# Large Spurious Cycle in Global Static Analyses and Its Algorithmic Mitigation

Hakjoo Oh
pronto@ropas.snu.ac.kr

School of Computer Science and Engineering
Seoul National University
Korea

Dec 14, 2009
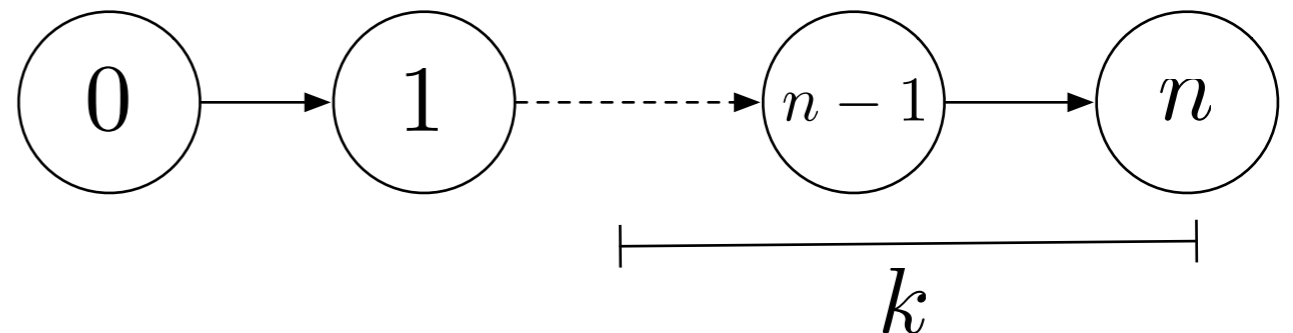Asian Symposium on Programming Language and Systems
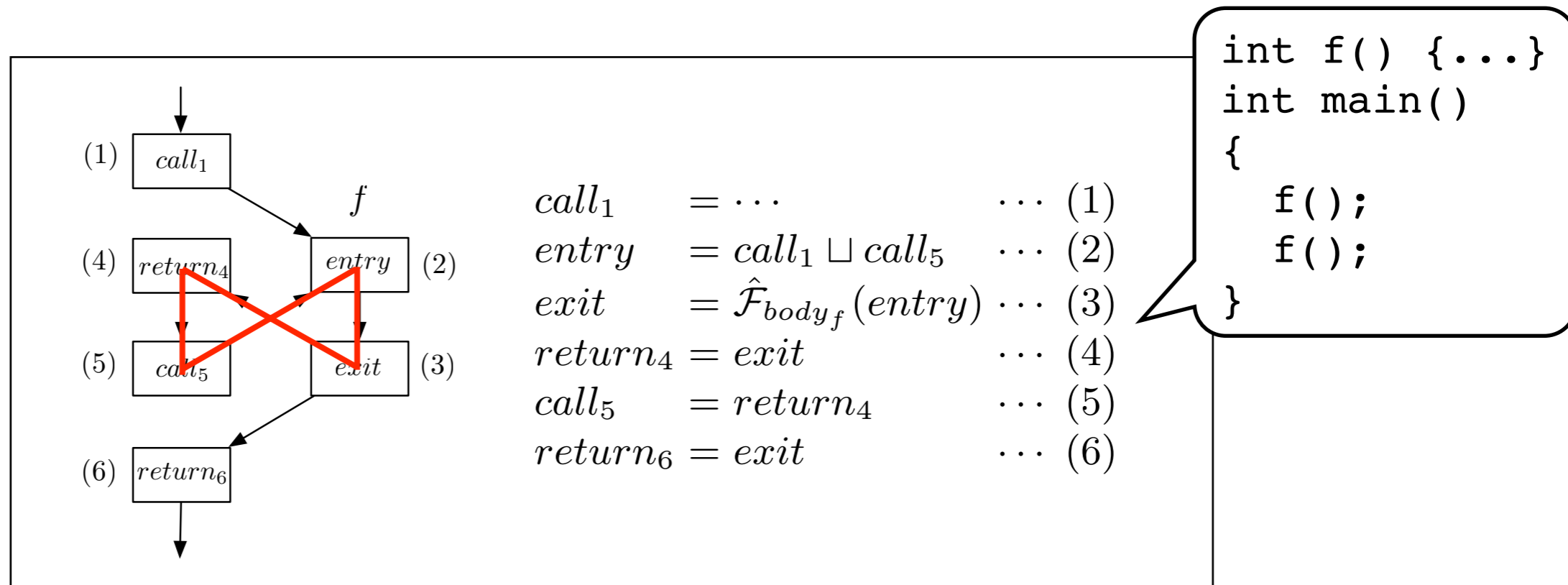
A performance problem is identified and solved.

An extension of the classical call-strings approach

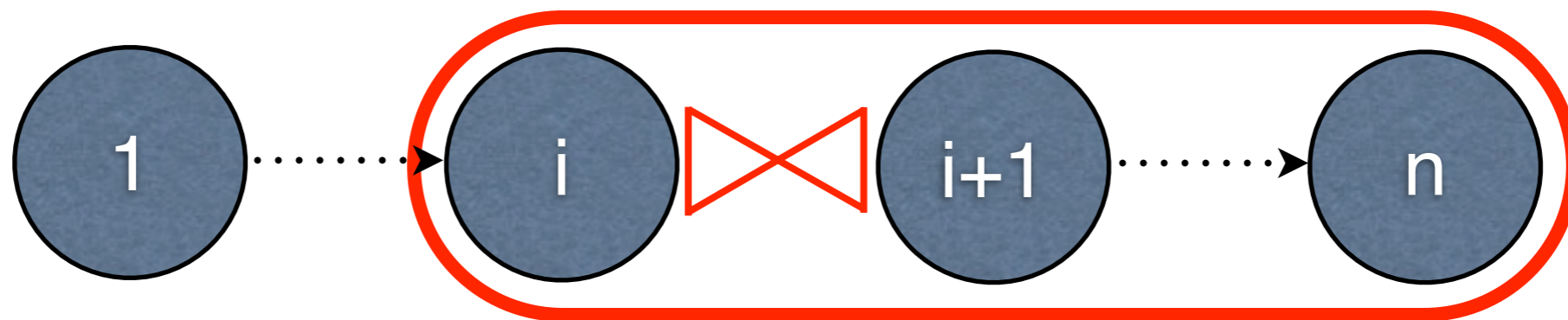Conventional context-sensitive analysis, distinguishing the last $k$ call-sites to each procedure (k-limiting).

$$time_\rho \Downarrow \quad precision \uparrow$$

$$Normal_k \longrightarrow RSS_k$$

$$0 \longrightarrow 1 \dashrightarrow n-1 \longrightarrow n$$

$k$

# Spurious Cycle in Static Analysis

```
int f() {...}
int main()
{
    f();
    f();
}
```

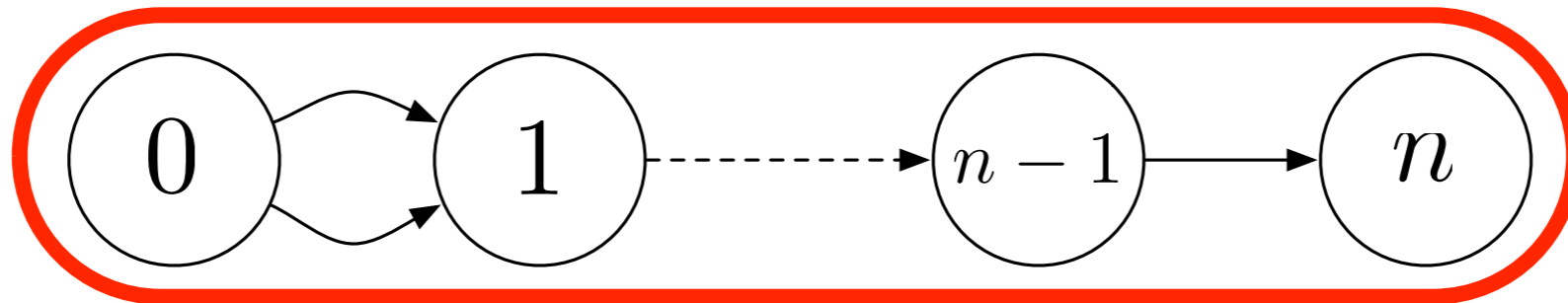- Spurious cycles degrade both precision and time
- Spurious information flows
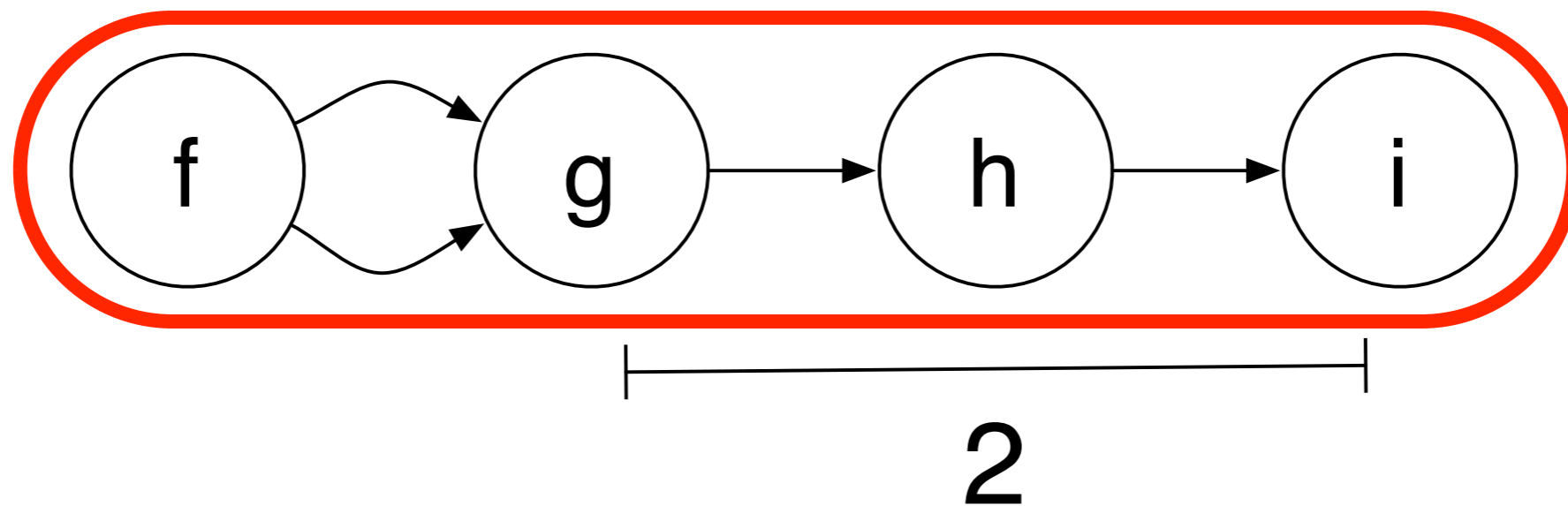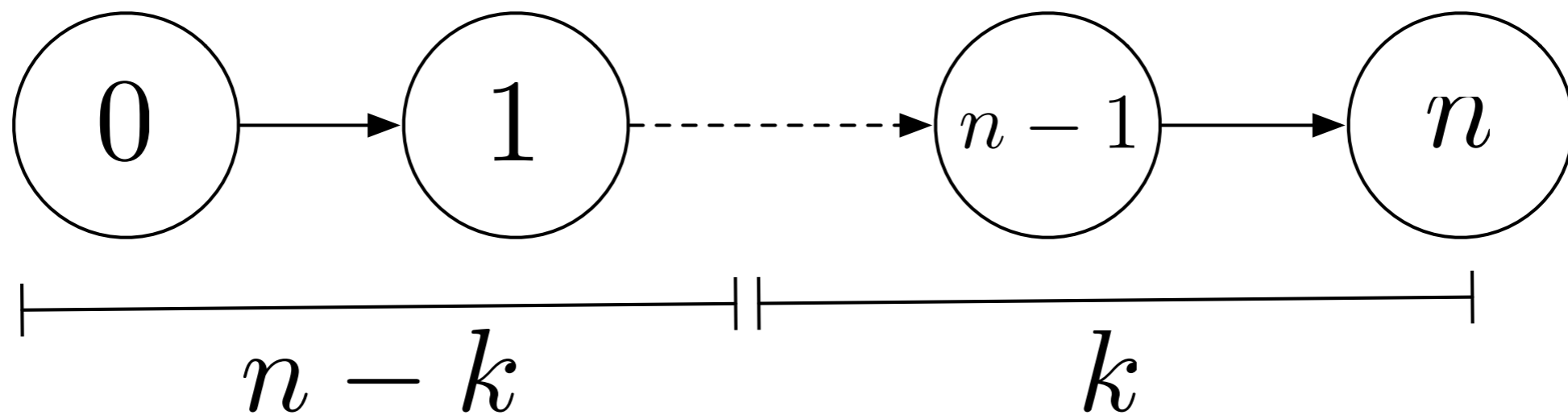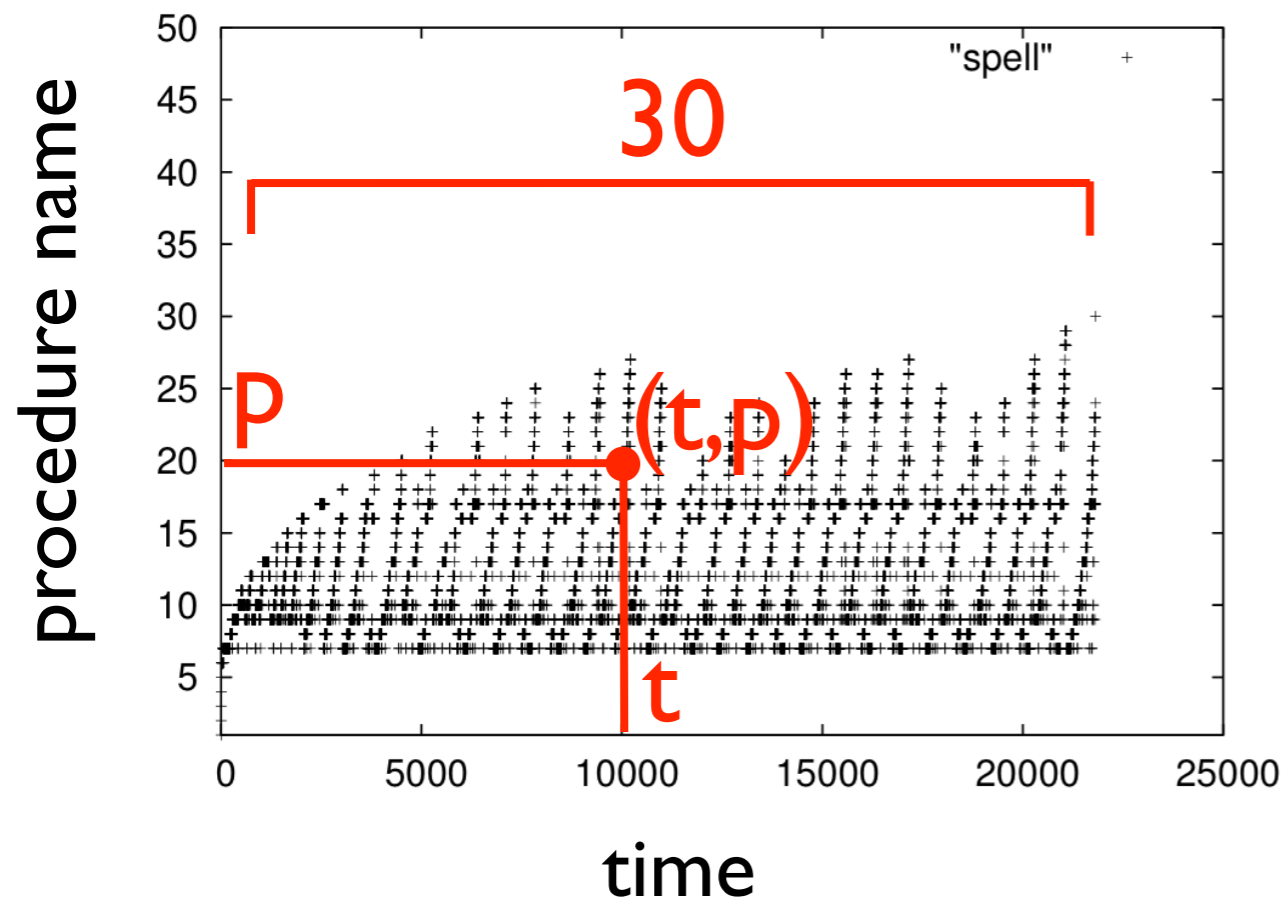- Fake cyclic dependence cycle



3 / 20

# Easy To Be Large

SEOUL NATIONAL UNIVERSITY

| | | |
|---|---|---|
| bison-1.875 | 410/832(49%) | 12,558/18,110(69%) |
| proftpd-1.3.1 | 940/1,096(85%) | 35,386/41,062(86%) |
| apache-2.2.2 | 1,364/2,075(66%) | 71,719/95,179(75%) |

$k$-limiting is not much effective to mitigate the problem.



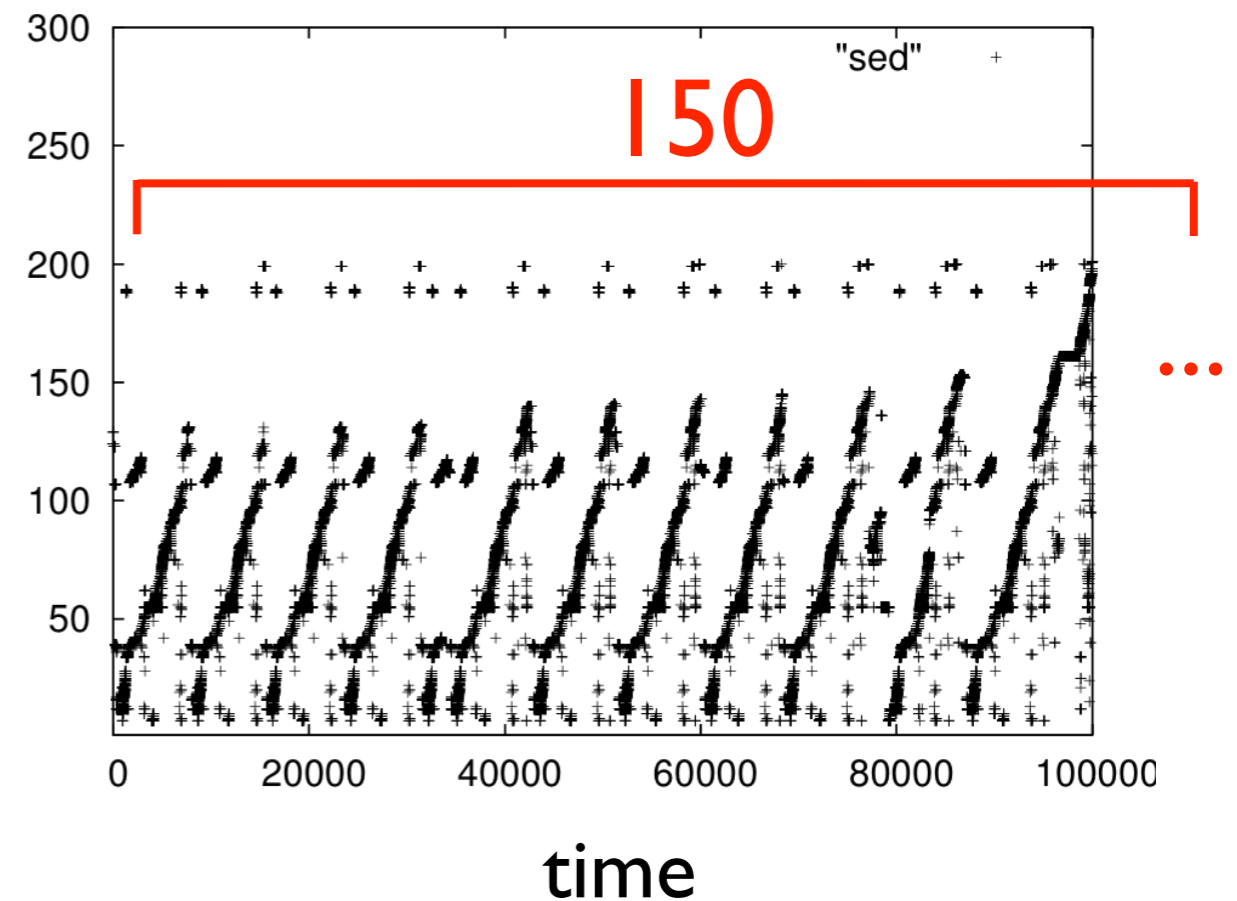| Program | Basic-blocks in the largest cycle |
|---|---|
| spell-1.0 | 751/782(95%) |
| gzip-1.2.4a | 5,988/6,271(95%) |
| sed-4.0.8 | 14,559/14,976(97%) |
| tar-1.13 | 10,194/10,800(94%) |
| wget-1.9 | 15,249/16,544(92%) |
| bison-1.875 | 12,558/18,110(69%) |
| proftpd-1.3.1 | 35,386/41,062(86%) |
| apache-2.2.2 | 71,719/95,179(75%) |

# Increasing $k$ Does Not Help

# Performance Problems due to Large Spurious Cycle
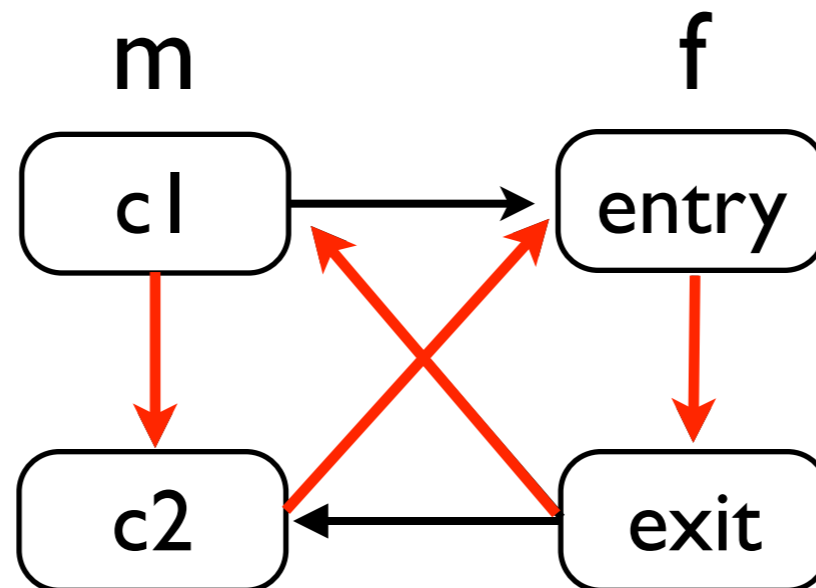
spell-1.0

sed-4.0.8

procedure name

"spell"

30

P

(t,p)

t

time

"sed"

150

...

time

Programming Research Laboratory

- In the paper,
- Generalization for k > 0
- Details on algorithm, correctness, and precision
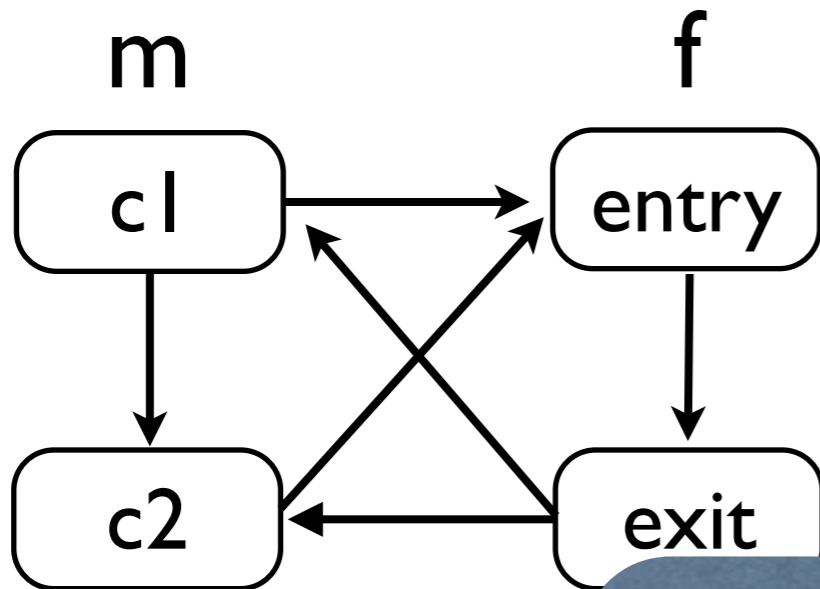
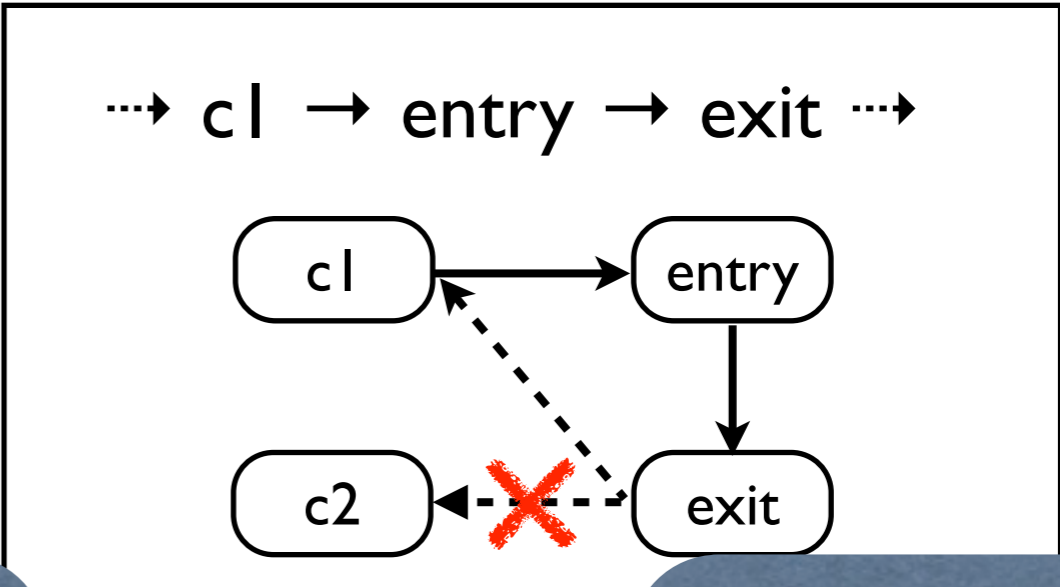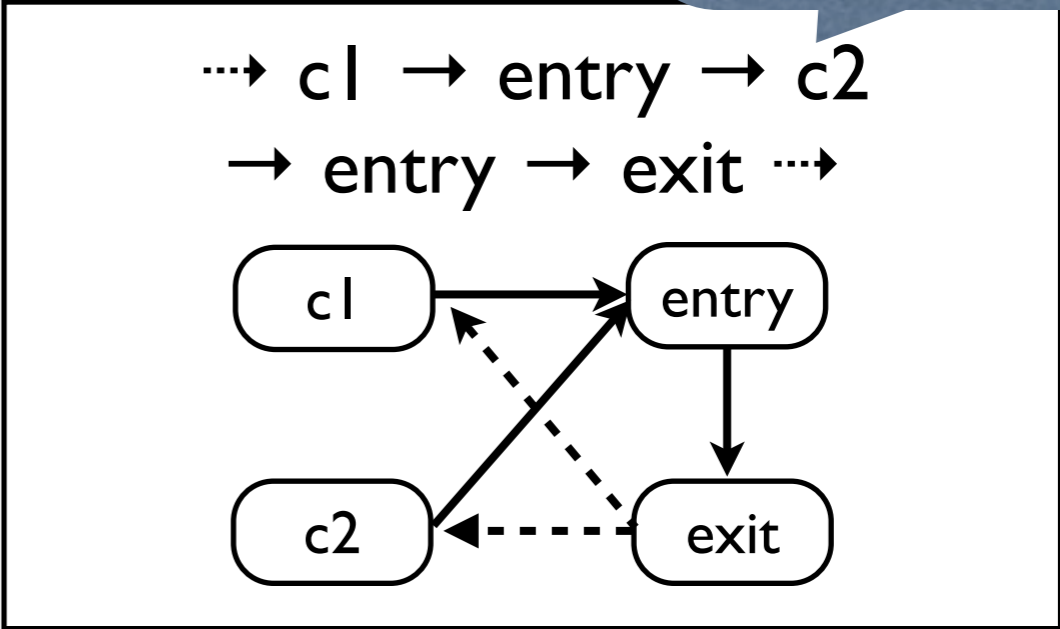# Worklist-based Normalo Algorithm

m

f



c1 → entry

c2 → exit

→ c2 → entry → exit → c1 → c2 → ...

cycle

# Observation

$\cdots\to$ cl $\to$ entry $\delta\mapsto$ exit $\rho\to$

$\rho\vee\epsilon$

cl $\to$ entry

c2 $\;\times\;$ exit

**single return sequence**

**multiple return sequence**

$\cdots\to$ cl $\to$ entry $\to$ c2
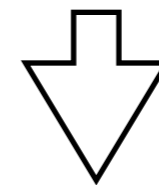$\to$ entry $\to$ exit $\cdots\to$

cl $\to$ entry

c2 $\;\times\;$ exit

# Basic Idea

1. Change M.R.S. into S.R.S
2. Enforce single return to the last called site
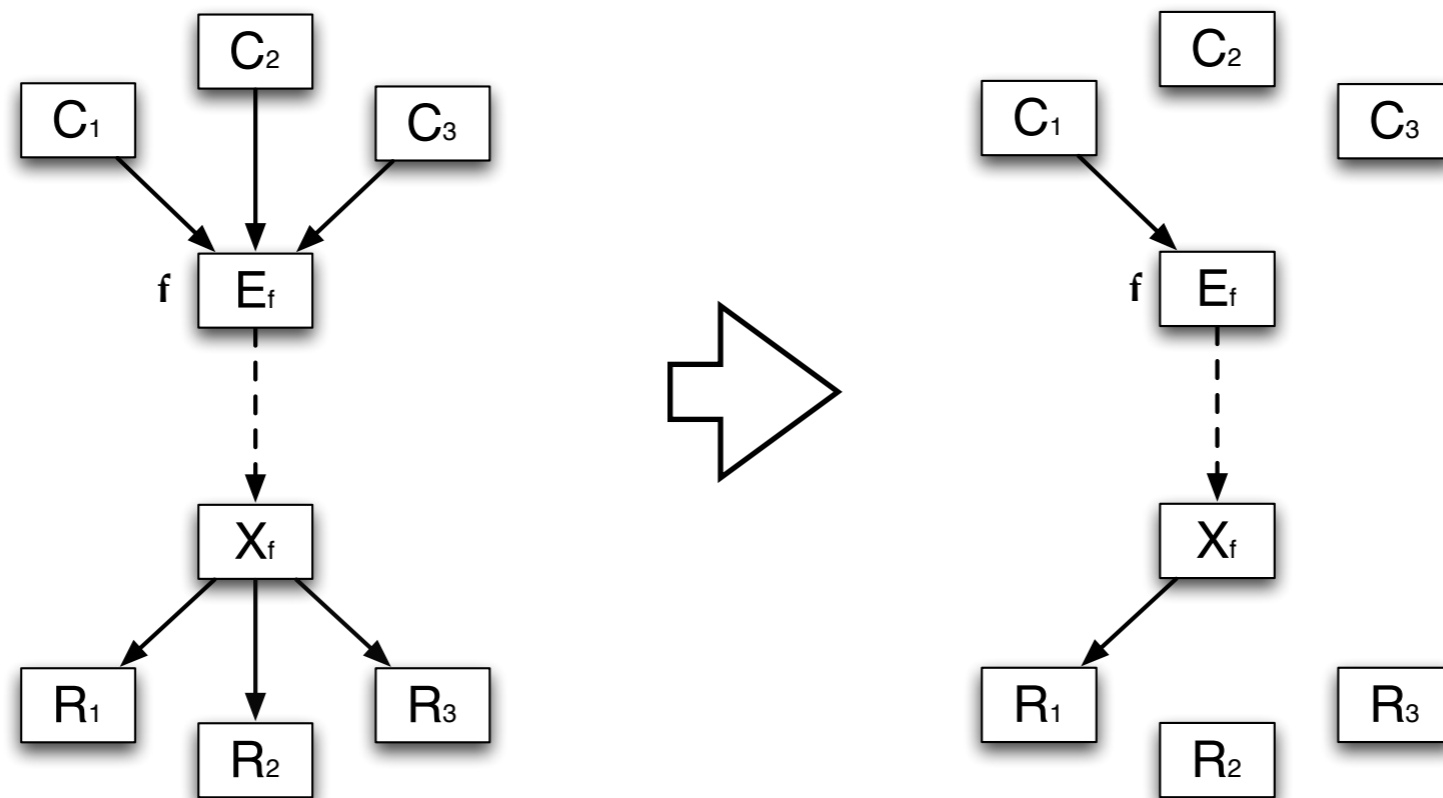
**multiple return sequence**

$c1 \rightarrow entry \rightarrow c2 \rightarrow entry \rightarrow exit \cdots$

$c1 \rightarrow entry \rightarrow exit \rightarrow c2 \rightarrow entry \rightarrow exit \cdots$

**single return sequence**

**single return sequence**

# One-call-per-procedure



- A procedure is analyzed exclusively for its one particular call-site.
- Each called procedure is locked.
- The other call-sites wait until the lock is released.

# Controlling Worklist

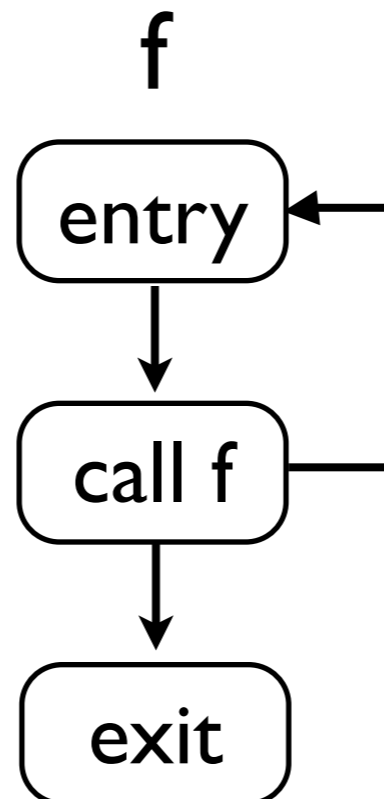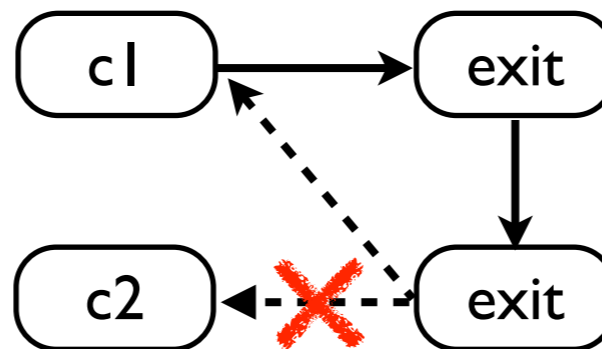Prioritize a callee procedure over its call-sites

For example,

$F \xrightarrow{c1} G \xrightarrow{c2} H$

$W = \{c2, n_H\}$ → select $n_H$ first

$W = \{c2, n_H\}$

# Recursion Handling

- Recursive procedures are handled in the same way as the normal algorithm.
- We cannot finish analyzing a recursive procedure without considering other calls in it.

f

entry

call f

exit

- The result is not a fixpoint of the given equation system.

- But still a sound approximation of the program semantics.

# Summary

Control worklist

Recursion handling

Enforce single return

# Experiments

- Sparrow is an
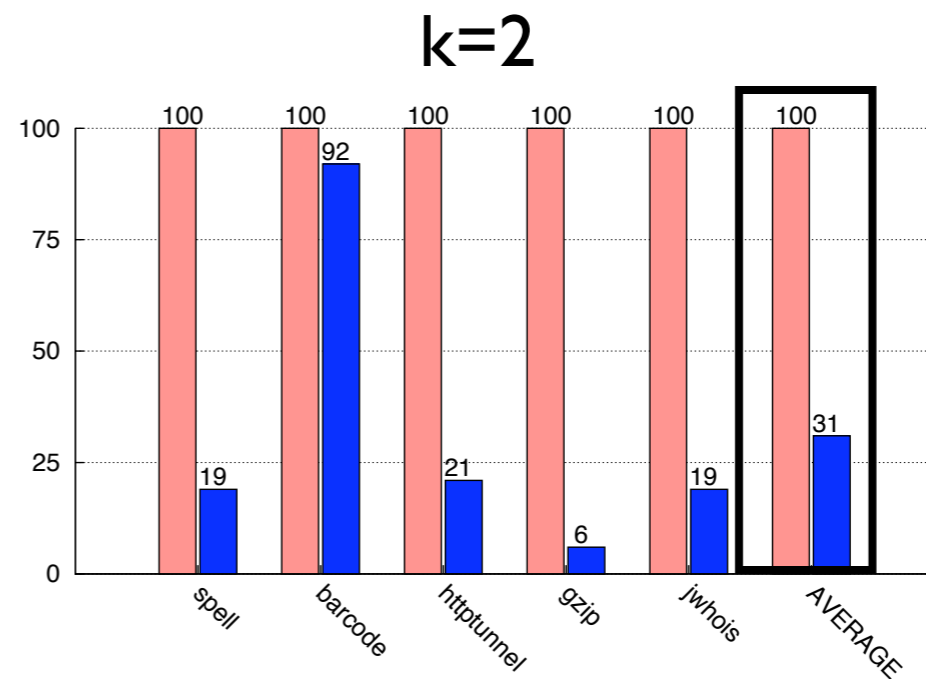- industrialized static analyzer
- interval-domain-based abstract interpreter
- open-source software packages

$\text{Normal}_k$ vs. $\text{RSS}_k$

$\text{Normal}_k$ vs. $\text{RSS}_{k+1}$

| Program | LOC | #Basic-Blocks |
|---|---|---|
| spell-1.0 | 2,213 | 782 |
| barcode-0.96 | 4,460 | 2,634 |
| httptunnel-3.3 | 6,174 | 2,757 |
| gzip-1.2.4a | 7,327 | 6,271 |
| jwhois-3.0.1 | 9,344 | 5,147 |
| parser | 10,900 | 9,298 |
| bc-1.06 | 13,093 | 4,924 |
| less-290 | 18,449 | 7,754 |
| twolf | 19,700 | 14,610 |
| tar-1.13 | 20,258 | 10,800 |
| make-3.76.1 | 27,304 | 11,061 |

# Normal$_k$ vs. RSS$_k$

## k=0

spell 100 / 14
barcode 100 / 91
httptunnel 100 / 9
gzip 100 / 14
jwhois 100 / 30
parser 100 / 5
bc 100 / 15
less 100 / 44
twolf 100 / 4
make 100 / 51
tar 100 / 13
AVERAGE 100 / 26

## k=1

spell 100 / 24
barcode 100 / 93
httptunnel 100 / 33
gzip 100 / 12
jwhois 100 / 50
bc 100 / 28
AVERAGE 100 / 40

## k=2

spell 100 / 19
barcode 100 / 92
httptunnel 100 / 21
gzip 100 / 6
jwhois 100 / 19
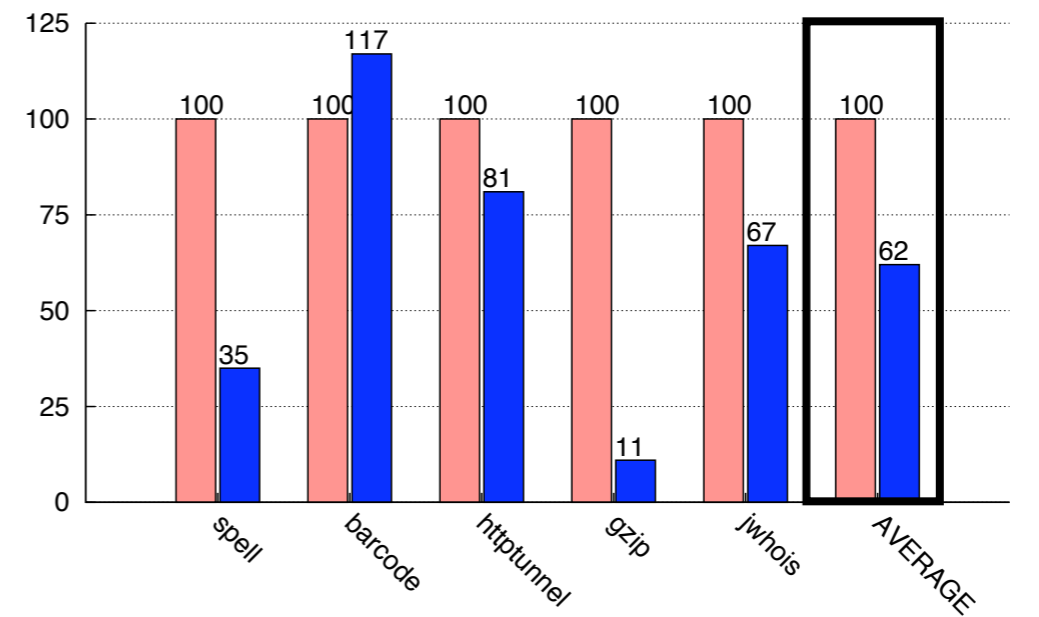AVERAGE 100 / 31

SEOUL NATIONAL UNIVERSITY
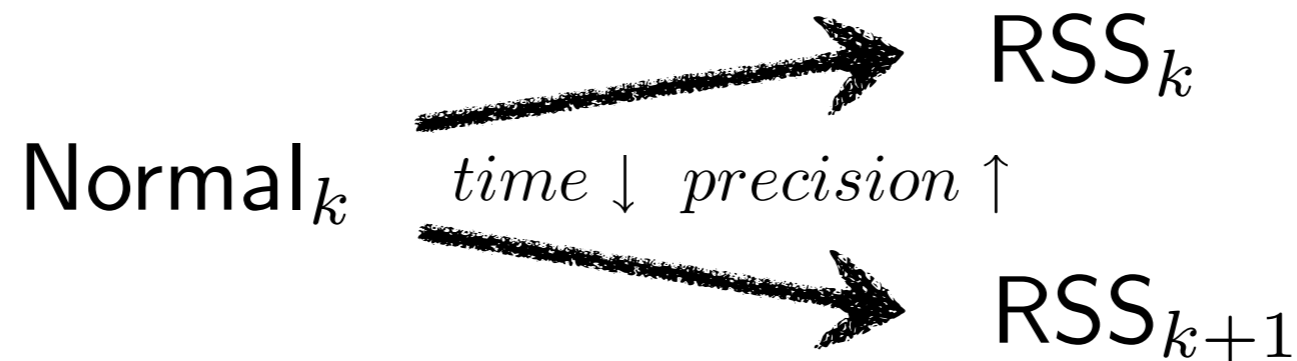
# Normal vs. RSS

## k=0



## k=1



Normal
RSS

# Conclusion

- One key reason why less accurate context-sensitivity makes the analysis very slow.

- A simple algorithm that mitigates the problem.

$RSS_k$

$RSS_{k+1}$

$Normal_k$

$time \downarrow \quad precision \uparrow$

# Thank you