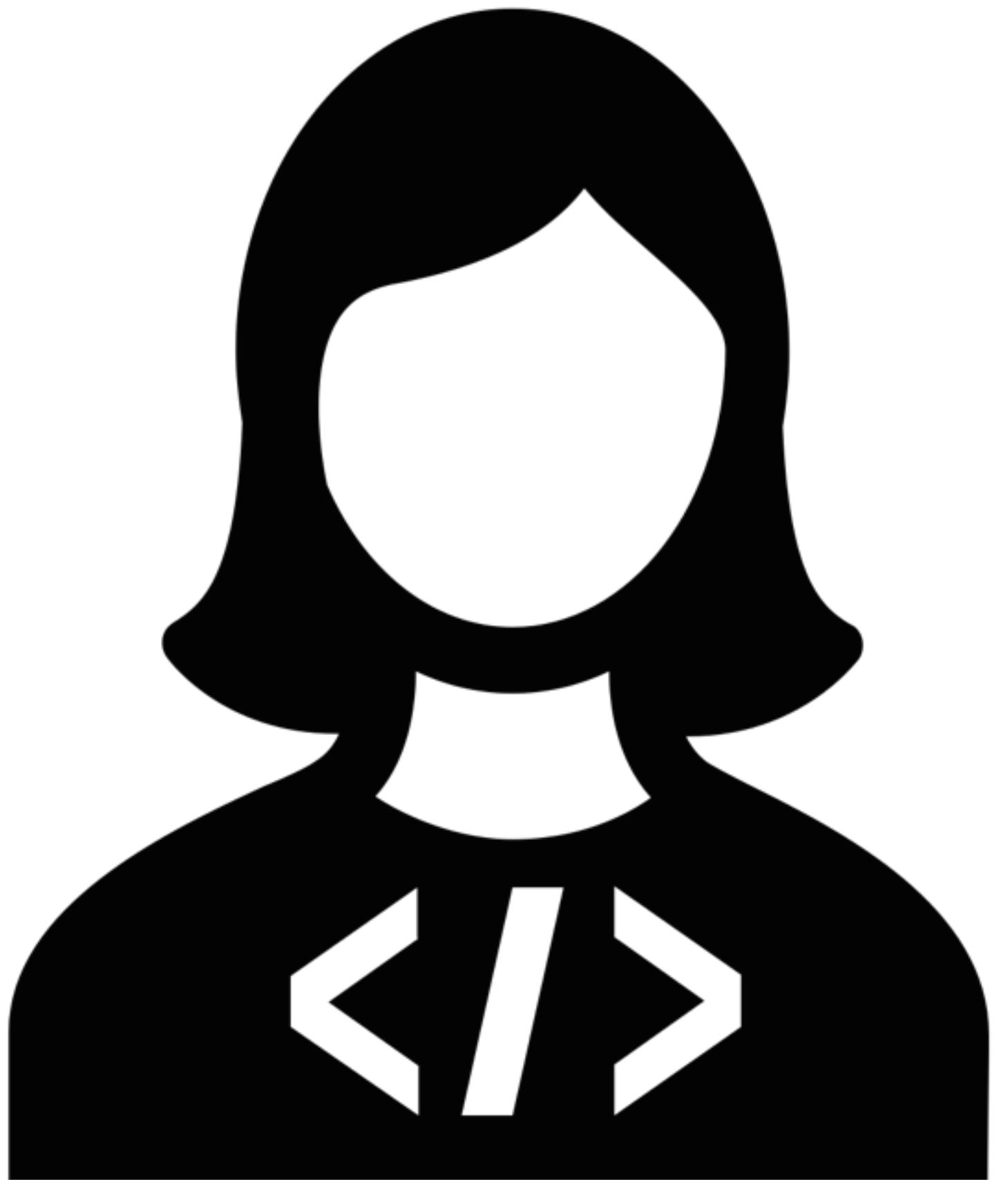
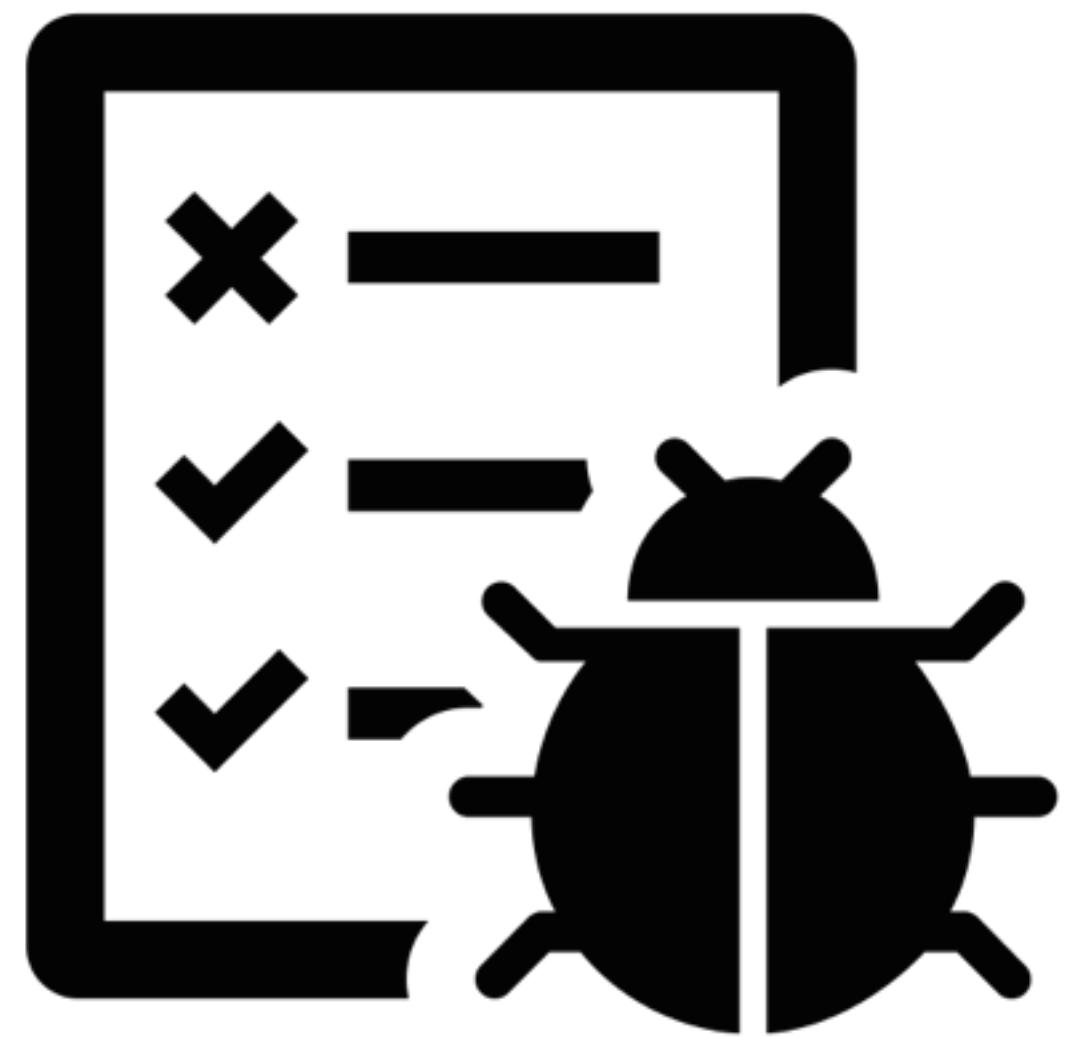


Minimizing Faulty Executions of Distributed Systems

Colin Scott, Aurojit Panda, Vjekoslav Brajkovic, George Necula,
Arvind Krishnamurthy, Scott Shenker

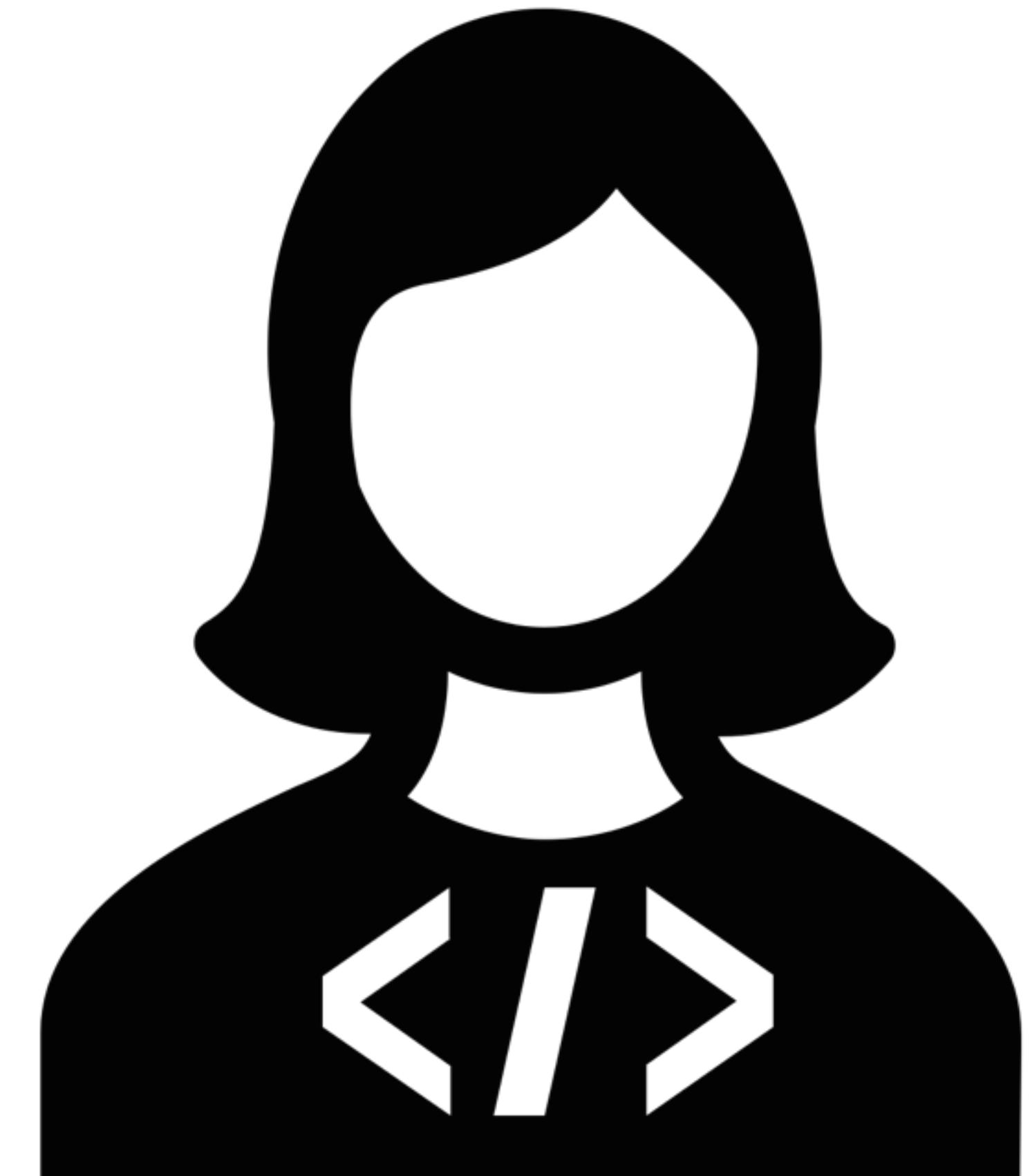




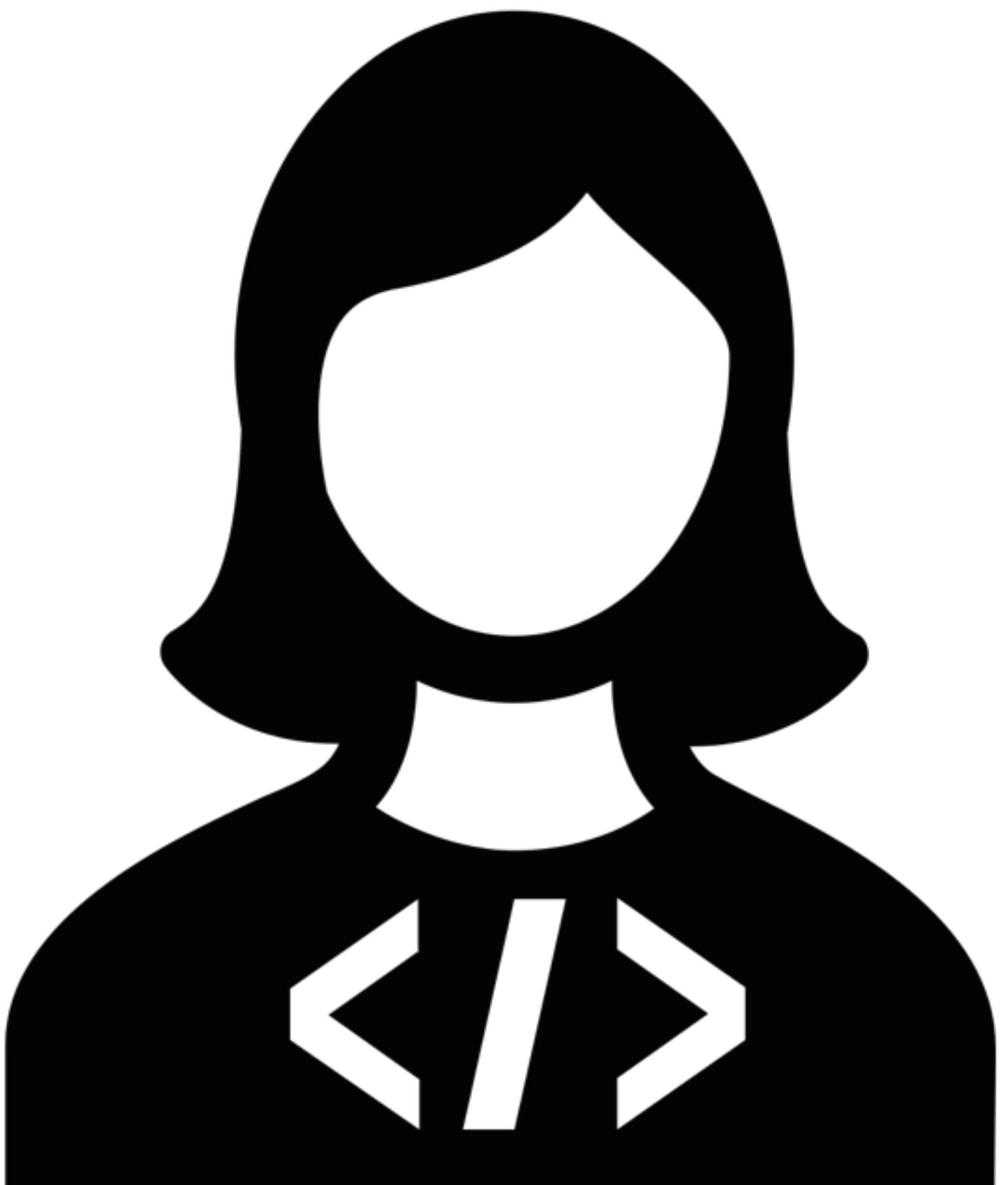
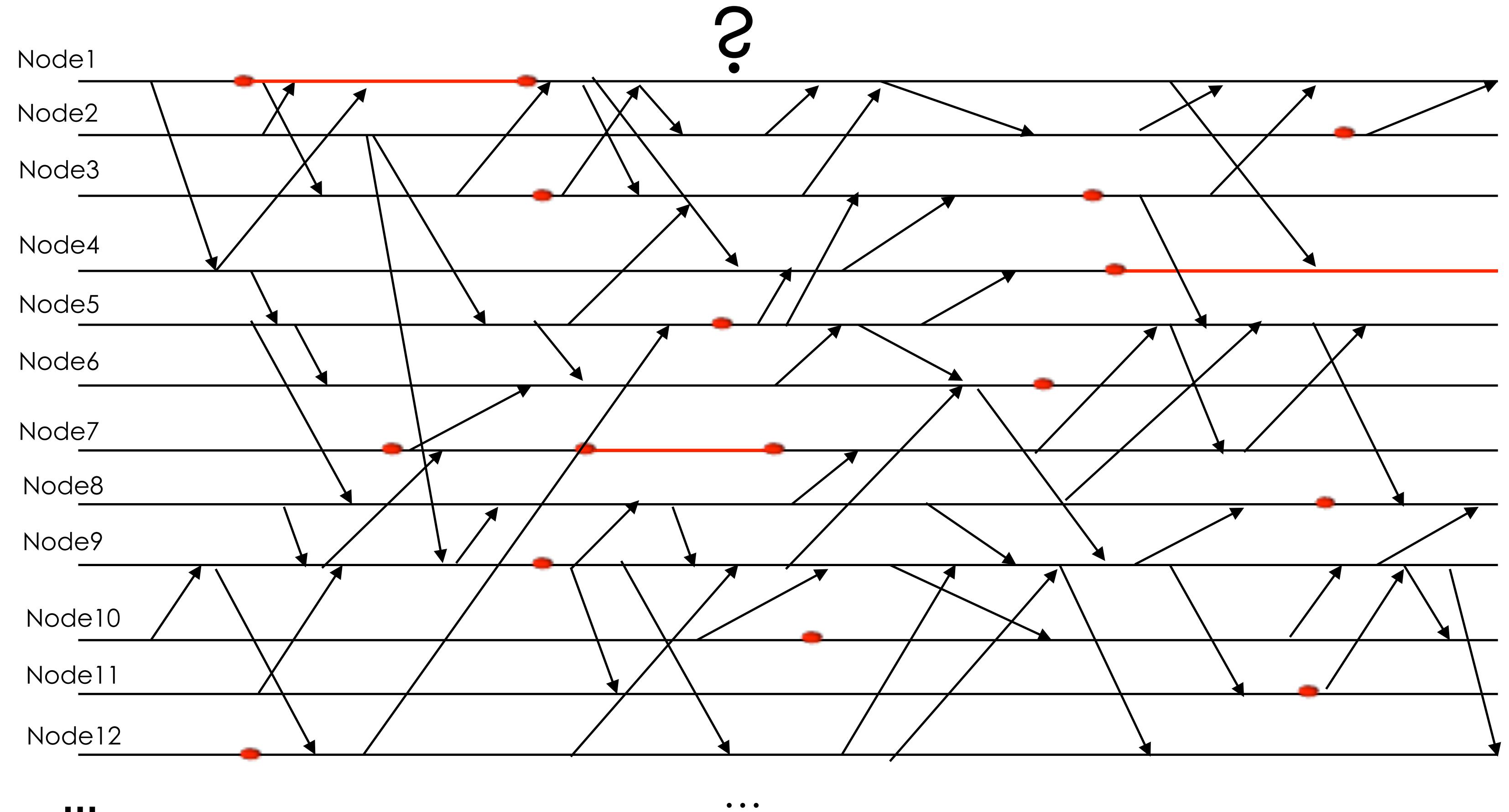
Software Developer



(GBs)

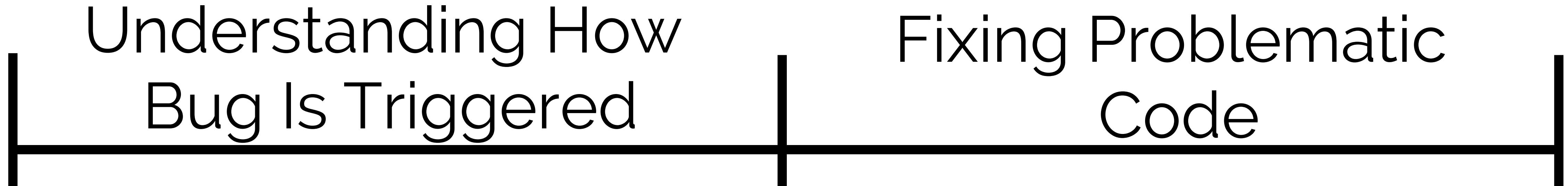


Software Developer



Software Developer

~49% of developers' time
spent on debugging!¹



¹ LaToza, Venolia, DeLine, ICSE' 06

Our Goal

Allow Developers To Focus on
Fixing the Underlying Bug

Problem Statement

Identify a minimal causal sequence of events that triggers the bug

Why Minimization?

Smaller event traces are
easier to understand

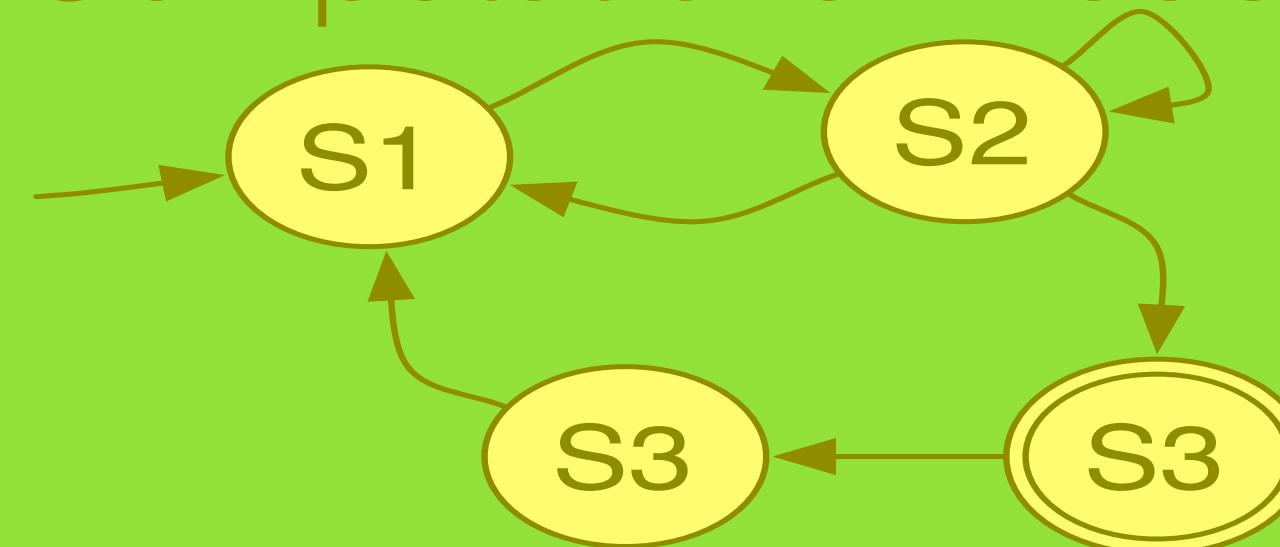
G. A. Miller. The Magical Number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information. Psychological Review '56.

Outline

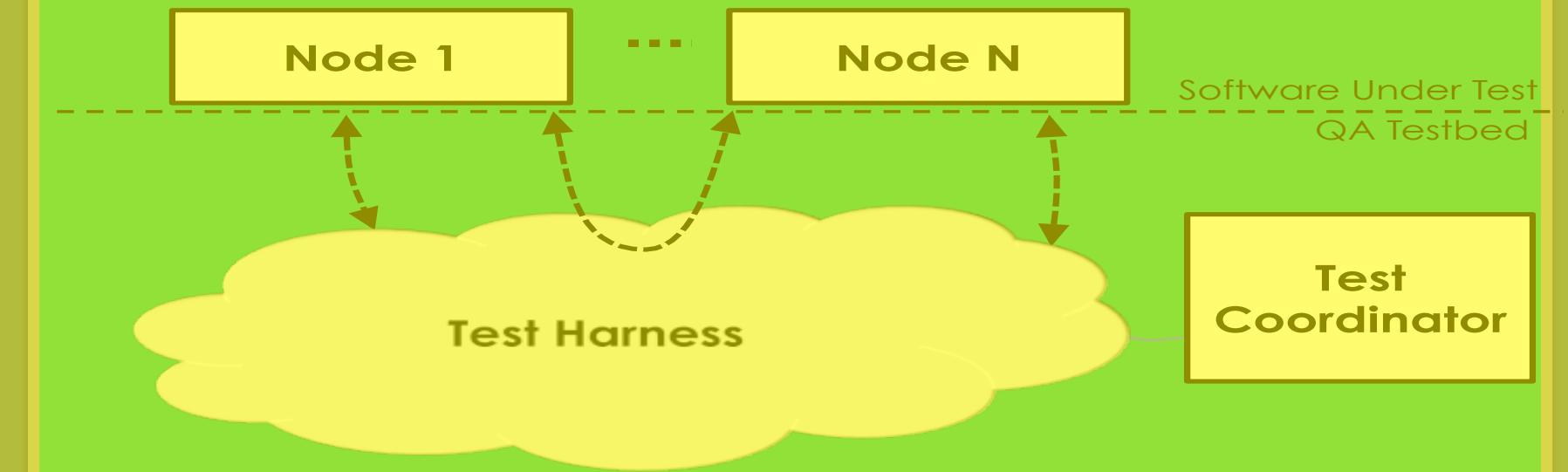
Introduction

Background

Computational Model



Fuzz Testing w/ DEMi



Minimization

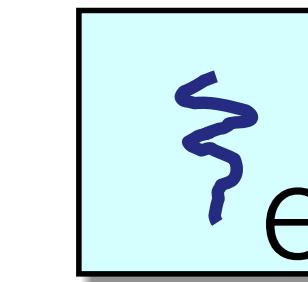
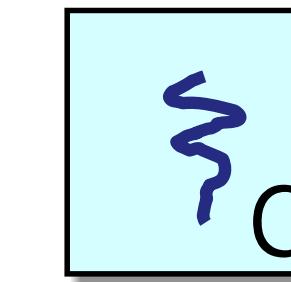
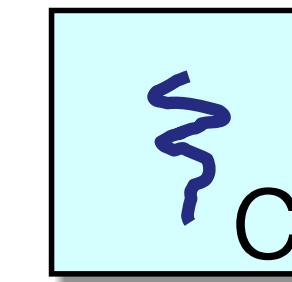
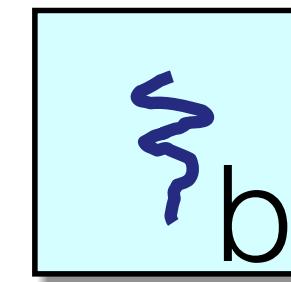
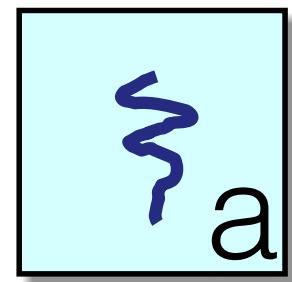


Evaluation



Conclusion

Computational Model

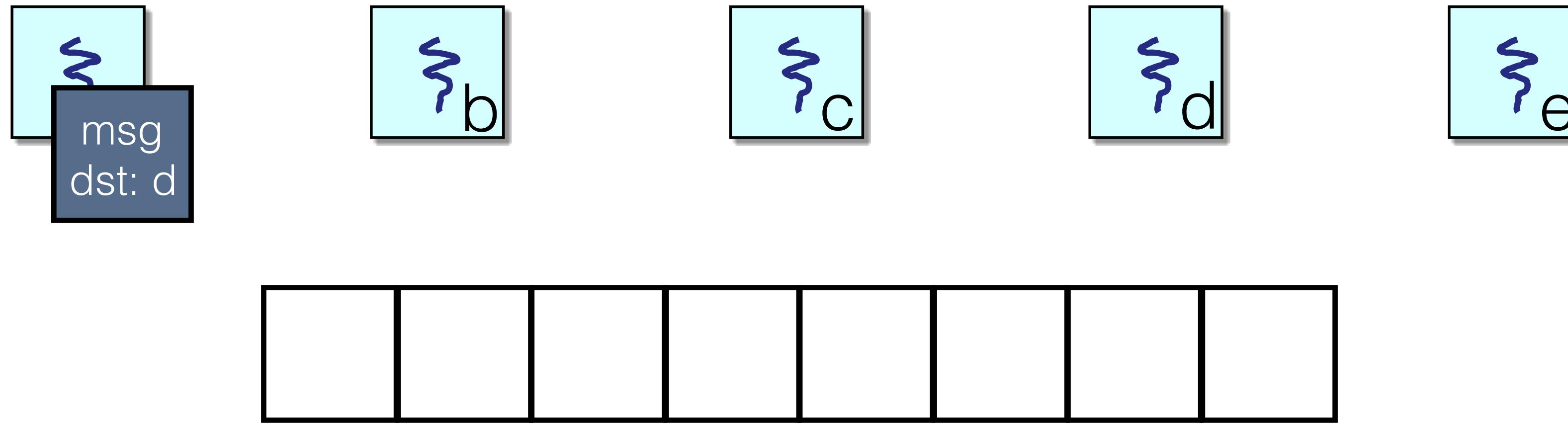


Distributed System: Collection of N processes

Each process p:

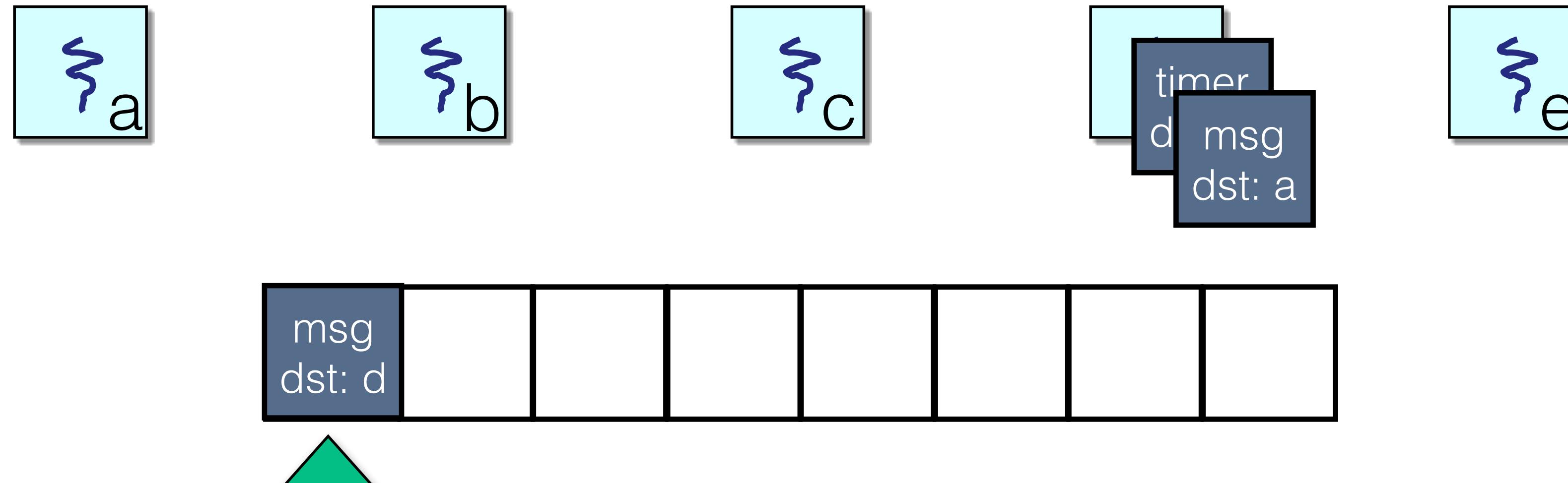
- ▶ Has unbounded memory
- ▶ Starts in a known initial state
- ▶ Changes states deterministically

Computational Model



The network maintains a buffer of sent but not yet delivered messages

Computational Model

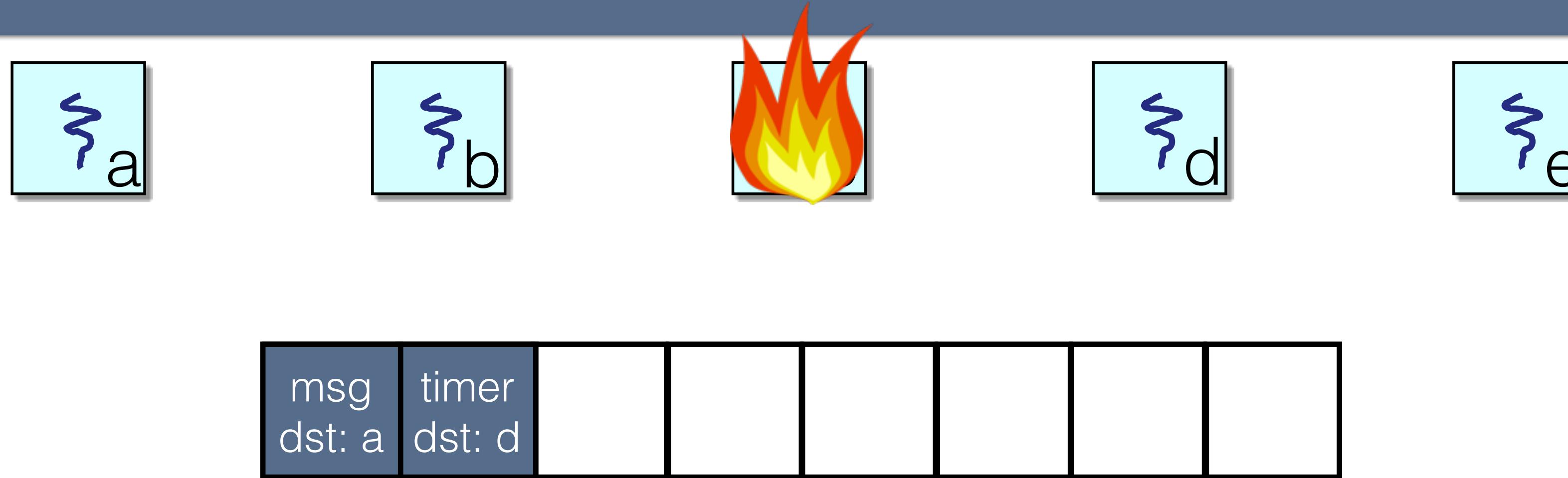


Message deliveries occur one at a time:

- destination enters a new state according to old state & message
- destination sends a finite set of messages to other processes*

*May include timer messages to be delivered to itself later

Computational Model

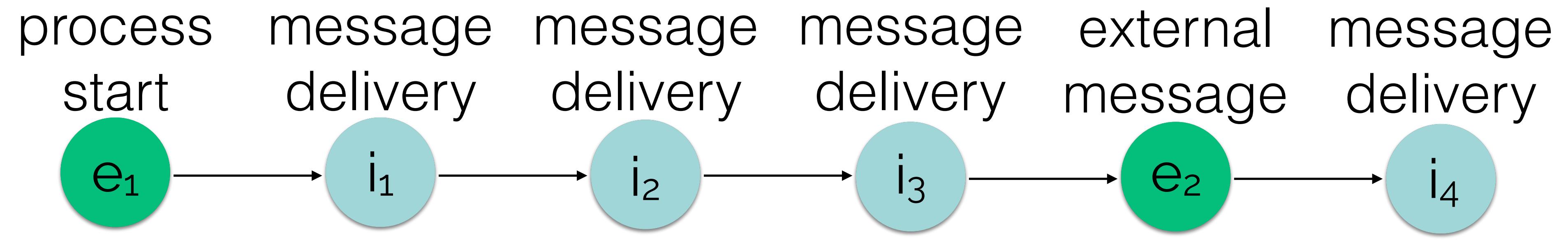


Steps may also be *external*:

- ▶ External message is sent
- ▶ Process is created
- ▶ Process crash-recovers

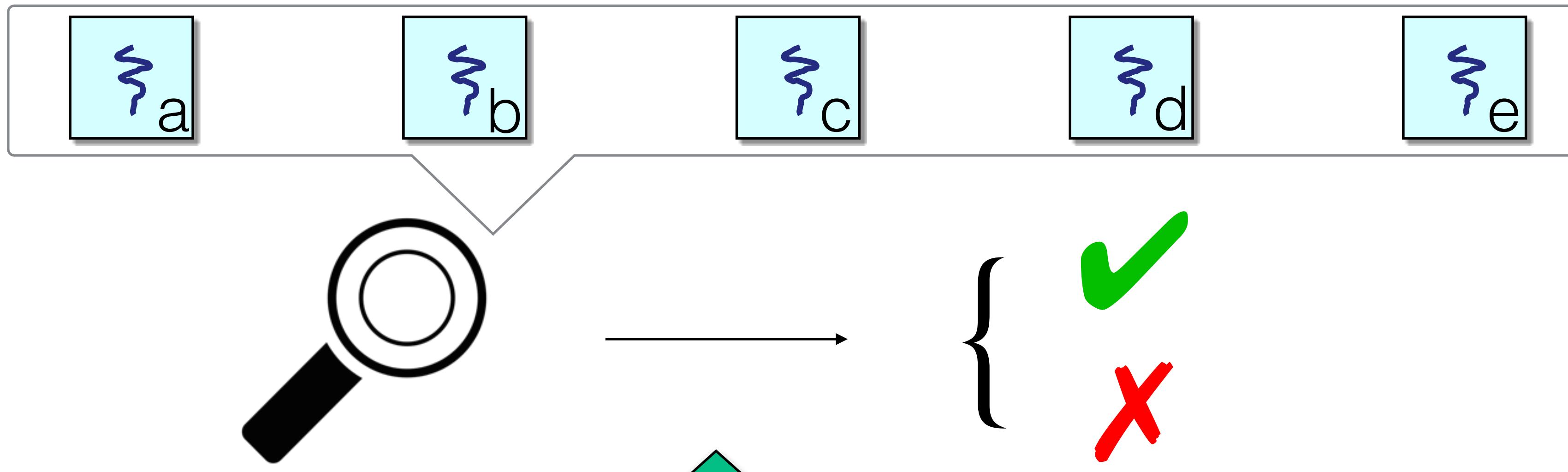
msg
dst: e

Computational Model



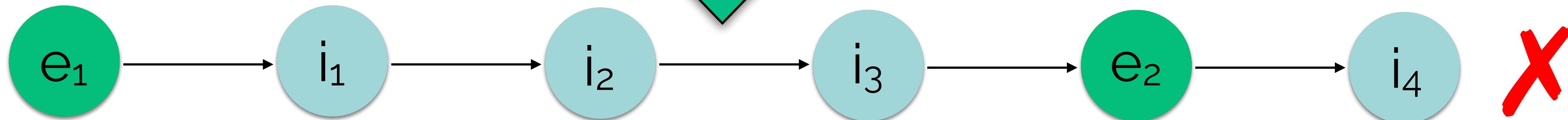
A schedule τ is a sequence of events (either external or internal message deliveries) that can be applied in turn starting from the initial configuration.

Invariant Checking



An invariant is a predicate P over the state of all processes.

A faulty execution is one that ends in an invariant violation.



Formal Problem Statement

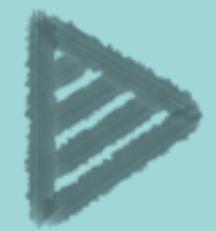
Given: schedule τ that results in violation of P

Find: schedule τ' s.t.

- ▶ τ' violates P, $|\tau'| \leq |\tau|$
- ▶ τ' contains a subsequence of the external events of τ
- ▶ if we remove any external event e from τ' ,
- ▶ $\neg \exists \tau''$ containing same external events - e, s.t. τ'' violates P

Formal Problem Statement

After finding τ' :



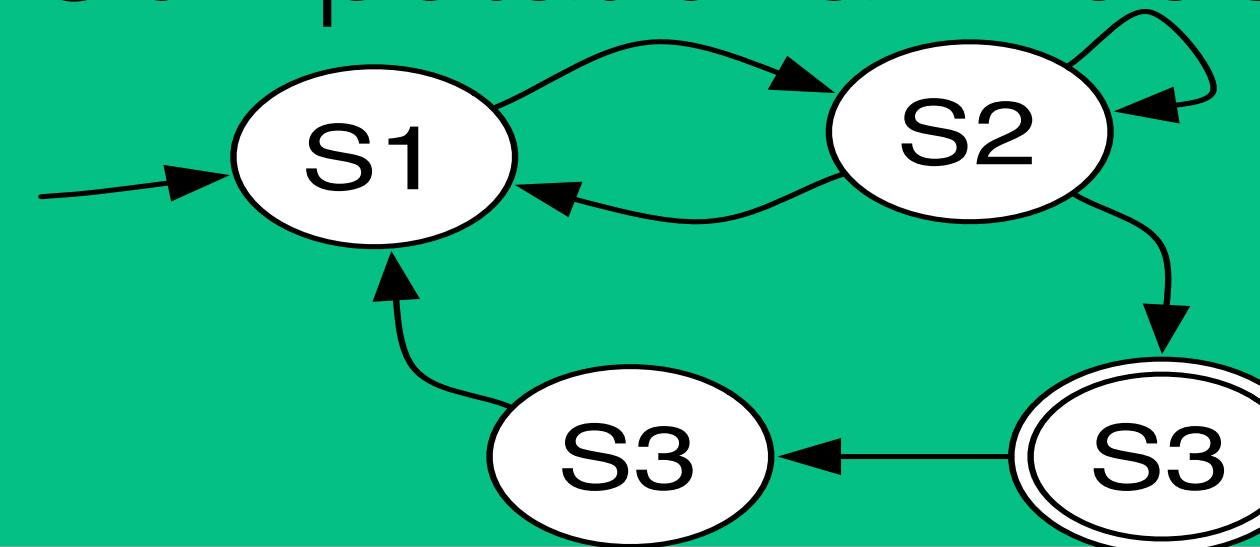
remove extraneous message deliveries from τ'

Outline

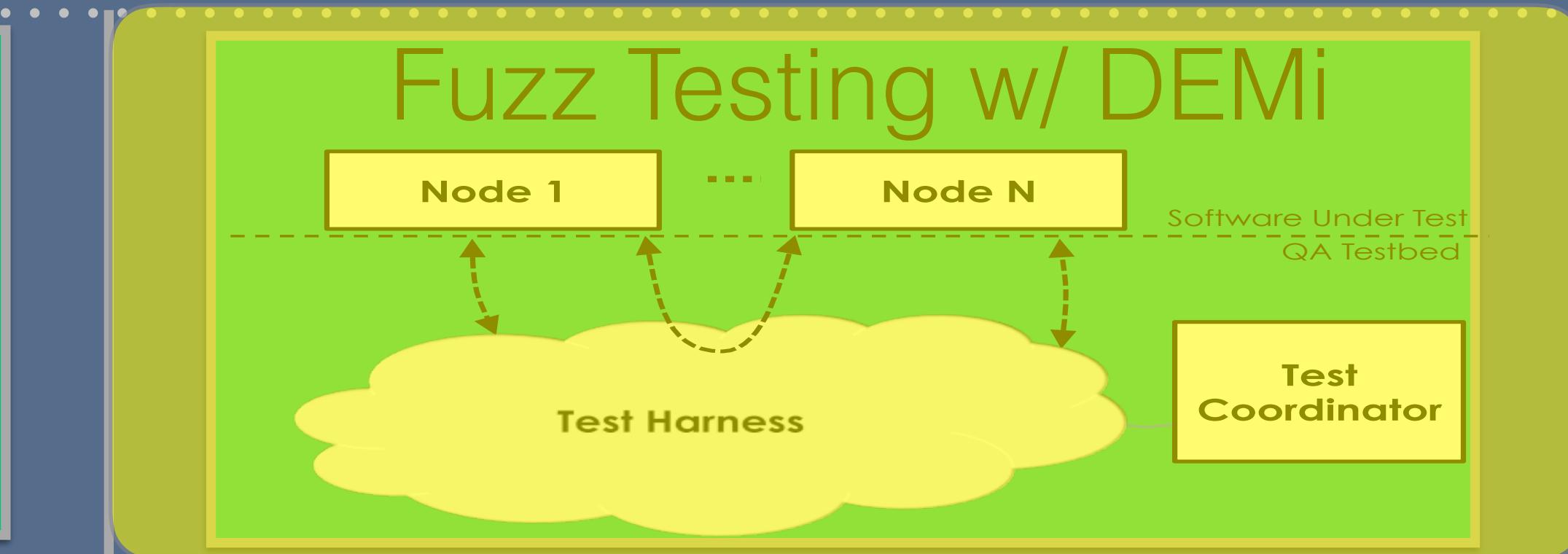
Introduction

Background

Computational Model



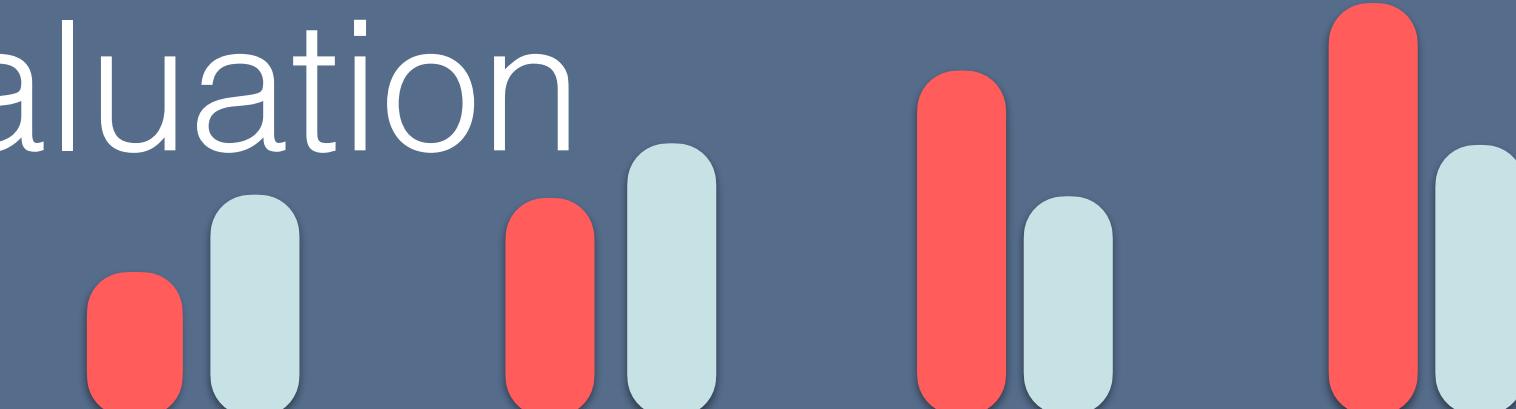
Fuzz Testing w/ DEMi



Minimization

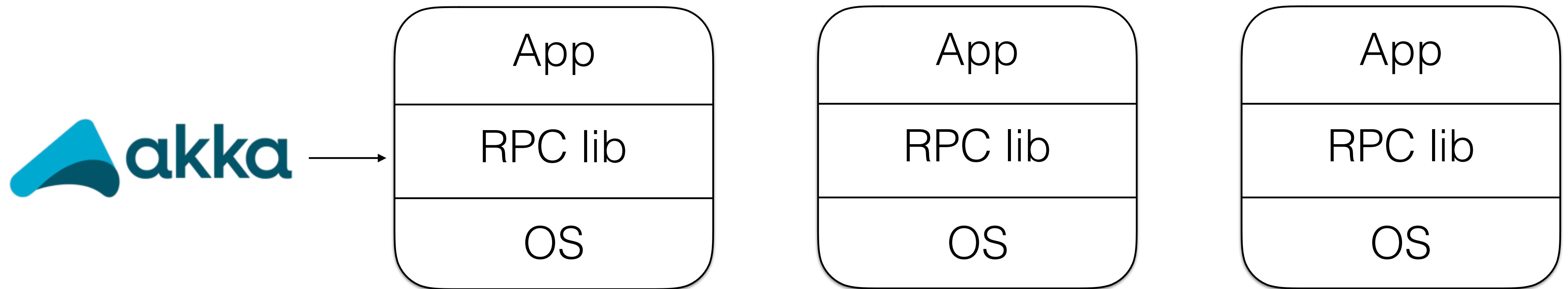


Evaluation

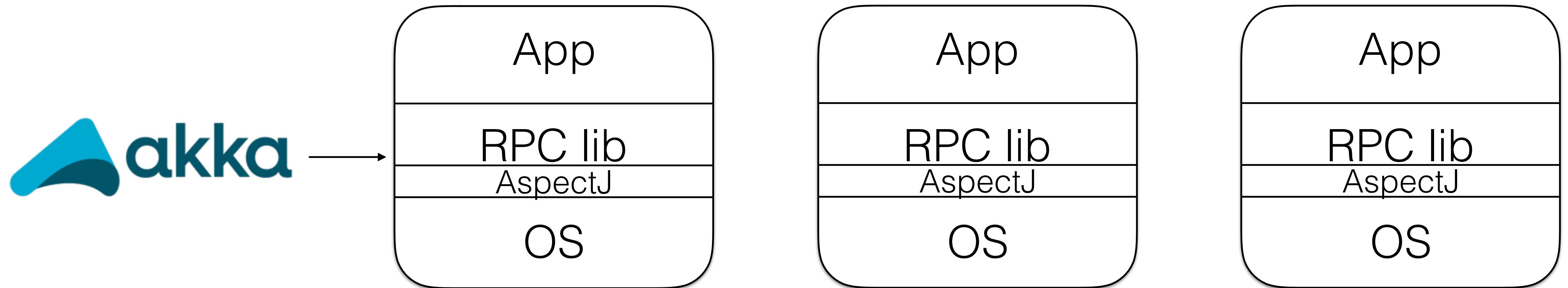


Conclusion

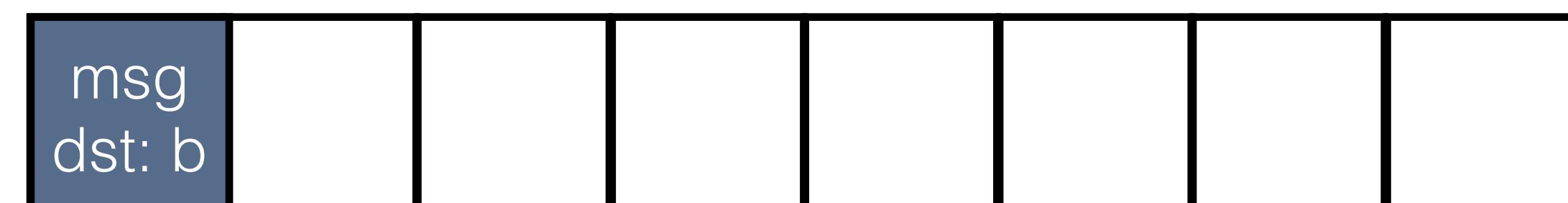
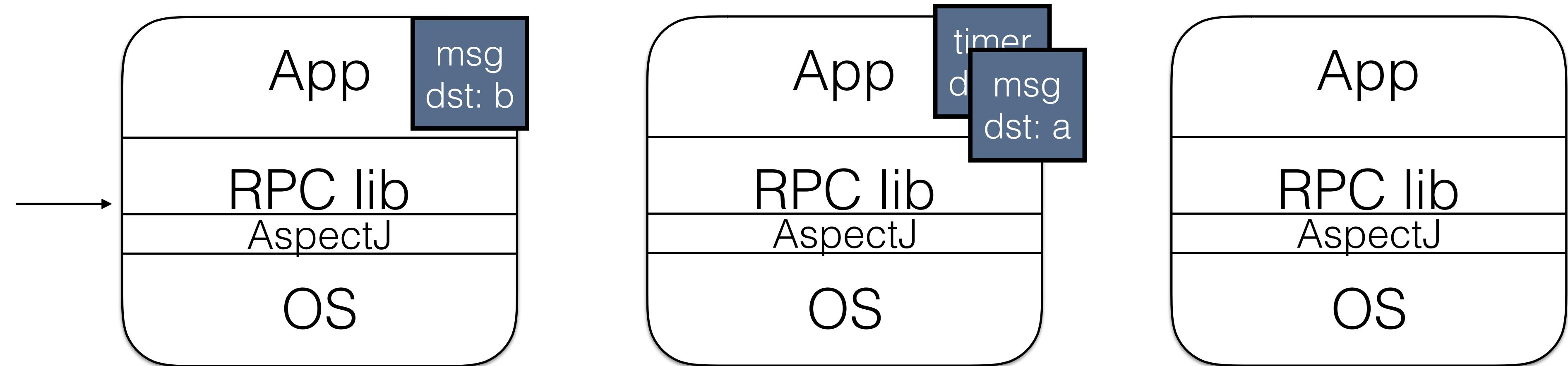
Fuzz Testing with DEMi



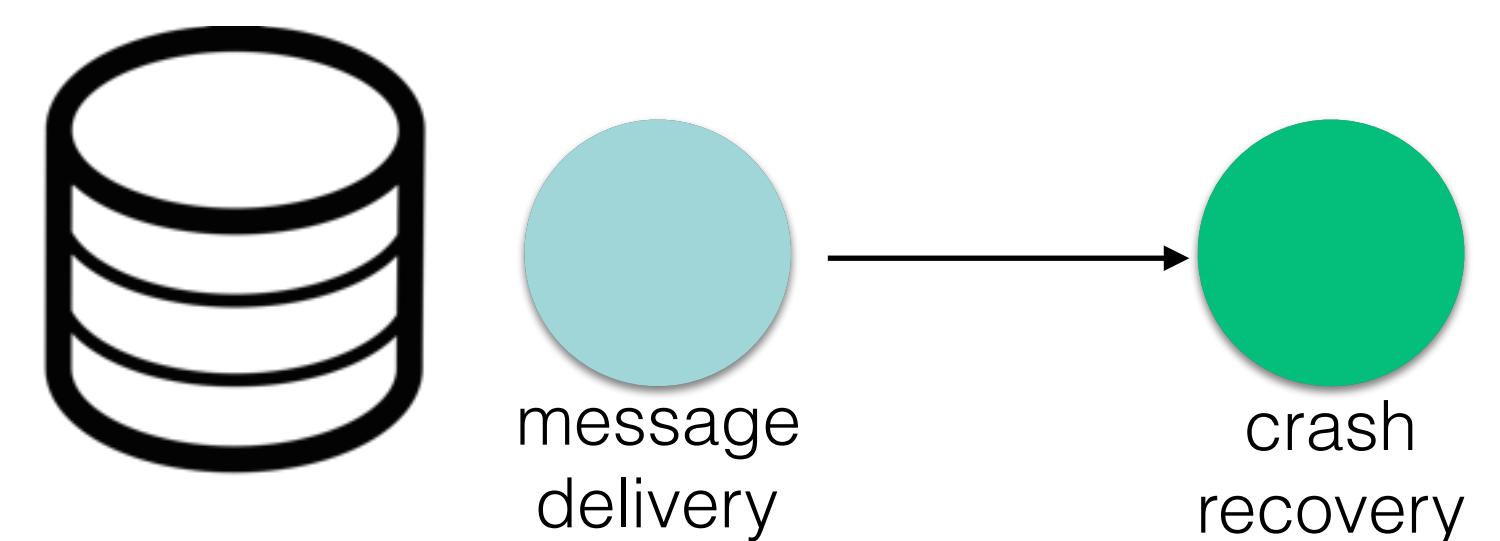
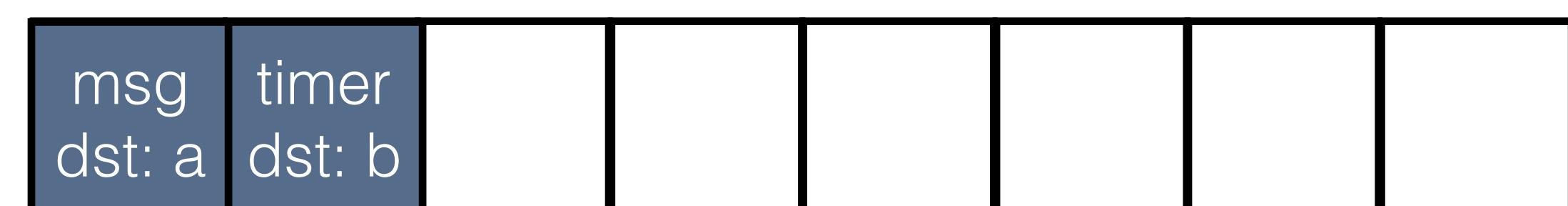
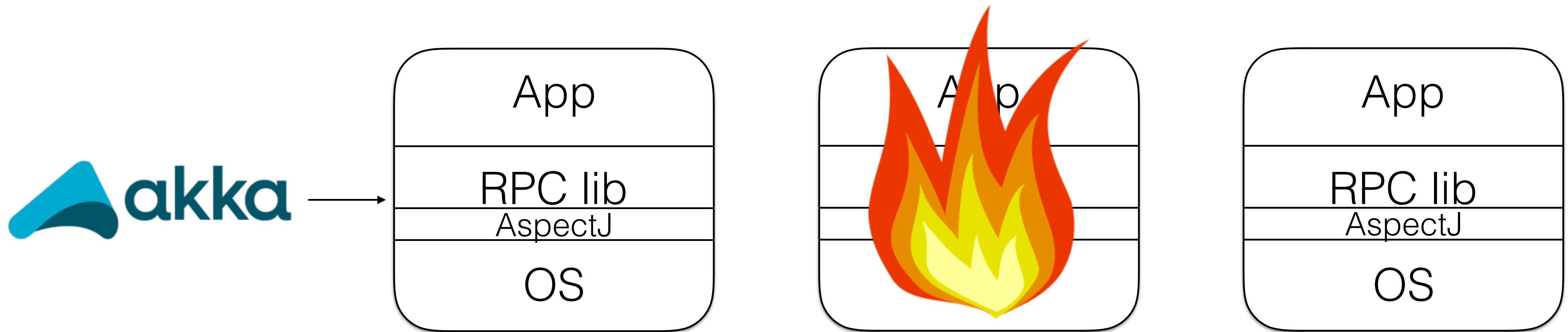
Fuzz Testing with DEMi



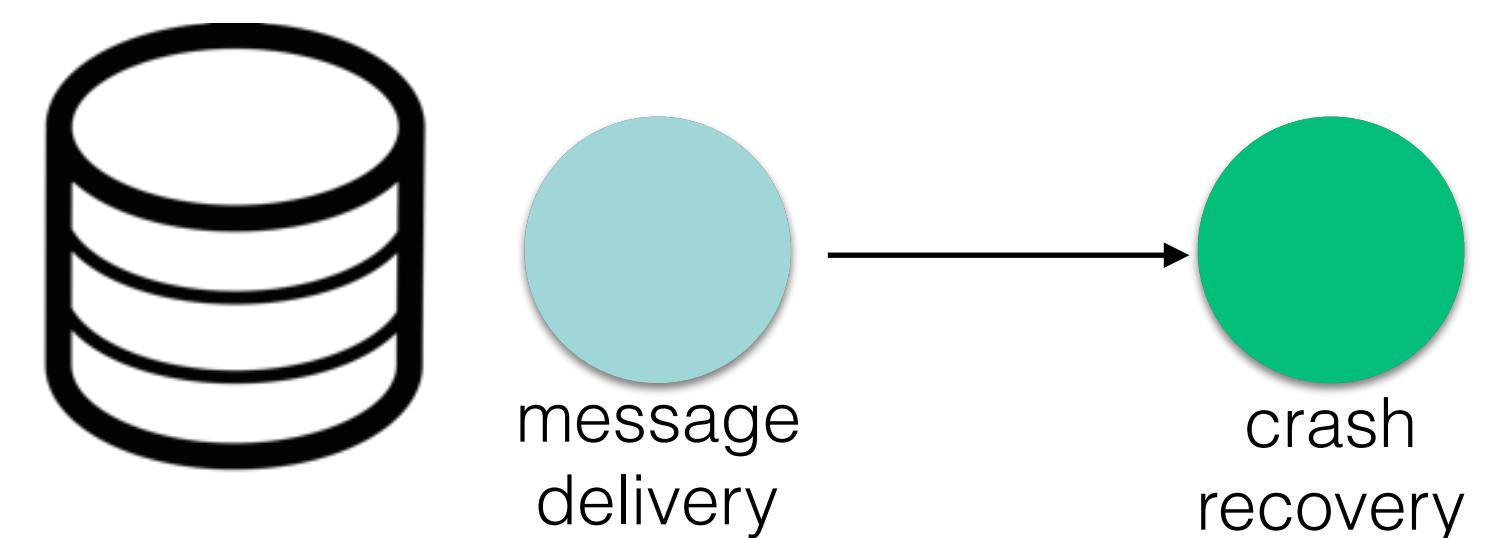
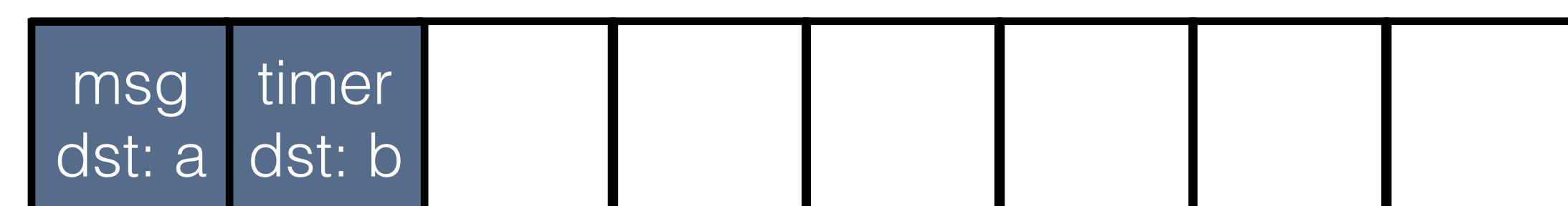
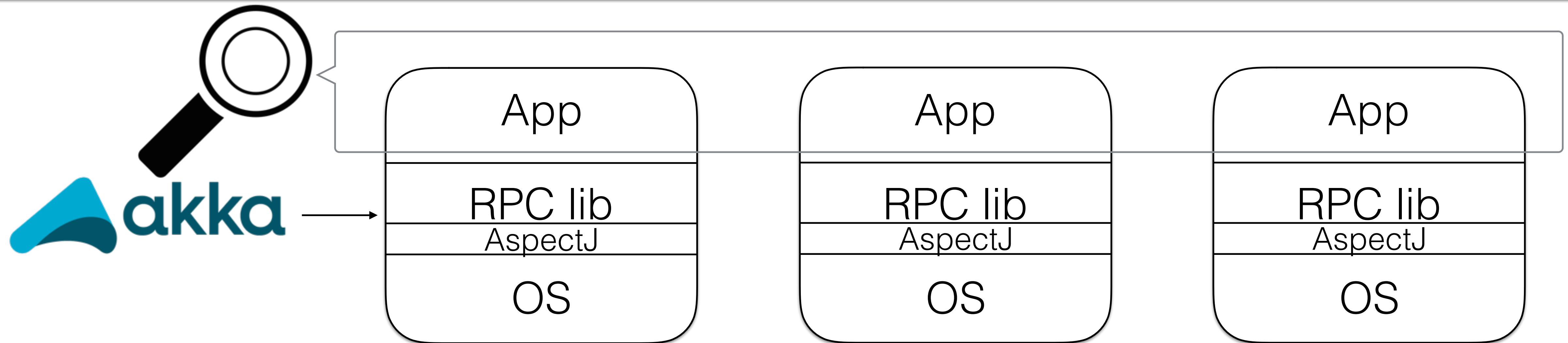
Fuzz Testing with DEMi



Fuzz Testing with DEMi



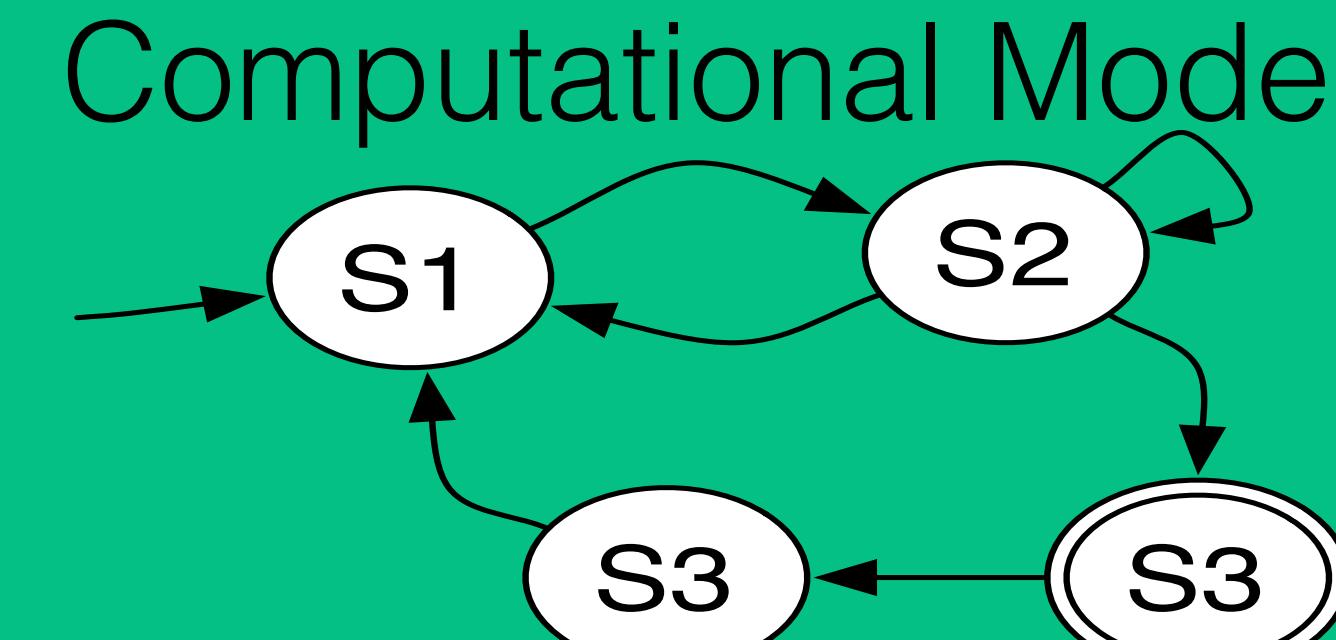
Fuzz Testing with DEMi



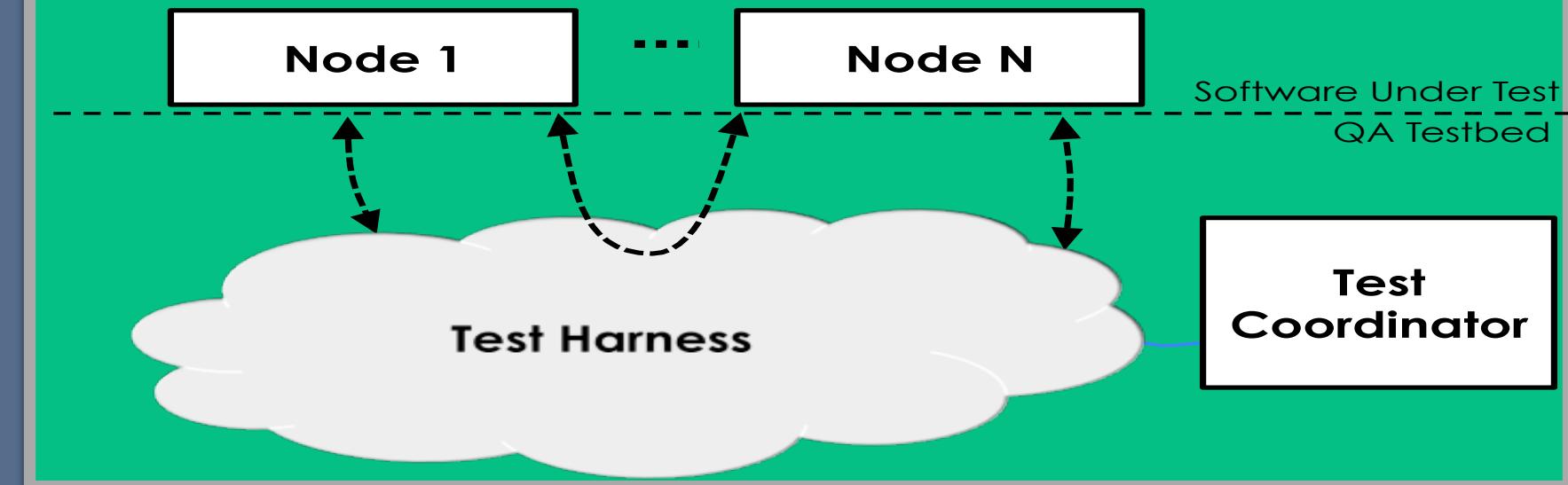
Outline

Introduction

Background



Fuzz Testing w/ DEMi



Minimization

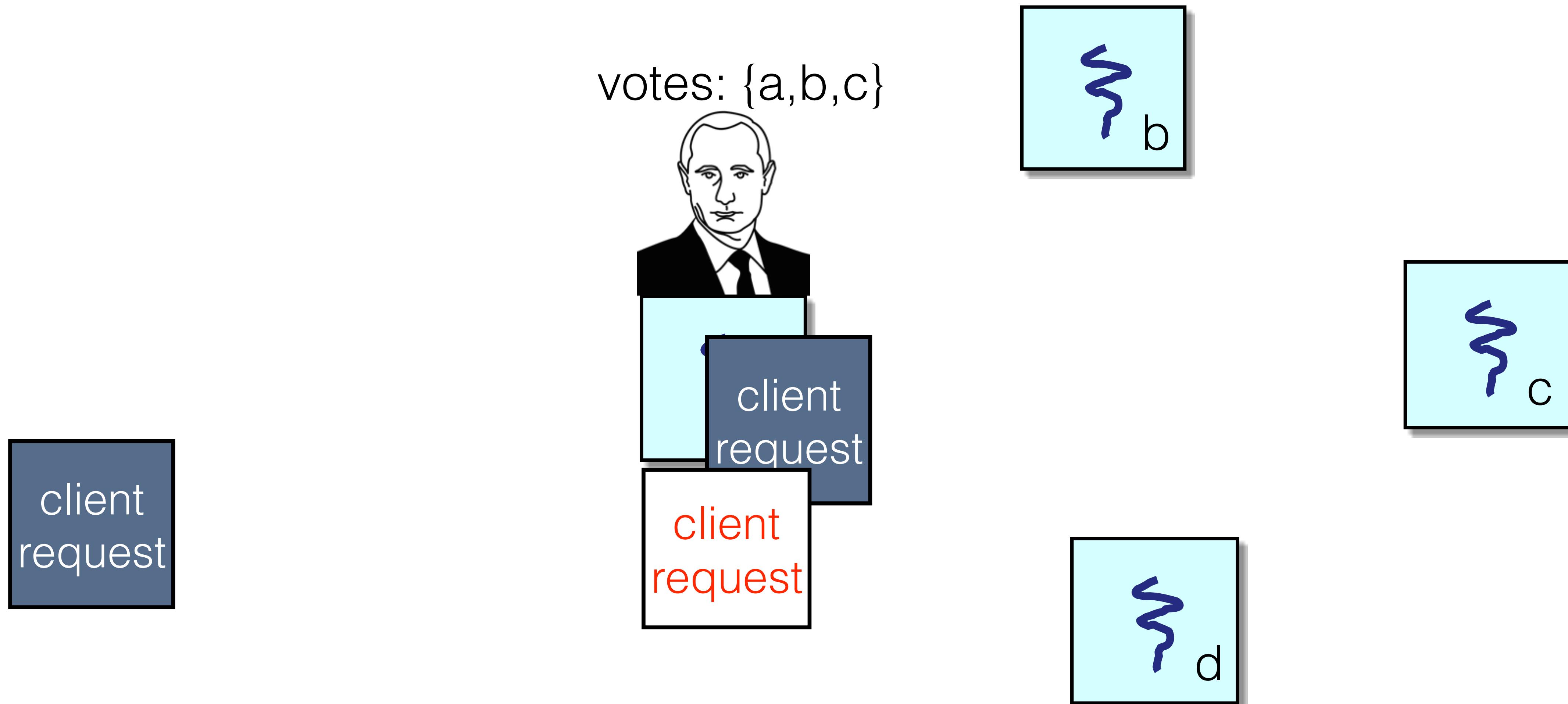


Evaluation

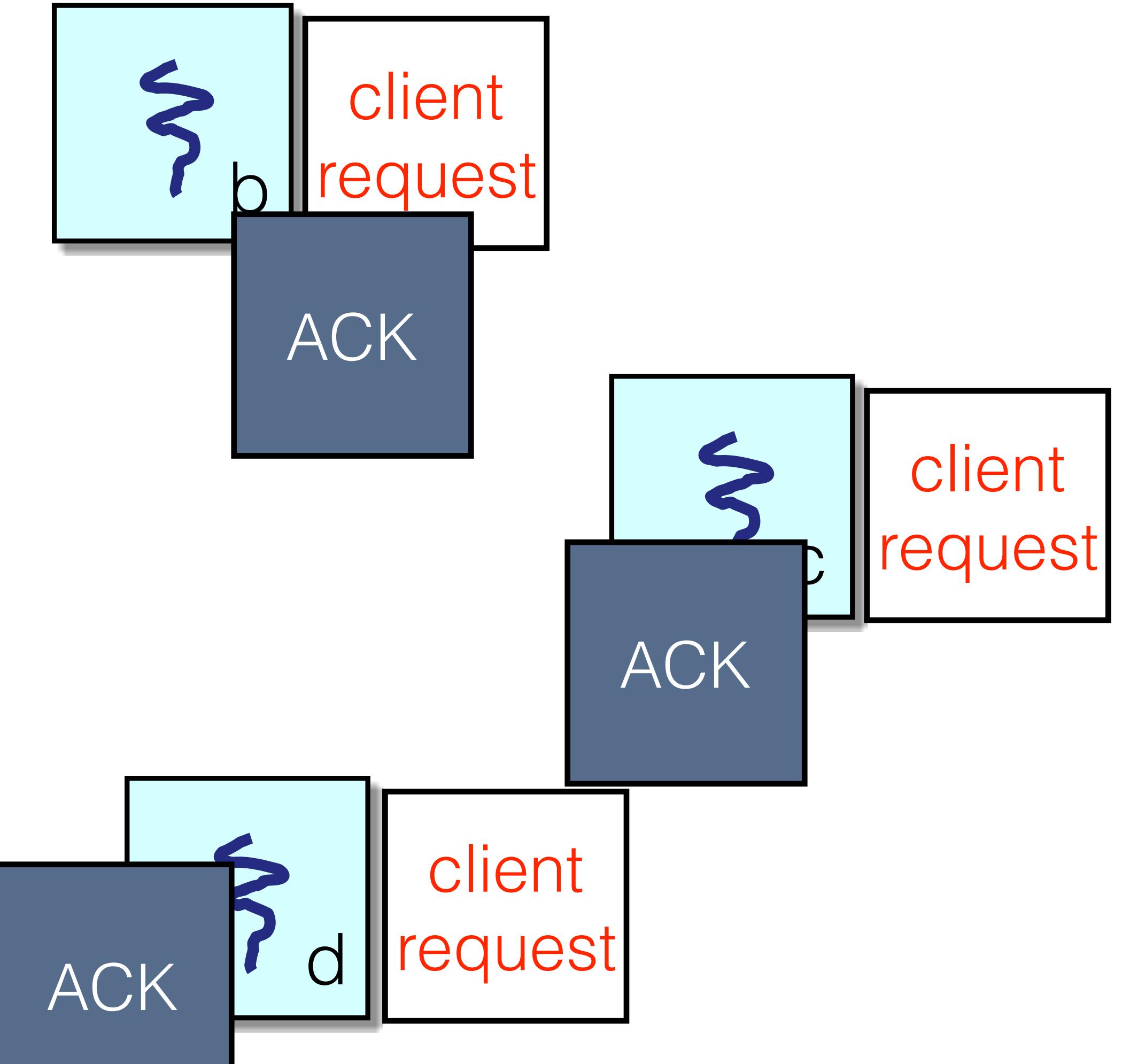
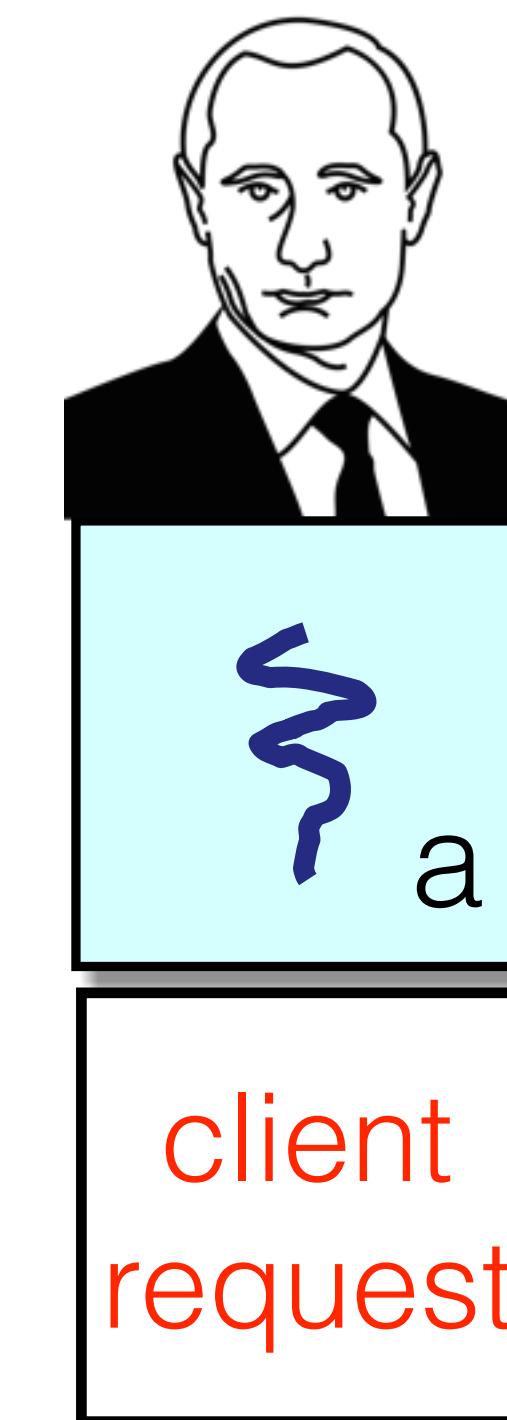


Conclusion

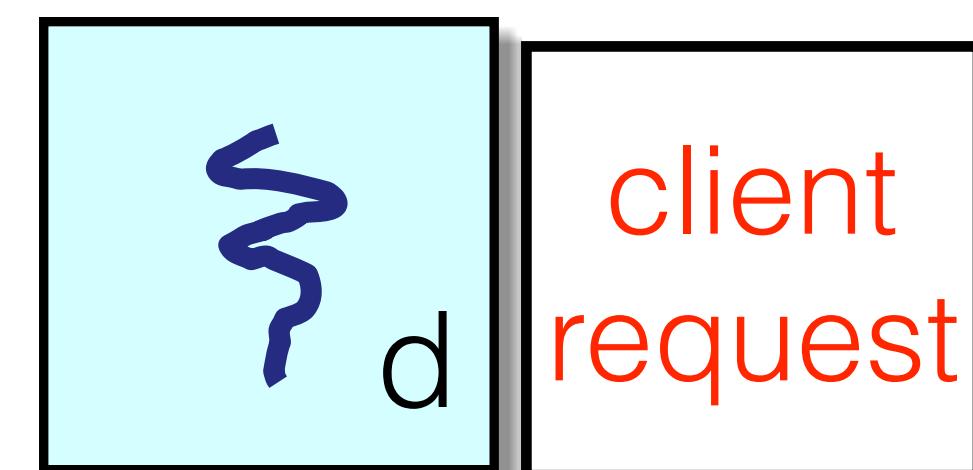
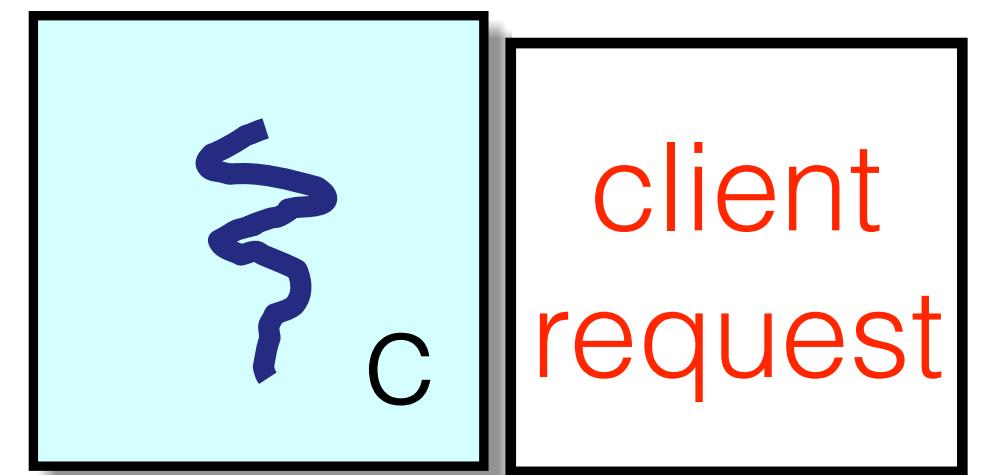
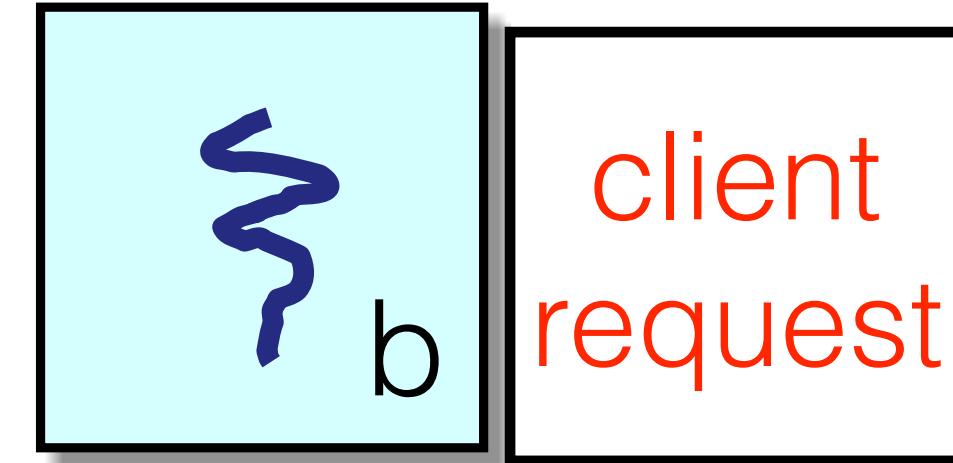
Running Example: Raft Consensus



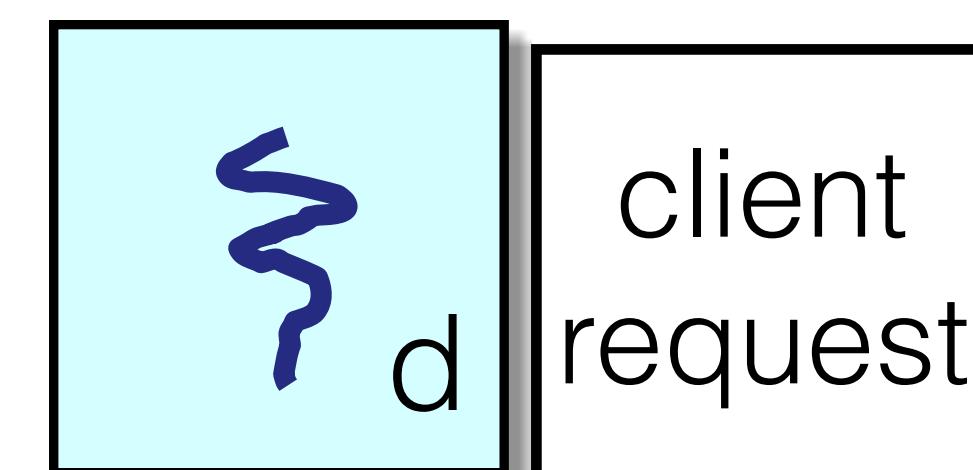
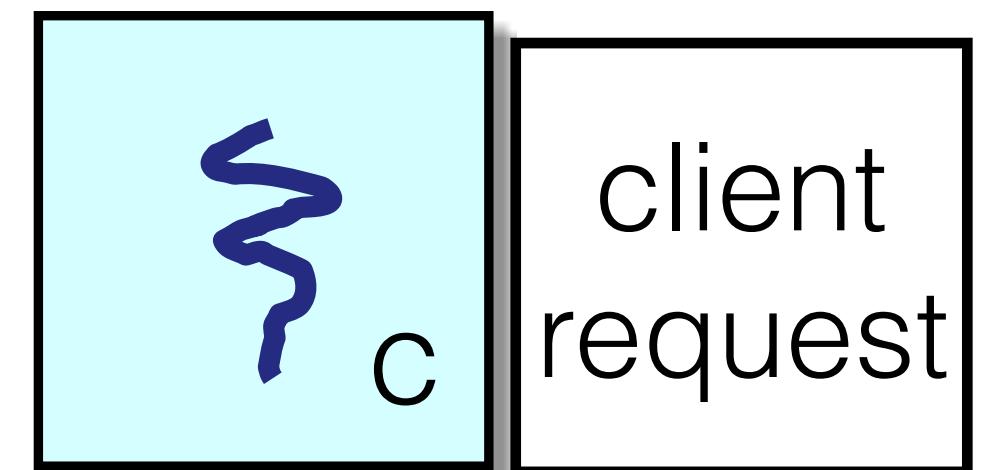
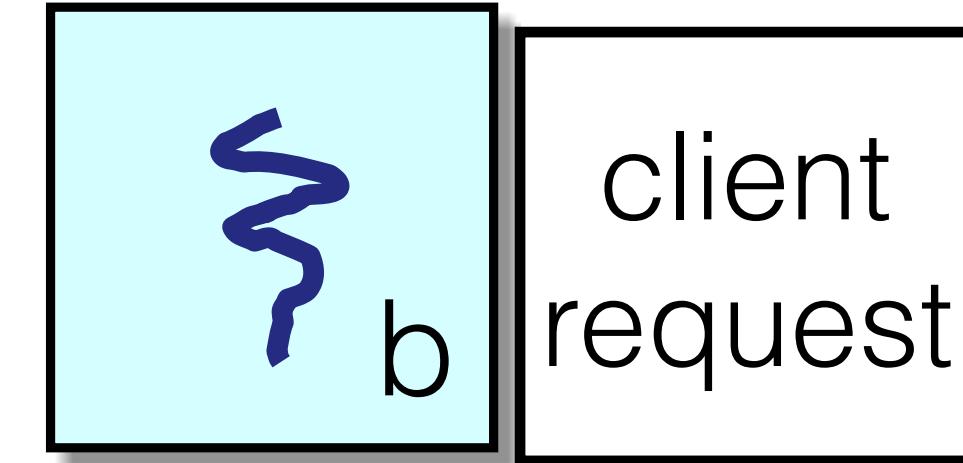
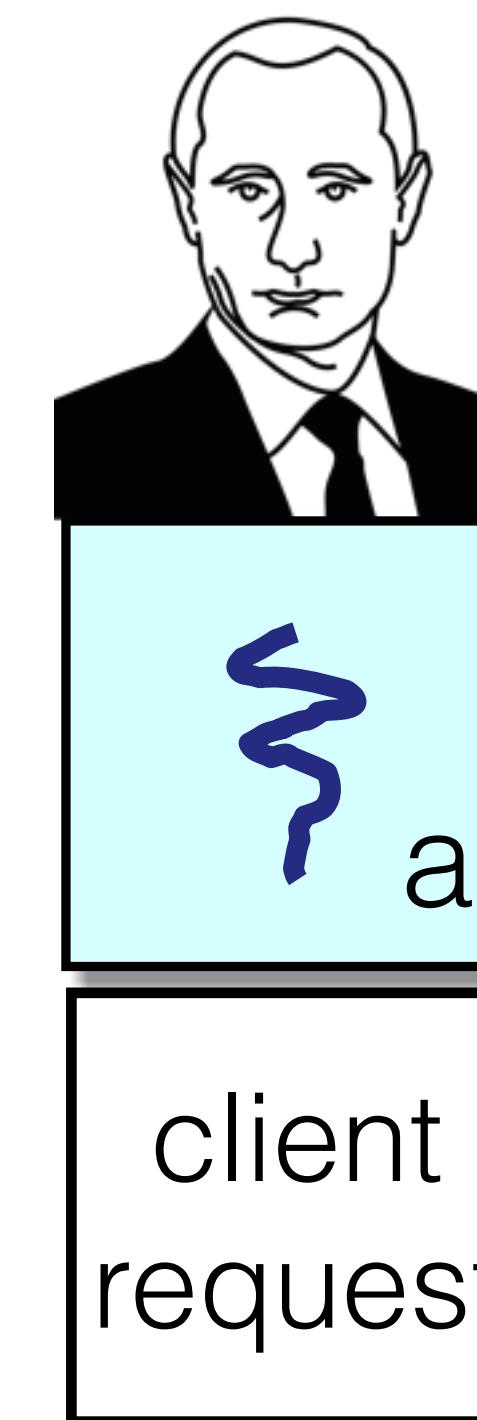
Running Example: Raft Consensus



Running Example: Raft Consensus



Running Example: Raft Consensus



Log lines

Motifs

Sending to raft-member-4: RequestVote(Te
Sending to raft-member-2: BeginElection

Received message from raft-member-3: Req
Received timer: Timer(election-timer,Ele
Received message from raft-member-2: Beg

Sending to raft-member-3: VoteCandidate(
Sending to raft-member-3: BeginElection
Sending to raft-member-1: RequestVote(Te

Received message from raft-member-4: Vot
25 Received message from raft-member-2: R
equestVote(Term(2),Actor[akka://new-sy
stem-0/user/raft-member-2#-69475157
Rec 3],Term(0),0)

Sending to raft-member-2: VoteCandidate(

Sending to raft-member-4: RequestVote(Te
Received message from raft-member-1: Vot
Received timer: Timer(election-timer,Ele

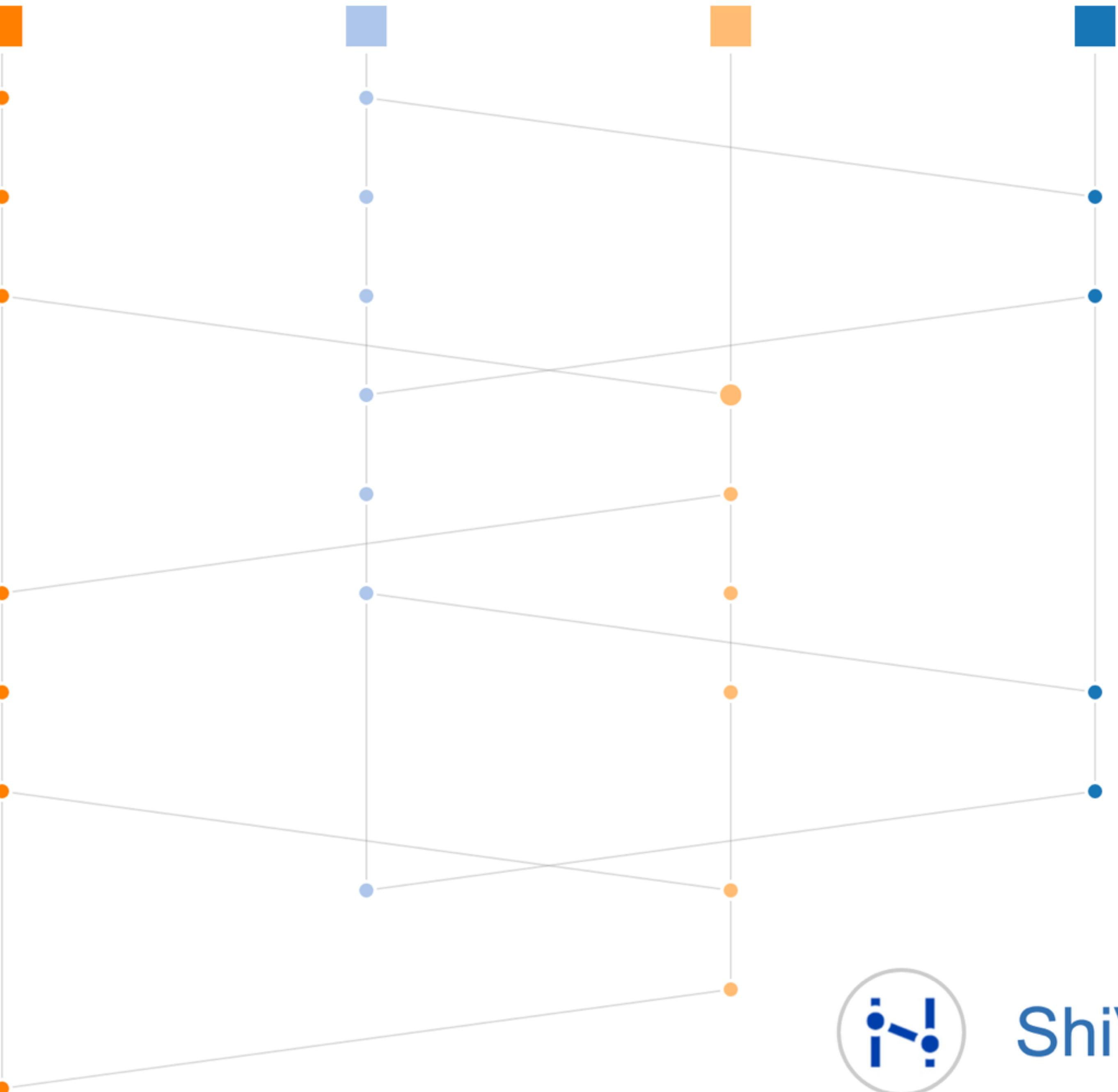
Received message from raft-member-3: Req
Received message from raft-member-2: Beg
Received timer: Timer(election-timer,Ele

Sending to raft-member-3: VoteCandidate(
Sending to raft-member-1: RequestVote(Te

Received message from raft-member-4: Vot
Received message from raft-member-2: Req

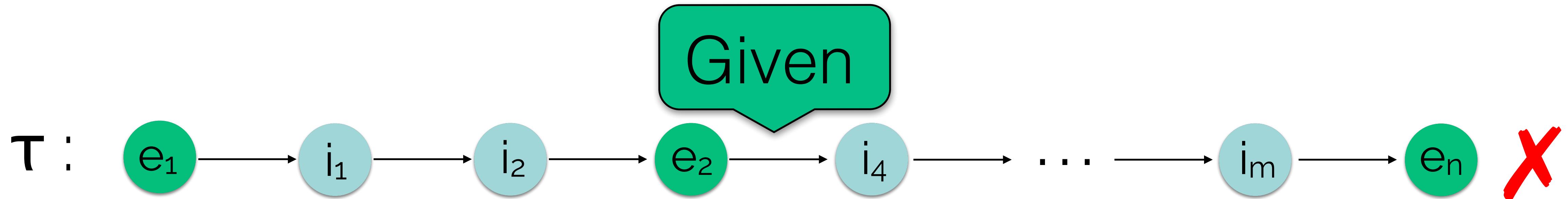
Sending to raft-member-2: VoteCandidate(

Received message from raft-member-1: Vot



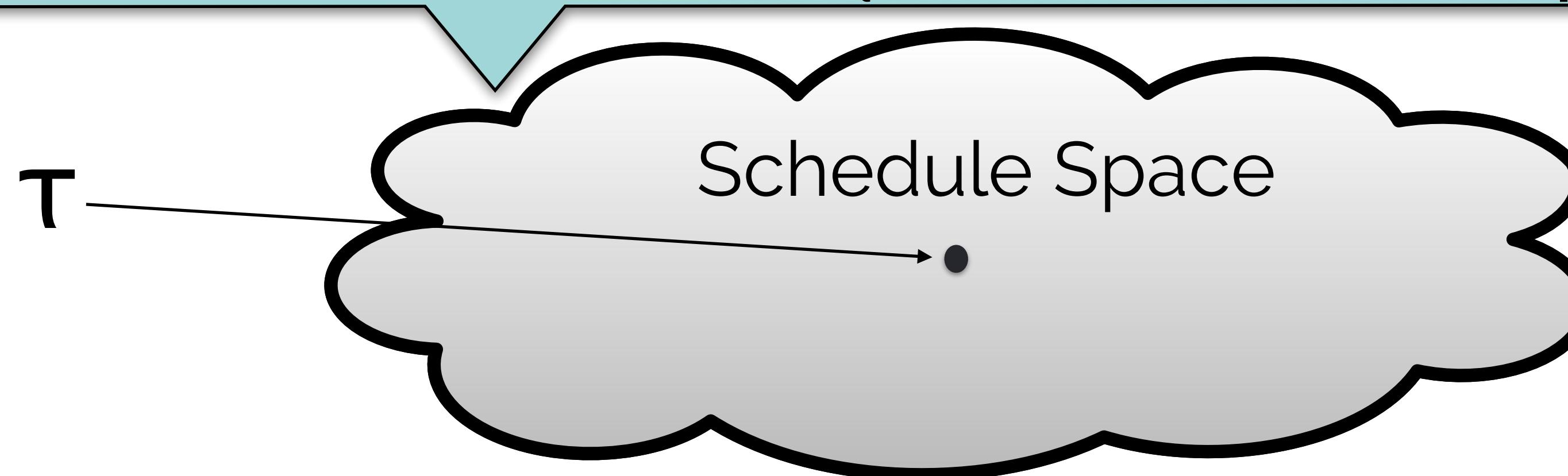
ShiViz

Minimization



Straightforward approach:

- ▶ Enumerate all schedules $|\tau'| \leq |\tau|$,
- ▶ Pick shortest sequence that reproduces X



on!

Observation #1: many schedules are commutative

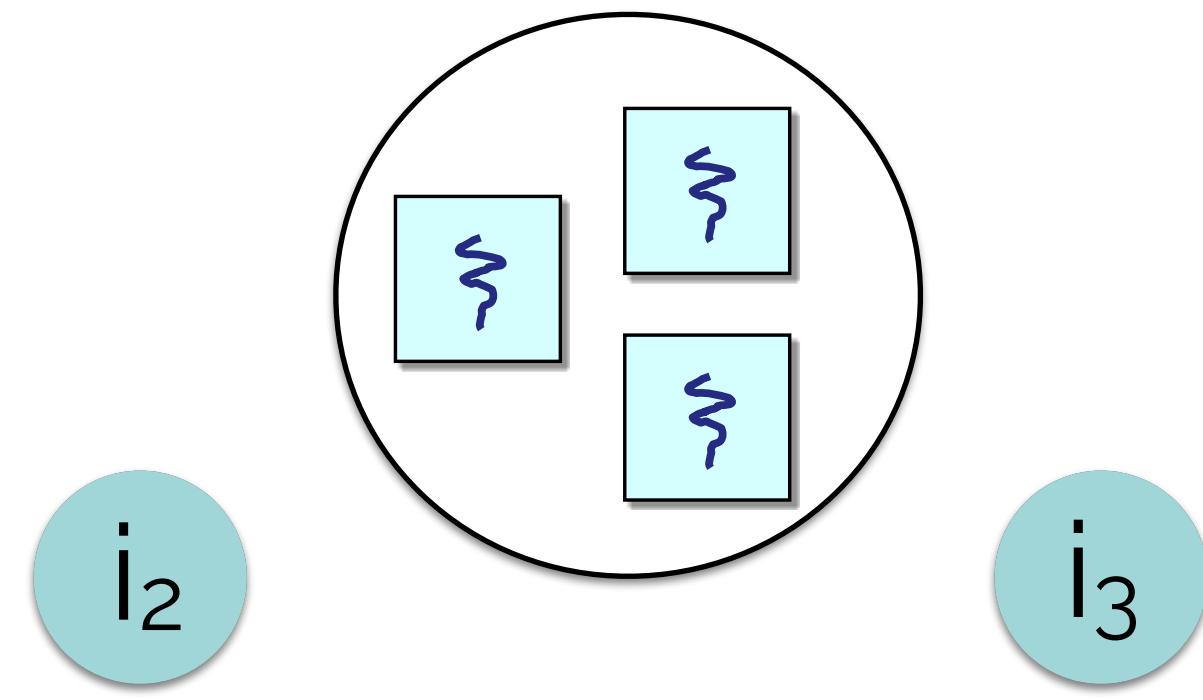
i_2

i_3

- ▶ $i_2 \rightarrow i_3$
- ▶ $i_3 \rightarrow i_2$
- ▶ $\text{dst}(i_2) \neq \text{dst}(i_3)$

Observation #1: many schedules are commutative

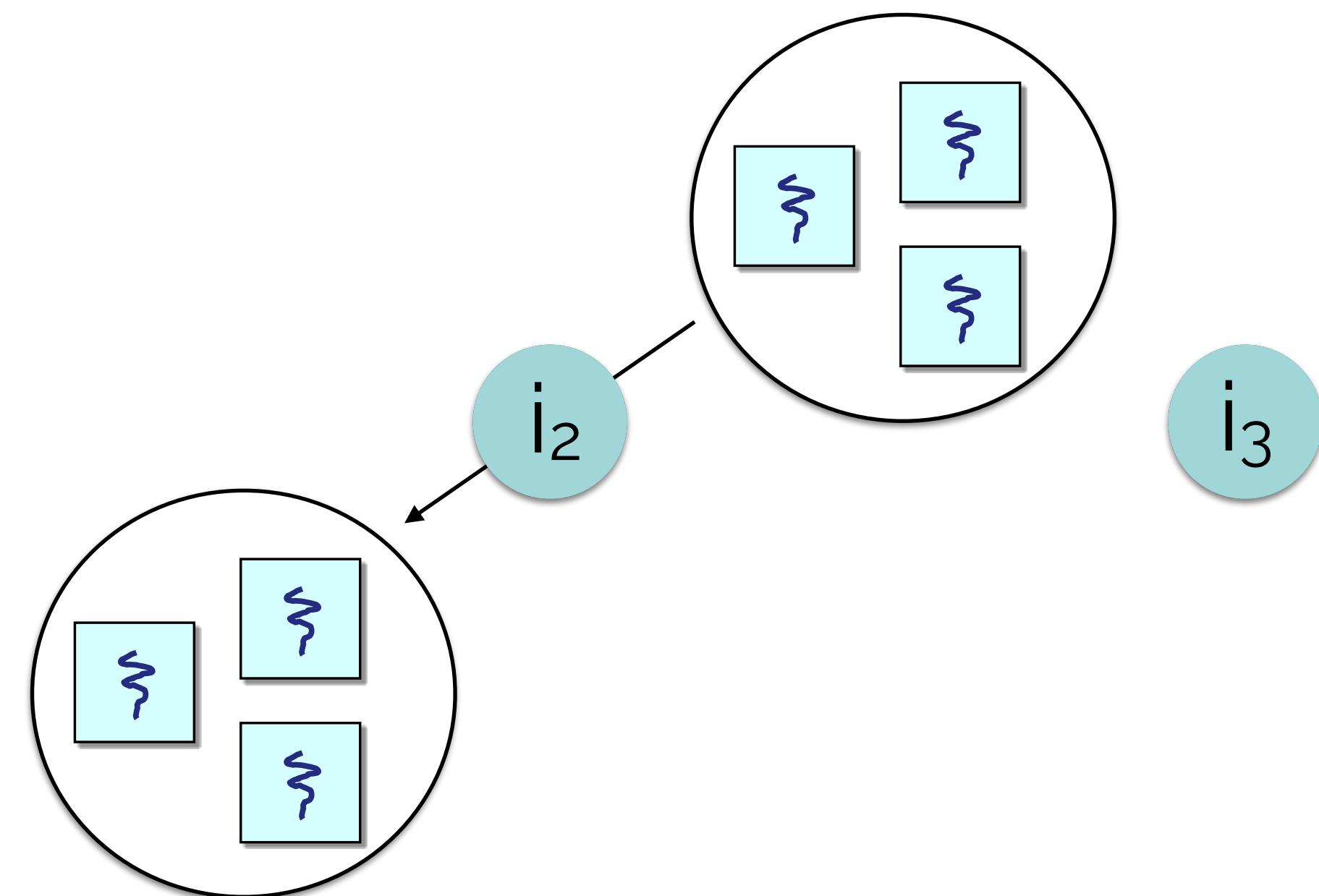
Step n:



- ▶ $i_2 \rightarrow i_3$
- ▶ $i_3 \rightarrow i_2$
- ▶ $\text{dst}(i_2) \neq \text{dst}(i_3)$

Observation #1: many schedules are commutative

Step n:

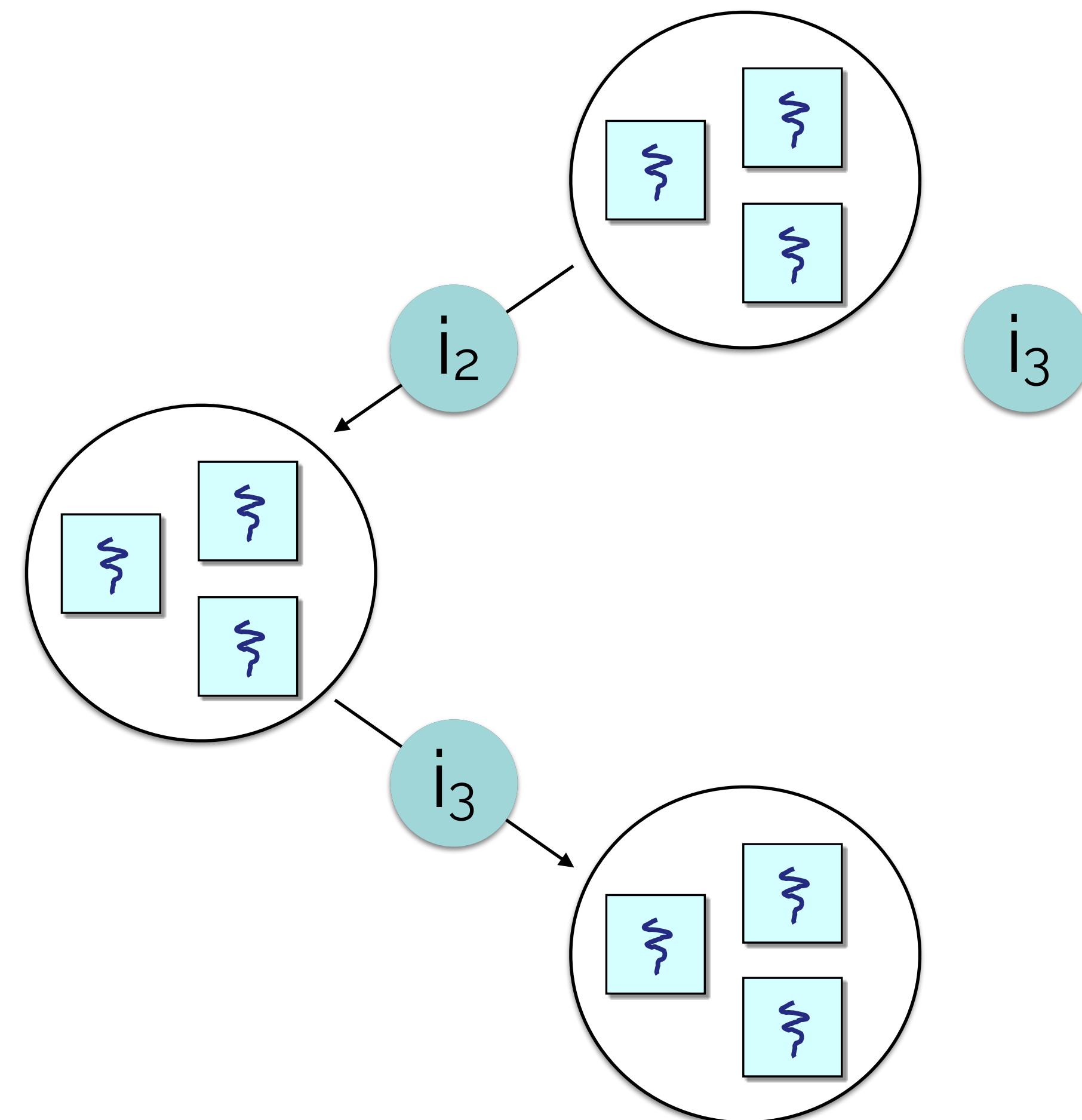


Step n+1:

- ▶ $i_2 \rightarrow i_3$
- ▶ $i_3 \rightarrow i_2$
- ▶ $\text{dst}(i_2) \neq \text{dst}(i_3)$

Observation #1: many schedules are commutative

Step n:



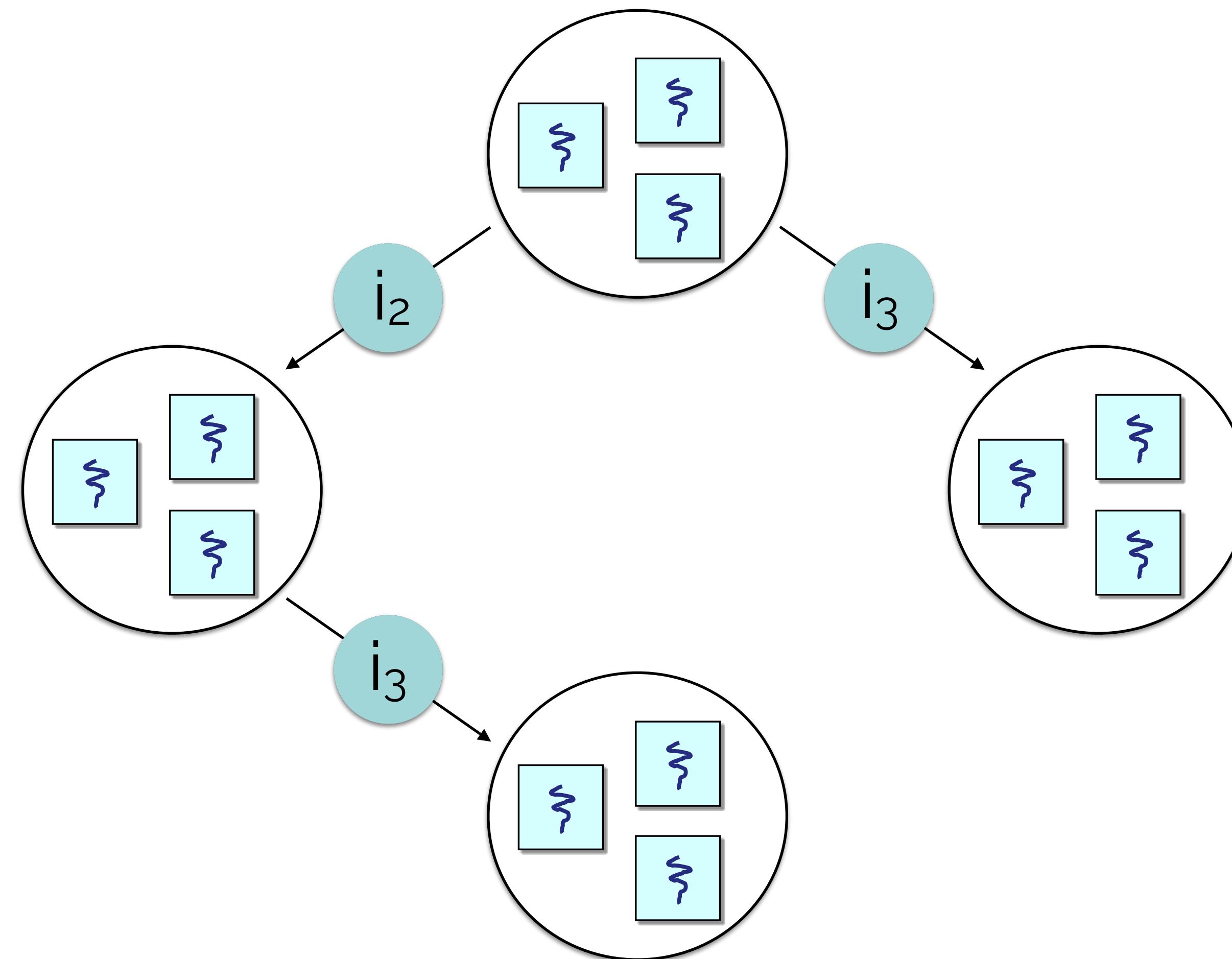
Step n+1:

- ▶ $i_2 \rightarrow i_3$
- ▶ $i_3 \rightarrow i_2$
- ▶ $\text{dst}(i_2) \neq \text{dst}(i_3)$

Step n+2:

Observation #1: many schedules are commutative

Step n:



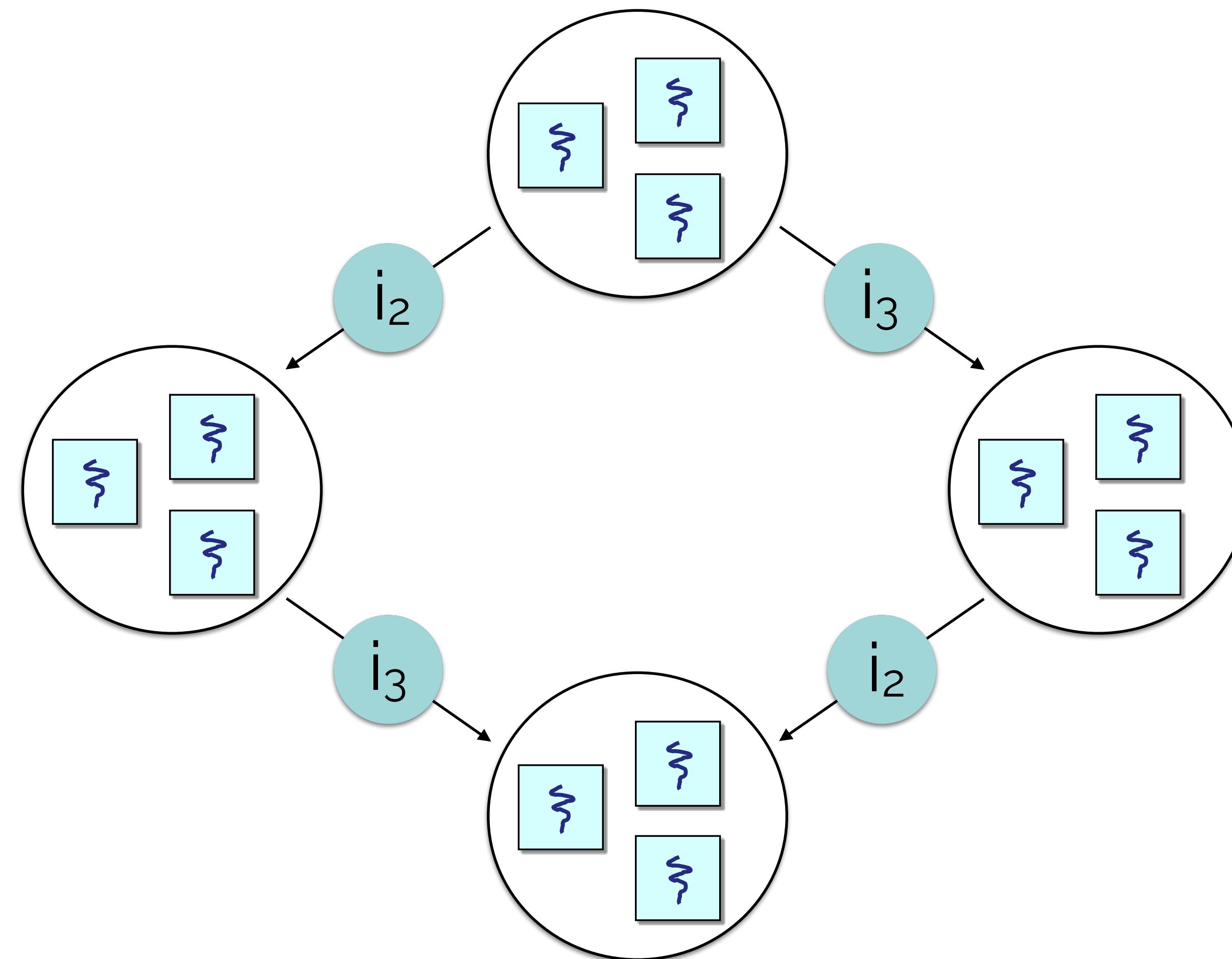
Step n+1:

Step n+2:

$\Rightarrow i_2 \rightarrow i_3$
 $\Rightarrow i_3 \rightarrow i_2$
 $\Rightarrow \text{dst}(i_2) \neq \text{dst}(i_3)$

Observation #1: many schedules are commutative

Step n:



Step n+1:

► $i_2 \rightarrow i_3$
► $i_3 \rightarrow i_2$
► $\text{dst}(i_2) \neq \text{dst}(i_3)$

Step n+2:

Observation #1: many schedules are commutative

Adopt DPOR:
Dynamic Partial Order Reduction

C. Flanagan, P. Godefroid, “Dynamic Partial-Order Reduction for Model Checking Software”, POPL ‘05

$$\Theta\left(\frac{n!}{k!}\right)$$

Approach: prioritize schedule space exploration

Assume: fixed time budget
Objective: quickly find small failing schedules



Given:

Prioritization function



Given:

▶ Prioritization function

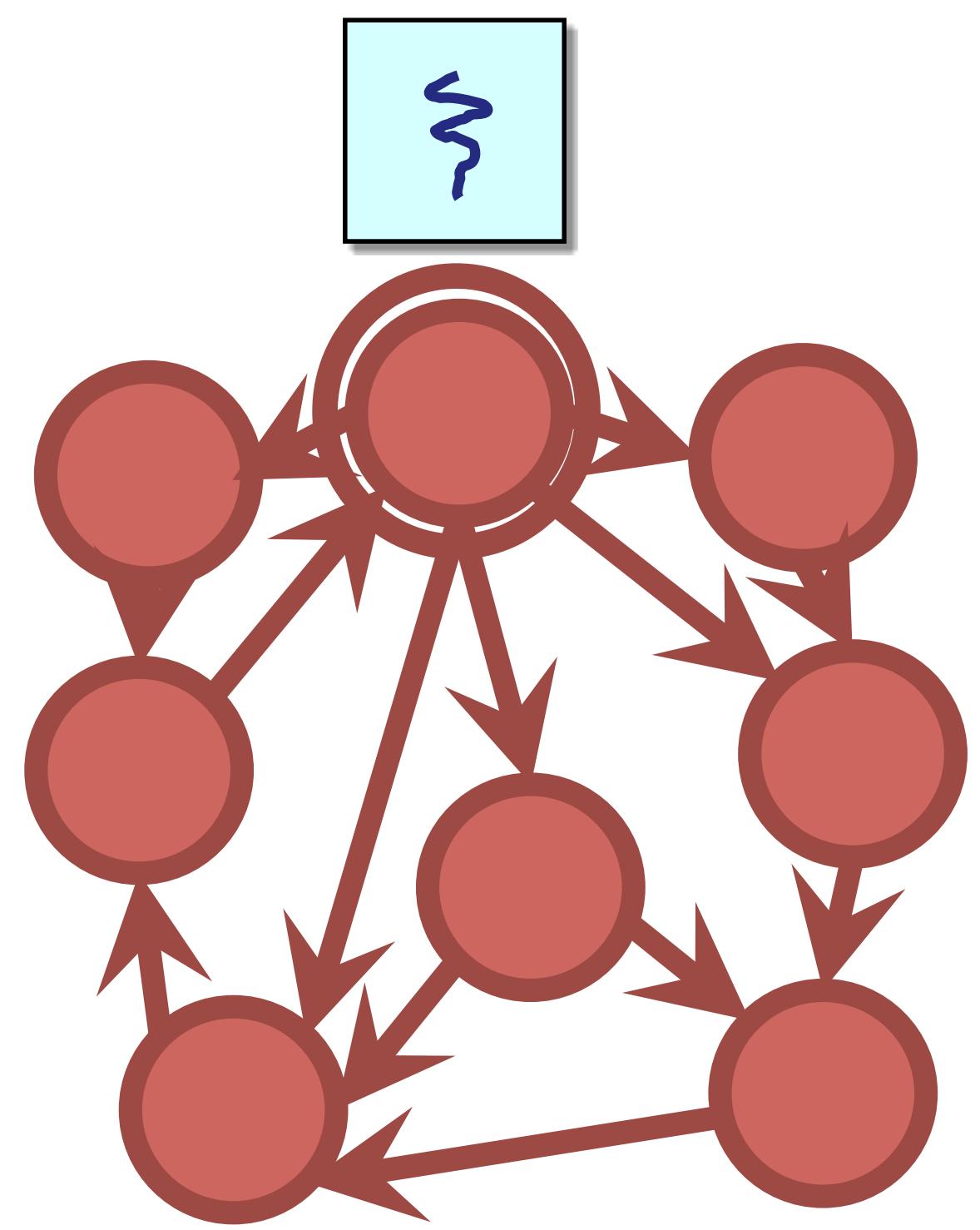
Produce:

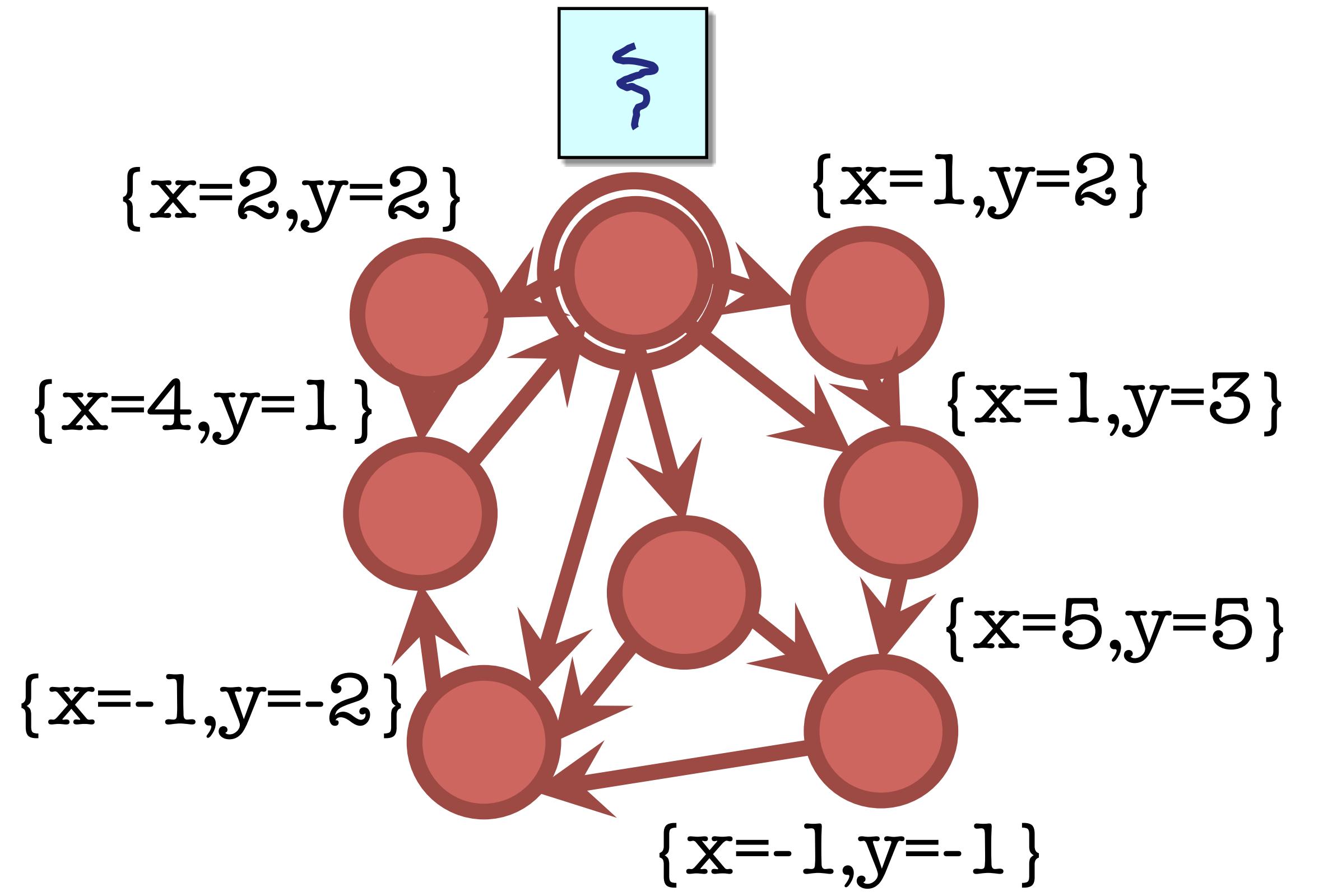
▶ Program under test
▶ Initial execution

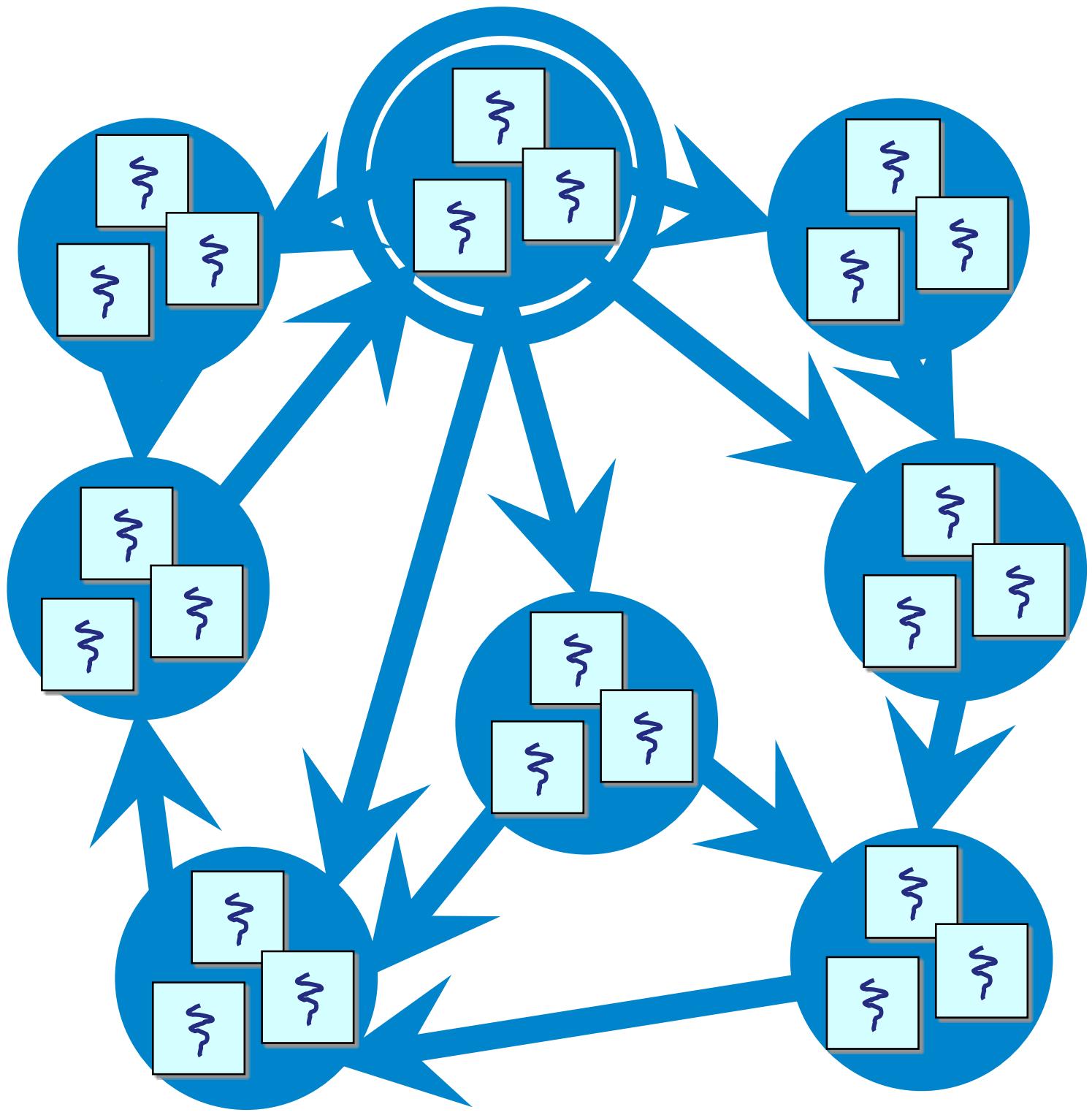
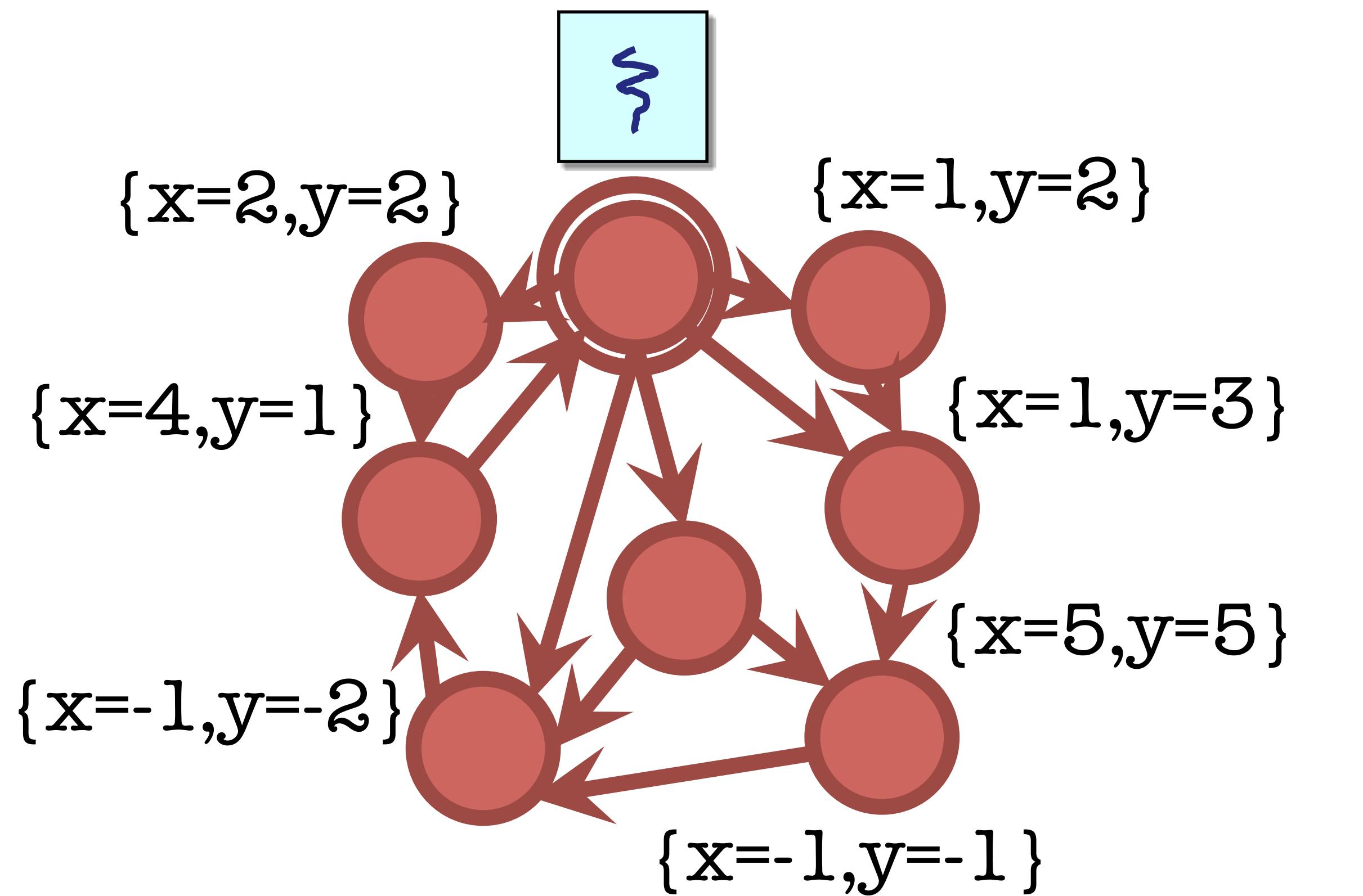
s.t. prioritization makes scant progress

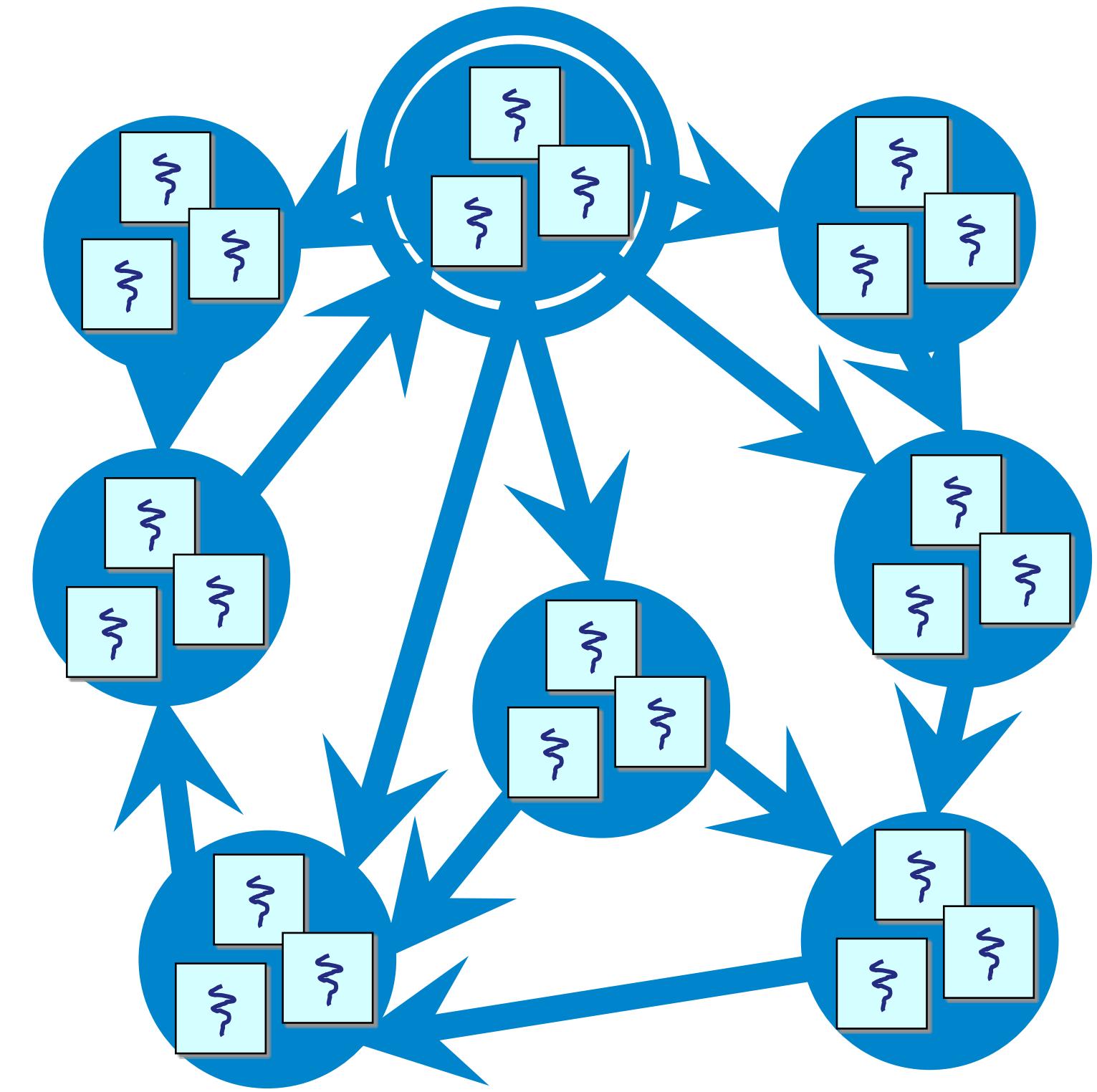
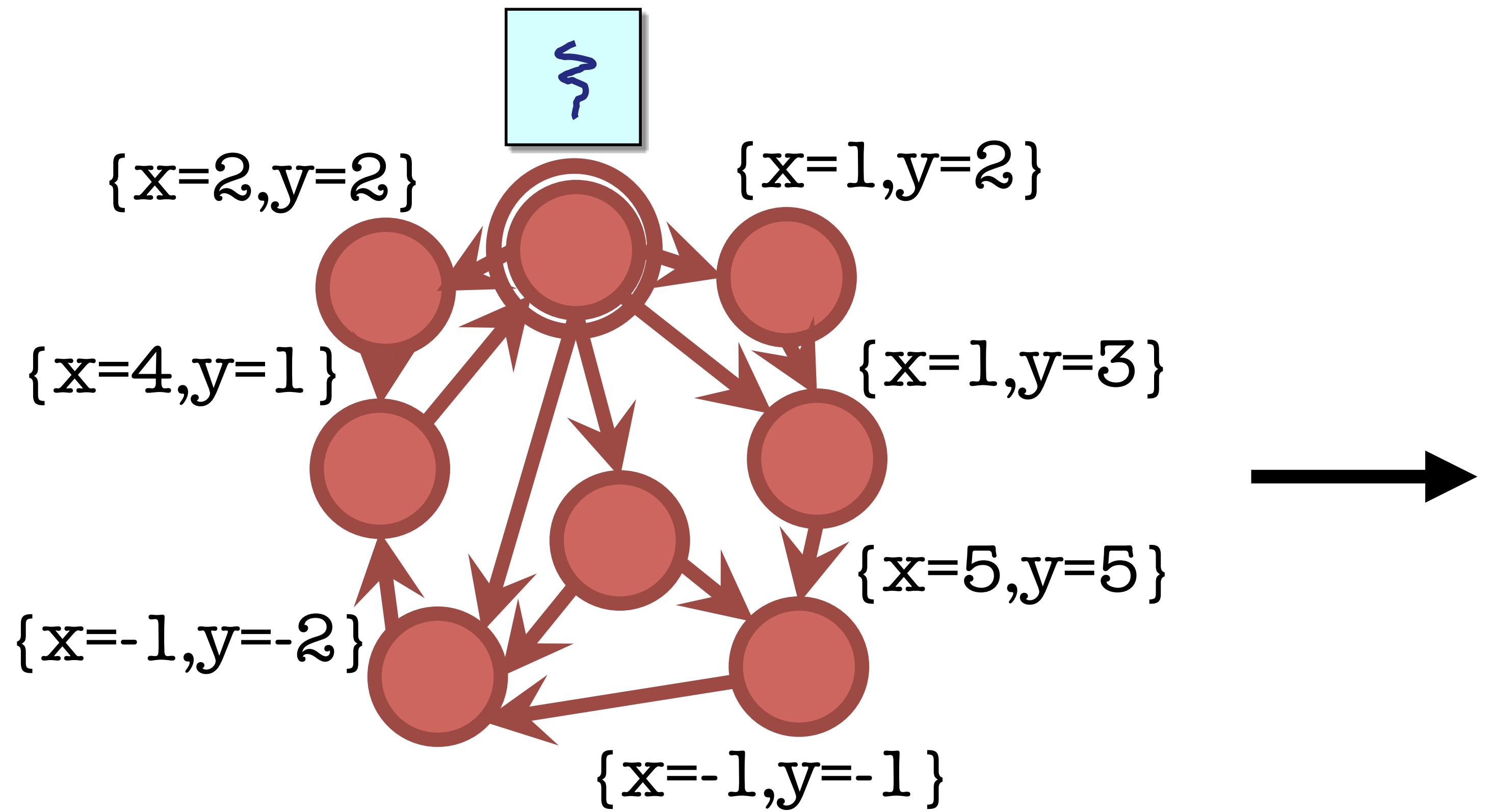
Conjecture:

Systems we care about
exhibit program properties
amenable with prioritization

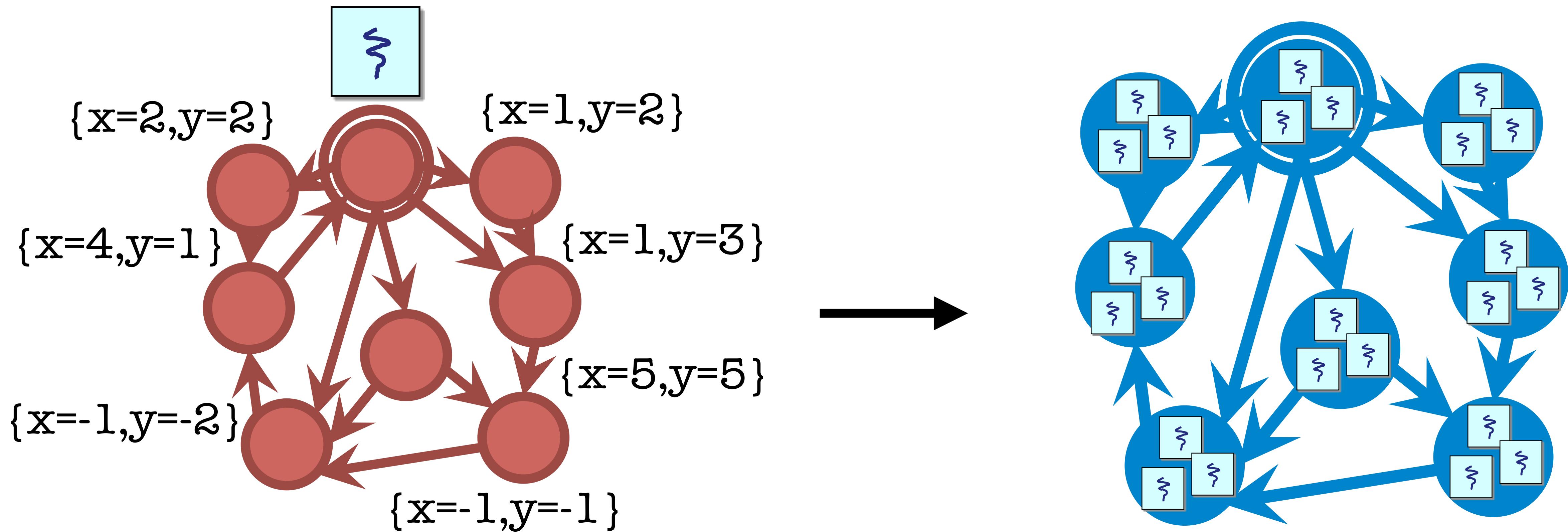




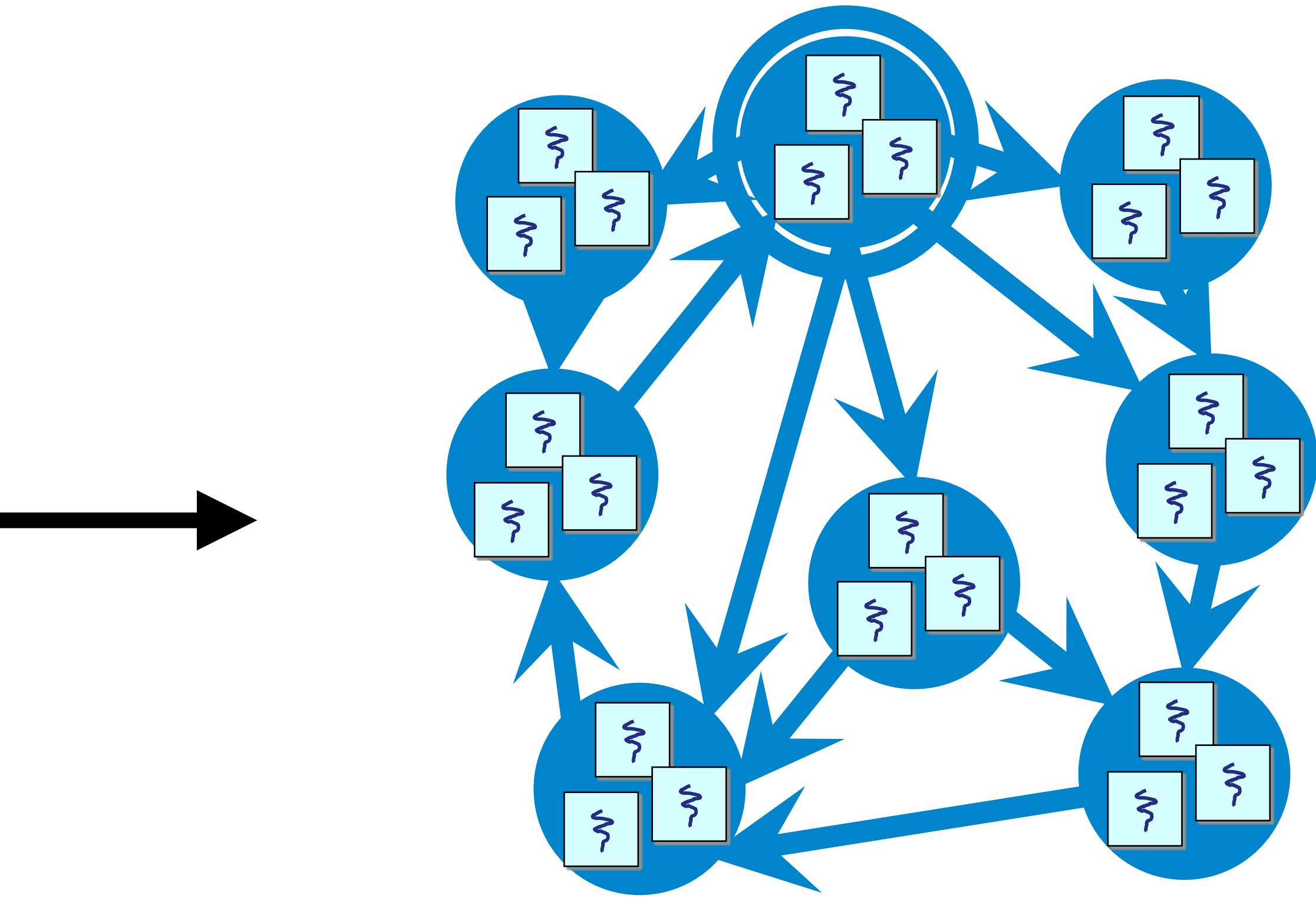
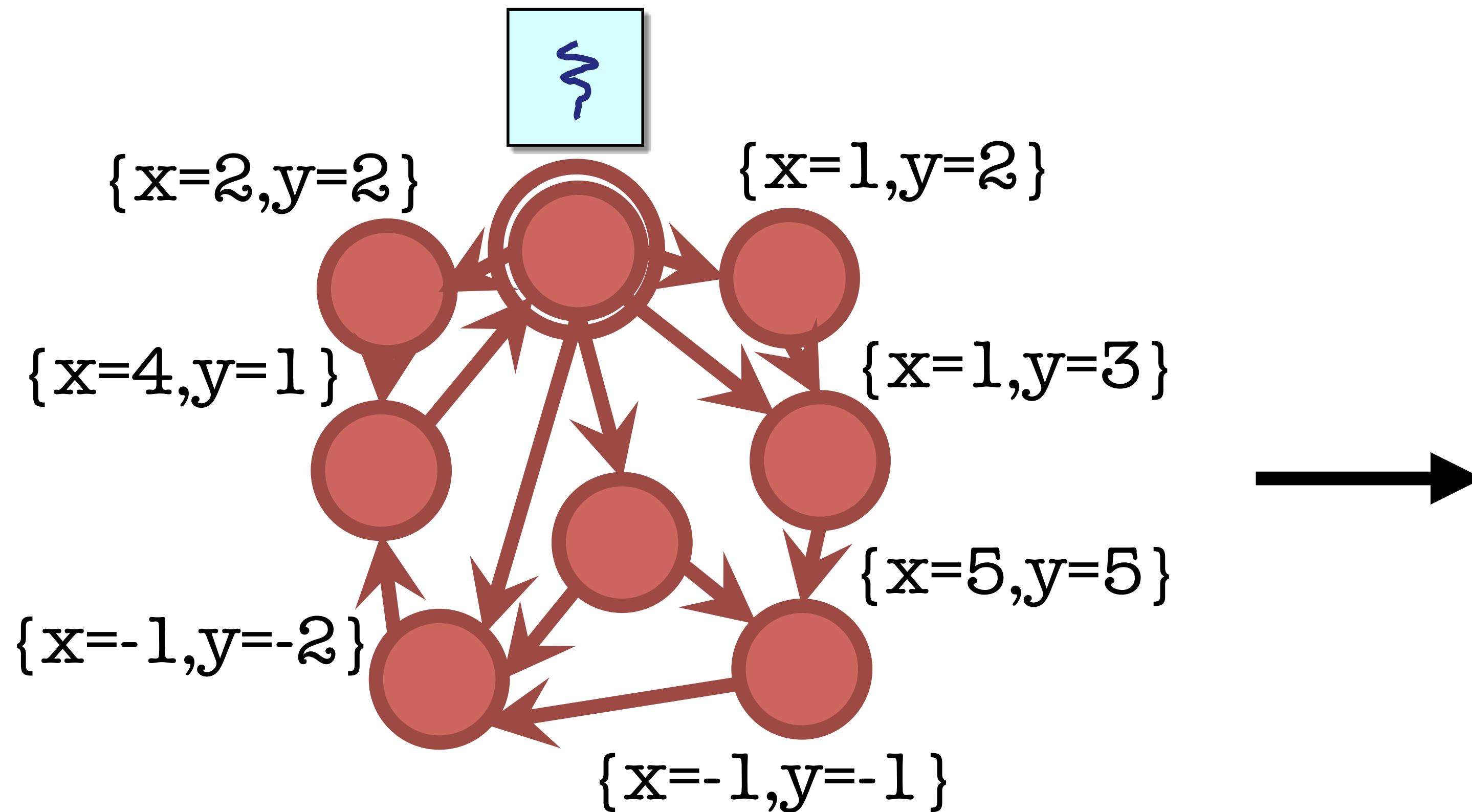




Invariant defined over small subset of processes' variables



- ▶ Invariant defined over small subset of processes' variables
- ▶ Each event affects a small subset of receiver's variables



- ▶ Invariant defined over small subset of processes' variables
- ▶ Each event affects a small subset of receiver's variables

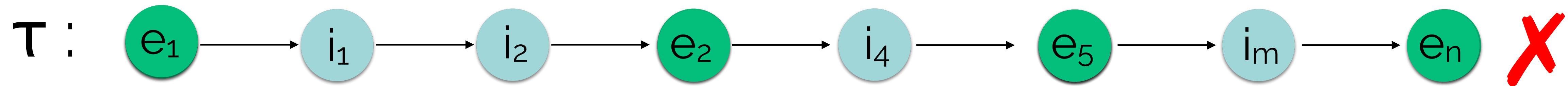
Initial execution contains events that don't affect invariant

Challenge:
Don't know which events are important

Approach:

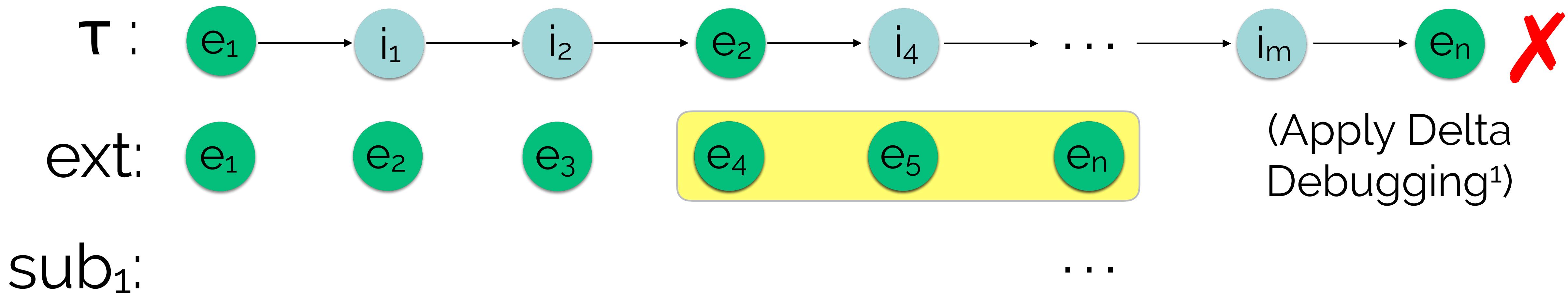
- ▶ experimentally “infer” important events
- ▶ stay close to the original execution

Observation #2: selectively mask original events



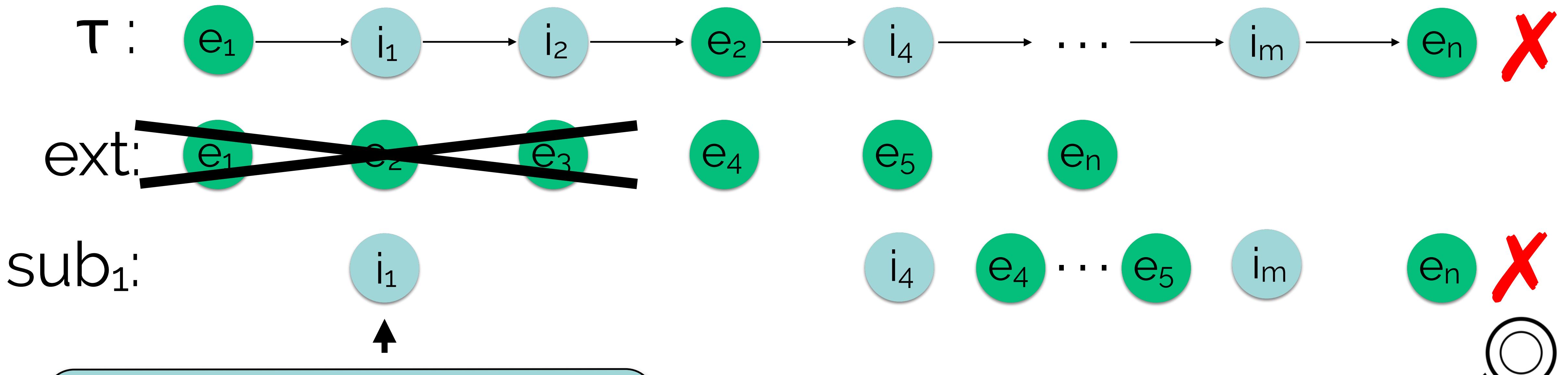
ext:

Observation #2: selectively mask original events



¹A Zeller, R. Hildebrandt, "Simplifying and Isolating Failure-Inducing Input", IEEE '02

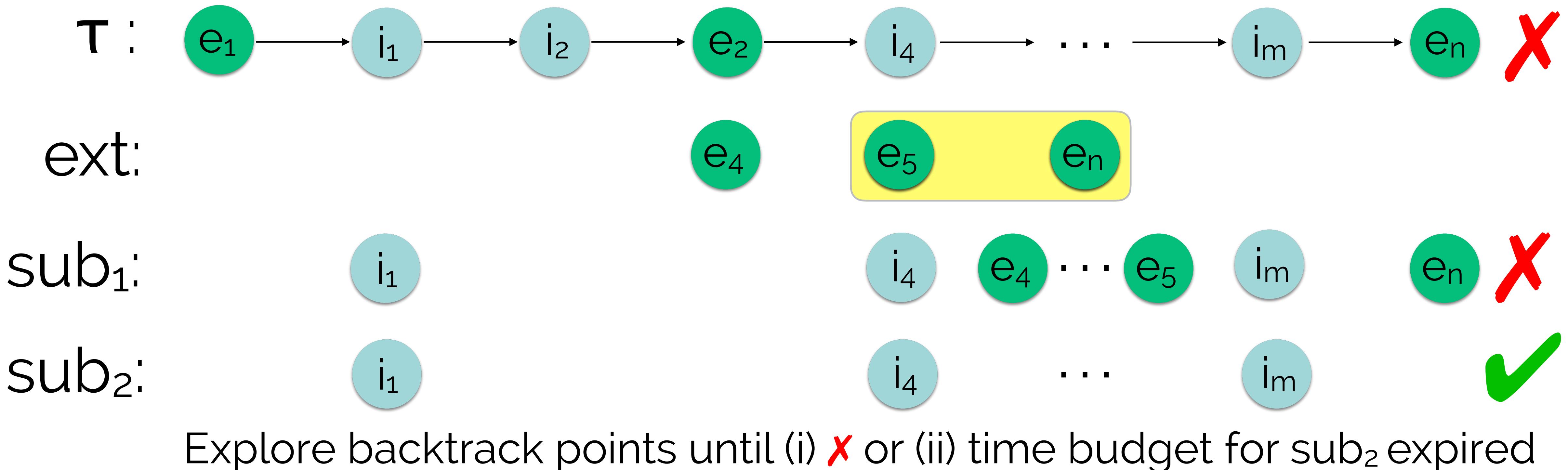
Observation #2: selectively mask original events



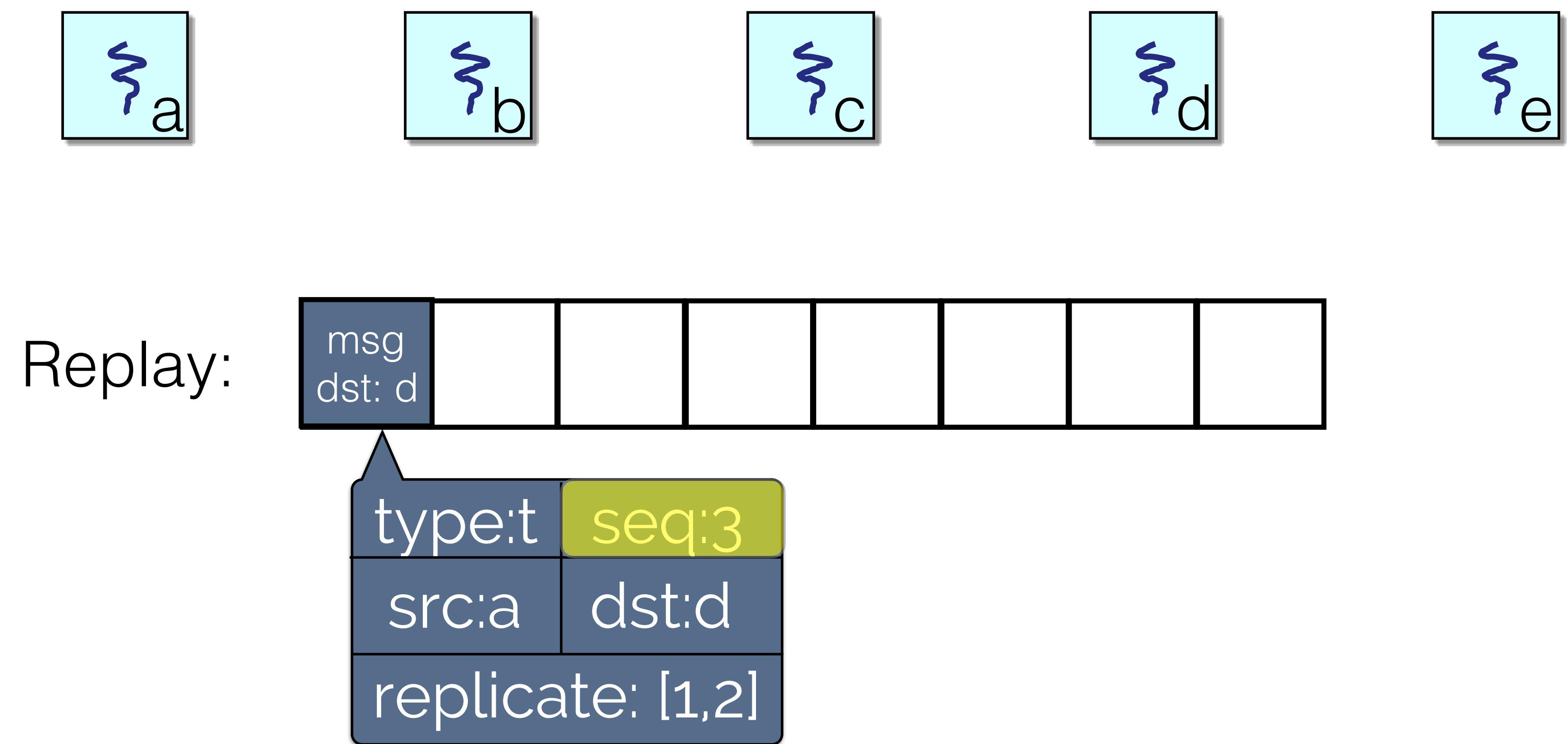
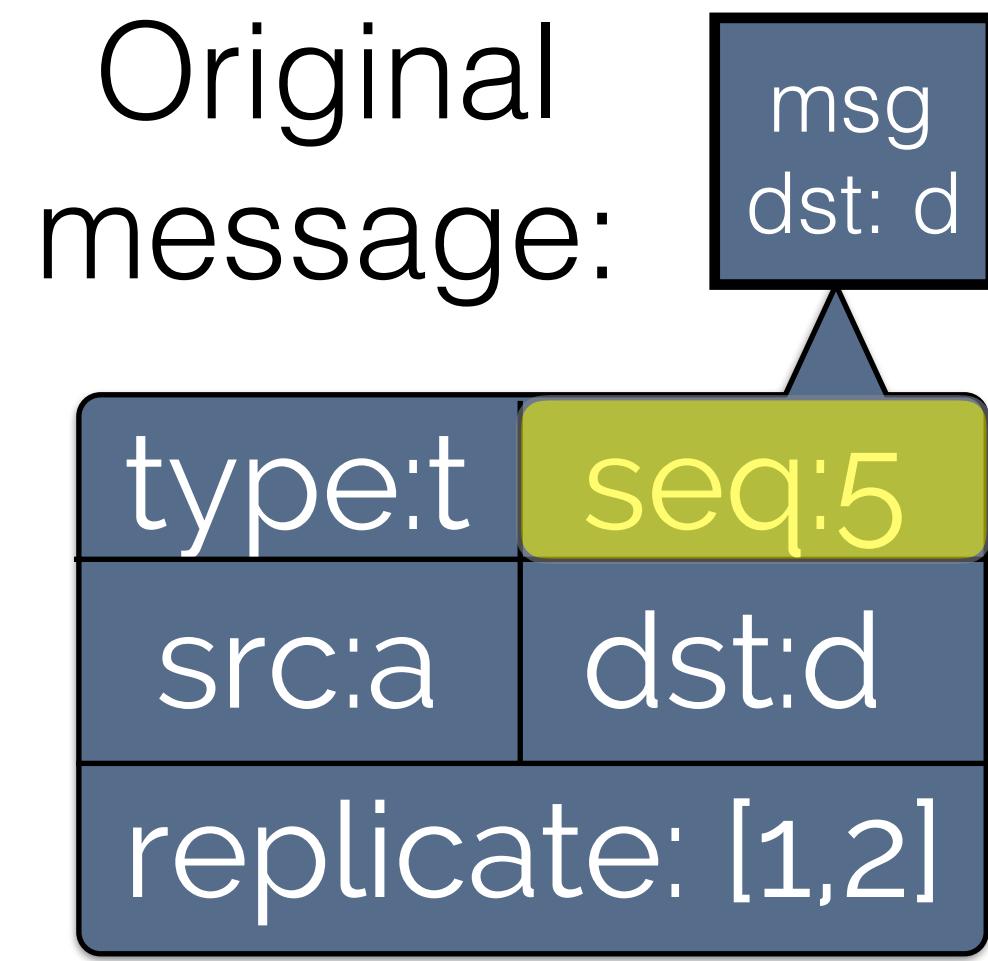
```
foreach i in  $\tau$ :  
    if i is pending:  
        deliver i  
    # ignore unexpected
```



Observation #2: selectively mask original events



Message contents
may differ across
executions!



Observation #3: some contents should be masked

Observation #3: some contents should be masked

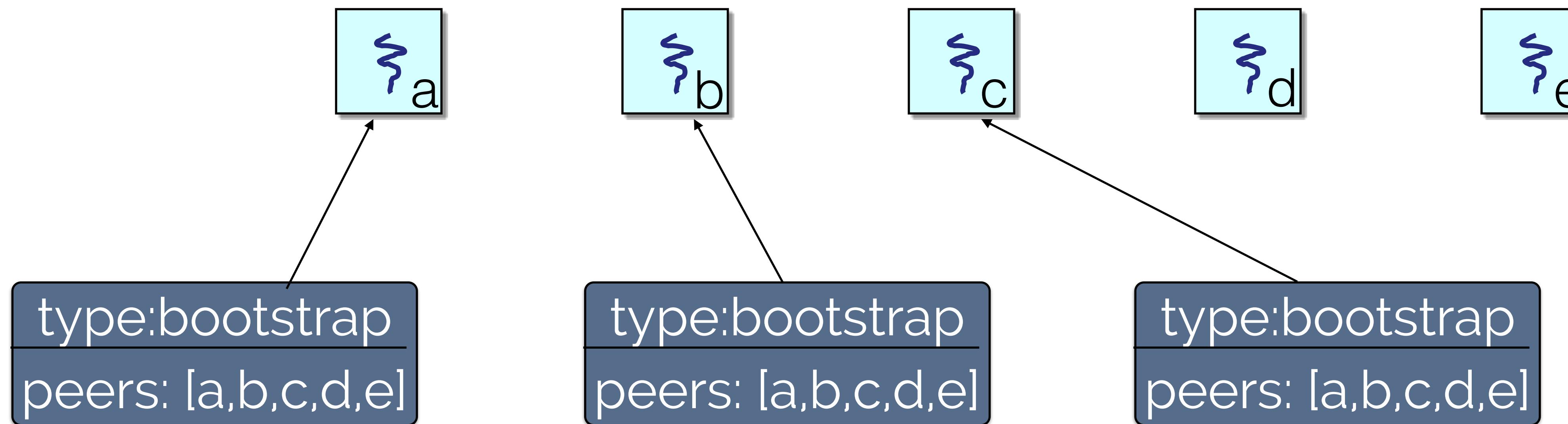
Phase 1: choose initial schedule

- ▶ Match messages by user-defined “fingerprint”

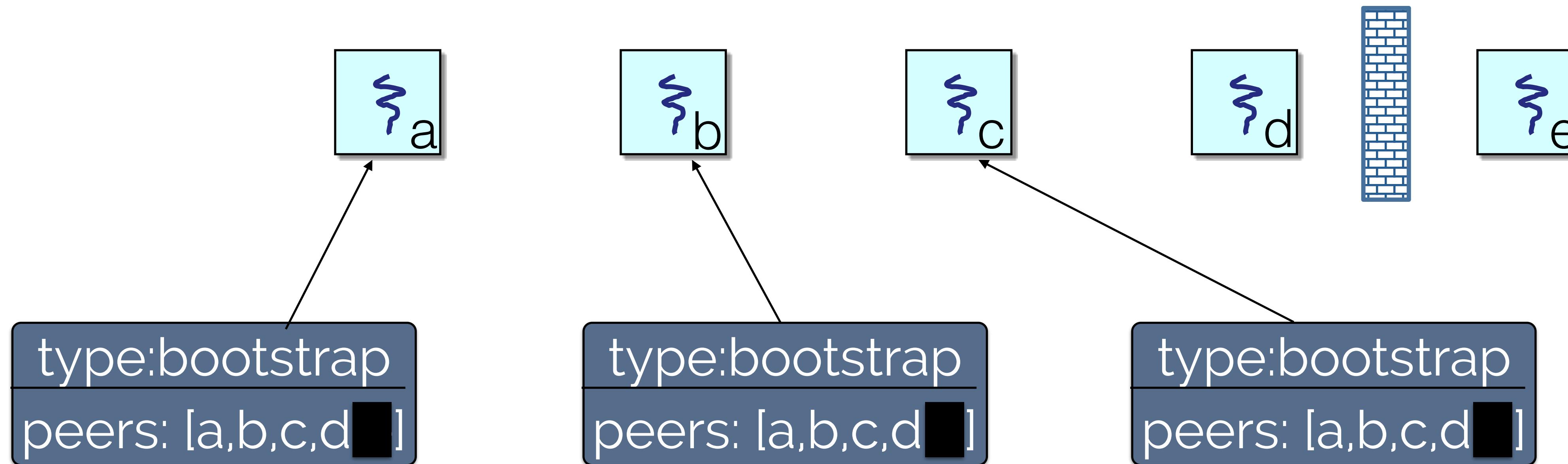
Phase 2: prioritize backtrack points

- ▶ Match messages by type only
- ▶ Backtrack whenever multiple pending messages match by type

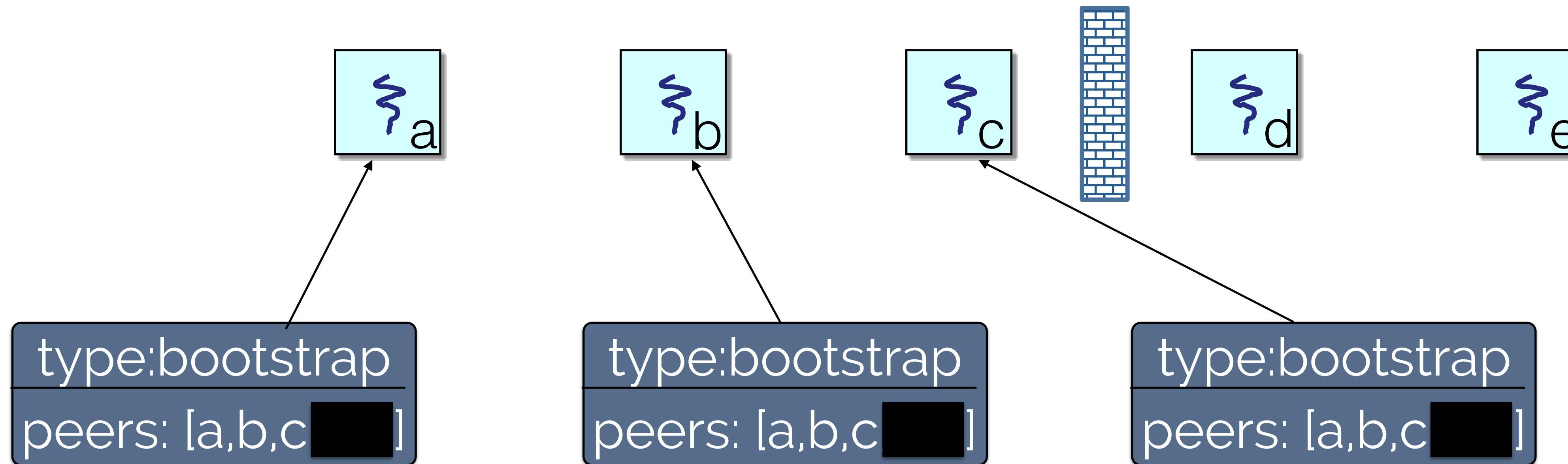
Observation #4: shrink external message contents



Observation #4: shrink external message contents



Observation #4: shrink external message contents



Goal: find minimal schedule that produces violation

Observation #1: many schedules are commutative

Approach: prioritize schedule space exploration

Observation #2: selectively mask original events

Observation #3: some contents should be masked

Observation #4: shrink external message contents

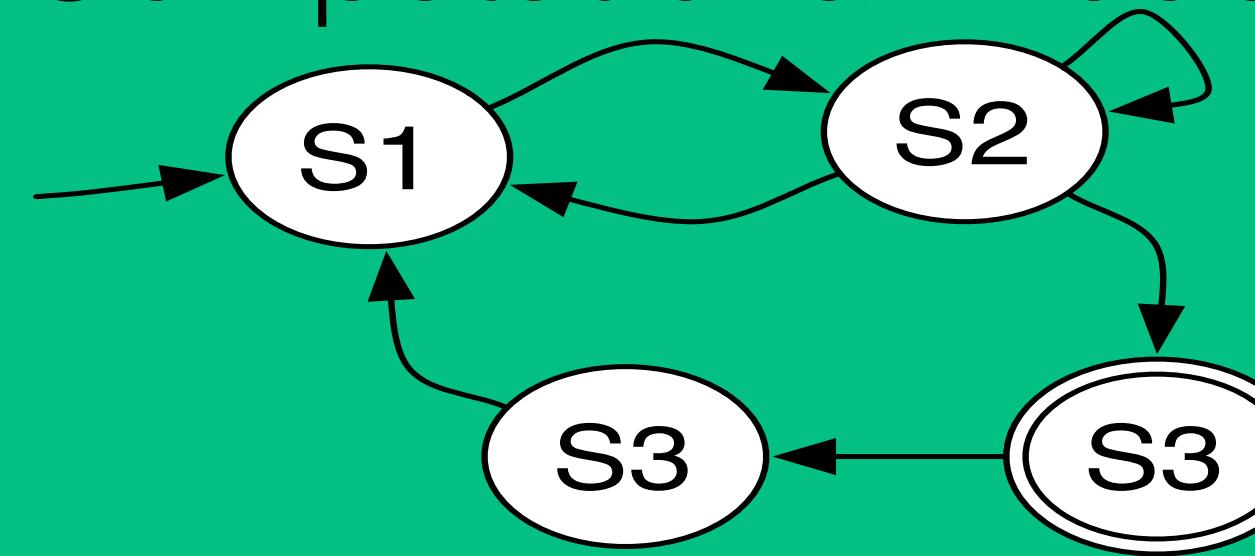
Minimize internal events after externals minimized

Outline

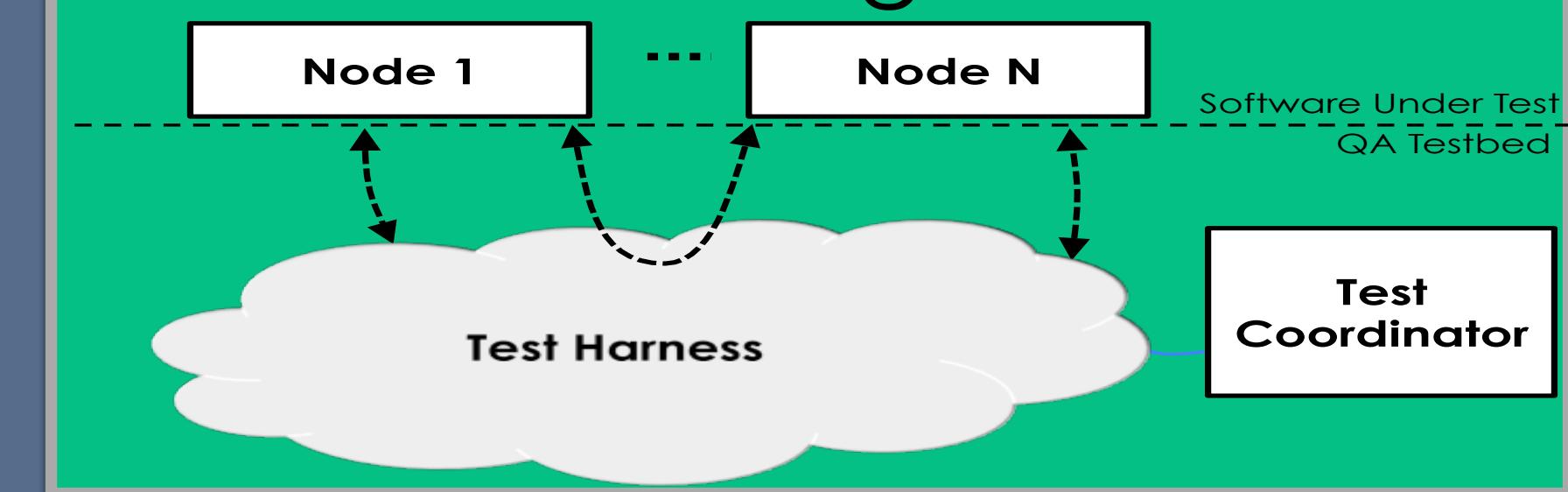
Introduction

Background

Computational Model



Fuzz Testing w/ DEMi



Minimization



Evaluation

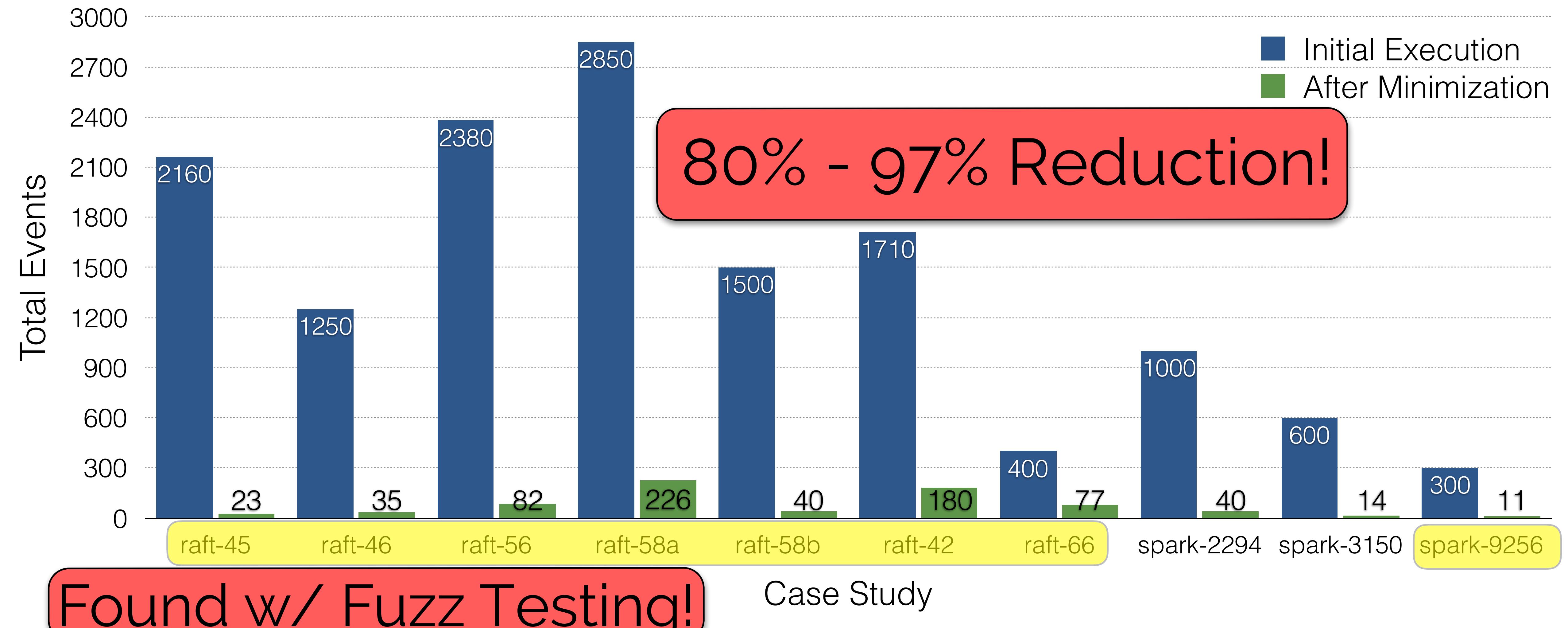


Conclusion

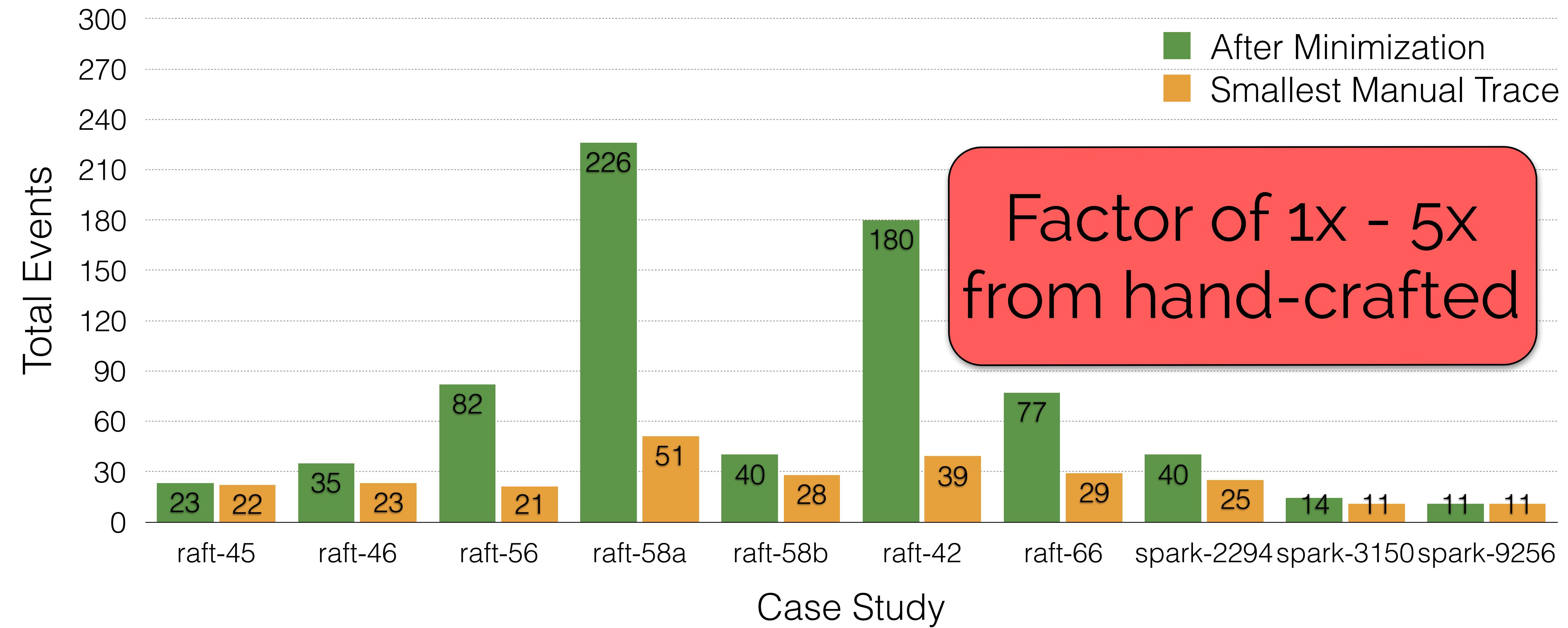
Target Systems



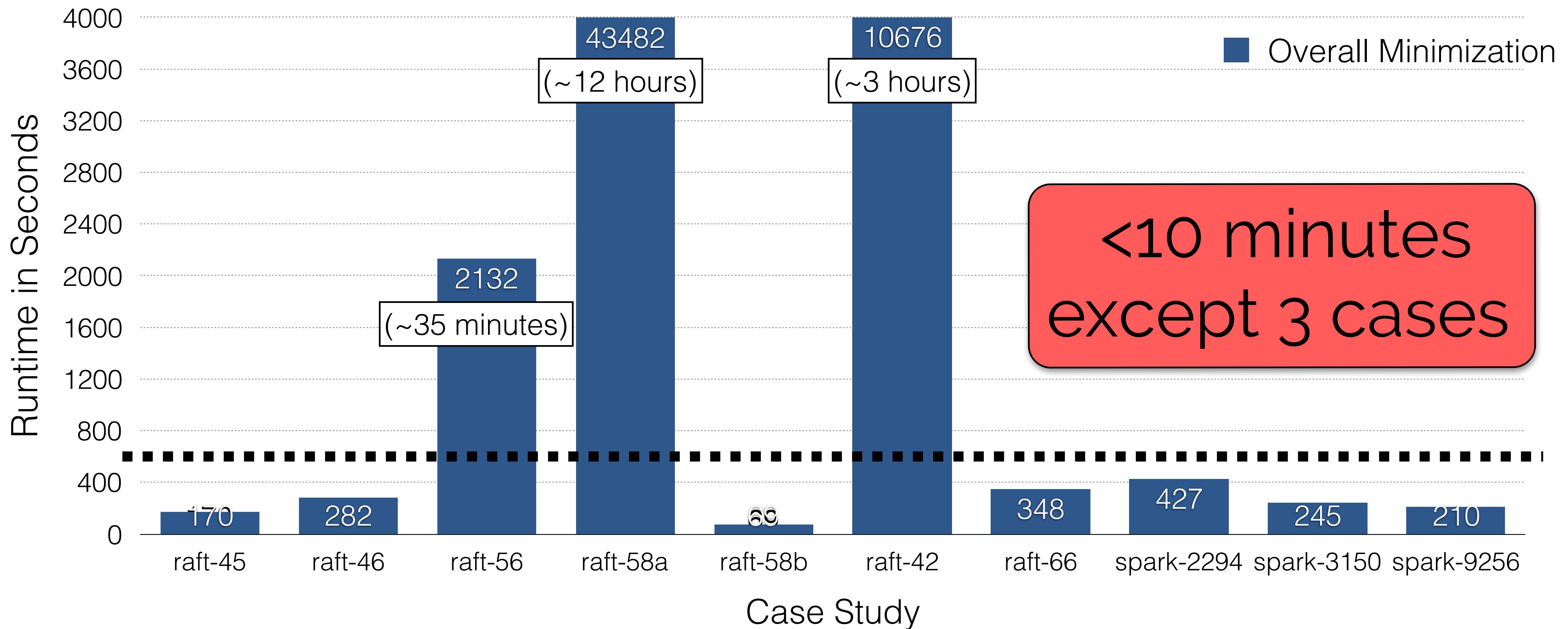
How well does DEMi work?



How well does DEMi work?



How quickly does DEMi work?



See the paper for...

- ▶ How we handle non-determinism
- ▶ Handling multithreaded processes
- ▶ Supporting other RPC libraries
- ▶ Sketch for minimizing production traces
- ▶ More in-depth evaluation
- ▶ Related work
- ▶ ...

Conclusion

Optimistic that these techniques can
be successfully applied more broadly

Open source tool: github.com/NetSys/demi

Read our paper!

eecs.berkeley.edu/~rcs/research/nsdi16.pdf

Contact me! cs@cs.berkeley.edu

Thanks for your time!

Attributions

Inspiration for slide design: Jay Lorch's IronFleet slides

Graphic Icons: thenounproject.org

logfile: mantisshrimpdesign

magnifying glass: Ricardo Moreira

disk: Anton Outkine

hook: Seb Cornelius

bug report: Lemon Liu

devil: Mourad Mokrane

Putin: Remi Mercier

Production Traces

Model: feed partially ordered log into single machine DEMi

Require:

- Partial ordering of all message deliveries
- All crash-recoveries logged to disk

Instrumentation Complexity

	akka-raft	Spark
Message Fingerprint	59	56
Non-Determinism	2	~250

Table 4: Instrumentation complexity (lines of code) needed to define message fingerprints, and to mitigate non-determinism.

Related Work

► Thread Schedule Minimization

- Isolating Failure-Inducing Thread Schedules. SIGSOFT '02.
- A Trace Simplification Technique for Effective Debugging of Concurrent Programs. FSE '10.

► Program Flow Analysis.

- Enabling Tracing of Long-Running Multithreaded Programs via Dynamic Execution Reduction. ISSTA '07.
- Toward Generating Reducible Replay Logs. PLDI '11.

► Best-Effort Replay of Field Failures

- A Technique for Enabling and Supporting Debugging of Field Failures. ICSE '07.
- Triage: Diagnosing Production Run Failures at the User's Site. SOSP '07.

DDmin in more detail

Input: $T_{\mathbf{x}}$ s.t. $T_{\mathbf{x}}$ is a trace and $\text{test}(T_{\mathbf{x}}) = \mathbf{x}$. Output: $T'_{\mathbf{x}} = ddmin(T_{\mathbf{x}})$ s.t. $T'_{\mathbf{x}} \subseteq T_{\mathbf{x}}$, $\text{test}(T'_{\mathbf{x}}) = \mathbf{x}$, and $T'_{\mathbf{x}}$ is minimal.

$$ddmin(T_{\mathbf{x}}) = ddmin_2(T_{\mathbf{x}}, \emptyset) \quad \text{where}$$

$$ddmin_2(T'_{\mathbf{x}}, R) = \begin{cases} T'_{\mathbf{x}} & \text{if } |T'_{\mathbf{x}}| = 1 \text{ ("base case")} \\ ddmin_2(T_1, R) \\ ddmin_2(T_2, R) & \text{else if } \text{test}(T_1 \cup R) = \mathbf{x} \text{ ("in } T_1\text{")} \\ ddmin_2(T_1, T_2 \cup R) \cup ddmin_2(T_2, T_1 \cup R) & \text{else if } \text{test}(T_2 \cup R) = \mathbf{x} \text{ ("in } T_2\text{")} \\ & \text{otherwise ("interference")} \end{cases}$$

where $\text{test}(T)$ denotes the state of the system after executing the trace T , \mathbf{x} denotes an invariant violation, $T_1 \subset T'_{\mathbf{x}}$, $T_2 \subset T'_{\mathbf{x}}$, $T_1 \cup T_2 = T'_{\mathbf{x}}$, $T_1 \cap T_2 = \emptyset$, and $|T_1| \approx |T_2| \approx |T'_{\mathbf{x}}|/2$ hold.

DDmin assumptions

- *Monotonic:*

$$P \oplus C = \chi \Rightarrow P \oplus (C \cup C') \neq \checkmark$$

- *Unambiguous:*

$$P \oplus C = \chi \wedge P \oplus C' = \chi \Rightarrow P \oplus (C \cap C') \neq \checkmark$$

- *Consistent*

$$P \oplus C \neq ?$$

Local vs. global minima

Definition 8 (Global minimum). A set $c \subseteq c_\chi$ is called the global minimum of c_χ if: $\forall c' \subseteq c_\chi \cdot (|c'| < |c| \Rightarrow \text{test}(c') \neq \chi)$ holds.

Definition 10 (n -minimal test case). A test case $c \subseteq c_\chi$ is n -minimal if: $\forall c' \subset c \cdot |c| - |c'| \leq n \Rightarrow (\text{test}(c') \neq \chi)$ holds. Consequently, c is 1-minimal if $\forall \delta_i \in c \cdot \text{test}(c - \{\delta_i\}) \neq \chi$ holds.

Minimization Pace



Figure 2: Minimization pace for raft-58b. Significant progress is made early on, then progress becomes rare.

Dealing With Threads

If you're lucky: threads are largely independent (Spark)

If you're unlucky: key insight:

A write to shared memory is equivalent to a message delivery

Approach:

- interpose on virtual memory, thread scheduler
- pause a thread whenever it writes to shared memory / disk

Cf. "Enabling Tracing Of Long-Running Multithreaded Programs Via Dynamic Execution Reduction", ISSTA '07

Dealing With Non-Determinism

Interpose on:

- Timers
- Random number generators
- Unordered hash values
- ID allocation

Stop-gap: replay each schedule multiple times

Complete Results

Bug Name	Bug Type	Initial	Provenance	STSSched	TFB	Optimal	NoDiverge
raft-45	Akka-FIFO, reproduced	2160 (E:108)	2138 (E:108)	1183 (E:8)	23 (E:8)	22 (E:8)	1826 (E:11)
raft-46	Akka-FIFO, reproduced	1250 (E:108)	1243 (E:108)	674 (E:8)	35 (E:8)	23 (E:6)	896 (E:9)
raft-56	Akka-FIFO, found	2380 (E:108)	2376 (E:108)	1427 (E:8)	82 (E:8)	21 (E:8)	2064 (E:9)
raft-58a	Akka-FIFO, found	2850 (E:108)	2824 (E:108)	953 (E:32)	226 (E:31)	51 (E:11)	2368 (E:35)
raft-58b	Akka-FIFO, found	1500 (E:208)	1496 (E:208)	164 (E:13)	40 (E:8)	28 (E:8)	1103 (E:13)
raft-42	Akka-FIFO, reproduced	1710 (E:208)	1695 (E:208)	1093 (E:39)	180 (E:21)	39 (E:16)	1264 (E:43)
raft-66	Akka-UDP, found	400 (E:68)	392 (E:68)	262 (E:23)	77 (E:15)	29 (E:10)	279 (E:23)
spark-2294	Akka-FIFO, reproduced	1000 (E:30)	886 (E:30)	43 (E:3)	40 (E:3)	25 (E:1)	43 (E:3)
spark-3150	Akka-FIFO, reproduced	600 (E:20)	536 (E:20)	18 (E:3)	14 (E:3)	11 (E:3)	18 (E:3)
spark-9256	Akka-FIFO, found (rare)	300 (E:20)	256 (E:20)	11 (E:1)	11 (E:1)	11 (E:1)	11 (E:1)

Runtime Breakdown

Bug Name	STSSched	TFB	
raft-45	56s (594)	114s	(2854)
raft-46	73s (384)	209s	(4518)
raft-56	54s (524)	2078s	(31149)
raft-58a	137s (624)	43345s	(834972)
raft-58b	23s (340)	31s	(1747)
raft-42	118s (568)	10558s	(176517)
raft-66	14s (192)	334s	(10334)
spark-2294	330s (248)	97s	(78)
spark-3150	219s (174)	26s	(21)
spark-9256	96s (73)	0s	(0)

Integrating with other RPC libs

