

Merge-Replay: Efficient IFDS-Based Taint Analysis by Consolidating Equivalent Value Flows

Yujiang Gui^{*}, Dongjie He^{*†} and Jingling Xue[†]

School of Computer Science and Engineering
University of New South Wales

Presenter: Yujiang Gui
38th ASE, September 2023



^{*}, [†] The first two authors contributed equally and are listed in alphabetical order by their last names, while the last two authors share corresponding authorship.

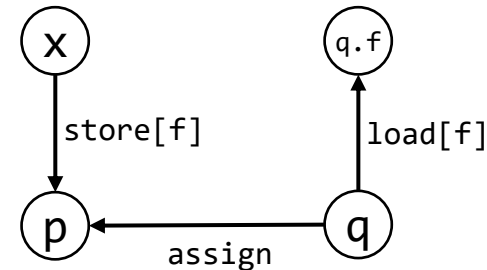
Static Taint Analysis

- Exposes potential sensitive data leaks ahead of time
 - Private messages
 - Location information
 - Financial details
 - ...

Static Taint Analysis

- Exposes potential sensitive data leaks ahead of time
 - Private messages
 - Location information
 - Financial details
 - ...
- Tracks the flows of tainted data to observe if they can reach sink

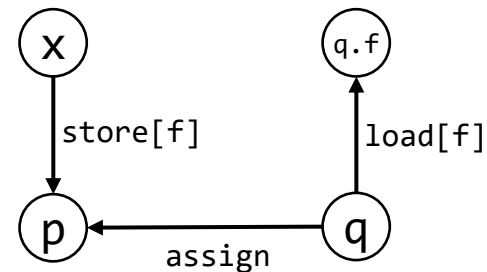
```
x = getPassword(); // source
p = q;
...
p.f = x;
sendMessage(q.f); // sink
```



Static Taint Analysis

- Exposes potential sensitive data leaks ahead of time
 - Private messages
 - Location information
 - Financial details
 - ...
- Tracks the flows of tainted data to observe if they can reach sink

```
x = getPassword(); // source
p = q;
...
p.f = x;
sendMessage(q.f); // sink
```



- An active research area
 - Security analysis for web apps (ASE'14, ISSTA'23)
 - Taint analysis for Android apps (PLDI'14, ISSTA'15, OOPSLA'18, ICSE'20)
 - Memory leak detection (CGO'21, ICSE'21, ISSTA'23)
 - ...

IFDS-Based Taint Analysis

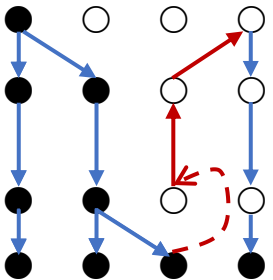
- Utilizes the IFDS framework*
 - Solves a class of interprocedural distributive analyses
 - Transforms an analysis into a graph-reachability problem
 - Context- and flow-sensitive

*Reps, Thomas, Susan Horwitz, and Mooly Sagiv. "Precise Interprocedural Dataflow Analysis via Graph Reachability.", POPL'95.

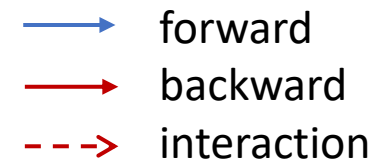
IFDS-Based Taint Analysis

- Utilizes the IFDS framework*
 - Solves a class of interprocedural distributive analyses
 - Transforms an analysis into a graph-reachability problem
 - Context- and flow-sensitive
- Employs two mutually iterative passes
 - Identify taints forwards
 - Detect aliases backwards

\emptyset x p.f q.f



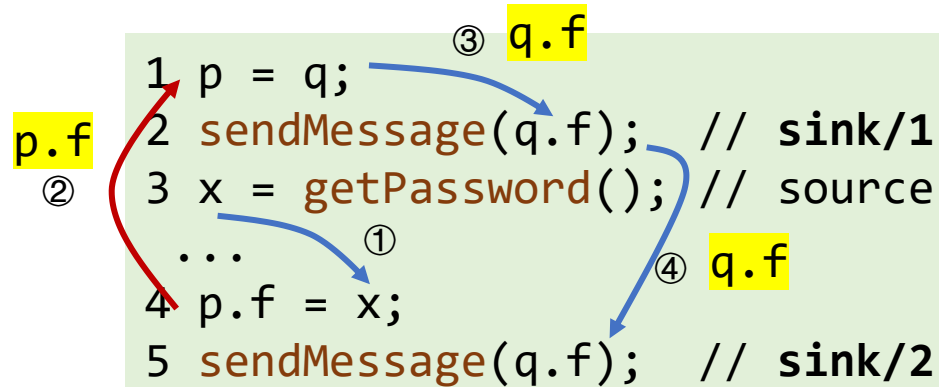
```
x = getPassword(); // source
p = q;
...
p.f = x;
sendMessage(q.f); // sink
```



*Reps, Thomas, Susan Horwitz, and Mooly Sagiv. "Precise Interprocedural Dataflow Analysis via Graph Reachability.", POPL'95.

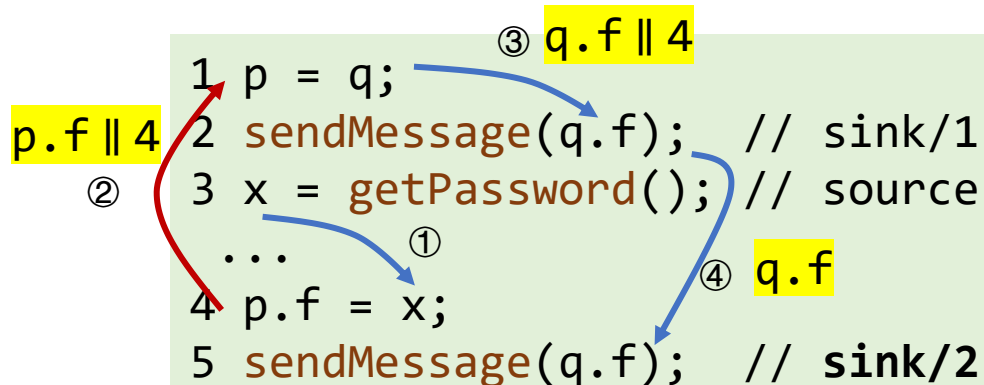
Maintaining Flow Sensitivity with Activation Statements

- Loss of flow sensitivity during the interactions of solvers



- FlowDroid*

- Recovers flow sensitivity with activation statements



*Arzt, Steven, Siegfried Rasthofer, Christian Fritz, Eric Bodden, Alexandre Bartel, Jacques Klein, Yves Le Traon, Damien Ocateau, and Patrick D. McDaniel. "FlowDroid: Precise Context, Flow, Field, Object-Sensitive and Lifecycle-Aware Taint Analysis for Android Apps.", PLDI'14.

Limitation 1: Redundant Propagation

```
1 class A {}
2 class B { A f; }

3 void main() {
4   B q = new B(); q.f || 14 q.f || 12
5   B p = q;
6   B r = new B(); ③ ② ③ ②
7   bar(q);
8   foo(p, q, r);
9 }
```

```
19 void bar(q) {
20   B p = new B();
21   B r = q;
22   foo(p, q, r);
23 }
```

p.f || 12

p.f || 14

```
10 void foo(p, q, r) {
11   if (...) {
12     p.f = <taint/1>;
13   } else {
14     p.f = <taint/2>;
15     sink(q.f);
16   }
17   sink(r.f);
18 }
```


Limitation 2: Overlooking Context

```
1 class A {}
2 class B { A f; }

3 void main() {
4     B q = new B();
5     B p = q;
6     B r = new B();
7     bar(q);
8     foo(p, q, r);
9 }
```

`q.f` || 14 `q.f` || 12

```
19 void bar(q) {
20     B p = new B();
21     B r = q;
22     foo(p, q, r);
23 }
```

`r.f` || 14

`r.f` || 12

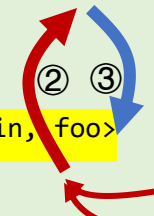
```
10 void foo(p, q, r) {
11     if (...) {
12         p.f = <taint/1>;
13     } else {
14         p.f = <taint/2>;
15         sink(q.f);
16     }
17     sink(r.f);
18 }
```

False alarm!

The Merge-Replay Strategy

```
1 class A {}
2 class B { A f; }

3 void main() {
4   B q = new B(); q.f || sym<main, foo>
5   B p = q;
6   B r = new B();
7   bar(q); p.f || sym<main, foo>
8   foo(p, q, r);
9 }
```



```
19 void bar(q) {
20   B p = new B();
21   B r = q;
22   foo(p, q, r);
23 }
```

p.f || 12

p.f || 14

```
10 void foo(p, q, r) {
11   if (...) {
12     p.f = <taint/1>;
13   } else {
14     p.f = <taint/2>;
15     sink(q.f);
16   }
17   sink(r.f);
18 }
```

Merge with a symbolic activation statement:

$\mathit{sym}_{p.f}^{\langle main, foo \rangle} \mapsto \{12, 14\}$

The Merge-Replay Strategy

```
1 class A {}
2 class B { A f; }

3 void main() {
4   B q = new B();
5   B p = q;
6   B r = new B();
7   bar(q);
8   foo(p, q, r);
9 }
```

```
19 void bar(q) {
20   B p = new B();
21   B r = q;
22   foo(p, q, r);
23 }
```

```
10 void foo(p, q, r) {
11   if (...) {
12     p.f = <taint/1>;
13   } else {
14     p.f = <taint/2>;
15     sink(q.f);
16   }
17   sink(r.f);
18 }
```

Prune when contexts are mismatched:

$$\text{sym}_{p.f}^{\langle \text{main}, \text{foo} \rangle} \mapsto \{12, 14\}$$

The Merge-Replay Strategy

```
1 class A {}
2 class B { A f; }

3 void main() {
4   B q = new B();
5   B p = q;
6   B r = new B();
7   bar(q);
8   foo(p, q, r);
9 }
```

① $q.f \parallel \text{sym}^{\langle \text{main}, \text{foo} \rangle}$

```
19 void bar(q) {
20   B p = new B();
21   B r = q;
22   foo(p, q, r);
23 }
```

②

```
10 void foo(p, q, r) {
11   if (...) { q.f || 12 q.f || 14
12     p.f = <taint/1>;
13   } else {
14     p.f = <taint/2>;
15     sink(q.f);
16   }
17   sink(r.f);
18 }
```

Replay using concrete activation statements:

$\text{sym}_{p.f}^{\langle \text{main}, \text{foo} \rangle} \mapsto \{12, 14\}$

Summary

- Reduce redundant propagation
- Enhance precision
- Boost overall performance
- Conceptually simple

The MergeDroid Algorithm

```

function ForwardAnalysis ()
11a for  $d_3$  such that  $\langle n, d_2 \rangle \rightarrow \langle s_{m'}, d_3 \rangle \in E_{FW}^\#$  do
11b   for  $d'_3 \in \text{Concretize}(\langle m, d_1 \rangle \rightarrow \langle n, d_2 \rangle, m', d_3)$  do
12a     Inject  $(\langle m, d_1 \rangle \rightarrow \langle n, d_2 \rangle, \langle m', d'_3 \rangle, E_{BW}^\#, \text{PathEdge}_{BW}, S_{BW}, W_{BW})$ 
13a     Prop  $(\langle m', d'_3 \rangle \rightarrow \langle s_{m'}, d_3 \rangle, W_{FW}, \text{PathEdge}_{FW})$ 
14a     for  $\langle m', d'_3 \rangle \rightarrow \langle e_{m'}, d_4 \rangle \in S_{FW} \wedge \langle e_{m'}, d_4 \rangle \rightarrow \langle r, d_5 \rangle \in E_{FW}^\#$  do
14b        $d'_5 = \text{AttachActivationStmt}(d_5, d_2)$ 
15a       Prop  $(\langle m, d_1 \rangle \rightarrow \langle r, d'_5 \rangle, W_{FW}, \text{PathEdge}_{FW})$ 
...
22a for  $\langle m'', d_3 \rangle \rightarrow \langle c, d_4 \rangle \in \text{PathEdge}_{FW} \wedge \langle c, d_4 \rangle \rightarrow \langle s_m, d_1 \rangle \in E_{FW}^\#$ 
     $\wedge \langle n, d_2 \rangle \rightarrow \langle r, d_5 \rangle \in E_{FW}^\#$  do
22b    $d'_5 = \text{AttachActivationStmt}(d_5, d_4)$ 
23a   Prop  $(\langle m'', d_3 \rangle \rightarrow \langle r, d'_5 \rangle, W_{FW}, \text{PathEdge}_{FW})$ 

function BackwardAnalysis ()
38a for  $d_3$  such that  $\langle n, d_2 \rangle \rightarrow \langle e_{m'}, d_3 \rangle \in E_{BW}^\#$  do
38b    $d'_3 = \text{DataAbstraction}(d_3) \parallel \text{GAS}$ 
39a   Inject  $(\langle m, d_1 \rangle \rightarrow \langle n, d_2 \rangle, \langle m', d'_3 \rangle, E_{FW}^\#, \text{PathEdge}_{FW}, S_{FW}, W_{FW})$ 
40a   Prop  $(\langle m', d'_3 \rangle \rightarrow \langle e_{m'}, d_3 \rangle, W_{BW}, \text{PathEdge}_{BW})$ 
41a   for  $\langle m', d'_3 \rangle \rightarrow \langle s_{m'}, d_4 \rangle \in S_{BW} \wedge \langle s_{m'}, d_4 \rangle \rightarrow \langle r, d_5 \rangle \in E_{BW}^\#$  do
41b      $d'_5 = \text{OnReturnFlow}(d_5, d_2, m, m')$ 
42a     Prop  $(\langle m, d_1 \rangle \rightarrow \langle r, d'_5 \rangle, W_{BW}, \text{PathEdge}_{BW})$ 
...
49a for  $\langle m'', d_3 \rangle \rightarrow \langle c, d_4 \rangle \in \text{PathEdge}_{BW} \wedge \langle c, d_4 \rangle \rightarrow \langle e_m, d_1 \rangle \in E_{BW}^\#$ 
     $\wedge \langle n, d_2 \rangle \rightarrow \langle r, d_5 \rangle \in E_{BW}^\#$  do
49b    $d'_5 = \text{OnReturnFlow}(d_5, d_4, m'', m)$ 
50a   Prop  $(\langle m'', d_3 \rangle \rightarrow \langle r, d'_5 \rangle, W_{BW}, \text{PathEdge}_{BW})$ 

61 function OnActivationStmtAdded ( $sym, as$ )
62    $\langle m, m' \rangle = \text{Context}(sym)$  //  $m$  is caller,  $m'$  is callee
63   for  $\langle \langle m, d_1 \rangle \rightarrow \langle c, d_2 \rangle, abs \rangle \in \text{SymbolIncoming}(sym)$  do
64      $d_3 = abs \parallel as$ 
65     Prop  $(\langle m', d_3 \rangle \rightarrow \langle s_{m'}, d_3 \rangle, W_{FW}, \text{PathEdge}_{FW})$ 
66     Let  $r$  be the return node of  $c$ 
67     for  $\langle m', d_3 \rangle \rightarrow \langle e_{m'}, d_4 \rangle \in S_{FW} \wedge \langle e_{m'}, d_4 \rangle \rightarrow \langle r, d_5 \rangle \in E_{FW}^\#$  do
68       Prop  $(\langle m, d_1 \rangle \rightarrow \langle r, d_5 \rangle, W_{FW}, \text{PathEdge}_{FW})$ 

69  $\text{SymbolIncoming} = \text{Symb2Reps} = \{\}$ 
70 function Symbolize ( $d_{ret}, caller, callee$ )
71    $as = \text{ActivationStmt}(d_{ret})$ 
72    $abs = \text{DataAbstraction}(d_{ret})$ 
73    $sym = \langle abs, caller, callee \rangle$  // Symbolic activation stmt
74   if  $as \notin \text{Symb2Reps}(sym)$  then
75      $\text{Symb2Reps}(sym) \ni as$ 
76     OnActivationStmtAdded ( $sym, as$ )
77   return  $abs \parallel sym$ 

78 function Concretize ( $\langle caller, d_1 \rangle \rightarrow \langle c, d_2 \rangle, callee, d_3$ )
79   if  $d_3$  is active then return  $\{d_3\}$ 
80    $as = \text{ActivationStmt}(d_3)$ 
81   if  $as$  is a concrete statement or  $as = \text{GAS}$  then return  $\{d_3\}$ 
82   assert  $as$  is a symbolic activation statement
83    $sym = as$  // Rename variable
84    $reps = \{\}$  // Set of represented facts
85    $abs = \text{DataAbstraction}(d_3)$ 
86   if  $\text{Context}(sym) = \langle caller, callee \rangle$  then
87      $\text{SymbolIncoming}(sym) \ni \langle \langle caller, d_1 \rangle \rightarrow \langle c, d_2 \rangle, abs \rangle$ 
88     for  $v \in \text{Symb2Reps}(sym)$  do
89        $reps \ni abs \parallel v$ 
90   else
91      $reps \ni abs \parallel \text{GAS}$ 
92   return  $reps$ 

93 function AttachActivationStmt ( $d_{ret}, d_{call}$ )
94    $u = \text{ActivationStmt}(d_{ret})$ 
95    $v = \text{ActivationStmt}(d_{call})$ 
96   if  $u = \text{GAS}$  then
97     return  $\text{DataAbstraction}(d_{ret}) \parallel v$ 
98   return  $d_{ret}$ 

99 function OnReturnFlow ( $d_{ret}, d_{call}, caller, callee$ )
100    $u = \text{ActivationStmt}(d_{ret})$ 
101   if  $u = \text{GAS}$  then
102     return  $\text{AttachActivationStmt}(d_{ret}, d_{call})$ 
103   return  $\text{Symbolize}(d_{ret}, caller, callee)$ 

```

Implementation

- Implemented on top of **FlowDroid** (PLDI'14)
- In about **400 lines** of Java code
- Available at <https://www.cse.unsw.edu.au/~corg/MergeDroid>

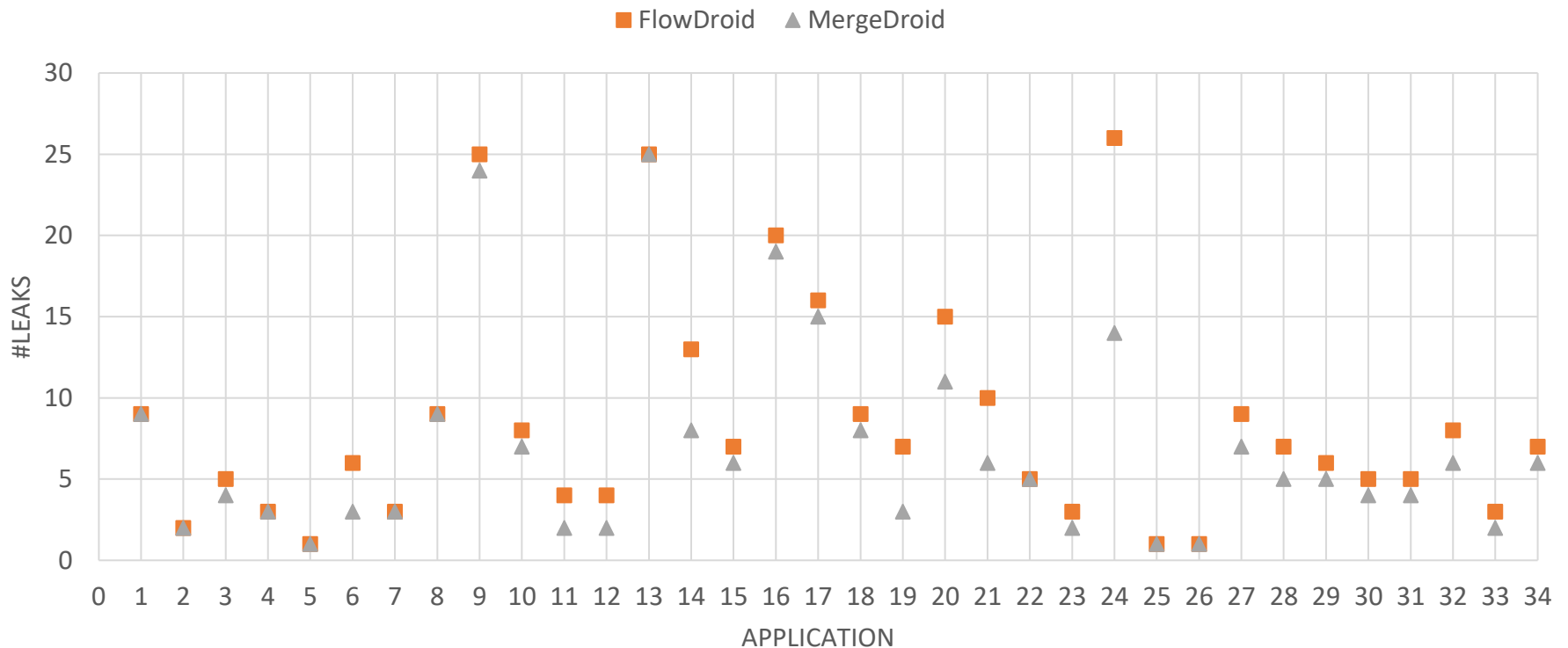
Evaluation

- Benchmark (40 apps)
 - From previous papers (ASE'19 and CGO'21)
 - From F-Droid
- Analysis budget
 - 3 hours and 256GB per app
- Metrics
 - Precision
 - Analysis time
 - Memory usage

RQ1: Precision

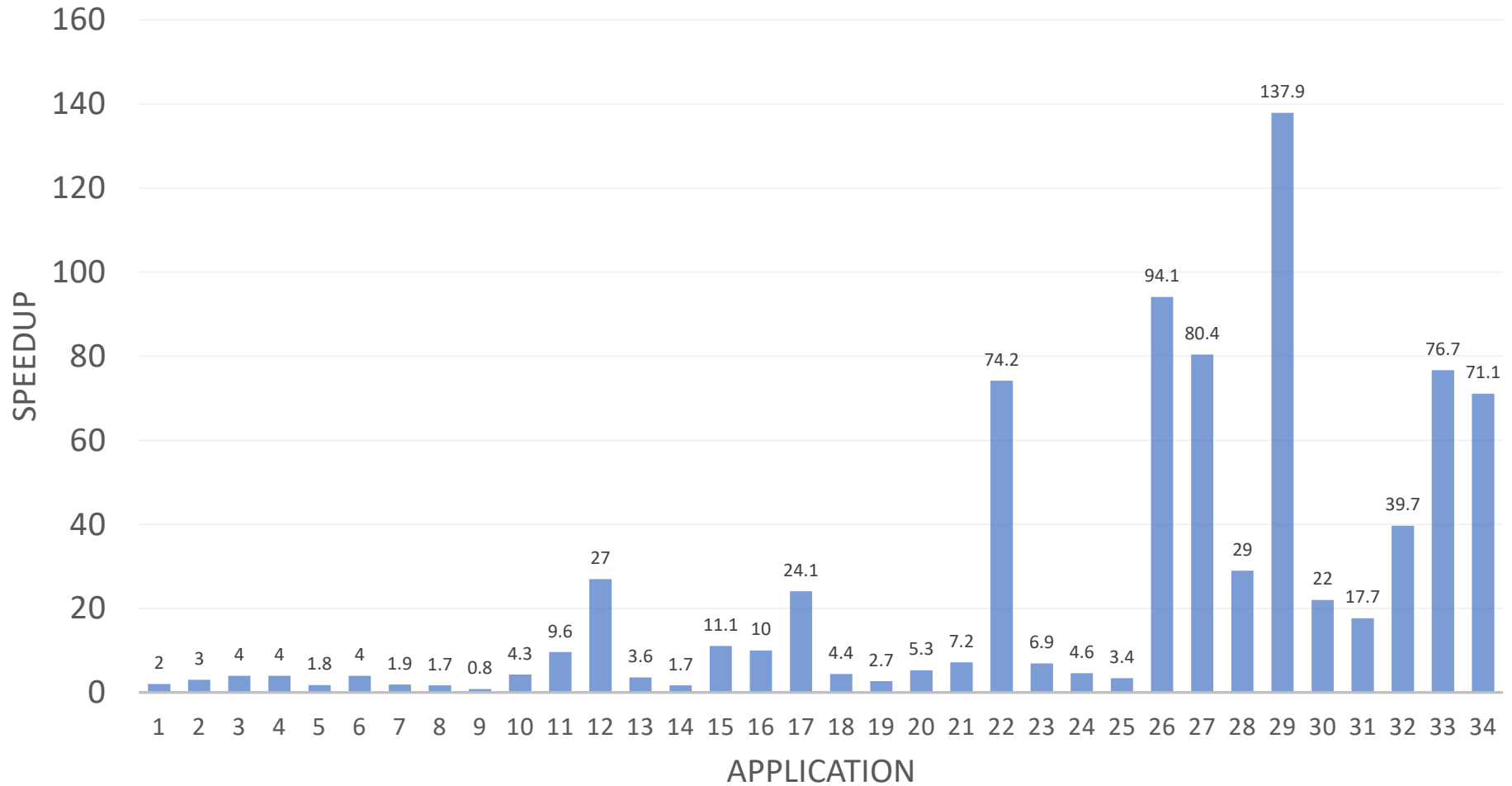
○ Validated correctness using

- DroidBench
- TaintBench



MergeDroid reduces false positives by decreasing reported leak warnings by FlowDroid by 19.2% on average.

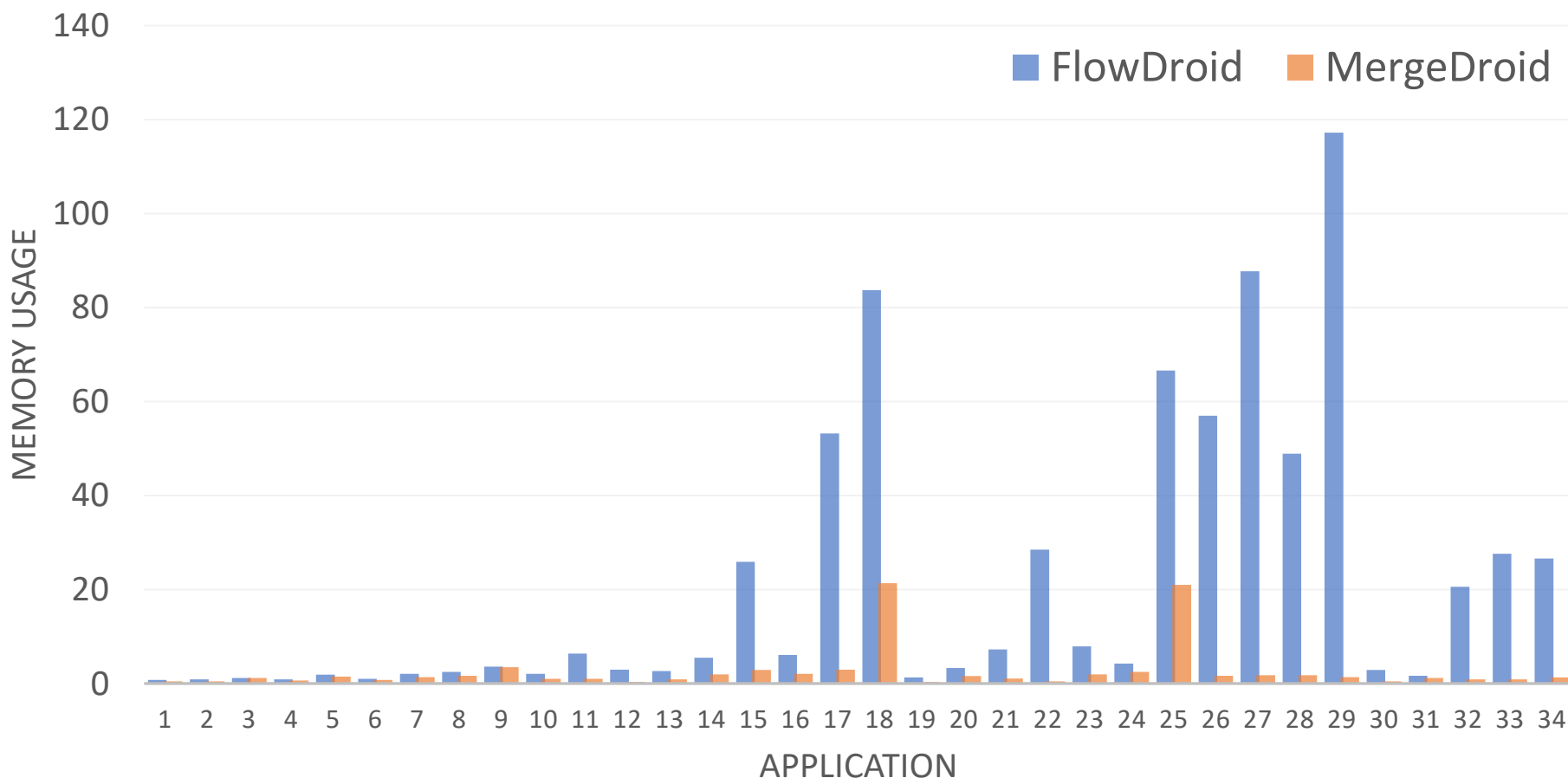
RQ2: Speedups



MergeDroid improves the efficiency of FlowDroid, and scales 6 more apps. The speedups range from 0.8 \times to 137.9 \times with an average of 9.0 \times .

RQ3: Memory Requirements

Maximum Memory Usage (GB)



MergeDroid uses less memory than FlowDroid for all apps analyzed, the ratio range from 1.0x to 83.6x with an average of 5.2x.

Thank you!