# Overcoming the Data-flow limit with **Structural Approximation**

Vignesh Balaji

Brandon Lucia

Radu Marculescu

Carnegie Mellon University

Electrical & Computer
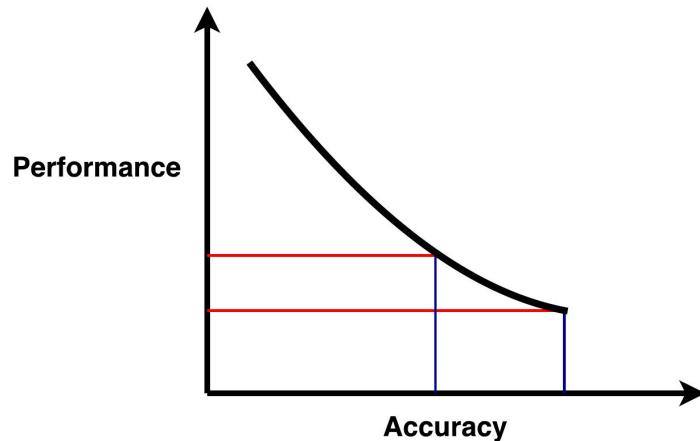ENGINEERING

**Carnegie Mellon University**

# Outline

❖ Trends in Approximate Computing

❖ Structural Approximation

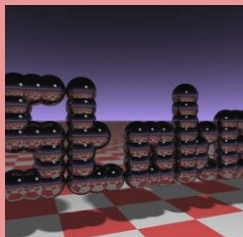❖ Making Structural Approximation viable
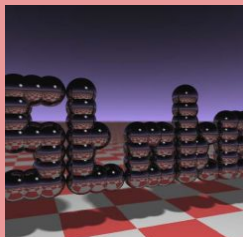
# Approximate Computing

Correctness is *costly*

AND

Precision is not always *required*

Performance

Accuracy

# Successful Domains for Approximation



*Precise Output*     *Approx. Output*     *Pixel error values*

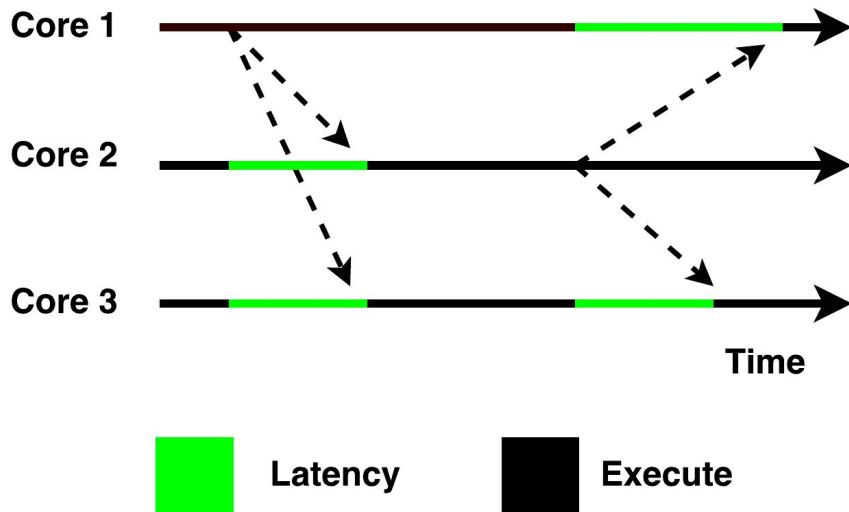Performance improves by 6.4% *



Multiple *correct* answers possible.

## Most existing approximations typically target only *numeric values*

**\*** Eric Schkufza, Rahul Sharma, and Alex Aiken. 2014. Stochastic optimization of floating-point programs with tunable precision. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation* (PLDI '14)
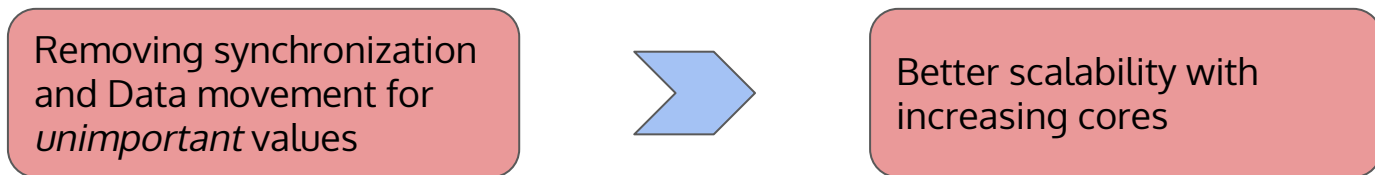
# Another Interesting Domain - Parallel Computing

- Parallel programs communicate to enforce *correctness*:
  - Synchronization primitives (locks, barriers, mutexes, etc.)
  - Cache Coherence
- Communication limits *scalability*



**Should we still communicate for values that are not *critical* for the app's correctness?**

# Removing Communication via Approximation

Removing synchronization and Data movement for *unimportant* values  ⟩  Better scalability with increasing cores
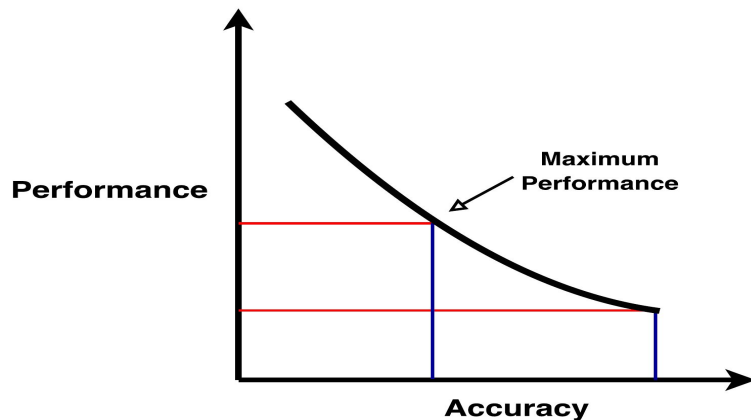
Prior work has successfully removed synchronization for such *values*

*Rely, Relaxing Synchronization for Performance and Insight, Dancing With Uncertainty, Unsynchronized Techniques for Approximate Parallel Computing, ...*

**Are such *value* approximations enough? Or is there a need for something more?**

# Quantifying the Performance Limit



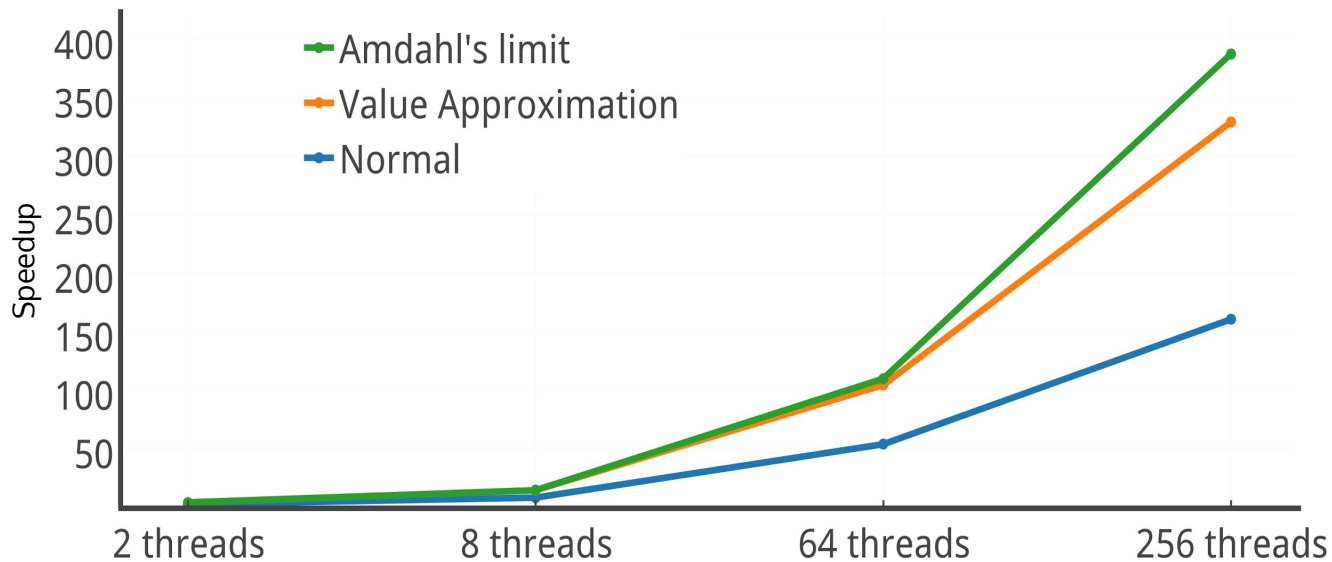Definition of correctness for the limit study -

The application should not *crash*

Simulated the maximum performance possible with Value Approximation :

➢ Removed *all* locks guarding values

➢ Removed all *effects* of cache coherence

➢ Removed the *cost* of performing reduction operations
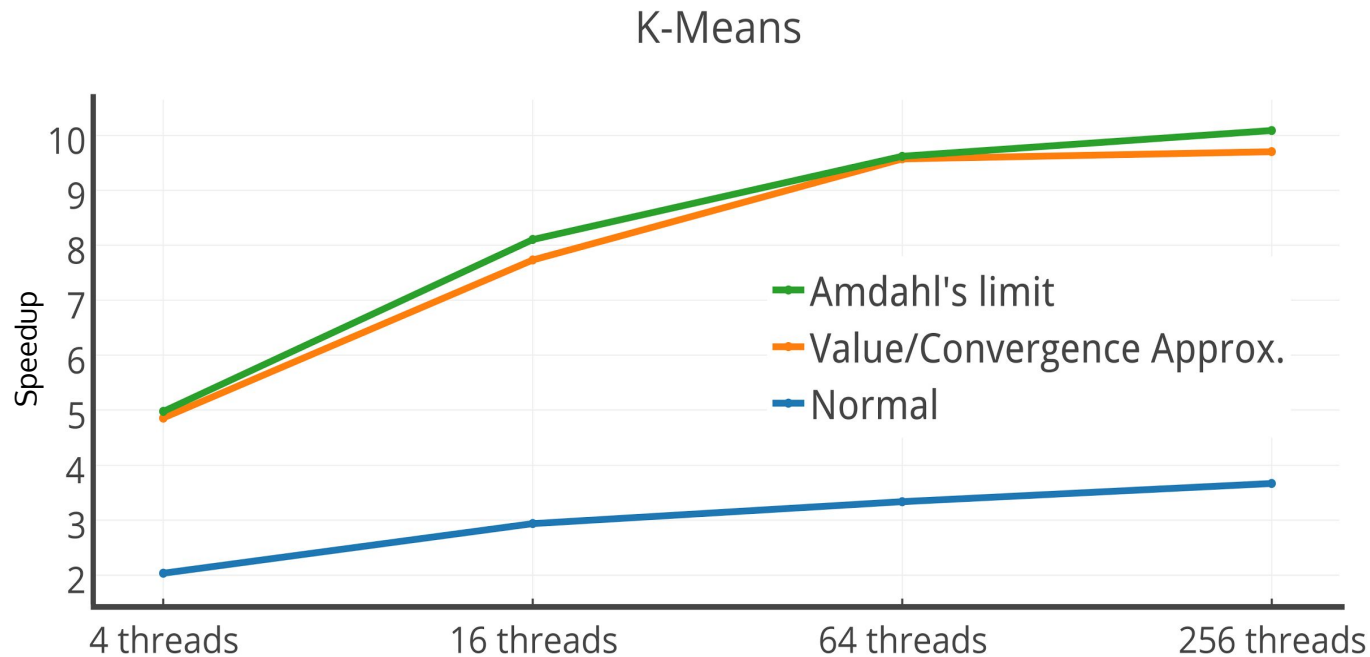
# Approximating Key-Value Store

## Key-Value Store



Speedup vs. threads chart with three lines:
- Amdahl's limit (green)
- Value Approximation (orange)
- Normal (blue)

X-axis: 2 threads, 8 threads, 64 threads, 256 threads
Y-axis: Speedup (50, 100, 150, 200, 250, 300, 350, 400)

```
Lock(table[x].mutex);
table[x].data++;
Unlock(table[x].mutex);
```

# Approximating K-Means

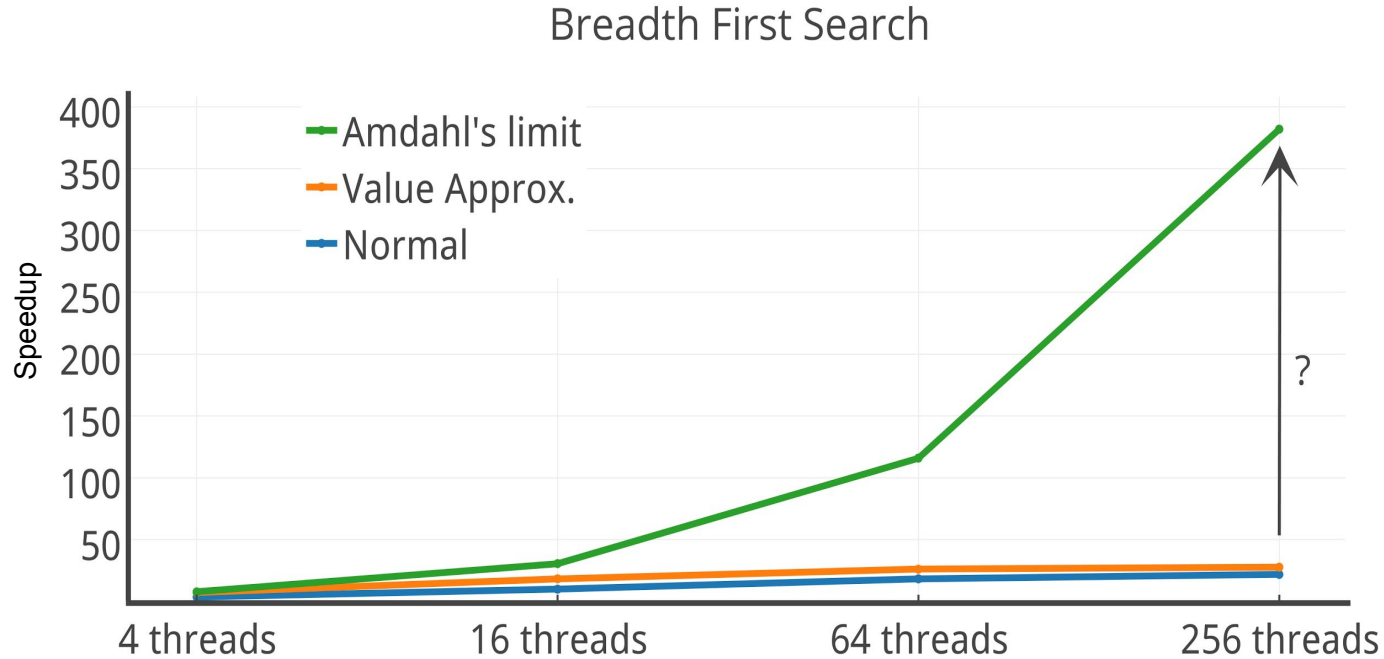K-Means



Step 1: Assign points
to the closest cluster

Step 2: Compute means
of the new clusters

# Value Approximation Works Well

- Approximating data *values* took performance close to Amdahl's limit

- Previous two applications only shared *numerical values*

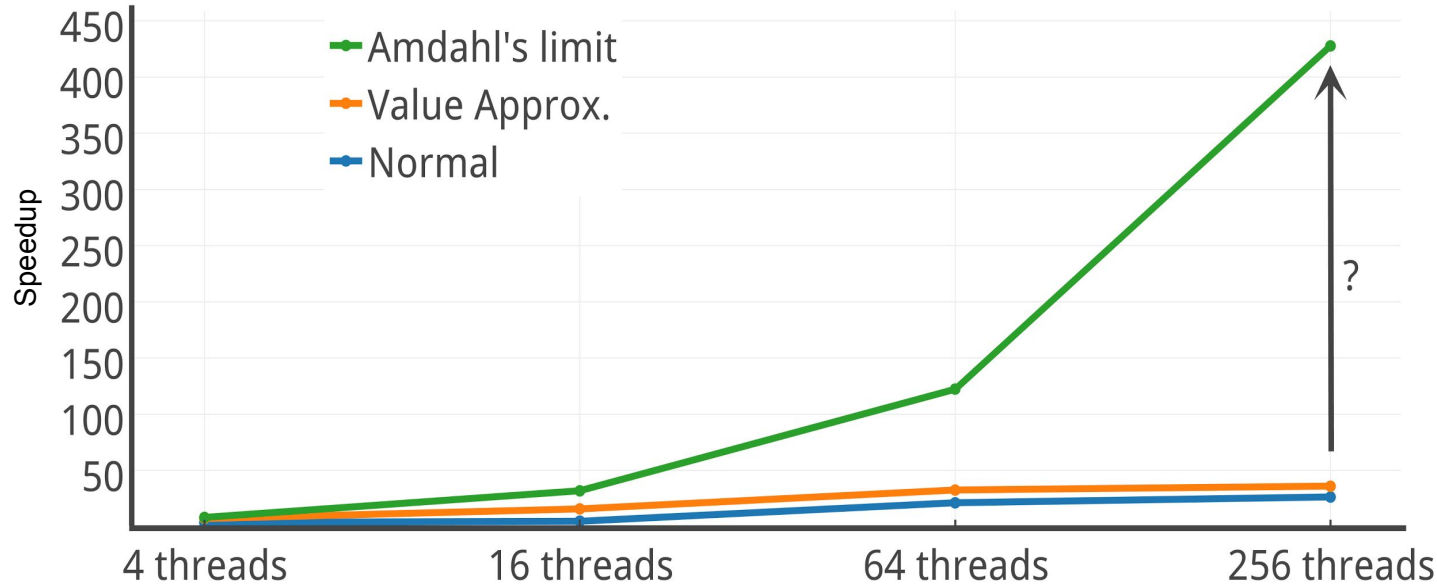**What if we have applications that share *more* than just numerical values?**

# Breadth First Search



Breadth First Search

Legend: Amdahl's limit, Value Approx., Normal

Y-axis: Speedup (50, 100, 150, 200, 250, 300, 350, 400)

X-axis: 4 threads, 16 threads, 64 threads, 256 threads

*Data structure* that holds the children nodes of a frontier is shared

# Single Source Shortest Path

SSSP



Speedup

450
400
350
300
250
200
150
100
50

Amdahl's limit
Value Approx.
Normal

4 threads     16 threads     64 threads     256 threads

?

*Data structure* maintaining nodes closest to the source is shared

# End of the Road for Value Approximation

**What's different in these applications**?

These application *share data structures* in addition to values

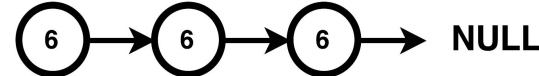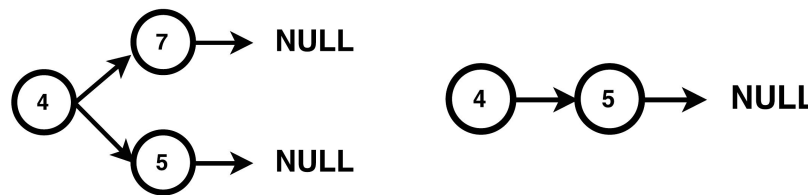**What if we remove synchronization for operations guarding data structures**?

Removing synchronization for data structures leaves them *inconsistent*

**How do we reason about output quality now**?
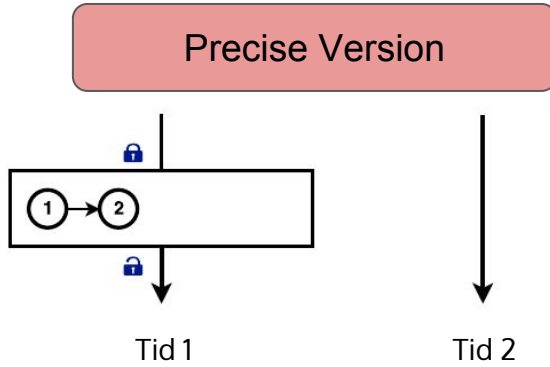
# Outline

❖    Trends in Approximate Computing

❖    Structural Approximation

❖    Making Structural Approximation viable
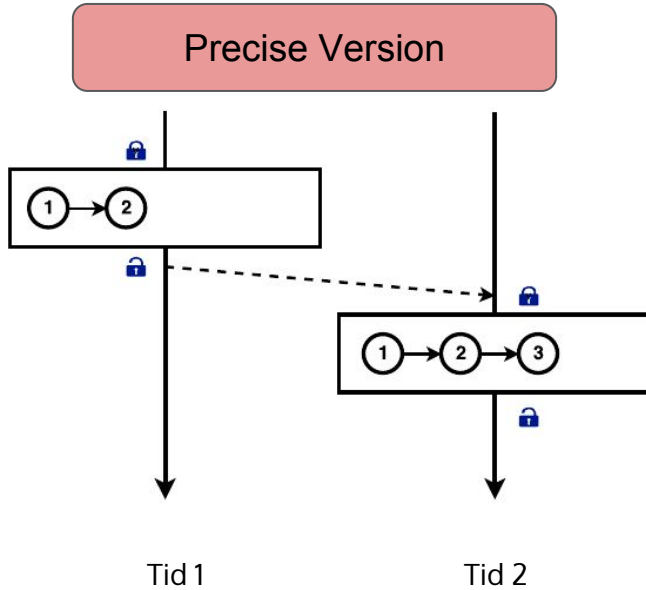
# Overview of Structural Approximation

| | |
|---|---|
| Precise Output | 4 → 7 → 5 → **NULL** |
| Output after Value Approximation | 6 → 6 → 6 → **NULL** |
| Output after Structural Approximation | 4 → 7 → **NULL**, 4 → 5 → **NULL**, 4 → 5 → **NULL** |

# Mechanism of Structural Approximation
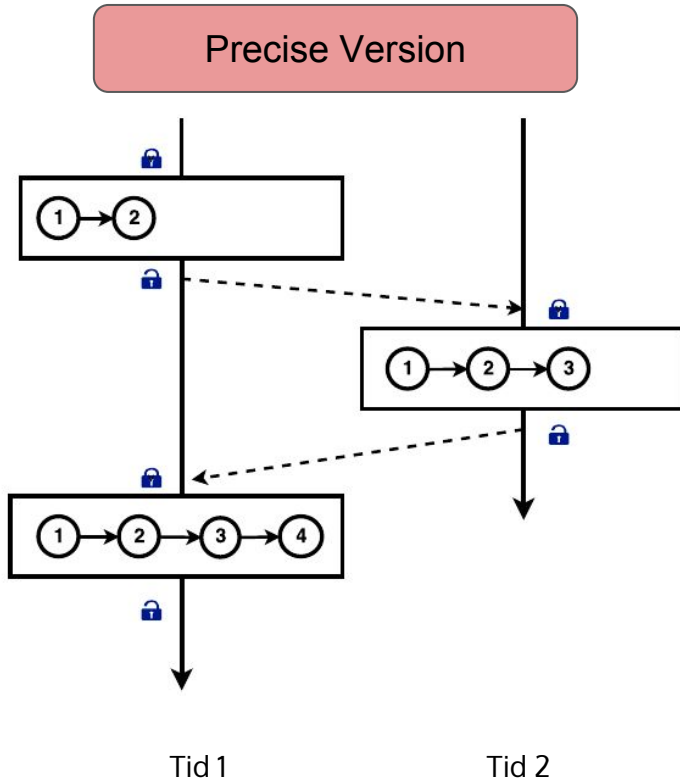


Precise Version

Tid 1          Tid 2
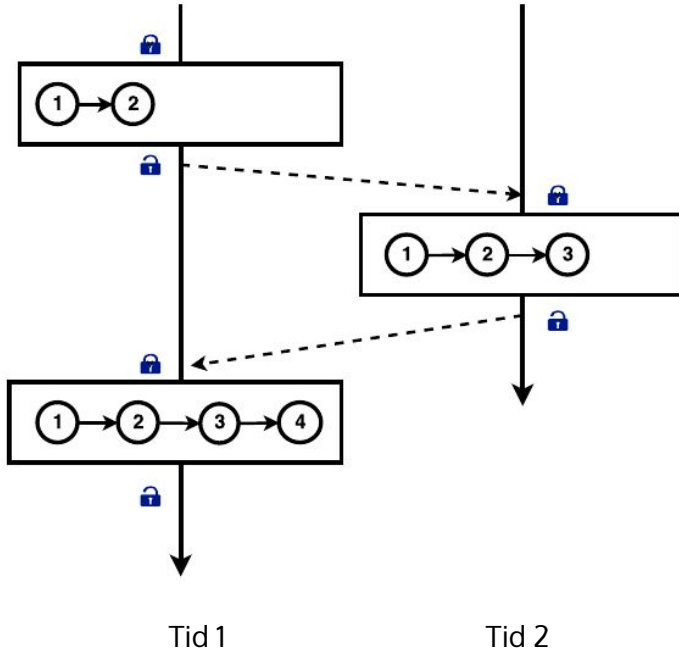
# Mechanism of Structural Approximation

# Mechanism of Structural Approximation



Tid 1　　　　　Tid 2

# Mechanism of Structural Approximation
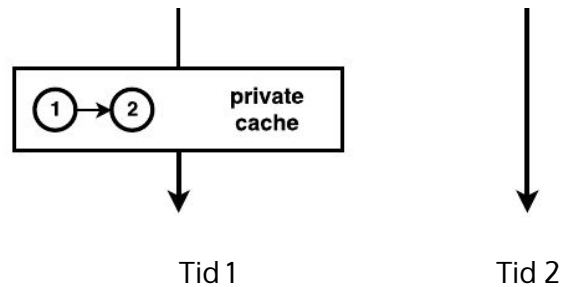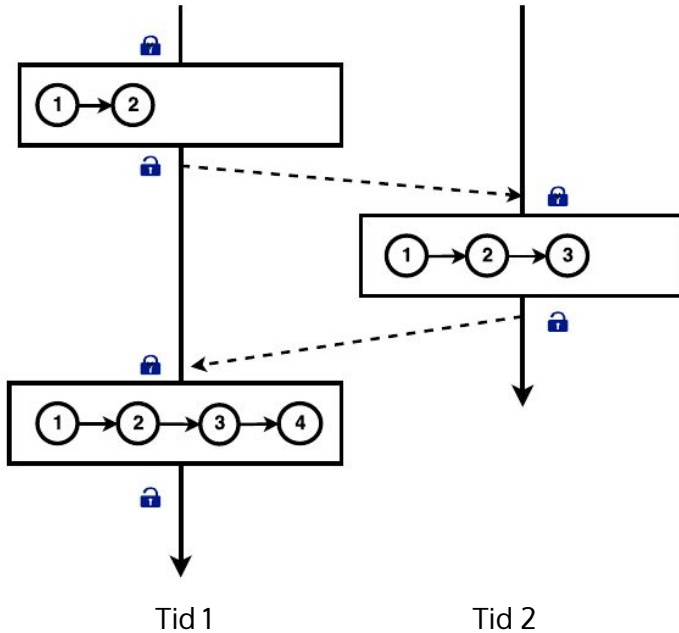


Precise Version

Approximate Version

private cache

Tid 1          Tid 2

Tid 1          Tid 2

# Mechanism of Structural Approximation



Precise Version

Approximate Version

Tid 1        Tid 2

Tid 1        Tid 2

# Mechanism of Structural Approximation



Precise Version

Approximate Version

Tid 1     Tid 2

private cache

private cache

private cache

Tid 1     Tid 2

How can we ensure *usable* results from Structural Approximation?

# Outline

❖ Trends in Approximate Computing

❖ Structural Approximation

❖ Making Structural Approximation viable

# Existing Definition of Correctness

- Value based approximations use *numerical distance* to measure error

$$Metric = (Output_{precise} - Output_{approximate})$$

- Makes sense for value approximation



| 3 | 4 | 6 | 8 | — | 2 | 3 | 5 | 7 | = | 1 | 1 | 1 | 1 |

Precise Output — Approx. Output = Error

**We cannot use this error definition for Structural Approximation!**

# New Definition of Correctness

The new definition should:

- ➢ Accommodate *temporary* inconsistency in data structures
- ➢ *Tolerate* the loss of a few elements in data structures

**Need mechanisms to *enforce* this correctness definition**
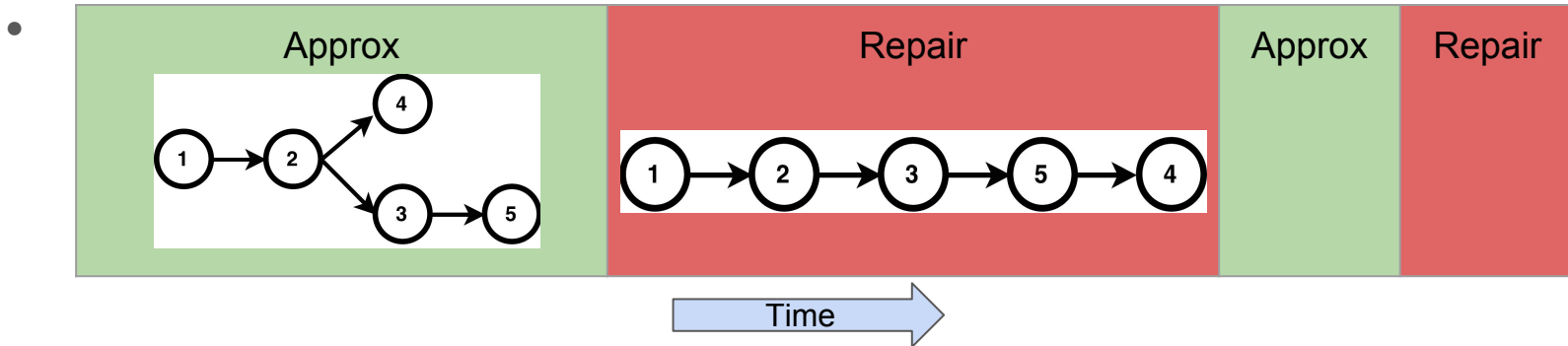
# Resilient Operators

- Operators that *successfully* service queries even on inconsistent data structures
- Similar to the concept of Defensive Programming

```
int risky_visitor(int size){
  while(i < size){
    do_work(node->value); //Fails if elements
                          //dropped from list
    node = node->next;
    i++;
}}
```

```
int secure_visitor(int size){
  while(i < size){
    if (node != NULL){
      do_work(node->value); //saved from
                            //NULL ptr exception
      node = node->next;
    i++;
    }
    else
      break;
}}
```
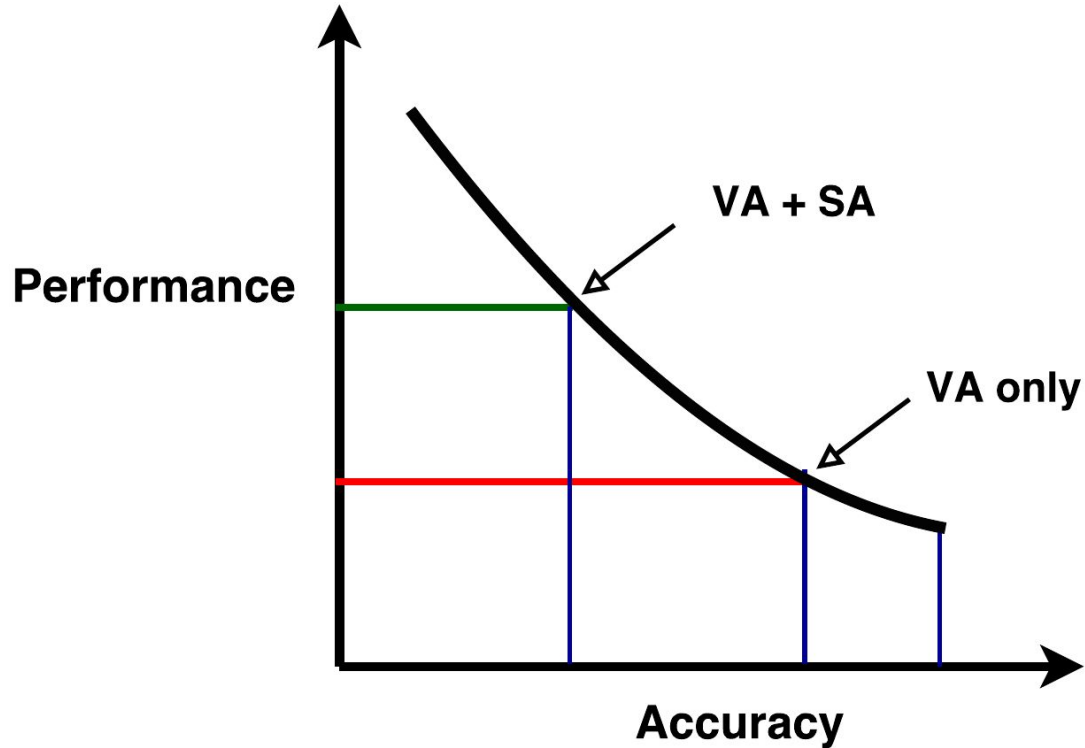
**Handle all failures that are possible due to Structural Approximation**

# Data Structure Repair



- 

- Major challenge - How to implement *low-cost* repairs?
  - ➢ Offload repair to cloud?

# Role of SA in the Future of Approximate Computing

# Thank You