

Optimistic Hybrid Analysis

David Devecsery,

Peter M. Chen, Satish Narayanasamy, Jason Flinn

University of Michigan



Motivation – Problem

- Dynamic analysis is essential
 - Memory errors
 - Buffer overflow, use-after-free, double free, memory leak...
 - Information flow
 - Taint tracking, information providence, program slicing...
 - Race freedom
 - Type safety
 - Many more
- We want to enforce these properties in programs
 - Need precise (no false positives) and sound (no false negatives) program analysis



Motivation – Problem

- Problem: Dynamic analysis too expensive
 - Information flow 1-2 orders of magnitude overhead
 - Race detection order of magnitude slowdown
- Question: Can we create a faster dynamic analysis?

Code

```
int a = my_malloc()
int b = my_malloc()

*a = read()
*b = read()

int c = *a + *b

Free(a)
Free(b)
```

Instrumenter

Binary +

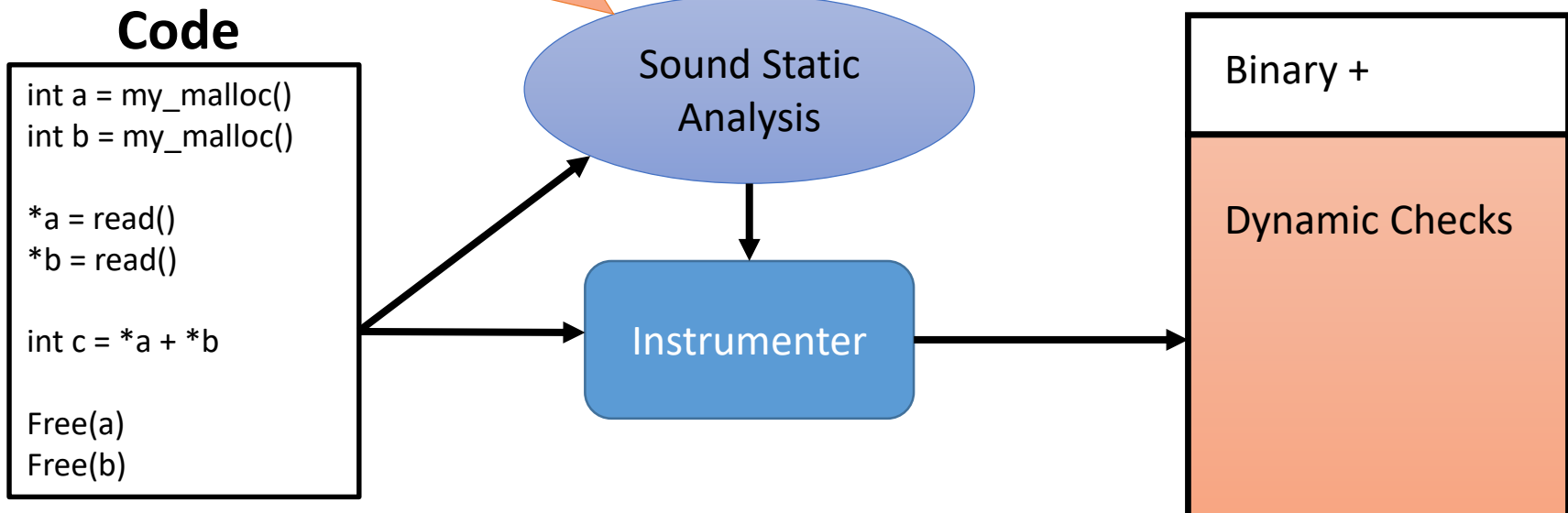
Dynamic Checks



Motivation – Hybrid Analysis

- Use static analysis to optimize dynamic analysis

If static analysis is sound dynamic analysis retains soundness



Sound Static Analysis is Inaccurate

- Makes conservative approximations for soundness
- Sound static analysis considers too many states to adequately optimize dynamic analysis!

Sound static analysis considers too many states to adequately optimize dynamic analysis!

States
used

Optimistic Hybrid Analysis

- Novel dynamic analysis optimization methodology
 - Combines profiling, **unsound static analysis**, and **speculative dynamic** execution
 - Unsound static analysis – better optimize analysis
 - Speculative execution – handles analysis unsoundness
 - No loss of soundness or precision
- Optimistic Hybrid Backward Slicer: **OptSlice**
 - Retains accuracy of traditional dynamic slicer
 - On average: **7.7x faster** than hybrid dynamic slicer



Overview

- ~~Motivation~~
- Optimistic Hybrid Analysis Overview
 - Likely Invariant Gathering
 - Optimistic Static Analysis
 - Speculative Dynamic Analysis
- OptSlice – Optimistic Dynamic Backward Slicer
- Evaluation
- Conclusion



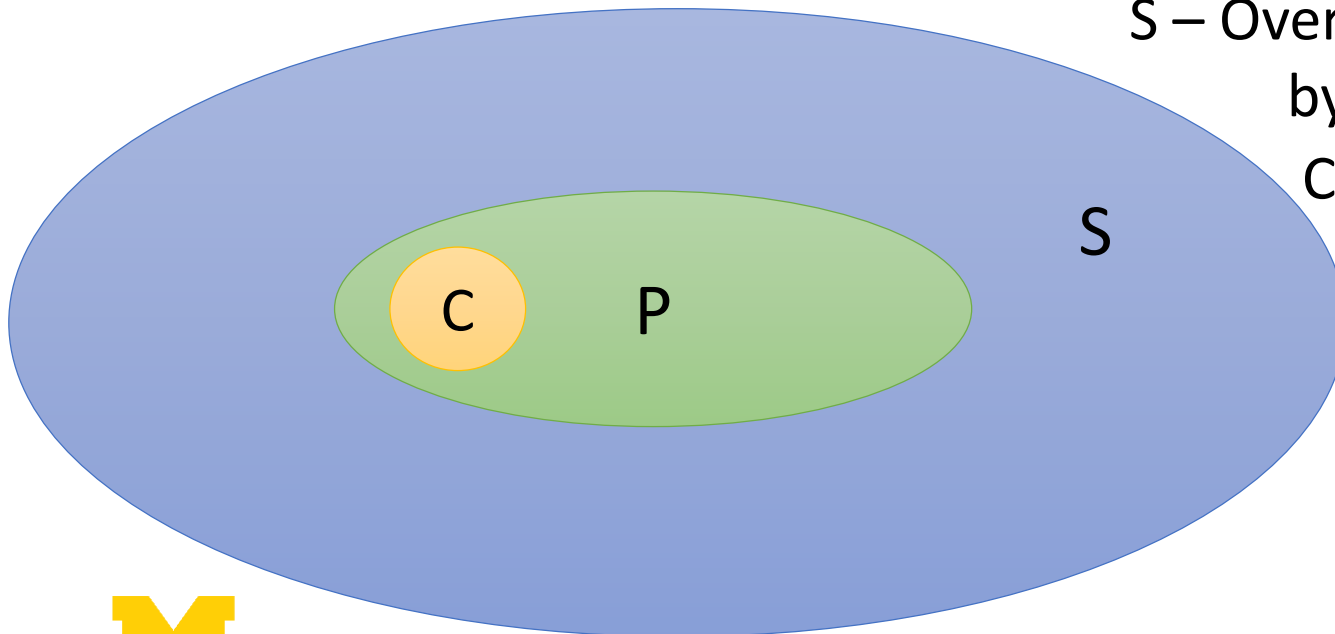
Idea: Optimistic Analysis

- **Observation:** Hybrid analysis need static analysis to explore **only** states executed by analyzed executions

P – Possible Program States

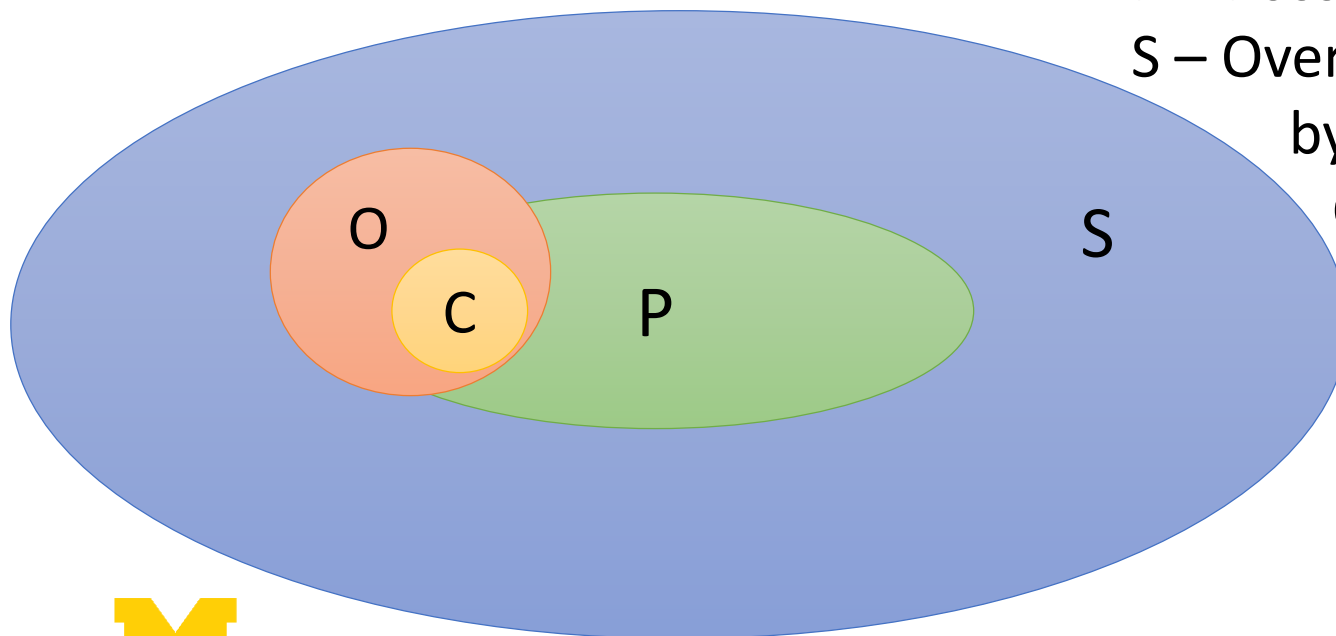
S – Over-approximation used
by analysis

C – Common Executed
States



Idea: Optimistic Analysis

- Idea: Instead use unsound static analysis
 - We call: **Optimistic Static Analysis**
- Have dynamic analysis detect & handle unsoundness



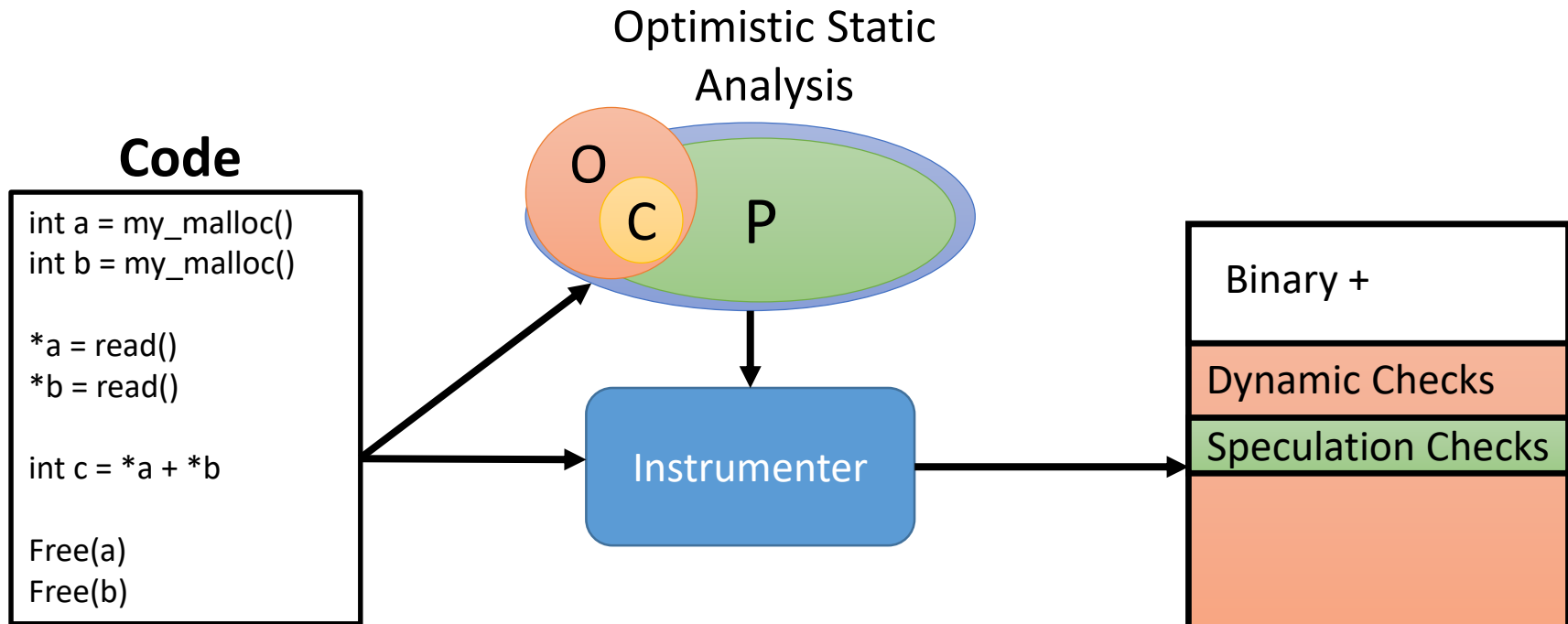
P – Possible Program States

S – Over-approximation used
by analysis

C – Common states

O - Optimistic
Static Analysis

Motivation – Hybrid Analysis

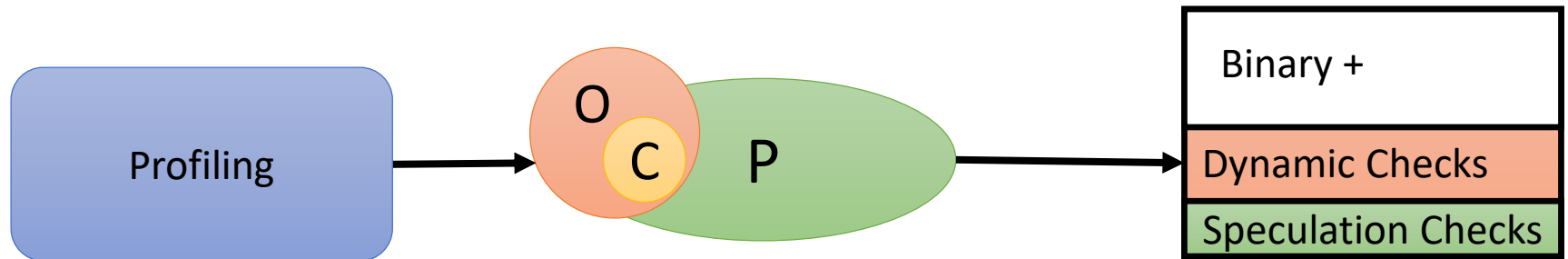


Optimistic Hybrid Analysis Phases

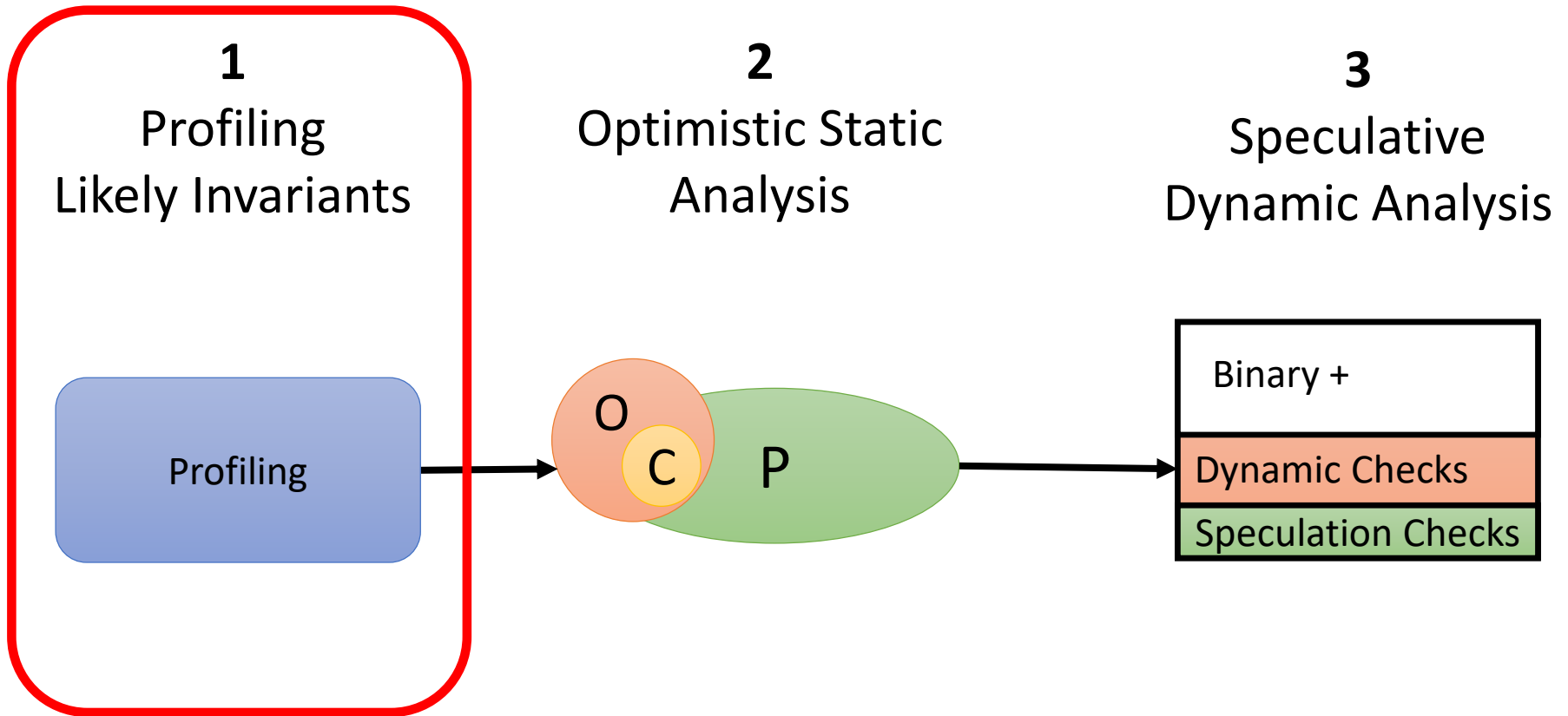
1
Profiling
Likely Invariants

2
Optimistic Static
Analysis


3
Speculative
Dynamic Analysis



Optimistic Hybrid Analysis



Likely Invariants

- **Goal:** Guide optimistic static analysis towards 
- Static analysis assumes likely invariants true
 - Reduce state space searched
 - Add **bounded** unsoundness – if invariant false
- Likely Invariants should be:
 - **Strong** - Help prune states considered by static analysis
 - **Cheap** - Be inexpensive to check dynamically
 - **Stable** - Hold true for most analyzed executions
- Gathered by profiling pass over program



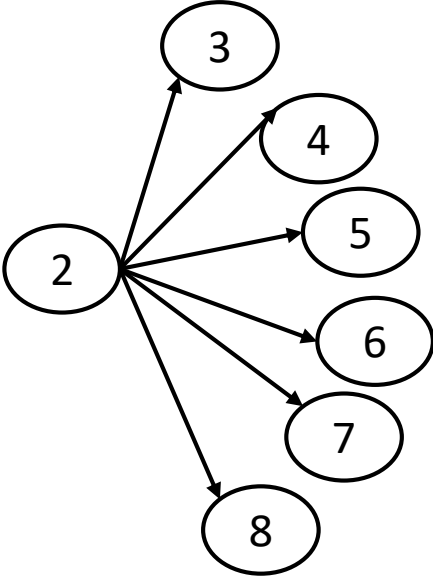
Our Invariants

- Likely Unreachable Code
 - Identifies regions of code which are not dynamically run
 - Stops static analysis from considering code
- Likely Callee Sets
 - Identifies targets of indirect function calls
 - Removes indirect function calls from analysis
- Likely Unreachable Call Contexts
 - Identifies call-stacks reached dynamically
 - Prunes call-graph in context-sensitive analysis



Likely Callee Sets

- Assume indirect calls can only call dynamically observed functions

Source Code	Static Analysis Graph
<pre>fcn() { 1: hash_fcn_t *fcn = get_hash_fcn(); 2: hash = *fcn(arg1, arg2); } 3: sha1(arg1, arg2) { ... } 4: sha256(arg1, arg2) {...} 5: sha3(arg1, arg2) {...} 6: md5 (arg1, arg2) {...} 7: ripemd(arg1, arg2) {...} 8: whirlpool(arg1, arg2) {...}</pre>	 <pre>graph LR; 2((2)) --> 3((3)); 2 --> 4((4)); 2 --> 5((5)); 2 --> 6((6)); 2 --> 7((7)); 2 --> 8((8));</pre>

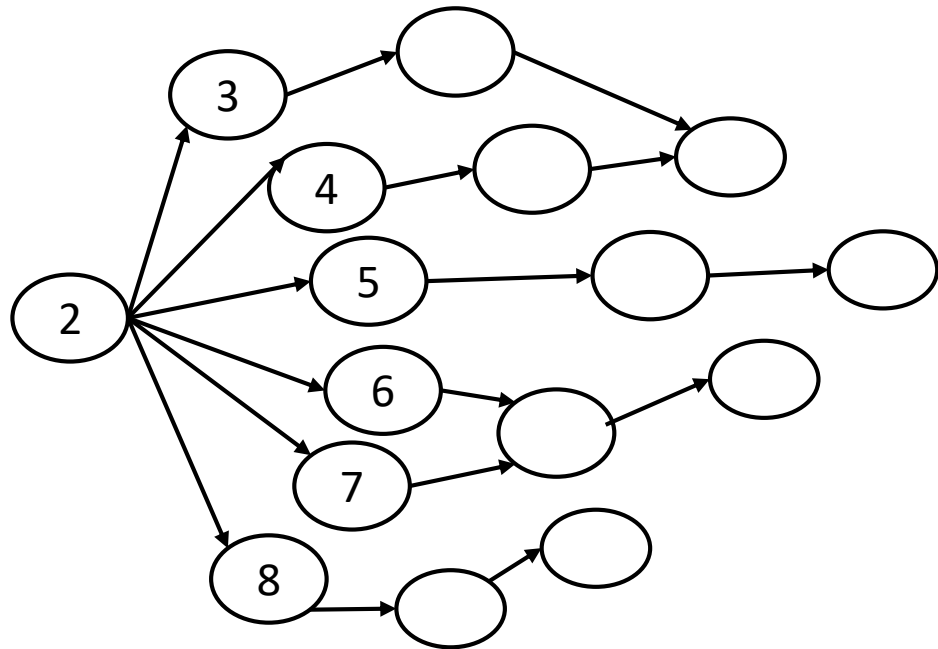


Likely Callee Sets

Source Code

```
fcn() {  
1: hash_fcn_t *fcn =  
   get_hash_fcn();  
2: hash = *fcn(arg1, arg2);  
}  
  
3: sha1(arg1, arg2) {...}  
4: sha256(arg1, arg2) {...}  
5: sha3(arg1, arg2) {...}  
6: md5 (arg1, arg2) {...}  
7: ripemd(arg1, arg2) {...}  
8: whirlpool(arg1, arg2) {...}
```

Static Analysis Graph



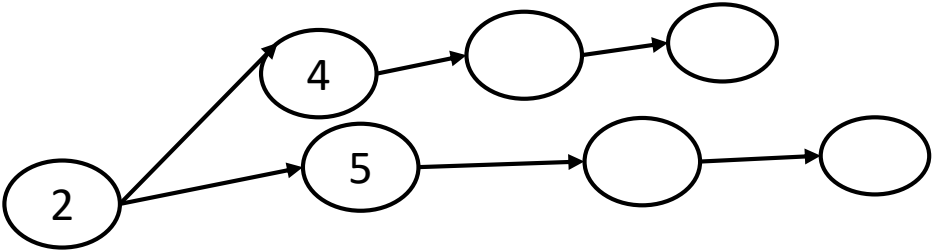
Likely Callee Sets

- Indirect Functions typically only call a small set of functions

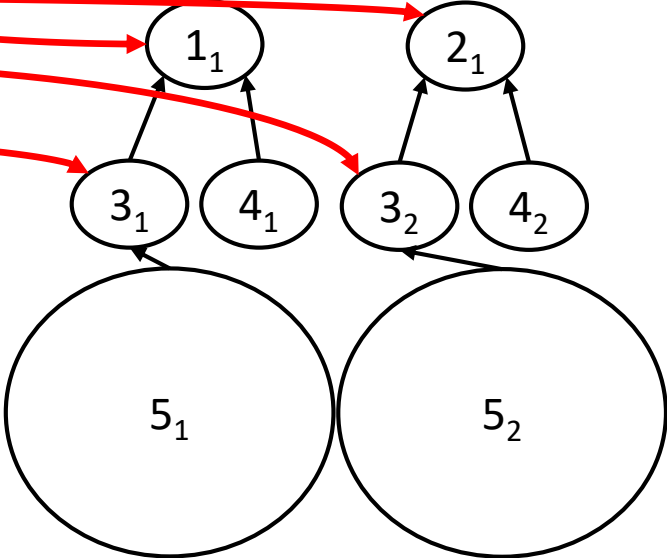
Source Code	Static Analysis Graph
<pre>fcn() { 1: hash_fcn_t *fcn = get_hash_fcn(); 2: hash = *fcn(arg1, arg2); 3: sha2(arg1, arg2) {...} 4: md5(arg1, arg2) {...} 5: ripemd(arg1, arg2) {...} 6: whirlpool(arg1, arg2) {...}</pre> <p>Likely Invariant: *fcn is sha2 or sha3</p>	<pre>graph LR 2((2)) --> 3((3)) 2 --> 4((4)) 2 --> 5((5)) 2 --> 6((6)) 3 --> 7((7)) 3 --> 8((8)) 3 --> 9((9)) 4 --> 10((10)) 4 --> 11((11)) 5 --> 12((12)) 5 --> 13((13)) 6 --> 14((14)) 6 --> 15((15))</pre>

Likely Callee Sets

- Indirect Functions typically only call a small set of functions

Source Code	Static Analysis Graph
<pre>fcn() { 1: hash_fcn_t *fcn = get_hash_fcn(); 2: hash = *fcn(arg1, arg2); 3: sha2(arg1, arg2) {...} 4: md5(arg1, arg2) {...} 5: ripemd(arg1, arg2) {...} 6: whirlpool(arg1, arg2) {...}</pre> <p>Likely Invariant: *fcn is sha2 or sha3</p>	 <pre>graph LR; 2((2)) --> 4((4)); 2 --> 5((5)); 4 --> 4a(()); 4a --> 4b(()); 5 --> 5a(()); 5a --> 5b(())</pre>

Likely Unused Call Contexts

Source Code	Analysis Graph
<pre>main() { 1: a = my_malloc(); 2: b = my_malloc();} my_malloc() { if (!g_init) 3: return do_init(); 4: return malloc(...);} do_init() { g_init = true; 5: // Long complex initialization code}</pre>	<p data-bbox="950 615 1375 696">Context-Sensitive</p>  <pre>graph TD 1_1((1_1)) --> 3_1((3_1)) 1_1 --> 2_1((2_1)) 3_1 --> 4_1((4_1)) 3_1 --> 5_1((5_1)) 2_1 --> 3_2((3_2)) 2_1 --> 4_2((4_2)) 3_2 --> 4_2 3_2 --> 5_2((5_2)) 5_1 --- 5_2</pre>

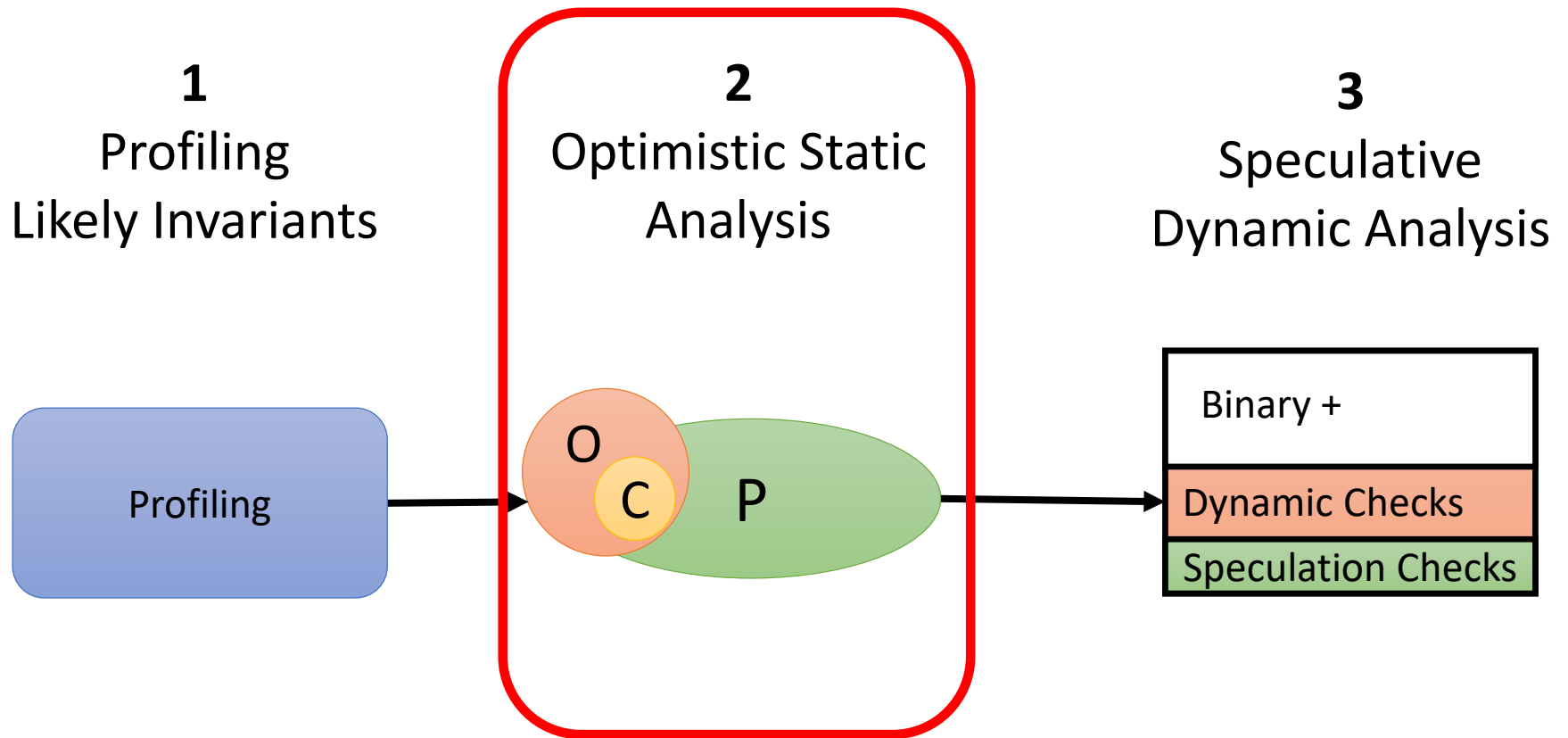


Likely Unused Call Contexts

- Remove call-stacks which never dynamically happen

Source Code	Analysis Graph
<pre>main() { 1: a = my_malloc(); 2: b = my_malloc(); my_malloc() { if (!g_init) 3: return do_init(); 4: return malloc(...); do_init() { g_init = true; 5: // Long complex initialization code }</pre>	<p>Context-Sensitive</p> <p>Dynamically Never called</p>

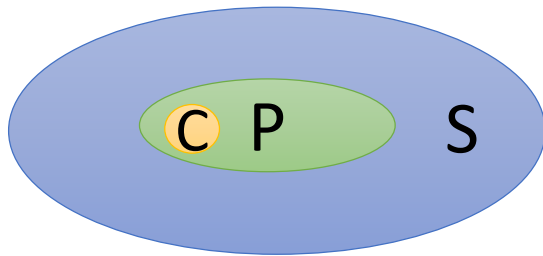
Optimistic Hybrid Analysis



Optimistic Static Analysis

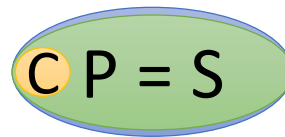
- Assumes likely invariants from profiling are true
- Likely Invariants allow analysis to approximate C
- If likely invariants hold, analysis is sound for execution

Sound Static Analysis



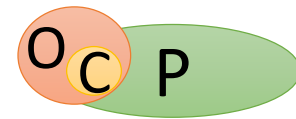
$$S > P$$

Perfect Static Analysis



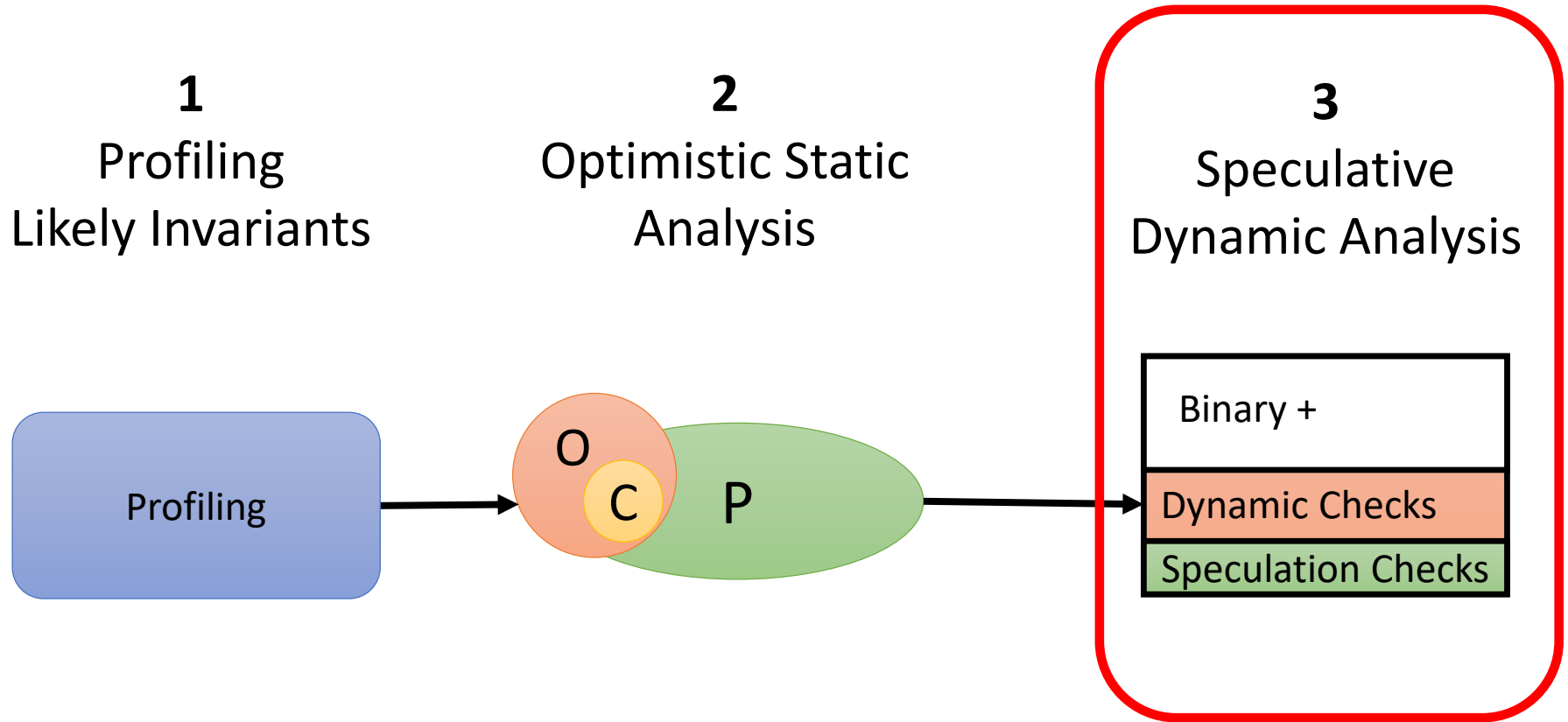
$$S = P$$

Optimistic Static Analysis



$$O < P$$

Optimistic Hybrid Analysis



Speculative Dynamic Analysis

- Runs optimized analysis
 - Fraction of original checks
 - Added likely-invariant speculation checks
 - Lightweight by design (by likely-invariant choices)
- May have to roll-back on mis-speculation
 - Mis-speculations are rare
 - Worst case: full-execution roll-back re-execute
 - On re-execute – use conservative dynamic analysis
 - Eg. A traditional hybrid analysis



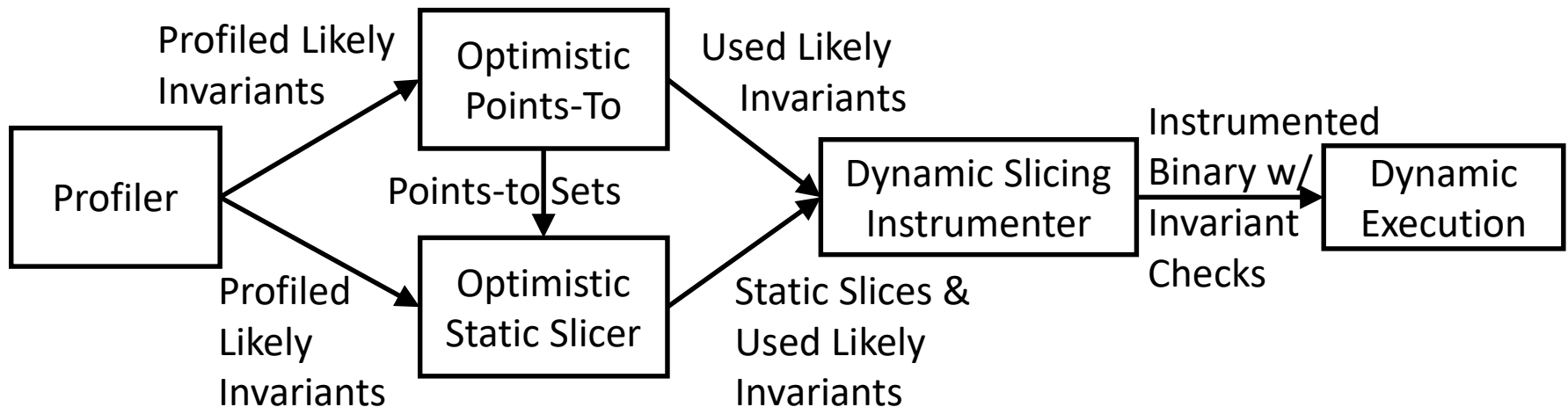
Overview

- ~~Motivation~~
- ~~Optimistic Hybrid Analysis Overview~~
 - ~~Likely Invariants~~
 - ~~Optimistic Static Analysis~~
 - ~~Speculative Dynamic Analysis~~
- OptSlice – Optimistic Dynamic Backward Slicer
- Evaluation
- Conclusion



OptSlice

- Optimistic Dynamic Backwards Slicer
- Uses three likely invariants discussed
- Uses optimistic static points-to and slicing analyses
- Heavyweight (full-execution) roll-back recovery



OptSlice – Misc. Details

- Static analyses have:
 - Call context-sensitive and insensitive versions
 - Use the most accurate version that will scale
 - Optimistic Static Analysis processes less states
 - Generally more scalable
- Built on-top of Giri dynamic slicer (ASPLOS 13)

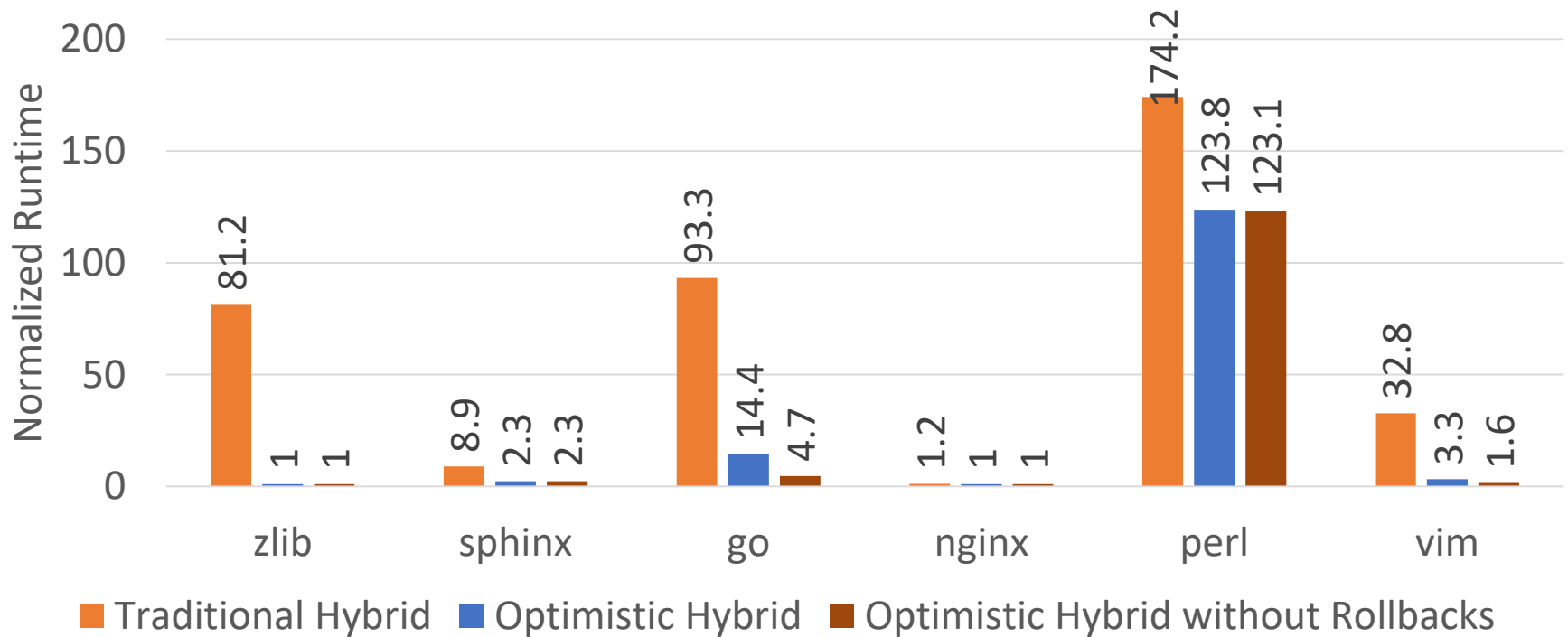


Evaluation

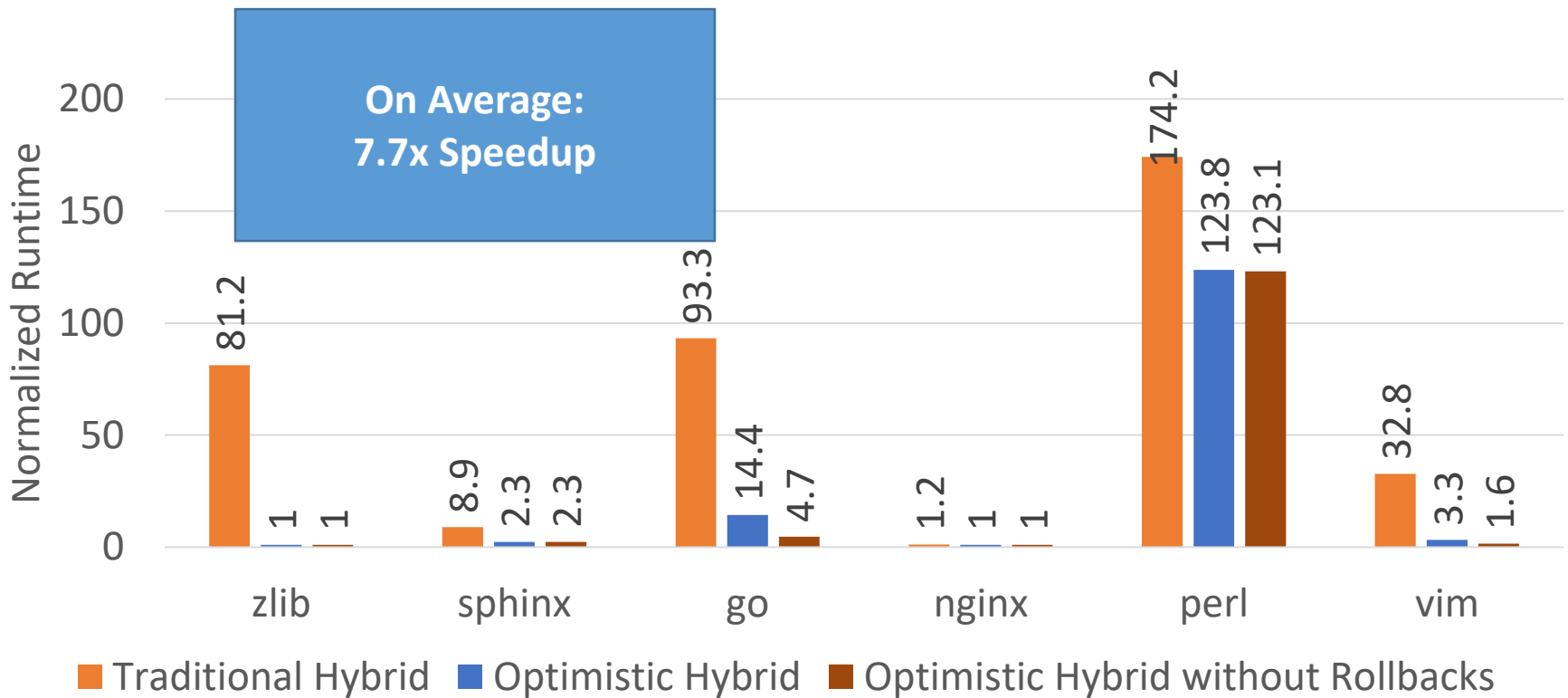
- All experiments bounded to 8hrs, 16GB RAM
 - We use the best static analysis that will complete
- Analysis built in LLVM-3.1
- **Goal:** Show how OHA techniques enable OptSlice to accelerate dynamic analysis
- Evaluate three portions
 - OptSlice speedup
 - Invariant profiling requirements
 - Static analysis improvements



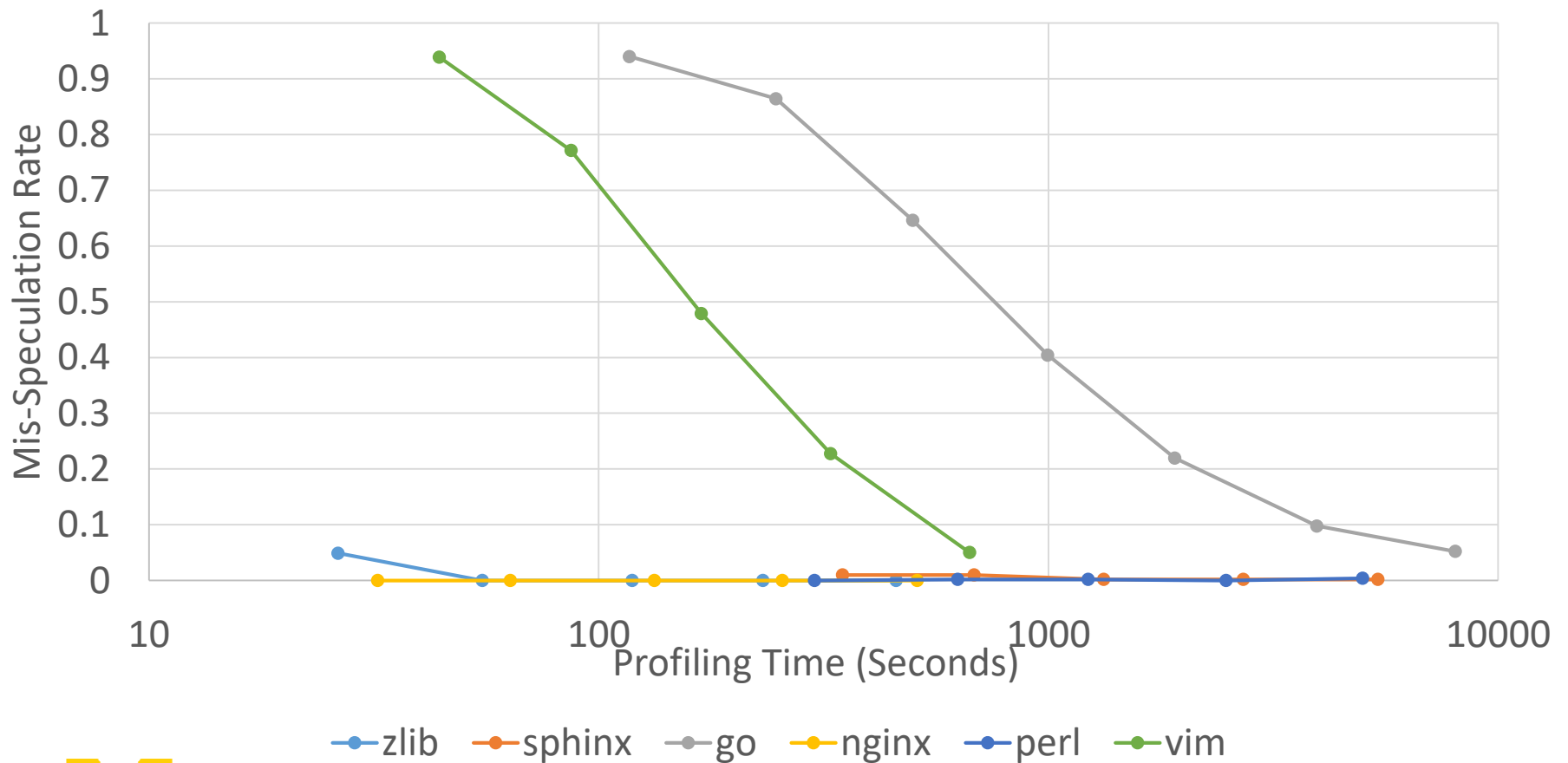
Runtime



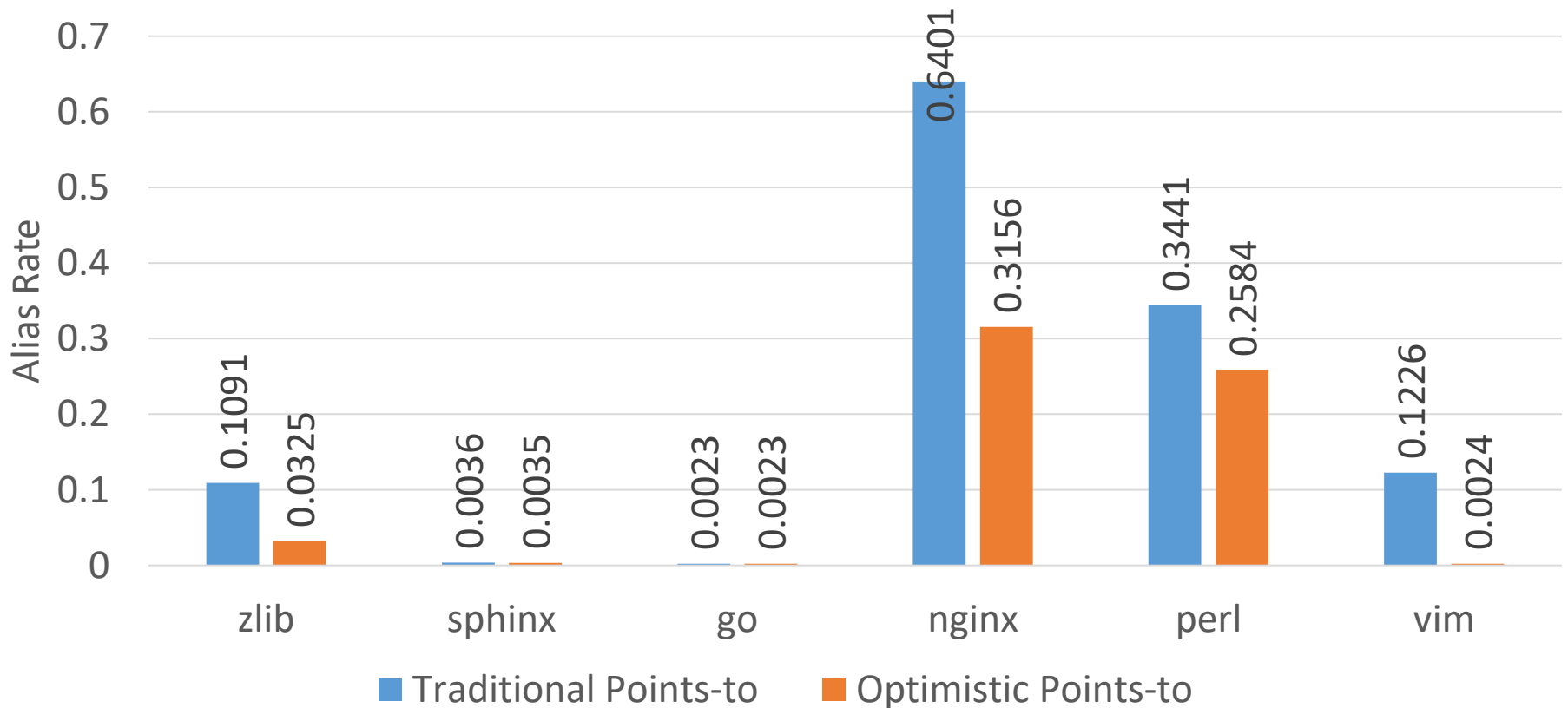
Runtime



Profiling



Static Points-to Analysis



Conclusion

- Optimistic Hybrid Analysis
 - Novel analysis methodology: combines **unsound** static analysis, and **speculative** dynamic analysis
 - Accelerates dynamic analysis without sacrificing accuracy or soundness
- OptSlice
 - Optimistic Hybrid Backward Slicer
 - 7.7x average speedup versus traditional hybrid slicing



Questions?

