30000   40000

# New Solution

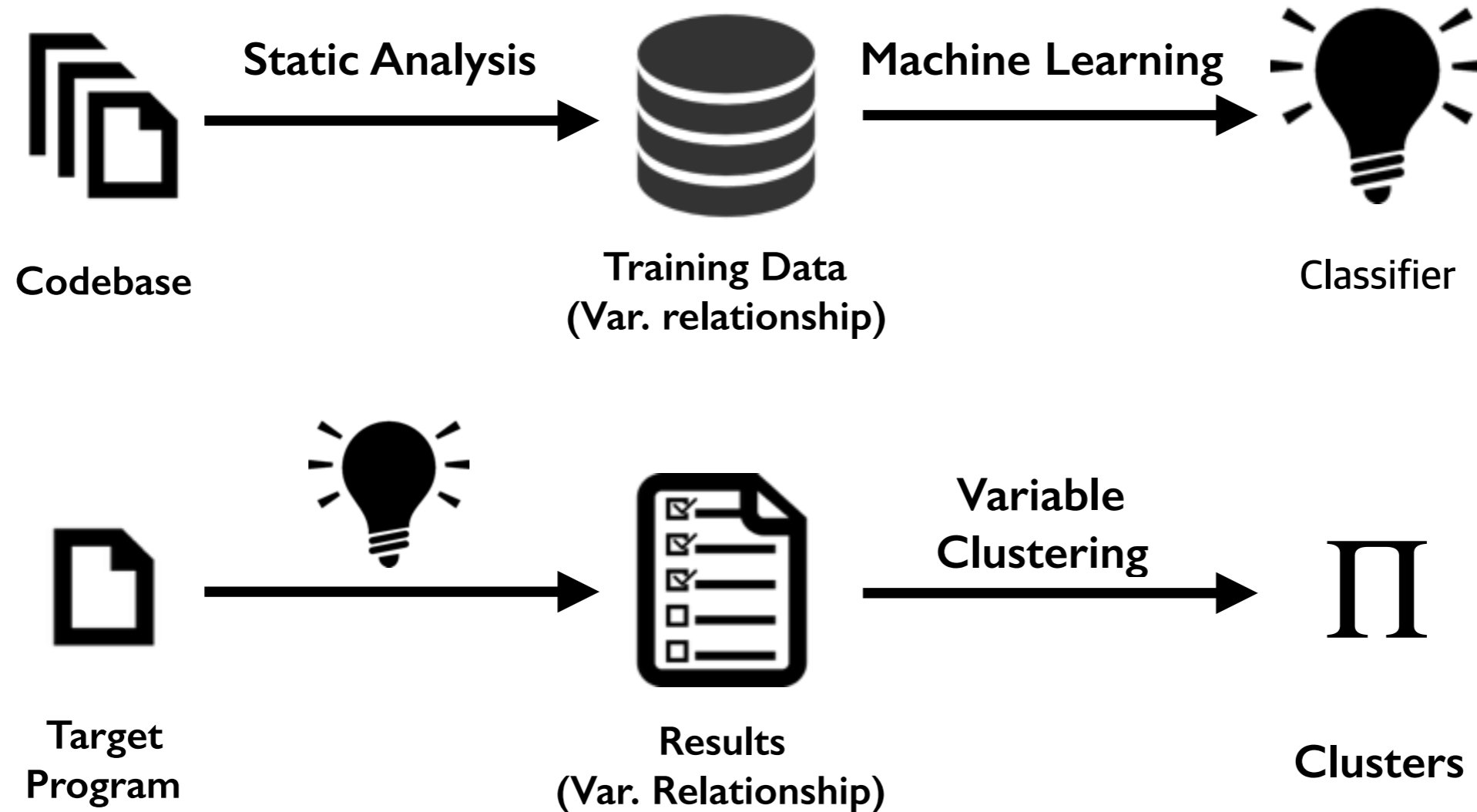- Learning a variable-clustering strategy from big data

  - fully relational pre-analysis <span style="color:red">as an offline teacher</span>

  - 33x faster yet similarly precise



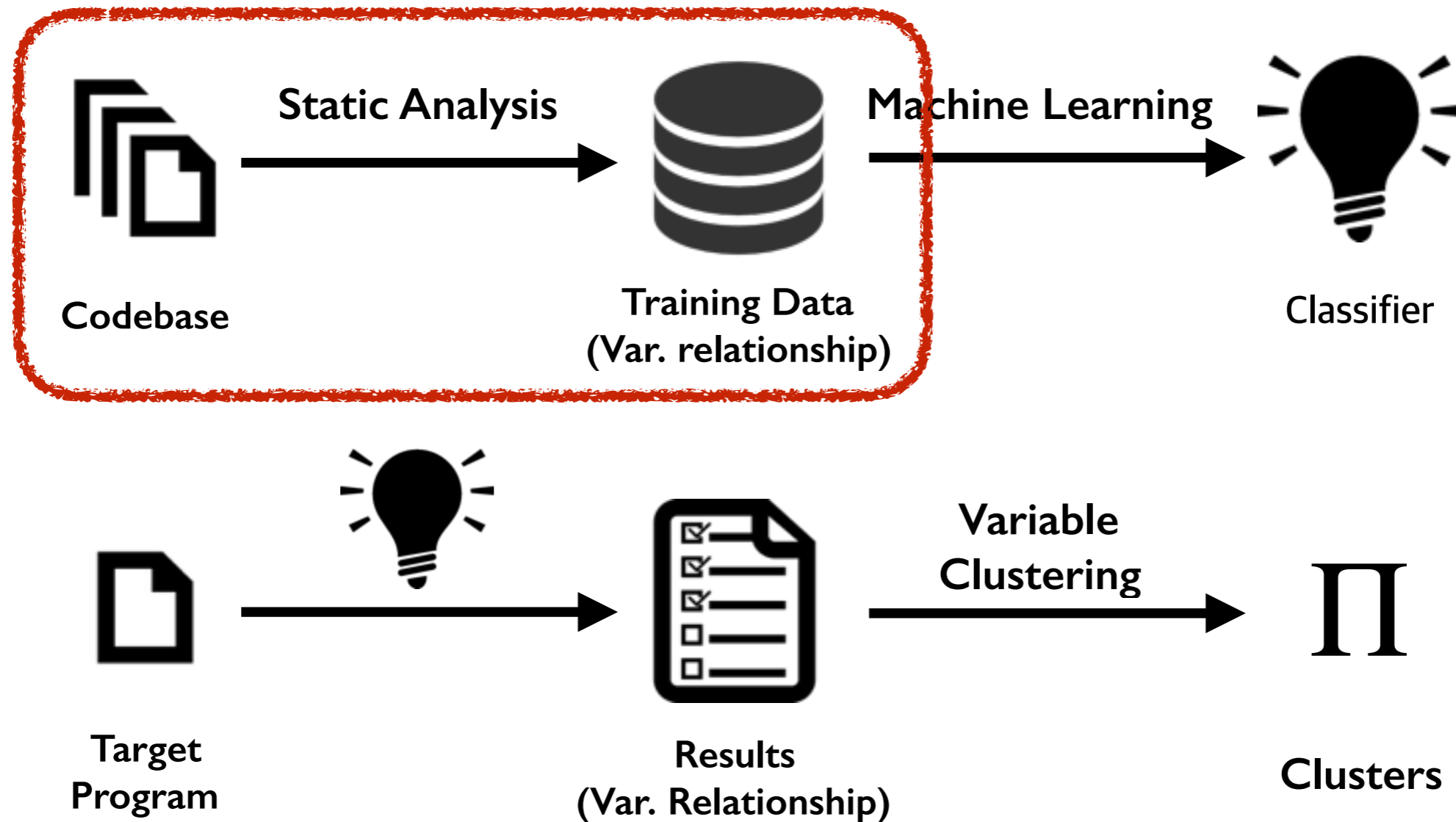Time — Var.Clustering / Main Analysis chart comparing [PLDI'14] and [ML-based]

# Big Picture

- Learning a variable-clustering strategy from big data

# Big Picture

- Learning a variable-clustering strategy from big data

# Training Data

- Pairs of two variables with label $\{\oplus, \ominus\}$

  - $\oplus$: precise $(< +\infty)$, $\ominus$: imprecise $(= +\infty)$

```
1  int a = b;
2  int c = input();              // User input
3  for (i = 0; i < b; i++) {
4    assert (i < a);             // Query 1
5    assert (i < c);             // Query 2
6  }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | 0 | ∞ | -1 |
| b | 0 | 0 | ∞ | -1 |
| c | ∞ | ∞ | 0 | ∞ |
| i | ∞ | ∞ | ∞ | 0 |

Octagon Analysis

$\oplus : \{(a,b), (a,i), (b,a) \dots\}$

$\ominus : \{(a,c), (b,c), (c,a) \dots\}$

# Training Data

- Automatically generated by impact pre-analysis[PLDI'14]

  - fully relational, yet more scalable than the full octagon

```
1   int a = b;
2   int c = input();              // User input
3   for (i = 0; i < b; i++) {
4     assert (i < a);             // Query 1
5     assert (i < c);             // Query 2
6   }
```

|   | a | b | c | i |
|---|---|---|---|---|
| a | 0 | 0 | $\infty$ | -1 |
| b | 0 | 0 | $\infty$ | -1 |
| c | $\infty$ | $\infty$ | 0 | $\infty$ |
| i | $\infty$ | $\infty$ | $\infty$ | 0 |

Octagon Analysis

$\gamma(\bigstar) = \mathbb{Z}$

$\gamma(\top) = \mathbb{Z} \cup \{+\infty\}$

$\oplus : \{(a,b), (a,i), (b,a) \dots\}$

$\ominus : \{(a,c), (b,c), (c,a) \dots\}$

|   | a | b | c | i |
|---|---|---|---|---|
| a | $\bigstar$ | $\bigstar$ | $\top$ | $\bigstar$ |
| b | $\bigstar$ | $\bigstar$ | $\top$ | $\bigstar$ |
| c | $\top$ | $\top$ | $\bigstar$ | $\top$ |
| i | $\top$ | $\top$ | $\top$ | $\bigstar$ |

Impact Pre-analysis

# Big Picture

- Learning a variable-clustering strategy from big data



Codebase → **Static Analysis** → Training Data (Var. relationship) → **Machine Learning** → Classifier

Target Program → → Results (Var. Relationship) → **Variable Clustering** → $\prod$ Clusters

# Features

- ## 30 Features of variable pairs

  - ## boolean predicate of (x,y) in program P

(Positive situations for Octagon)
- x=y+k or y=x+k
- x<=y+k or y<=x+k
- x=malloc(y) or y=malloc(x)
- x[y] or y[x]
- ...

(Negative situations for Octagon)
- x=cy or y=cx (c != 1)
- x=yz or y=xz
- x=y/z or y=x/z
- ...

(General syntactic features)
- x or y is a field
- x and y represent sizes of arrays
- x or y is the size of a const string
- x or y is a global variable
- ...

(General semantic features)
- x or y has a finite interval
- x or y is a local var in a recursive function
- x, y are not accessed in the same function
- ...

# Features

- Importance of features by Gini Index

  - negative & general > positive & domain-specific

(Positive situations for Octagon)
- x=y+k or y=x+k
- x<=y+k or y<=x+k
- x=malloc(y) or y=malloc(x)
- x[y] or y[x]
- …

(Negative situations for Octagon)
- x=cy or y=cx (c != 1)
- x=yz or y=xz
- x=y/z or y=x/z
- …

(General syntactic features)
- x or y is a field
- x and y represent sizes of arrays
- x or y is the size of a const string
- x or y is a global variable
- …

(General semantic features)
- x or y has a finite interval
- x or y is a local var in a recursive function
- x, y are not accessed in the same function
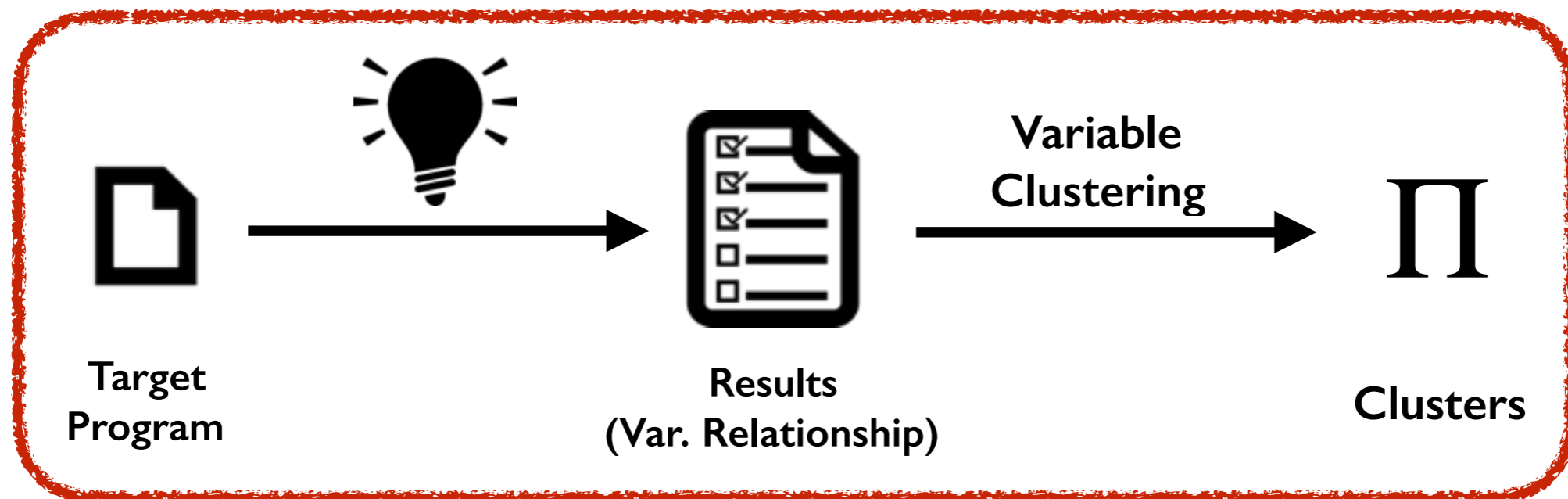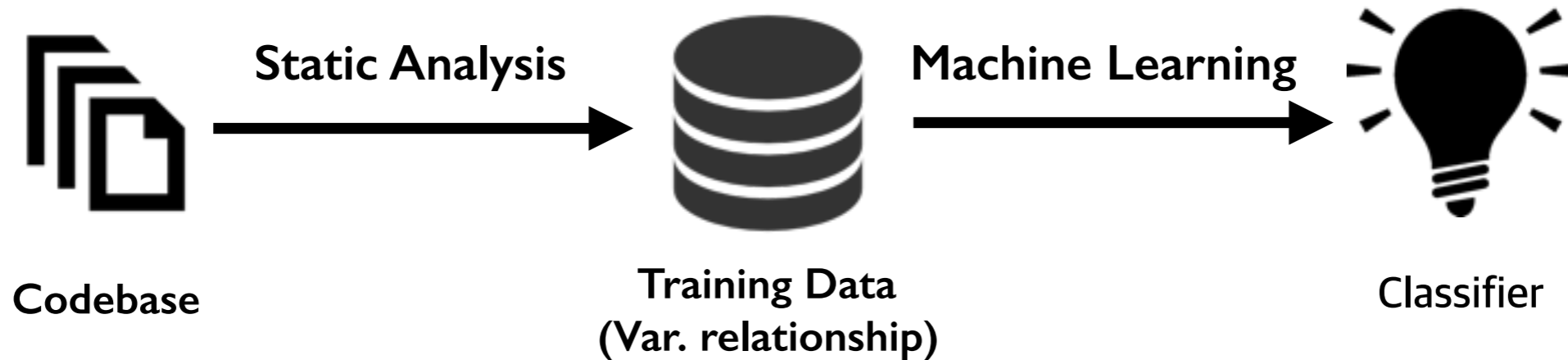- …

*Top 5 most important features

# Classifier

- Learning a binary classifier $\mathcal{C} : Var \times Var \rightarrow \{\oplus, \ominus\}$

  - using an off-the-shelf ML algorithm: decision tree

- Why decision tree?

  - more expressive than linear models

  - e.g.) Octagon with logistic regression : 10~12x slower

# Big Picture

- Learning a variable-clustering strategy from big data

# Clustering Strategy

- ⊕-marked variable pairs in the same cluster

- naturally covers transitive relationships

```
1  int a = b;
2  int c = input();              // User input
3  for (i = 0; i < b; i++) {
4    assert (i < a);             // Query 1
5    assert (i < c);             // Query 2
6  }
```

| | C(x,y) |
|---|---|
| (a,b) | ⊕ |
| (a,i) | ⊖ |
| (b,i) | ⊕ |
| (a,c) | ⊖ |
| ... | ... |

# Experiments

- Implemented on top of  Sparrow — The Early Bird

  - sound & global analyzer

  - a buffer overrun detector for full C

- 17 open source benchmarks (~100KLOC)

# Experimental Results

- Effectiveness (leave-one-out cross validation)

| Program | LOC | #Abs.Loc. | # Alarms | | | Time(s) | | |
|---|---|---|---|---|---|---|---|---|
| | | | Itv | Impt | ML | Itv | Impt | ML |
| brutefir | 103 | 54 | 4 | 0 | 0 | 0 | 0 | 0 |
| consol | 298 | 165 | 20 | 10 | 10 | 0 | 0 | 0 |
| id3 | 512 | 527 | 15 | 6 | 6 | 0 | 0 | 1 |
| spell | 2,213 | 450 | 20 | 8 | 17 | 0 | 1 | 1 |
| mp3rename | 2,466 | 332 | 33 | 3 | 3 | 0 | 1 | 1 |
| irmp3 | 3,797 | 523 | 2 | 0 | 0 | 1 | 2 | 3 |
| barcode | 4,460 | 1,738 | 235 | 215 | 215 | 2 | 9 | 6 |
| httptunnel | 6,174 | 1,622 | 52 | 29 | 27 | 3 | 35 | 5 |
| e2ps | 6,222 | 1,437 | 119 | 58 | 58 | 3 | 6 | 3 |
| bc | 13,093 | 1,891 | 371 | 364 | 364 | 14 | 252 | 16 |
| less | 23,822 | 3,682 | 625 | 620 | 625 | 83 | 2,354 | 87 |
| bison | 56,361 | 14,610 | 1,988 | 1,955 | 1,955 | 137 | 4,827 | 237 |
| pies | 66,196 | 9,472 | 795 | 785 | 785 | 49 | 14,942 | 95 |
| icecast-server | 68,564 | 6,183 | 239 | 232 | 232 | 51 | 109 | 107 |
| raptor | 76,378 | 8,889 | 2,156 | 2,148 | 2,148 | 242 | 17,844 | 345 |
| dico | 84,333 | 4,349 | 402 | 396 | 396 | 38 | 156 | 51 |
| lsh | 110,898 | 18,880 | 330 | 325 | 325 | 33 | 139 | 251 |
| Total | | | 7,406 | 7,154 | 7,166 | 656 | 40,677 | 1,207 |

# Experimental Results

- Effectiveness (leave-one-out cross validation)

| Program | LOC | #Abs.Loc. | # Alarms | | | Time(s) | | |
|---|---|---|---|---|---|---|---|---|
| | | | Itv | Impt | ML | Itv | Impt | ML |
| brutefir | 103 | 54 | 4 | 0 | 0 | 0 | 0 | 0 |
| consol | 298 | 165 | 20 | 10 | 10 | 0 | 0 | 0 |
| id3 | 512 | 527 | 15 | 6 | 6 | 0 | 0 | 1 |
| spell | 2,213 | 450 | 20 | 8 | 17 | 0 | 1 | 1 |
| mp3rename | 2,466 | 332 | 33 | 3 | 3 | 0 | 1 | 1 |
| irmp3 | 3,797 | 523 | 2 | 0 | 0 | 1 | 2 | 3 |
| barcode | 4,460 | 1,738 | 235 | 215 | 215 | 2 | 9 | 6 |
| httptunnel | 6,174 | 1,622 | 52 | 29 | 27 | 3 | 35 | 5 |
| e2ps | 6,222 | 1,437 | 119 | 58 | 58 | 3 | 6 | 3 |
| bc | 13,093 | 1,891 | 371 | 364 | 364 | 14 | 252 | 16 |
| less | 23,822 | 3,682 | 625 | 620 | 625 | 83 | 2,354 | 87 |
| bison | 56,361 | 14,610 | 1,988 | 1,955 | 1,955 | 137 | 4,827 | 237 |
| pies | 66,196 | 9,472 | 795 | 785 | 785 | 49 | 14,942 | 95 |
| icecast-server | 68,564 | 6,183 | 239 | 232 | 232 | 51 | 109 | 107 |
| raptor | 76,378 | 8,889 | 2,156 | 2,148 | 2,148 | 242 | 17,844 | 345 |
| dico | 84,333 | 4,349 | 402 | 396 | 396 | 38 | 156 | 51 |
| lsh | 110,898 | 18,880 | 330 | 325 | 325 | 33 | 139 | 251 |
| Total | | | 7,406 | 7,154 | 7,166 | 656 | 40,677 | 1,207 |

−252 −240

# Experimental Results

- Effectiveness (leave-one-out cross validation)

| Program | LOC | #Abs.Loc. | # Alarms | | | Time(s) | | |
|---|---|---|---|---|---|---|---|---|
| | | | Itv | Impt | ML | Itv | Impt | ML |
| brutefir | 103 | 54 | 4 | 0 | 0 | 0 | 0 | 0 |
| consol | 298 | 165 | 20 | 10 | 10 | 0 | 0 | 0 |
| id3 | 512 | 527 | 15 | 6 | 6 | 0 | 0 | 1 |
| spell | 2,213 | 450 | 20 | 8 | 17 | 0 | 1 | 1 |
| mp3rename | 2,466 | 332 | 33 | 3 | 3 | 0 | 1 | 1 |
| irmp3 | 3,797 | 523 | 2 | 0 | 0 | 1 | 2 | 3 |
| barcode | 4,460 | 1,738 | 235 | 215 | 215 | 2 | 9 | 6 |
| httptunnel | 6,174 | 1,622 | 52 | 29 | 27 | 3 | 35 | 5 |
| e2ps | 6,222 | 1,437 | 119 | 58 | 58 | 3 | 6 | 3 |
| bc | 13,093 | 1,891 | 371 | 364 | 364 | 14 | 252 | 16 |
| less | 23,822 | 3,682 | 625 | 620 | 625 | 83 | 2,354 | 87 |
| bison | 56,361 | 14,610 | 1,988 | 1,955 | 1,955 | 137 | 4,827 | 237 |
| pies | 66,196 | 9,472 | 795 | 785 | 785 | 49 | 14,942 | 95 |
| icecast-server | 68,564 | 6,183 | 239 | 232 | 232 | 51 | 109 | 107 |
| raptor | 76,378 | 8,889 | 2,156 | 2,148 | 2,148 | 242 | 17,844 | 345 |
| dico | 84,333 | 4,349 | 402 | 396 | 396 | 38 | 156 | 51 |
| lsh | 110,898 | 18,880 | 330 | 325 | 325 | 33 | 139 | 251 |
| Total | | | 7,406 | 7,154 | 7,166 | 656 | 40,677 | 1,207 |

-252    -240

# Experimental Results

- Effectiveness (leave-one-out cross validation)

| Program | LOC | #Abs.Loc. | # Alarms | | | Time(s) | | |
|---|---|---|---|---|---|---|---|---|
| | | | Itv | Impt | ML | Itv | Impt | ML |
| brutefir | 103 | 54 | 4 | 0 | 0 | 0 | 0 | 0 |
| consol | 298 | 165 | 20 | 10 | 10 | 0 | 0 | 0 |
| id3 | 512 | 527 | 15 | 6 | 6 | 0 | 0 | 1 |
| spell | 2,213 | 450 | 20 | 8 | 17 | 0 | 1 | 1 |
| mp3rename | 2,466 | 332 | 33 | 3 | 3 | 0 | 1 | 1 |
| irmp3 | 3,797 | 523 | 2 | 0 | 0 | 1 | 2 | 3 |
| barcode | 4,460 | 1,738 | 235 | 215 | 215 | 2 | 9 | 6 |
| httptunnel | 6,174 | 1,622 | 52 | 29 | 27 | 3 | 35 | 5 |
| e2ps | 6,222 | 1,437 | 119 | 58 | 58 | 3 | 6 | 3 |
| bc | 13,093 | 1,891 | 371 | 364 | 364 | 14 | 252 | 16 |
| less | 23,822 | 3,682 | 625 | 620 | 625 | 83 | 2,354 | 87 |
| bison | 56,361 | 14,610 | 1,988 | 1,955 | 1,955 | 137 | 4,827 | 237 |
| pies | 66,196 | 9,472 | 795 | 785 | 785 | 49 | 14,942 | 95 |
| icecast-server | 68,564 | 6,183 | 239 | 232 | 232 | 51 | 109 | 107 |
| raptor | 76,378 | 8,889 | 2,156 | 2,148 | 2,148 | 242 | 17,844 | 345 |
| dico | 84,333 | 4,349 | 402 | 396 | 396 | 38 | 156 | 51 |
| lsh | 110,898 | 18,880 | 330 | 325 | 325 | 33 | 139 | 251 |
| Total | | | 7,406 | 7,154 | 7,166 | 656 | 40,677 | 1,207 |
| | | | | | | | x62 | x2 |

# Experimental Results

- Generalization : training only with small (<60KLOC) pgms

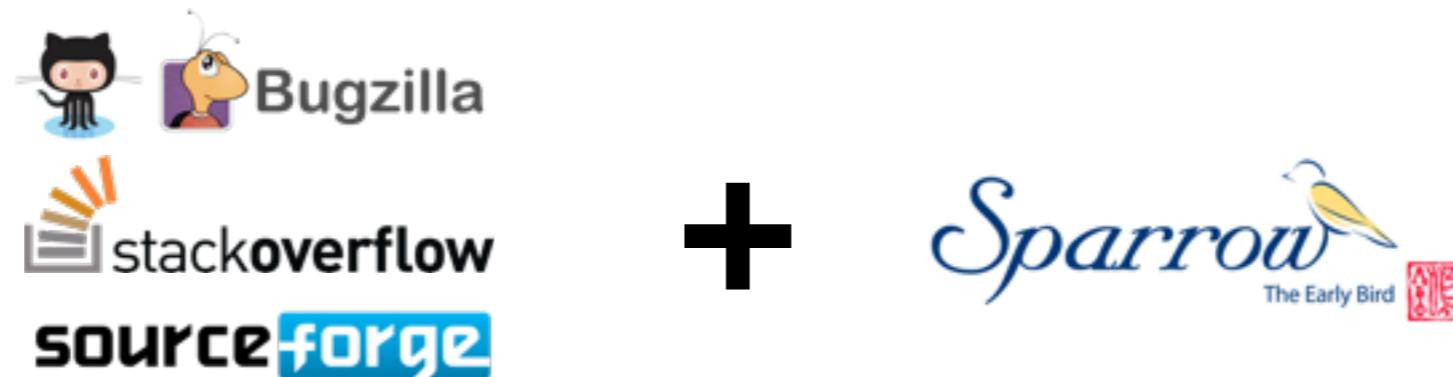| Program | LOC | Abs. Loc. | # Alarms | | | Time(s) | | |
|---|---|---|---|---|---|---|---|---|
| | | | Itv | All | Small | Itv | All | Small |
| pies | 66,196 | 9,472 | 795 | 785 | 785 | 49 | 95 | 98 |
| icecast-server | 68,564 | 6,183 | 239 | 232 | 232 | 51 | 113 | 99 |
| raptor | 76,378 | 8,889 | 2,156 | 2,148 | 2,148 | 242 | 345 | 388 |
| dico | 84,333 | 4,349 | 402 | 396 | 396 | 38 | 61 | 62 |
| lsh | 110,898 | 18,880 | 330 | 325 | 325 | 33 | 251 | 251 |
| Total | | | 7,406 | 3,886 | 3,886 | 413 | 865 | 898 |

+4%

# Summary



- Adaptive variable-clustering strategy for Octagon

  - **Machine Learning** (learner) + **Static Analysis** (teacher)

- 33x faster than a static-analysis-only approach

# Summary



- Adaptive variable-clustering strategy for Octagon

  - **Machine Learning** (learner) + **Static Analysis** (teacher)

- 33x faster than a static-analysis-only approach

Thank You