

Approximate Computing: (Old) Hype or New Frontier?

Sarita Adve

University of Illinois, EPFL

**Acks: Vikram Adve, Siva Hari, Man-Lap Li, Abdulrahman Mahmoud,
Helia Naemi, Pradeep Ramachandran, Swarup Sahoo, Radha
Venkatagiri, Yuanyuan Zhou**



What is Approximate Computing?

- Trading output quality for resource management?
- Exploiting applications' inherent ability to tolerate errors?

Old?

GRACE [2000-05] rsim.cs.illinois.edu/grace

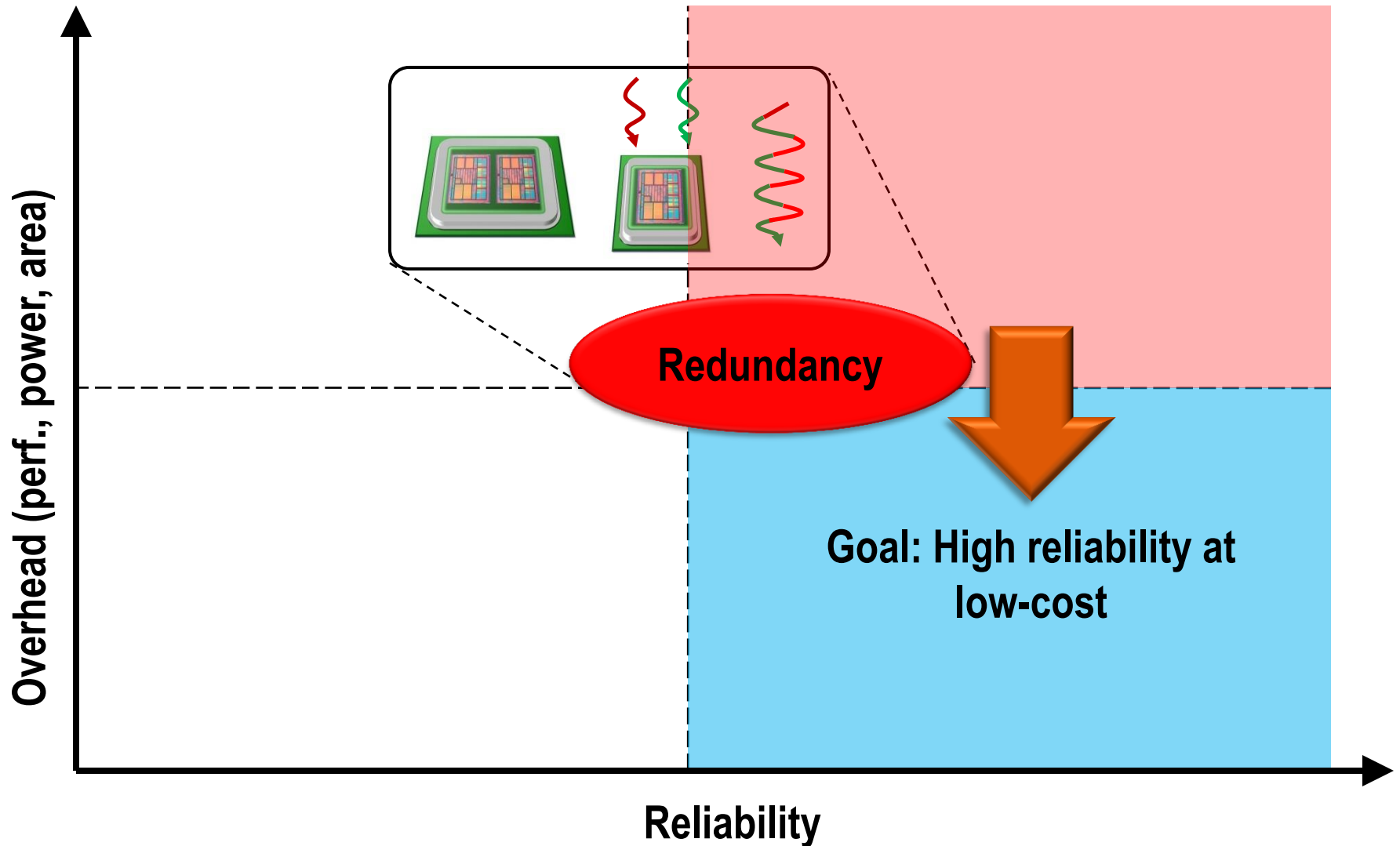
SWAT [2006-] rsim.cs.illinois.edu/swat

Many others

Or something new?

Approximate computing through the lens of hardware resiliency ...

Motivation



SWAT: A Low-Cost Reliability Solution

- Need handle only hardware faults that affect software

⇒ Watch for software anomalies (symptoms)

- Zero to low overhead “always-on” monitors

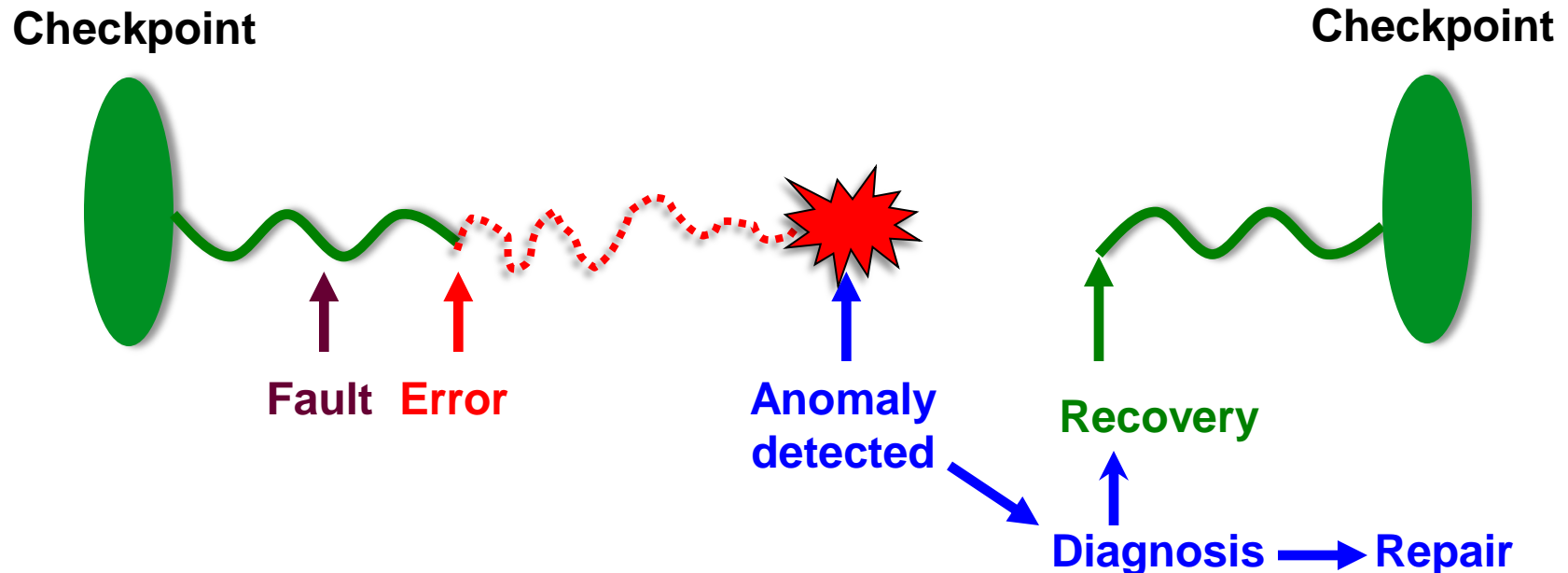
Diagnose cause after anomaly detected and recover

- May incur high overhead, but invoked infrequently

⇒ SWAT: SoftWare Anomaly Treatment

SWAT Framework Components

- **Detection:** Monitor symptoms of software misbehavior
- **Diagnosis:** Rollback/replay on multicore
- **Recovery:** Checkpoint/rollback or app-specific action on acceptable errors
- **Repair/reconfiguration:** Redundant, reconfigurable hardware
- **Flexible control through firmware**



SWAT Framework Components

- **Detection:** Monitor symptoms of software misbehavior
- **Diagnosis:** Rollback/replay on multicore
- **Recovery:** Checkpoint/rollback or app-specific action on acceptable errors
- **Repair/reconfiguration:** Redundant, reconfigurable hardware
- **Flexible control through firmware**

How to bound silent (escaped) errors?

When is an error acceptable (at some utility) or unacceptable?

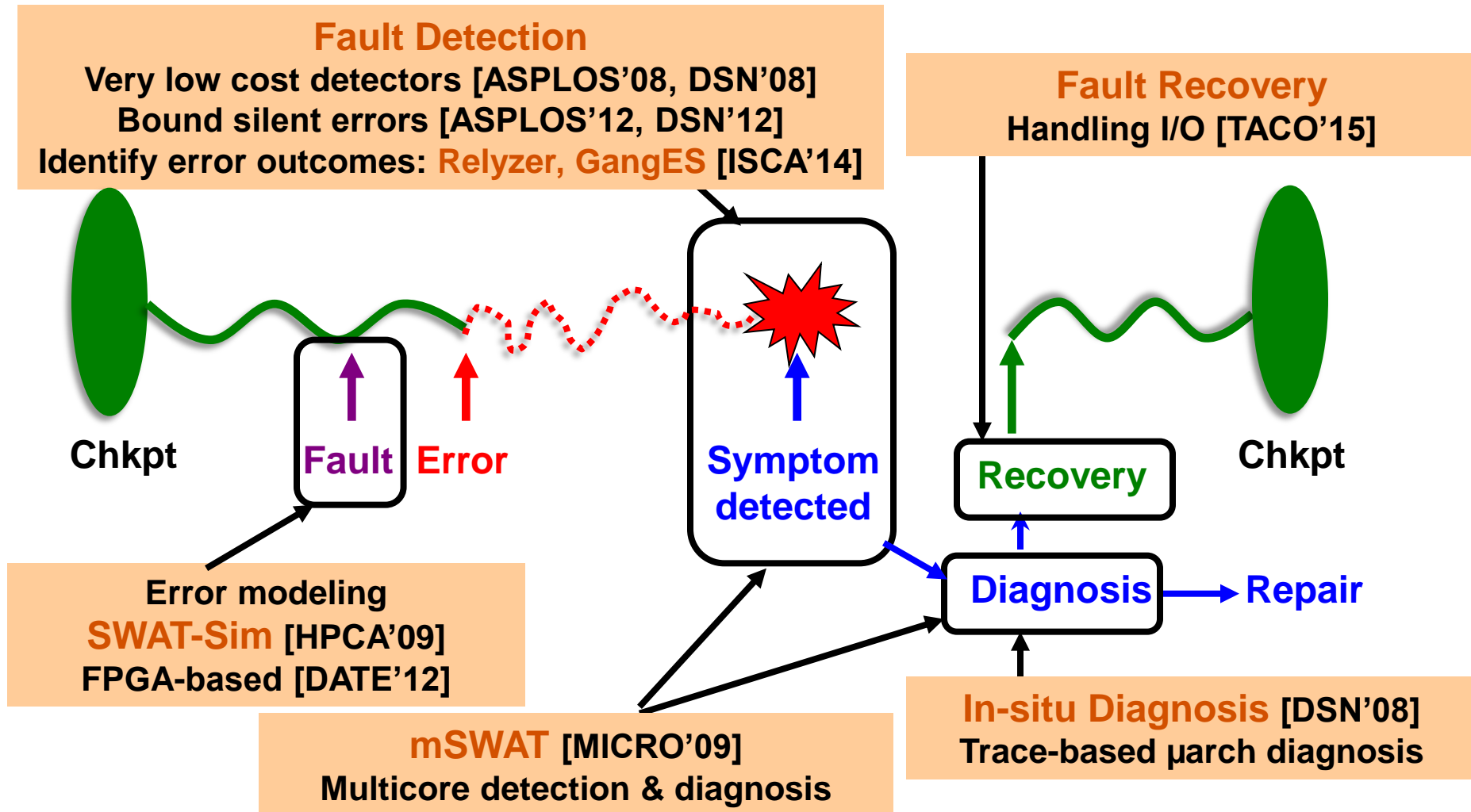
How to trade silent errors with resources?

How to associate recovery actions with errors?

Advantages of SWAT

- **Handles only errors that matter**
 - Oblivious to masked error, low-level failure modes, software-tolerated errors
- **Low, amortized overheads**
 - Optimize for common case, exploit software reliability solutions
- **Customizable and flexible**
 - Firmware control can adapt to specific reliability needs
- **Holistic systems view enables novel solutions**
 - Software-centric synergistic detection, diagnosis, recovery solutions
- **Beyond hardware reliability**
 - Long term goal: unified system (HW+SW) reliability
 - Systematic trade off between resource usage, quality, reliability

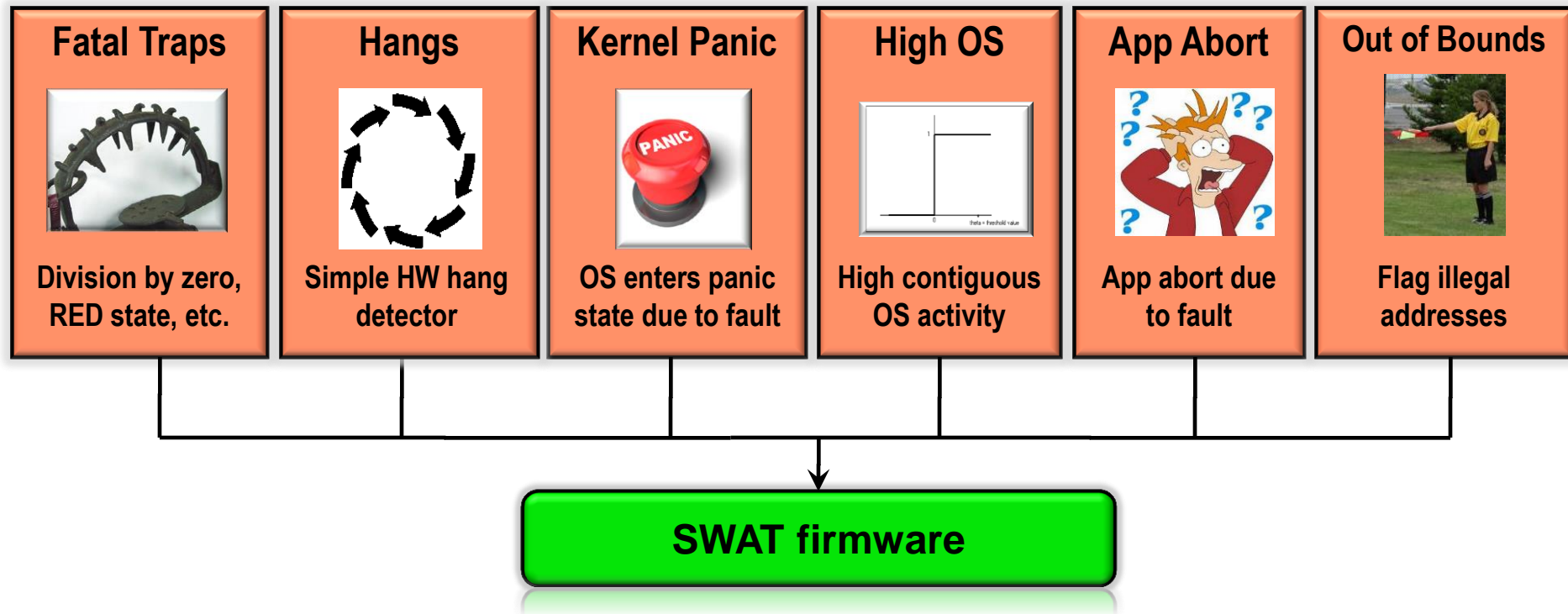
SWAT Contributions (for in-core hw faults)



Complete solution for in-core faults, evaluated for variety of workloads

SWAT Fault Detection

- Simple monitors observe **anomalous SW behavior** [ASPLOS'08, MICRO'09]

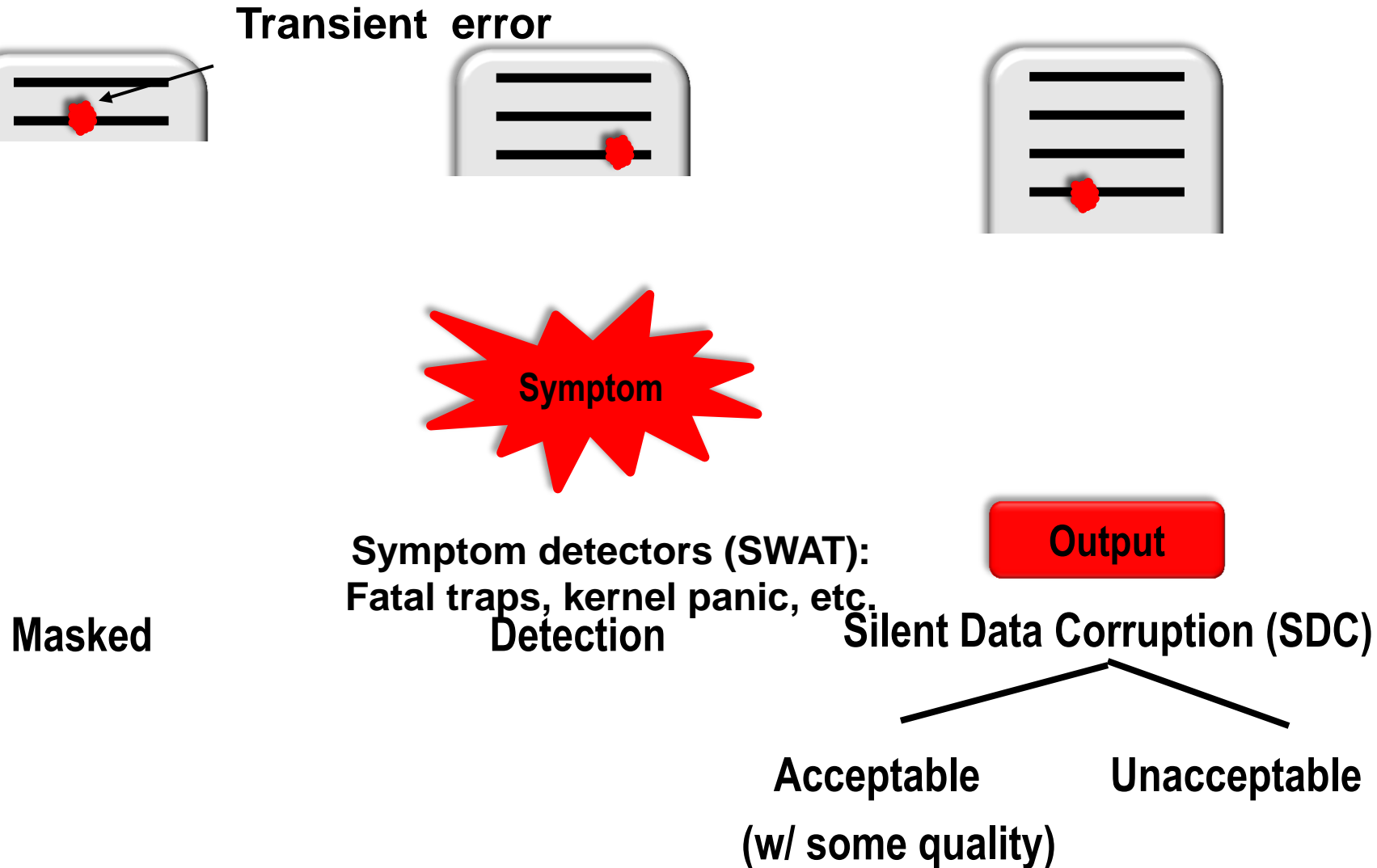


- Very low hardware area, performance overhead

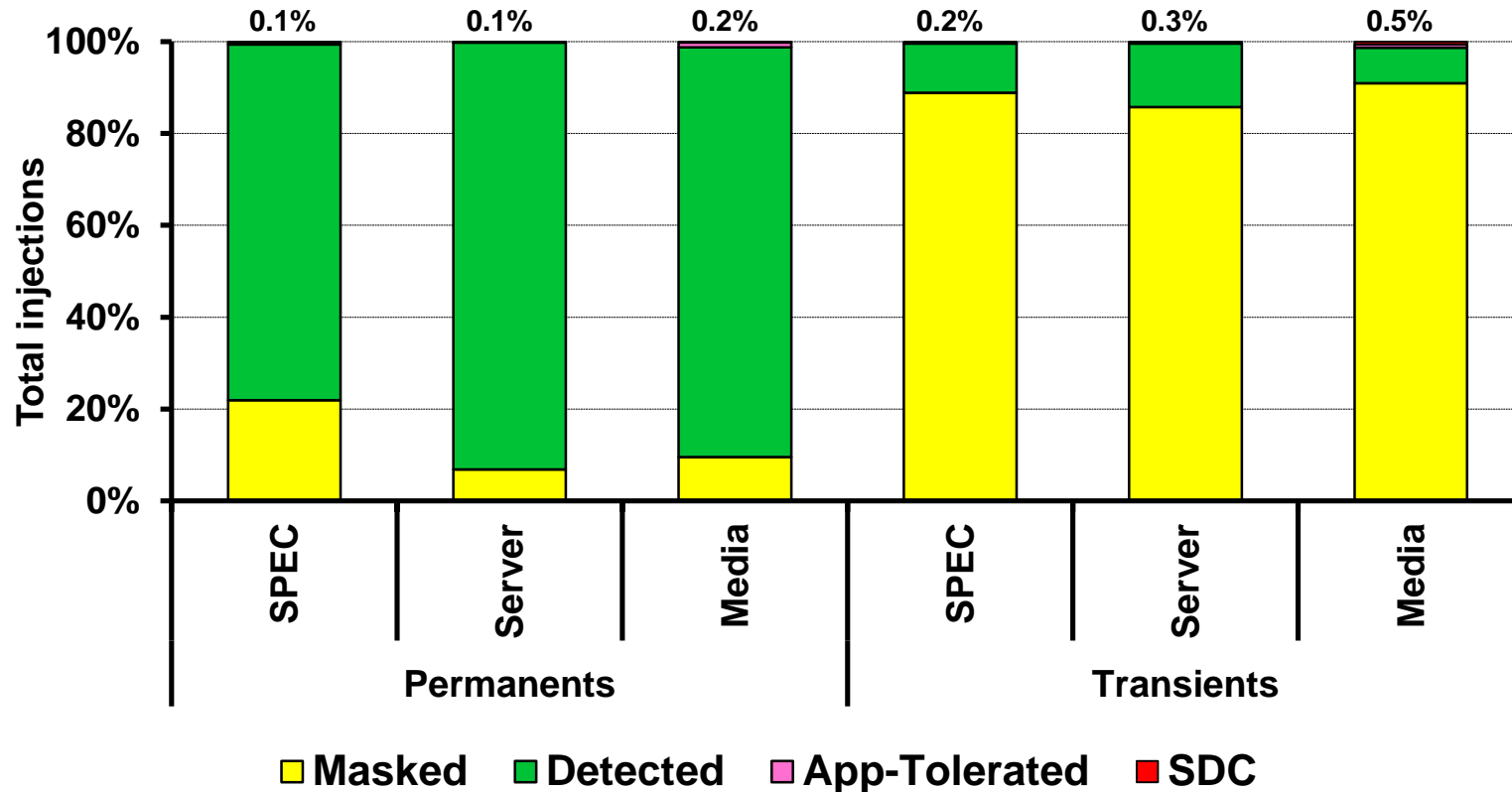
Evaluating SWAT Detectors So Far

- Full-system simulation with out-of-order processor
 - Simics functional + GEMS timing simulator
 - Single core, multicore, distributed client-server
- Apps: Multimedia, I/O intensive, & compute intensive
 - Errors injected at different points in app execution
- μ arch-level error injections (single error model)
 - Stuck-at, transient errors in latches of 8 μ arch units
 - ~48,000 total errors

Error Outcomes

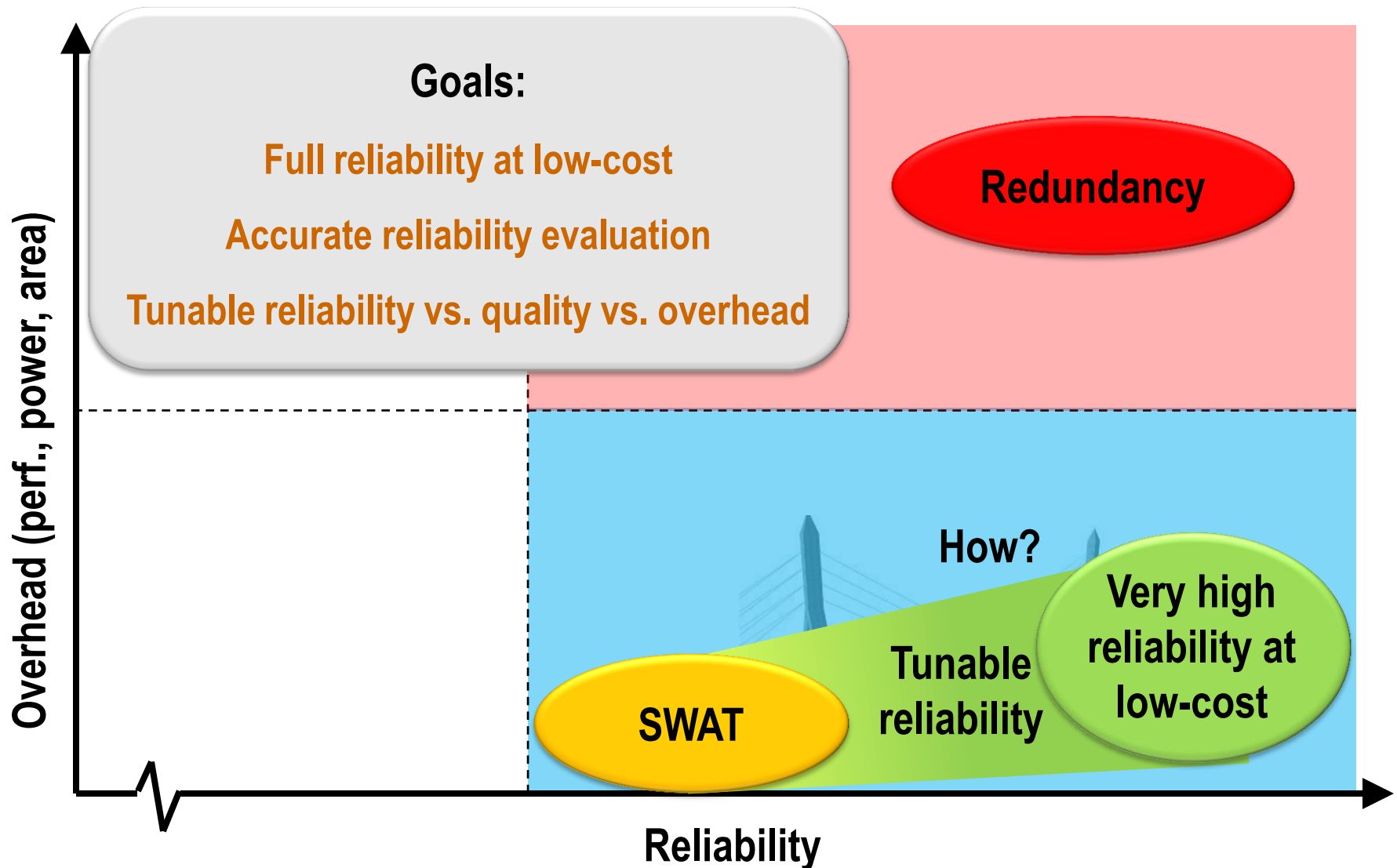


SWAT SDC Rates



- **SWAT detectors effective**
- **<0.2% of injected march errors give unacceptable SDCs (109/48,000 total)**
- **BUT hard to sell empirically validated, rarely incorrect systems**

Challenges



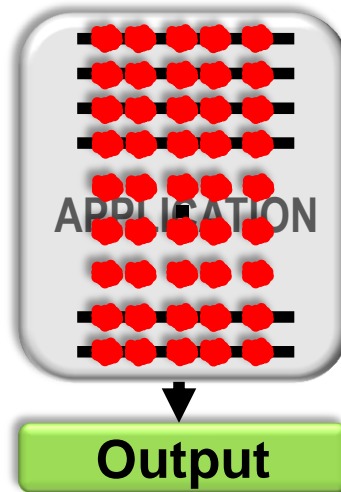
Research Strategy

- **Towards an Application Resiliency Profile**
 - For a given instruction what is the outcome of an error?
 - For now, focus on a transient error in single bit of register
- **Convert SDCs into (low-cost) detections for full reliability and quality**
- **OR let some acceptable or unacceptable SDCs escape**
 - **Quantitative tuning** of reliability vs. overhead vs. quality

Challenges and Approach

- Determine error outcomes for *all* application-sites

How?
Complete app
resiliency evaluation



Impractical, too many injections
>1,000 compute-years for one app

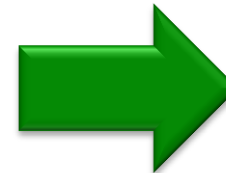
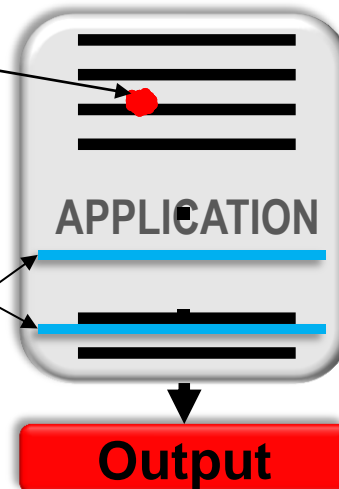
Challenge: Analyze all errors with few injections

- Cost-effectively convert (some) SDCs to Detections

How?

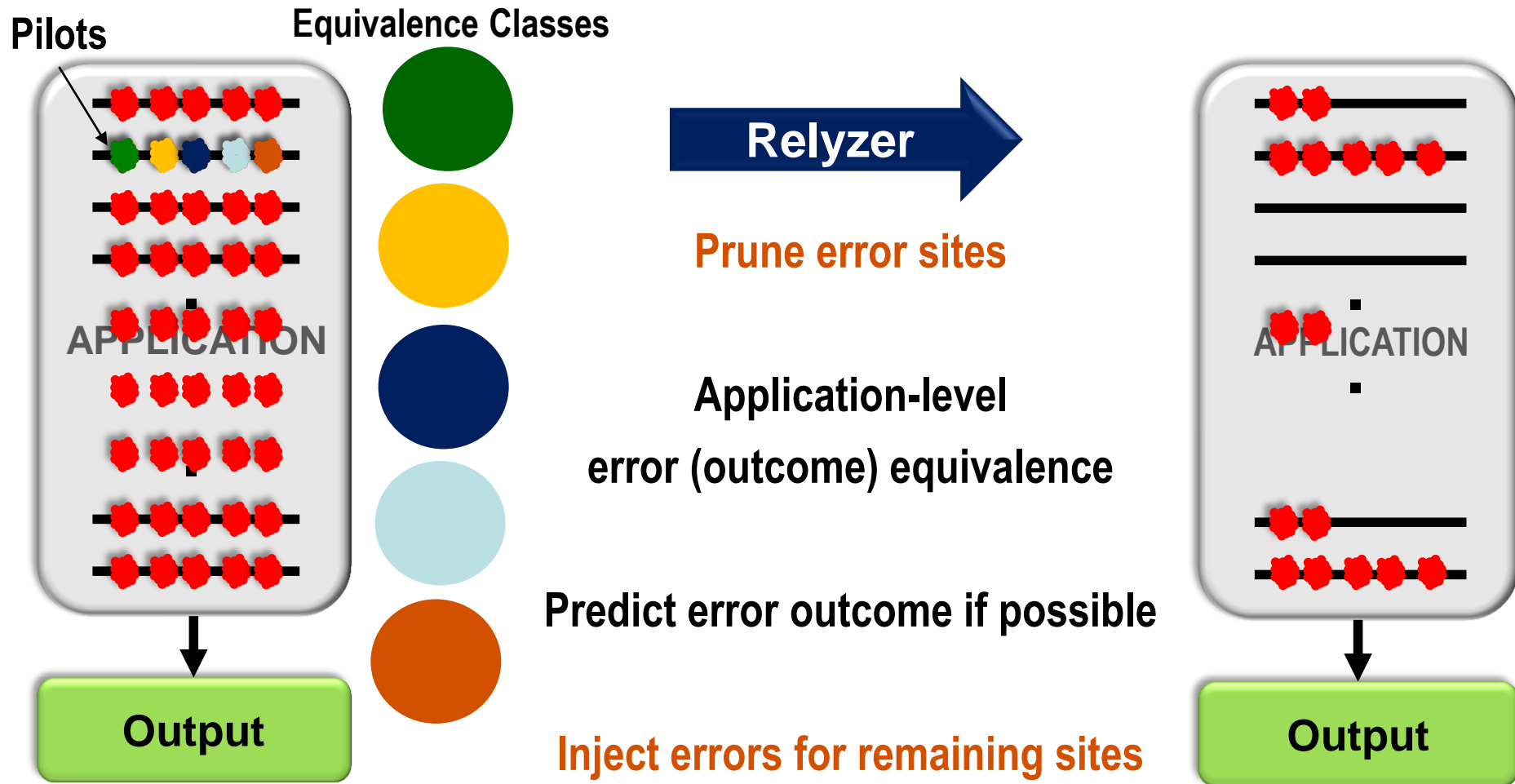
SDC-causing
error

Error
Detectors



Challenges:
What to use?
Where to place?
How to tune?

Relyzer: Application Resiliency Analyzer [ASPLOS'12]



Can list virtually all SDC-causing instructions

Def to First-Use Equivalence

- Fault in first use is equivalent to fault in def \Rightarrow prune def

Def \longrightarrow ~~r1~~ = r2 + r3

r4 = r1 + r5



First use

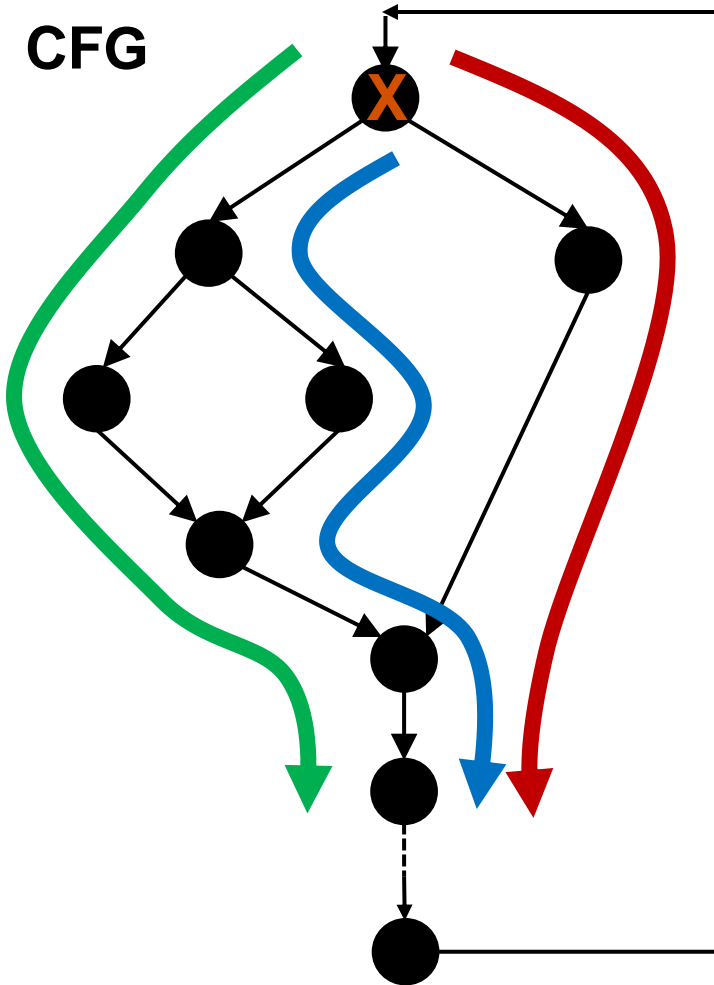
...

- If there is no first use, then def is dead \Rightarrow prune def

Control Flow Equivalence

Insight: Errors flowing through similar control paths may behave similarly

CFG

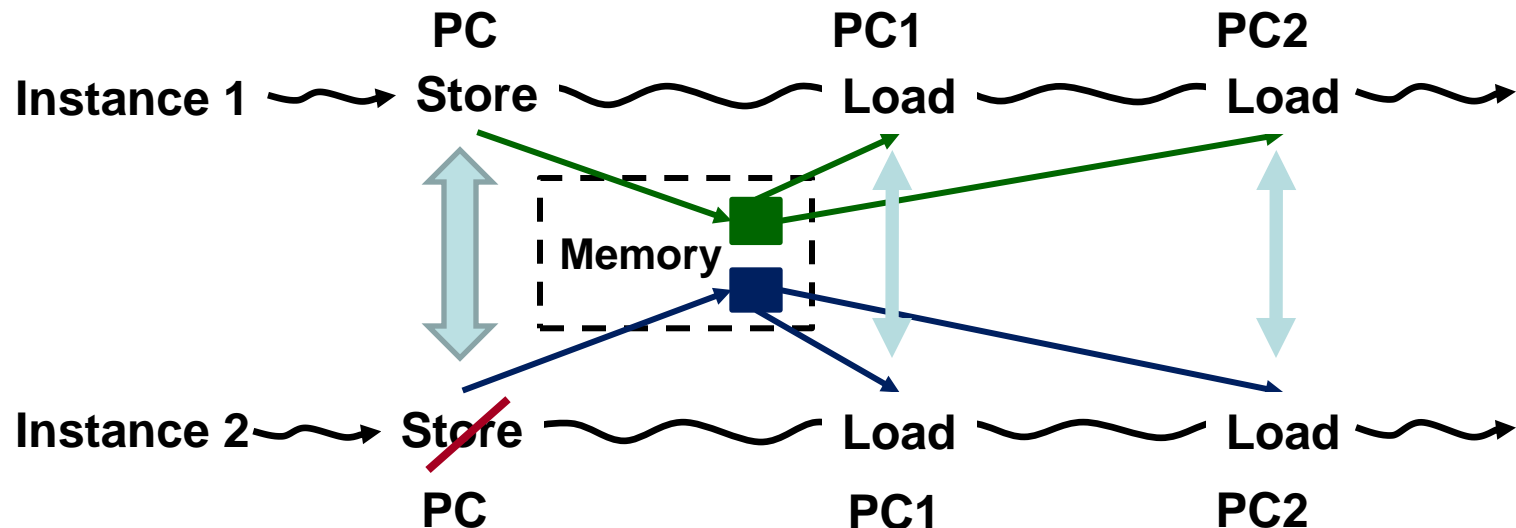


Errors in **X** that take ■ path behave similarly

Heuristic: Use direction of next 5 branches

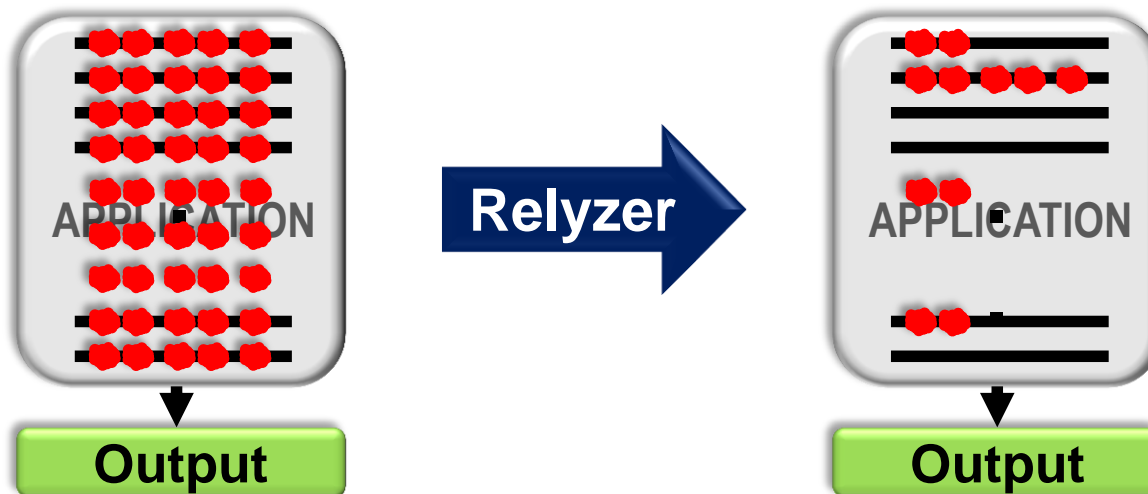
Store Equivalence

- **Insight: Errors in stores may be similar if stored values are used similarly**
- Heuristic to determine similar use of values:
 - Same number of loads use the value
 - Loads are from same PCs



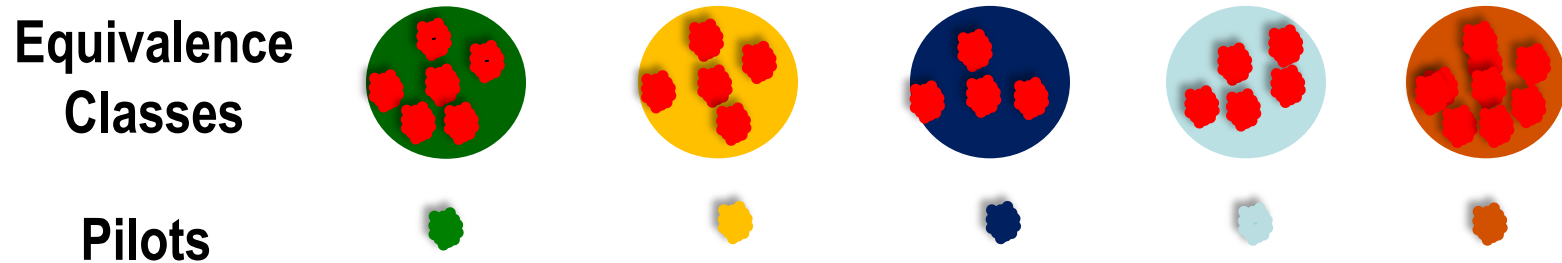
Relyzer

- Relyzer: A tool for complete application resiliency analysis
 - Employs systematic error pruning using static & dynamic analysis
 - Currently lists outcomes of masked, detection, SDC



- 3 to 6 orders of magnitude fewer error injections for most apps
 - 99.78% error sites pruned, **only 0.004% sites represent 99% of all sites**
 - GangES speeds up error simulations even further [ISCA'14]
- **Can identify virtually all SDC causing error sites**
- **What about unacceptable vs. acceptable SDC? Ongoing**

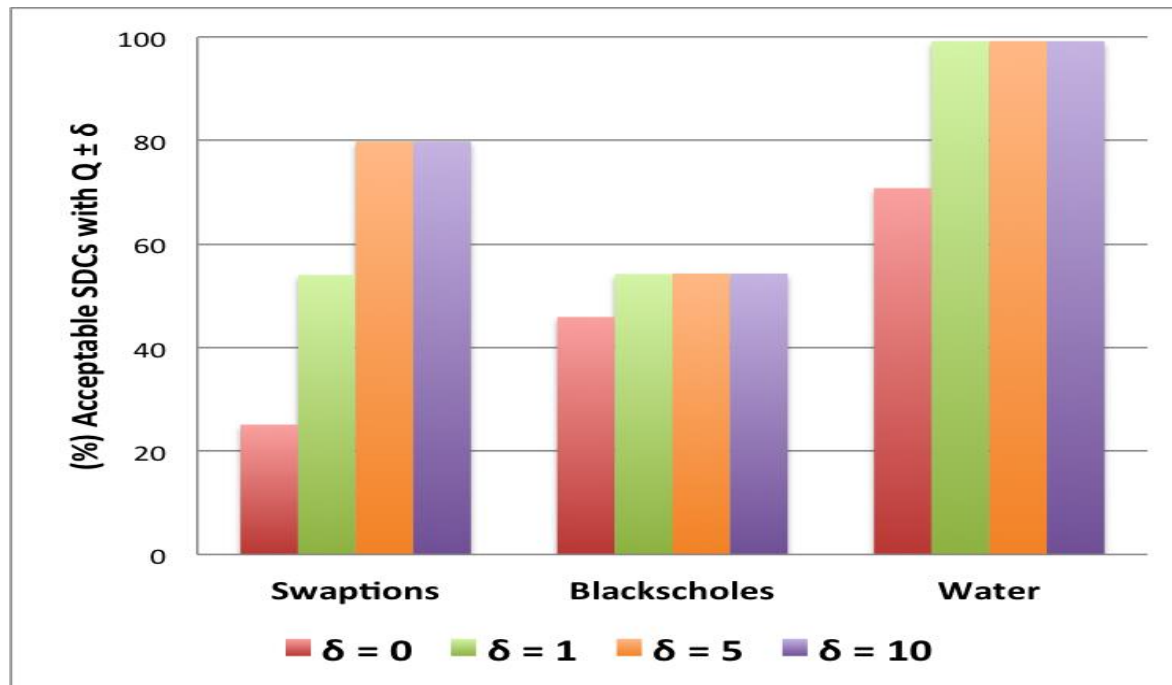
Can Relyzer Predict if SDC is Acceptable?



- Pilot = Acceptable SDC \Rightarrow
All faults in class = Acceptable SDCs ???
- Pilot = Acceptable SDC with quality $Q \Rightarrow$
All faults in class = Acceptable SDCs with quality Q ???

PRELIMINARY Results for Utility Validation

- Pilot = Acceptable SDC with quality $Q \Rightarrow$
All faults in class = Acceptable SDCs with quality $Q \pm \delta$???
- Studied several quality metrics
 - E.g., $E = \text{abs}(\text{percentage average relative error in output components})$,
capped to 100%
 $Q = 100 - \lceil E \rceil$ ($> 100\%$ error gives $Q=0$, 1% gives $Q=99$)



Research Strategy

- **Towards an Application Resiliency Profile**
 - **For a given instruction what is the outcome of a fault?**
 - **For now, focus on a transient fault in single bit of register**
- **Convert SDCs into (low-cost) detections for full reliability and quality**
- **OR let some acceptable or unacceptable SDCs escape**
 - **Quantitative tuning of reliability vs. overhead vs. quality**

SDCs → Detections [DSN'12]

What to protect?	SDC-causing fault sites	
How to Protect?	Low-cost detectors	
	Where to place?	Many errors propagate to few program values
	What detectors?	Program-level properties tests
	Uncovered fault-sites?	Selective instruction-level duplication

Insights for Program-level Detectors

- Goal: Identify *where* to place the detectors and *what* detectors to use
- Placement of detectors (*where*)
 - Many errors propagate to few program values
 - End of loops and function calls
- Detectors (*what*)
 - Test program-level properties
 - E.g., comparing similar computations and checking value equality

Loop Incrementalization

C Code

```
Array a, b;  
For (i=0 to n) {  
    ...  
    a[i] = b[i] + a[i]  
    ...  
}
```

ASM Code

```
A = base addr. of a  
B = base addr. of b
```

```
L: load  r1 ← [A]  
    ...  
    load  r2 ← [B]  
    ...  
    store r3 → [A]  
    ...
```

```
add  A = A + 0x8
```

```
add  B = B + 0x8
```

```
add  i = i + 1
```

```
branch (i < n) L
```

Loop Incrementalization

C Code

```
Array a, b;  
For (i=0 to n) {  
    ...  
    a[i] = b[i] + a[i]  
    ...  
}
```

SDC-hot app sites

Where: **Errors from all iterations** propagate here in **few quantities**

ASM Code

A = base addr. of a
B = base addr. of b

L: load r1 ← [A]
...
load r2 ← [B]
...
store r3 → [A]
...

add A = A + 0x8

add B = B + 0x8

add i = i + 1

branch (i < n) L

Collect initial values
of A, B, and i

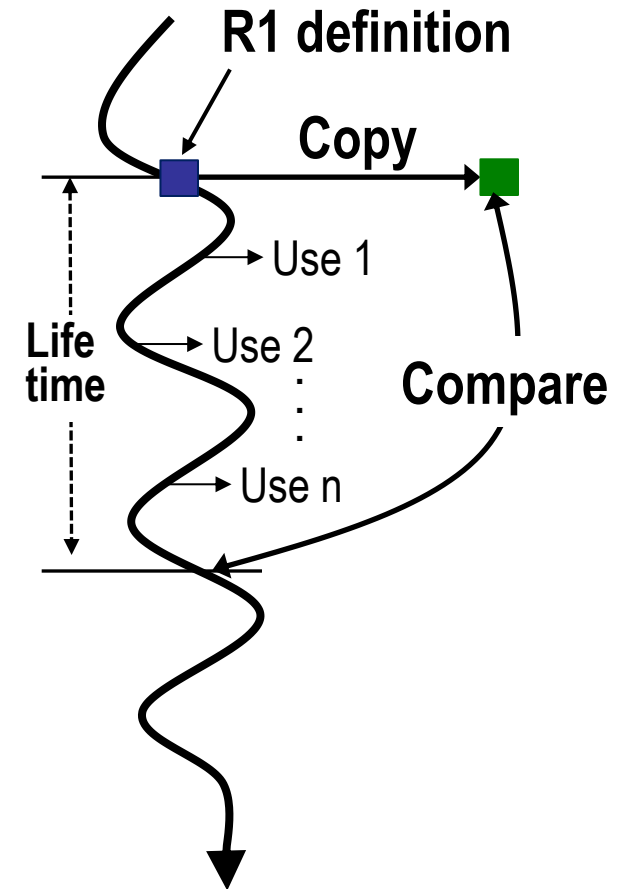
What: **Property checks**
on A, B, and i

Diff in A = Diff in B
Diff in A = 8 × Diff in i

No loss in coverage - *lossless*

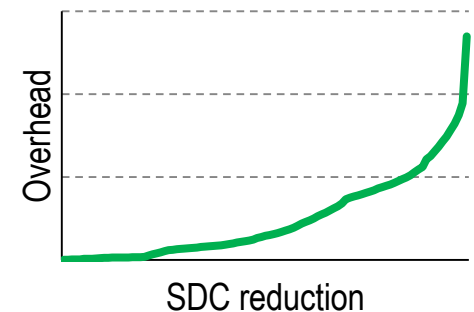
Registers with Long Life

- Some long lived registers are prone to SDCs
- For detection
 - Duplicate the register value at its definition
 - Compare its value at the end of its life
- No loss in coverage - *lossless*



Converting SDCs to Detections: Results

- Discovered common **program properties** around most SDC-causing sites
- Deviced **low-cost program-level detectors**
 - Average SDC reduction of 84%
 - Average cost of 10%
- New detectors + selective duplication = **Tunable resiliency at low-cost**
 - Found near optimal detectors for any SDC target



Identifying Near Optimal Detectors: Naïve Approach

Example: Target SDC coverage = 60%

Overhead = 10%

Sample 1

SFI

SDC coverage

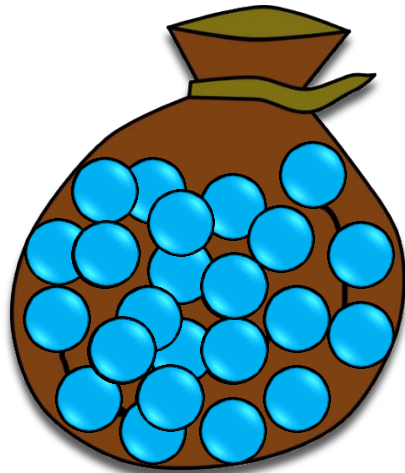
50%

Overhead = 20%

Sample 2

SFI

65%

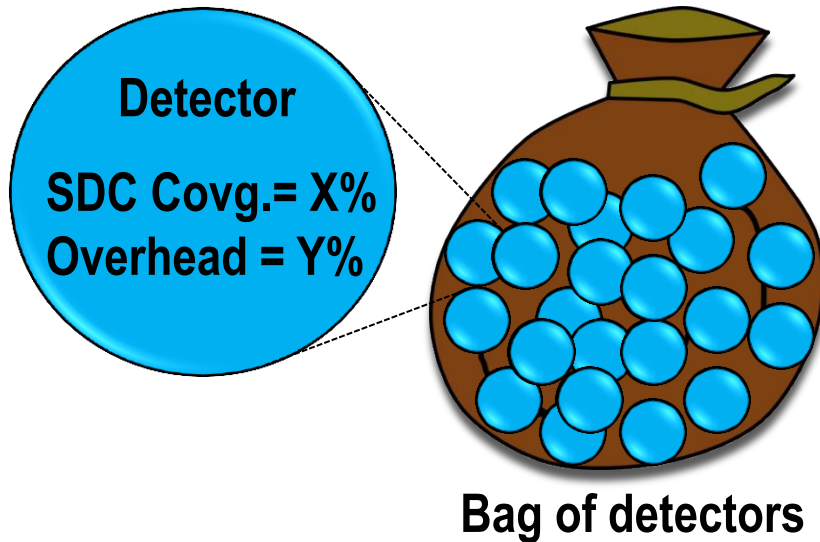


Bag of detectors

Tedious and time consuming

Identifying Near Optimal Detectors: Our Approach

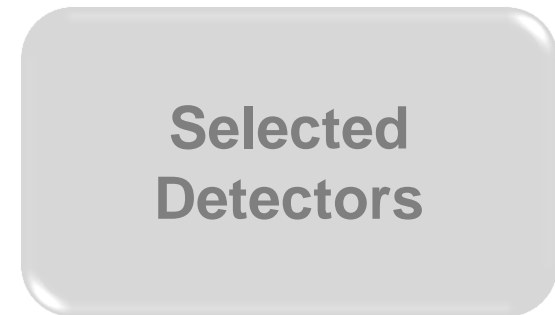
1. Set attributes, enabled by Relyzer



2. Dynamic programming

Constraint: Total SDC covg. $\geq 60\%$

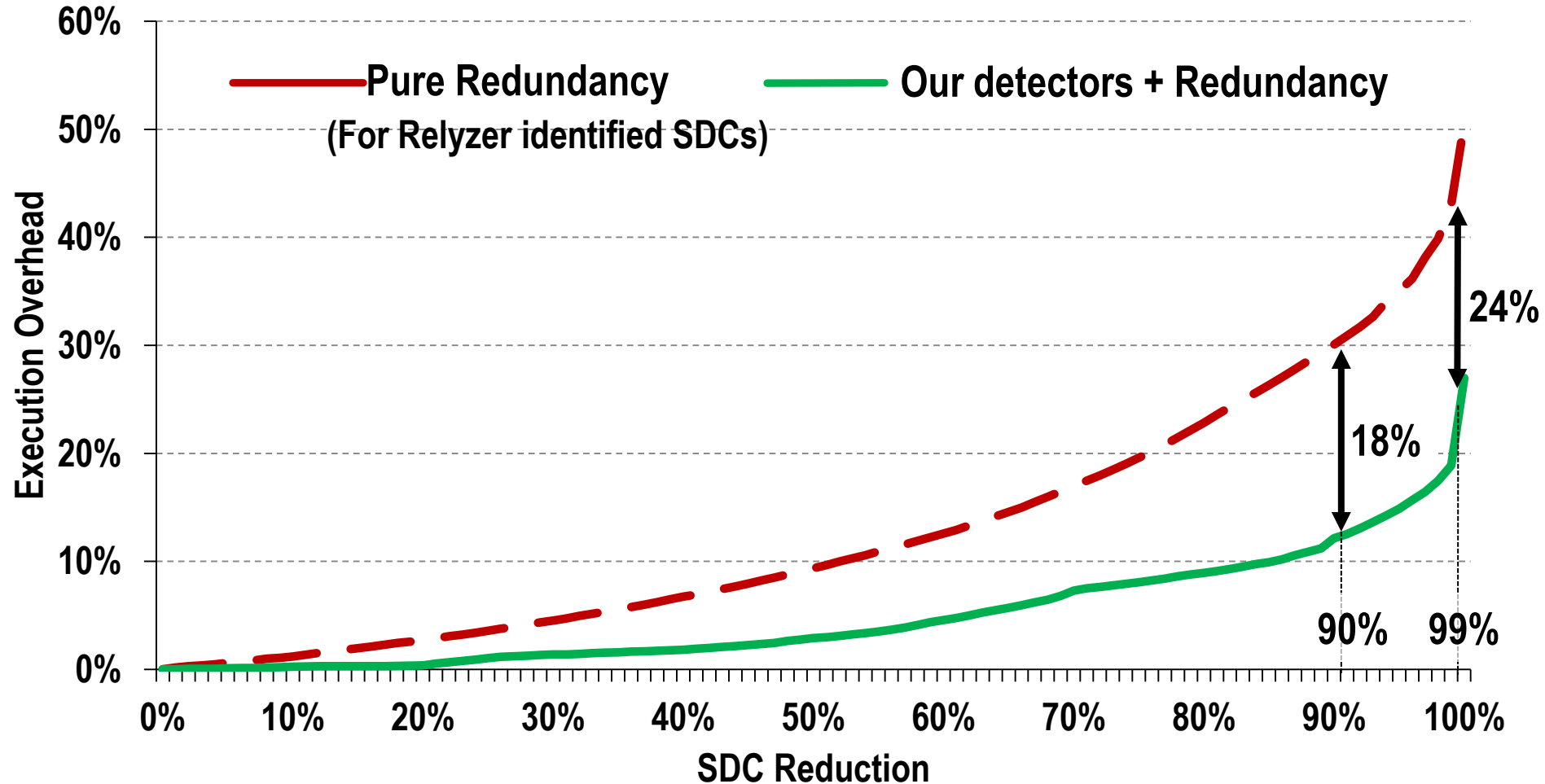
Objective: Minimize overhead



Overhead = 9%

Obtained SDC coverage vs. Performance trade-off curves [DSN'12]

SDC Reduction vs. Overhead Trade-off Curve



Consistently better over pure instruction-level duplication (w/ Relyzer)

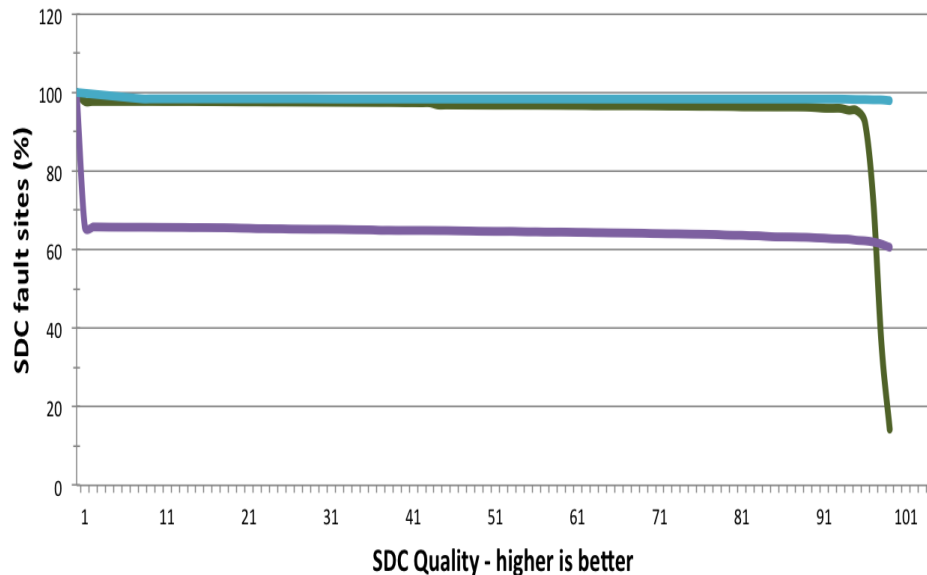
But overhead still significant for very high resilience

Remove protection overhead for acceptable (bounded) quality loss?

Understanding Quality Outcomes with Relyzer (Prelim Results)

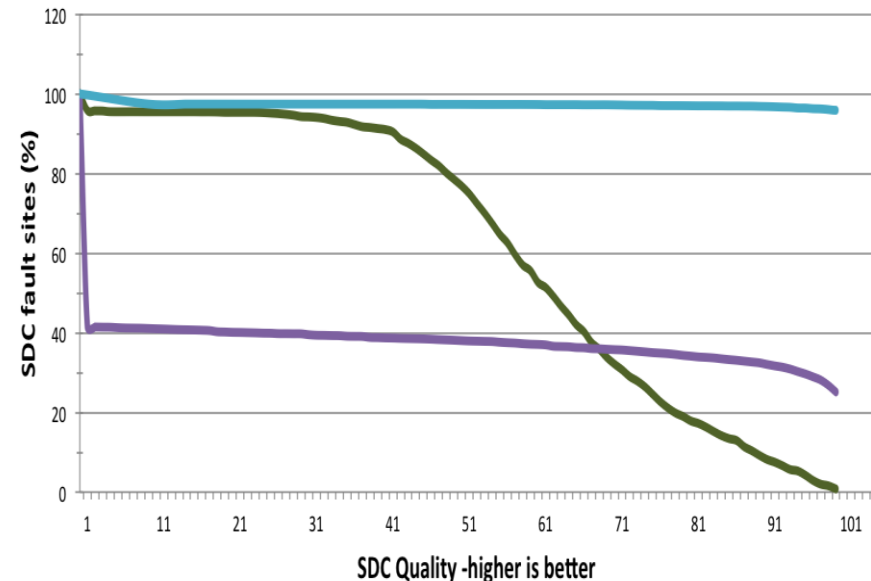
SDCs categorized by Quality
Average Relative Error (%)

Swaptions Blackscholes Water



SDCs categorized by Quality
Max Relative Error (%)

Swaptions Blackscholes Water

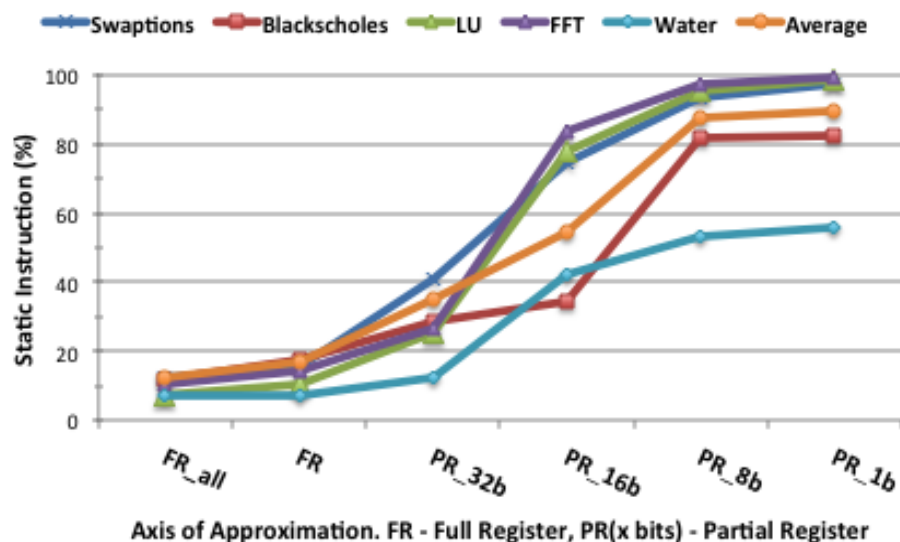


- Promising potential, but quality metric and application dependent
- Can use to quantitatively tune quality vs. reliability vs. resources

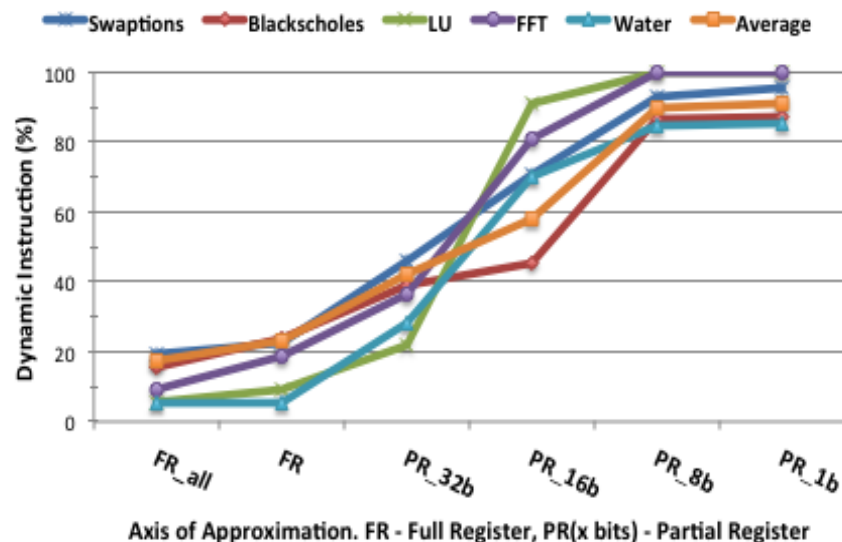
Quality Outcomes – An Instruction-Centric View

- Systems operate at higher granularity than error sites: **Instruction?**
- When is an instruction approximable?
 - Which errors in instruction should result in acceptable outcome?
 - * All? Single-bit errors in all register bits? **Single-bit errors in subset of bits?**
 - When is an outcome acceptable?
 - * **Best case** : all errors that are masked or produce SDCs are acceptable

PC Best case Approx - Static Instructions



PC Best case Approx - Dynamic Instructions



So - (Old) Hype or New Frontier?

- **SWAT inherently implies quality/reliability relaxation aka approx. computing**
- **Key: Bound quality/reliability loss? (Subject to resource constraint)**
 - Can serve as enabler for widespread adoption of resilience approach
 - Especially if automated
- **Some other open questions**
 - Instruction vs. data-centric?
 - Resiliency profiles and analysis at higher granularity?
 - How to compose?
 - Impact on recovery?
 - Other fault models?

Software doesn't ship with 100% test coverage, why should hardware?

Summary

