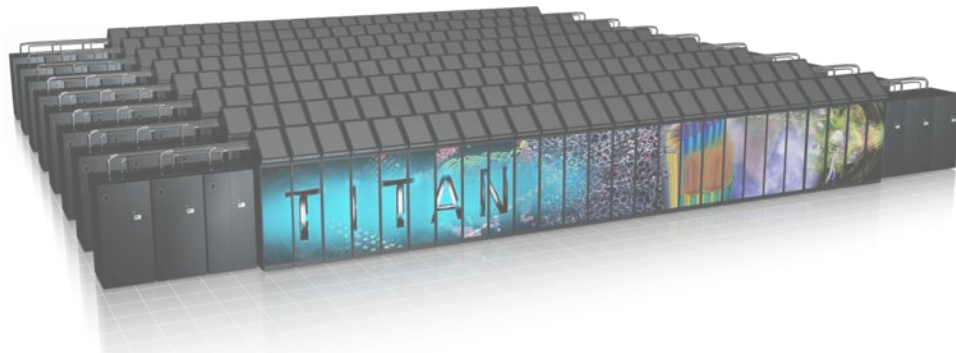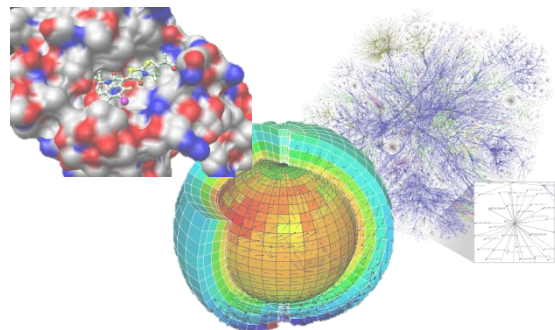# **Computing on the Lunatic Fringe:**
## Exascale Computers
## and Why You Should Care

### Mattan Erez

### The University of Texas at Austin

# Arch-focused whole-system approach

Efficiency requirements require crossing layers

Algorithms are key

– Compute less, move less, store less

Proportional systems

– Minimize waste

Utilize and improve emerging technologies

Explore (and act) up and down

– Programming model to circuits

Preferably implement at micro-arch/system

THE UNIVERSITY OF
**TEXAS**
— AT AUSTIN —

# Big problems and emerging platforms

Memory systems
- Capacity, bandwidth, efficiency – impossible to balance
- Adaptive and dynamic management helps
- New technologies to the rescue?
- Opportunities for in-memory computing
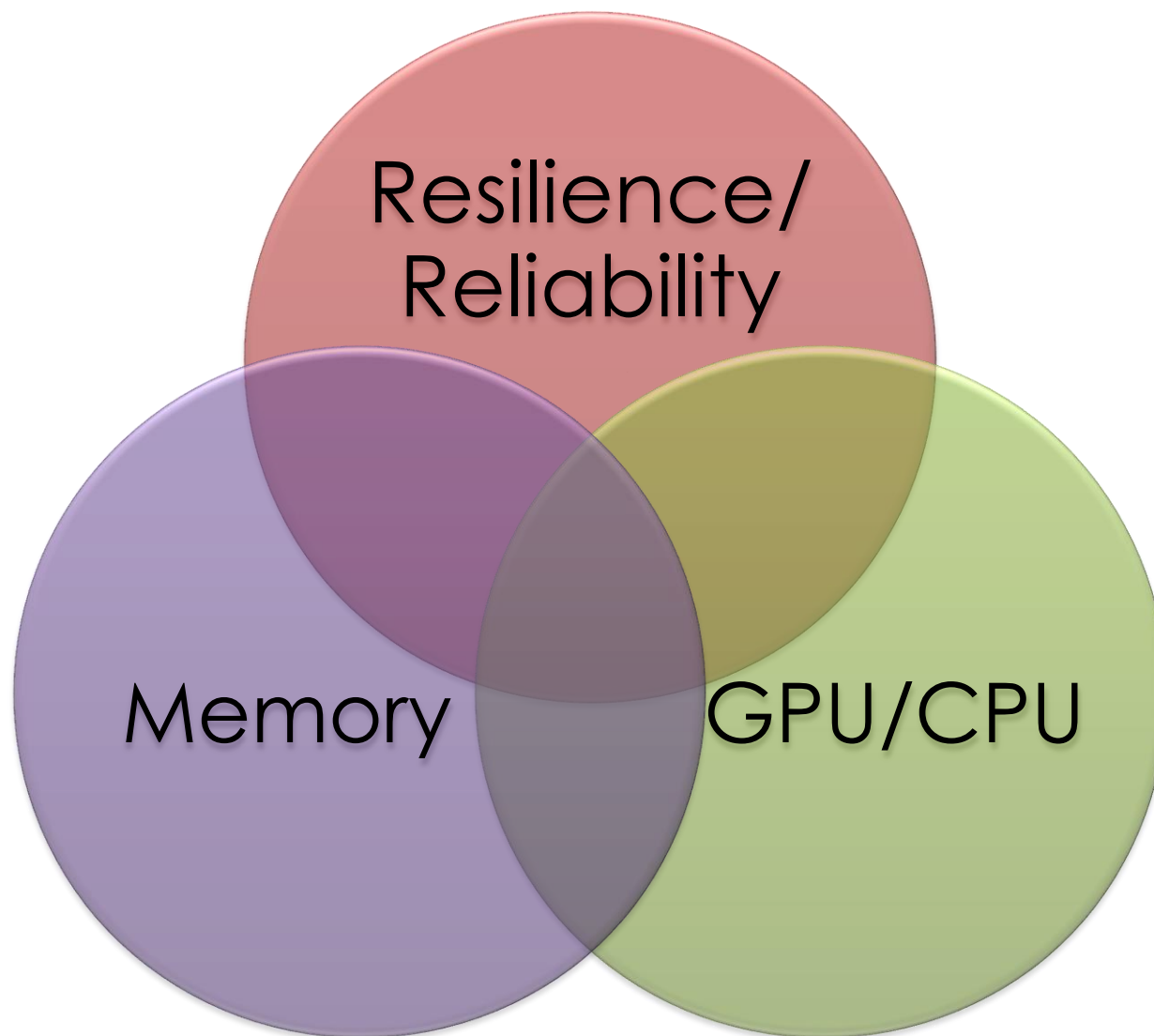
GPUs, supercomputers, clouds, and more
- Throughput oriented designs are a must
- Centralization trend is interesting

Reliability and resilience
- More, smaller devices – danger of poor reliability
- Trimmed margins – less room for error
- Hard constraints – efficiency is a must
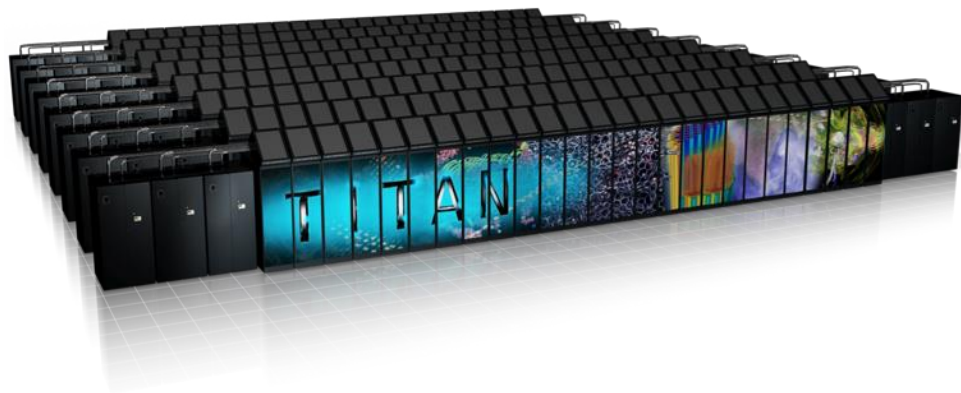- Scale exacerbates reliability and resilience concerns

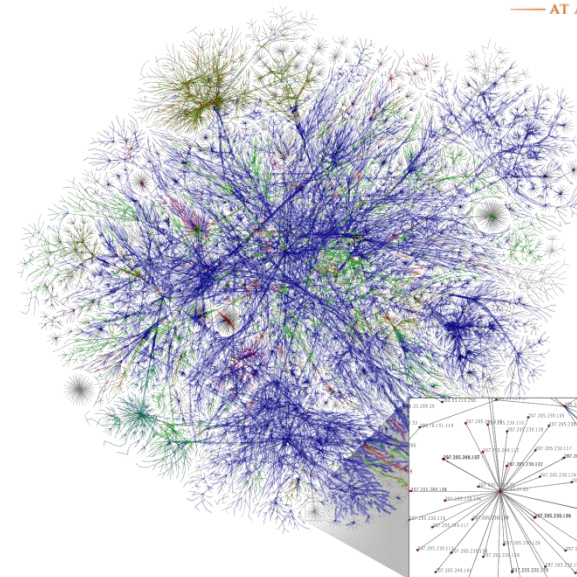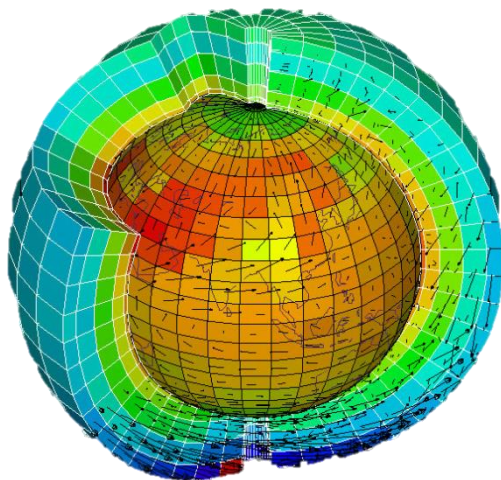# Lots of interesting multi-level projects
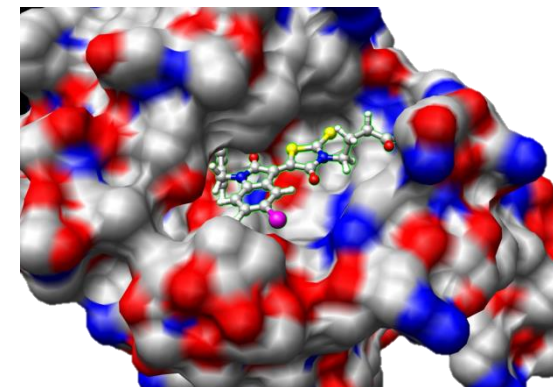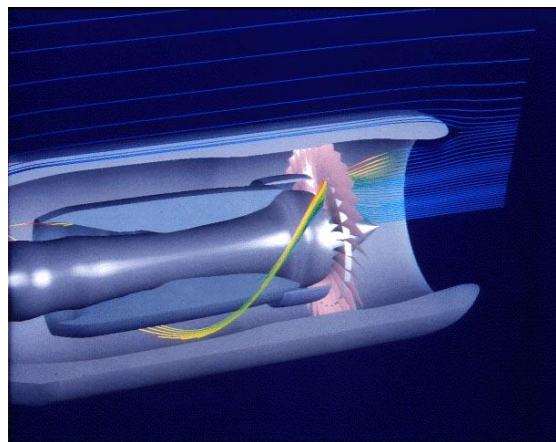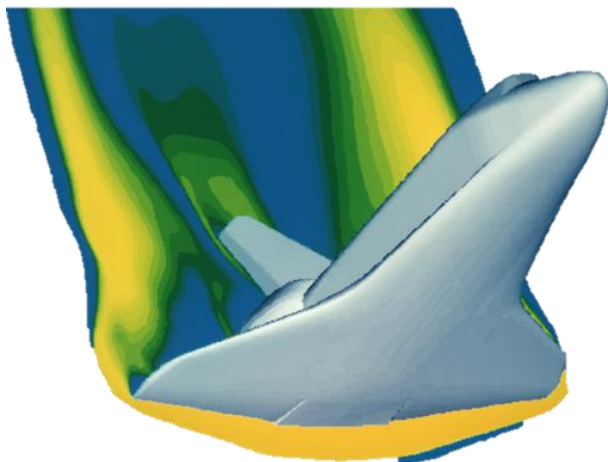
A **superscale computer** is
a high-performance system for
solving **big problems**

# A **supercomputer** is a high-performance system for solving **big cohesive problems**

THE UNIVERSITY OF
**TEXAS**
AT AUSTIN



A **supercomputer** is a high-performance system for solving **big cohesive problems**

# Simulate

# Analyze

# Predict

# Simulate

## Analyze

## Predict

# Supercomputers are general
– Not domain specific

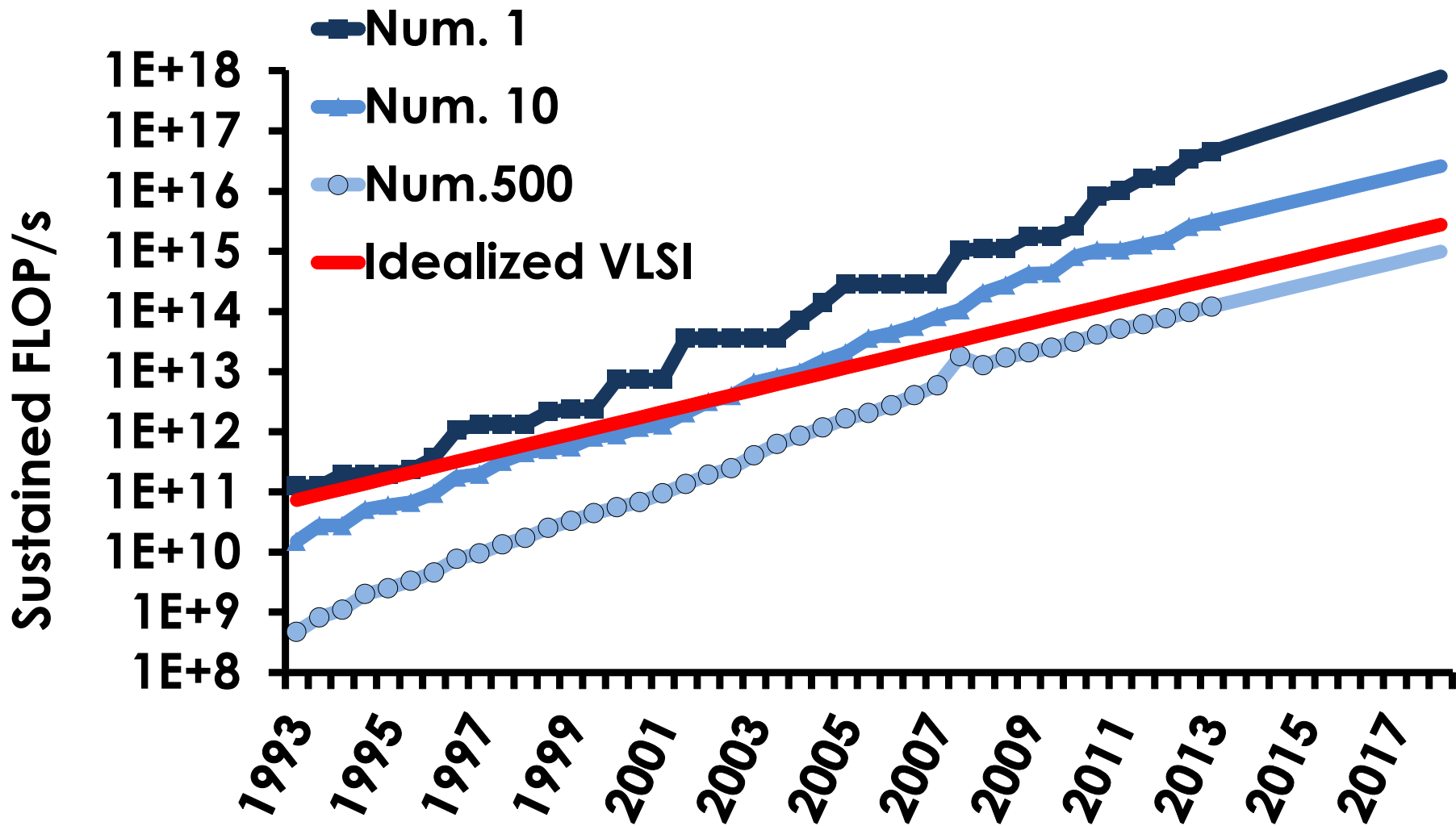# Cloud computers are general
– Algorithms keep changing

# Superscale computers must balance

– Compute

– Storage

– Communication

– Constraints

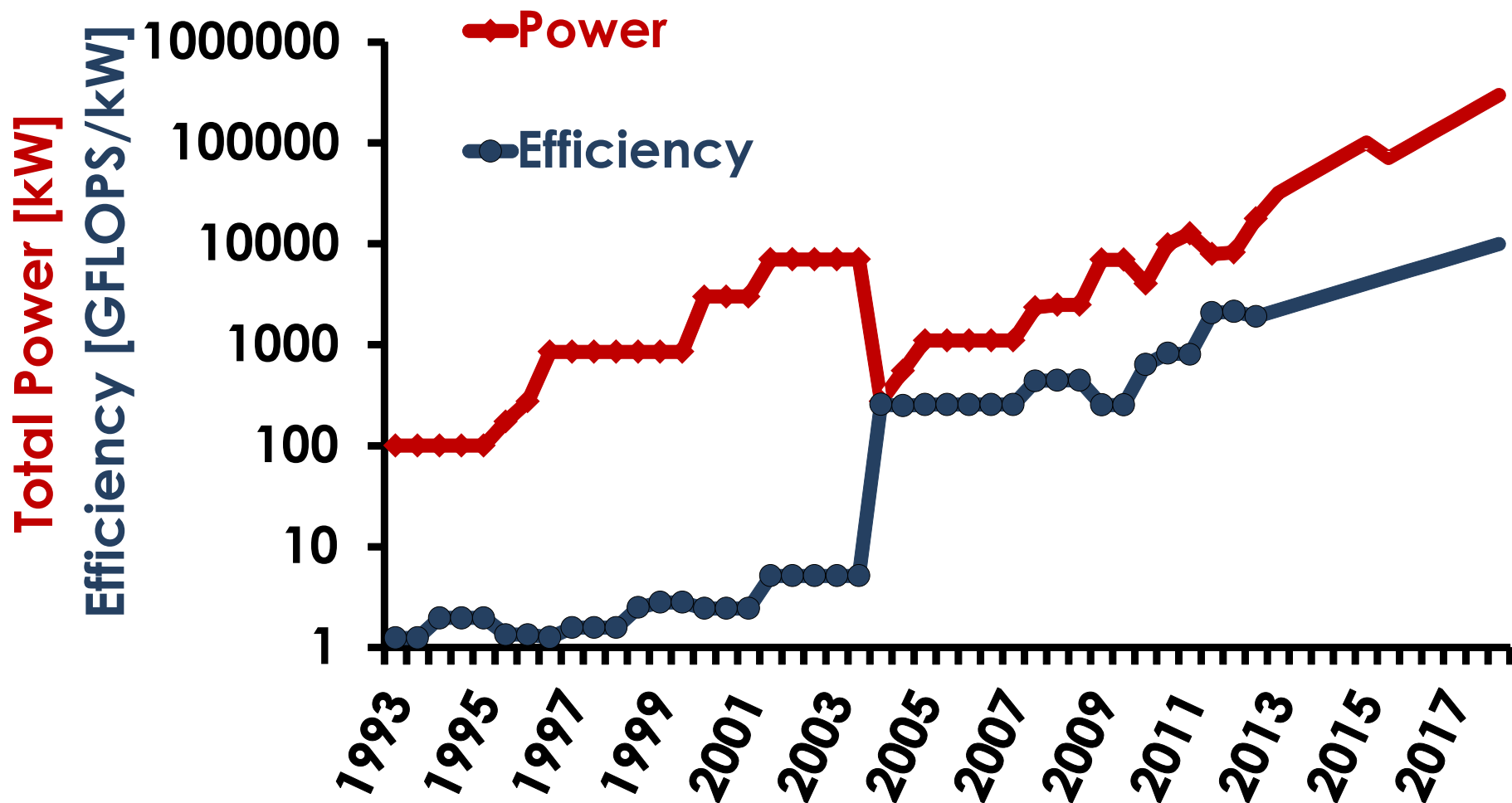  • OpEx and CapEx

# Super is never super enough

# What's keeping us from **exascale?**

# The **power** problem

# Solution: build more **efficient systems**

# Exascale computers are a lunacy:
– Crazy big (1M processors, >10MW)
– Need efficiency of embedded devices
– Effective for many domains
– Fully programmable

# Why should you care?

Harbingers of things to come
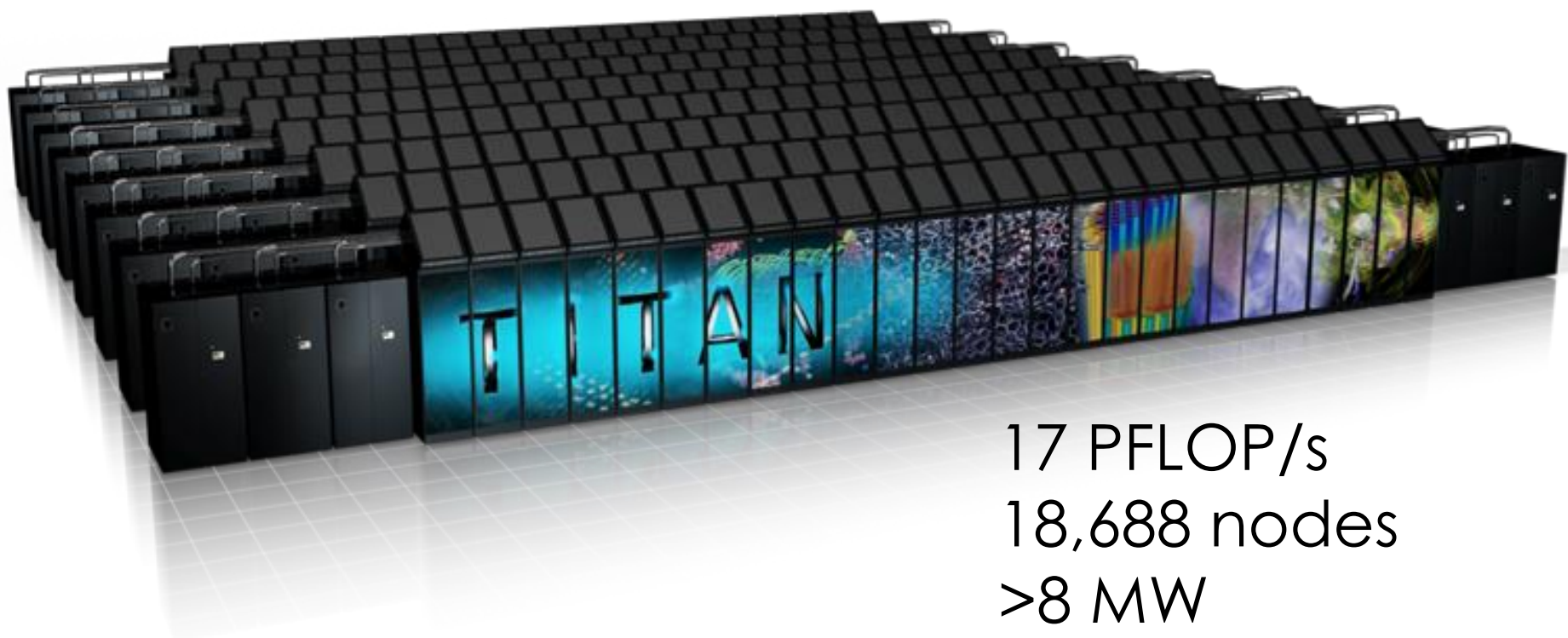Discovery awaits
**Enable** and promote **open innovation**

THE UNIVERSITY OF
**TEXAS**
— AT AUSTIN —

# Where does the power go?

# Cooling and infrastructure
– Amazing mechanical and electrical systems
– Getting close to optimal efficiency

17 PFLOP/s
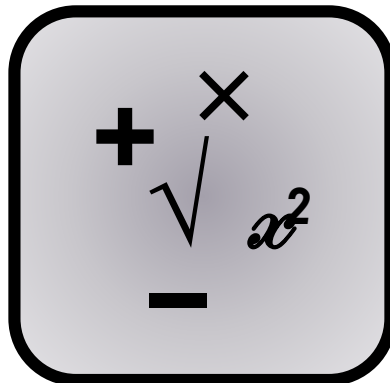18,688 nodes
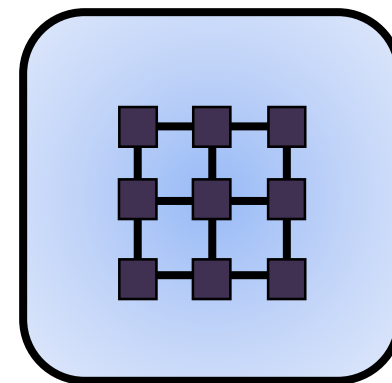>8 MW
~200 cabinets
~400 m² floorspace

# Actual processing
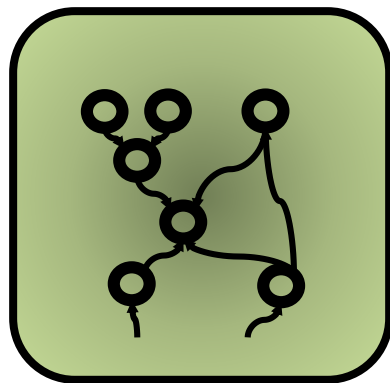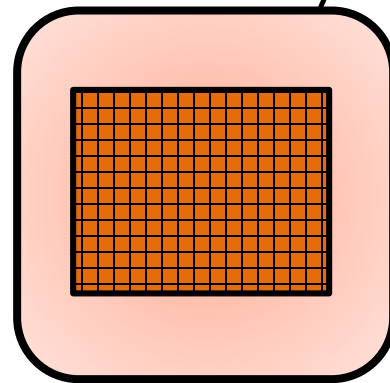
**I/O**

...100101001
0100...

**Arithmetic**

$+$ $\times$
$\sqrt{}$ $x^2$
$-$

Comm.

Control

Memory

Idle/margin

# How much of each component?

# The budget: **Power ≤ 20MW**

$$\text{Energy} \leq \frac{\text{20MW / 1 exa-FLOP/s}}{\textbf{20pJ/op}}$$
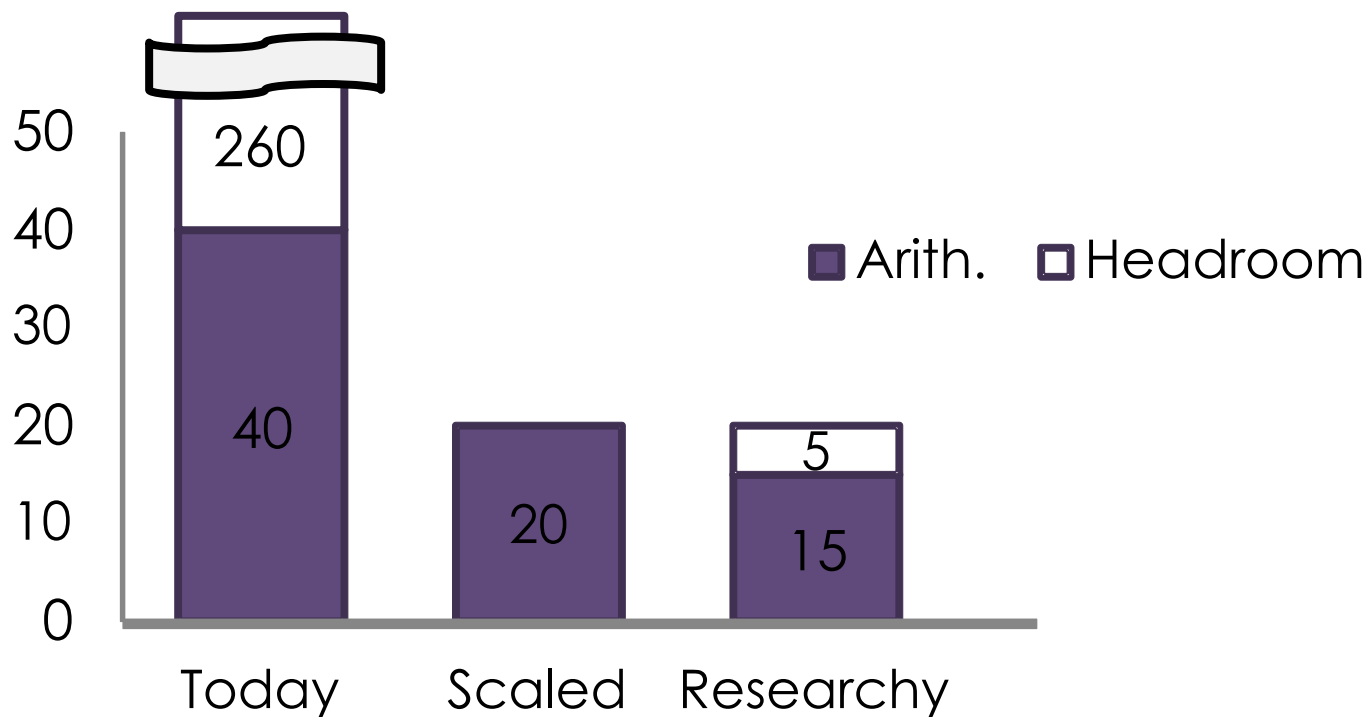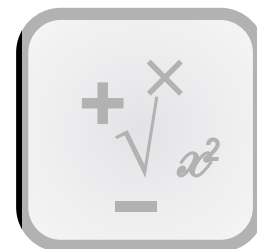
50 GFLOPs/W sustained

**Energy ≤**   20MW / 1 exa-FLOP/s

**20pJ/op**

50 GFLOPs/W sustained

Best supercomputer today: **~300pJ/op**
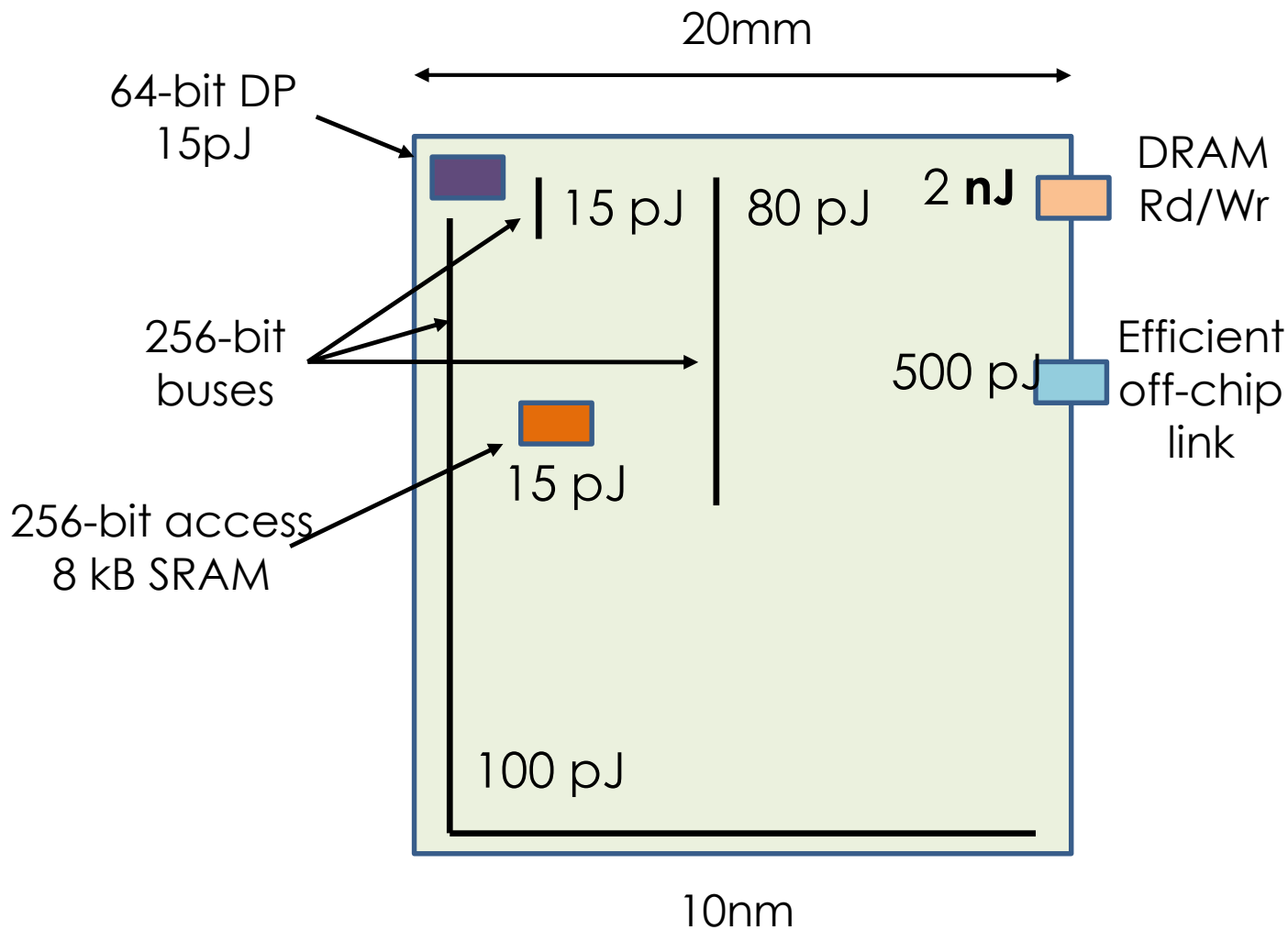
# Arithmetic

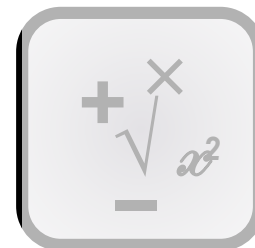## 64-bit floating-point operation
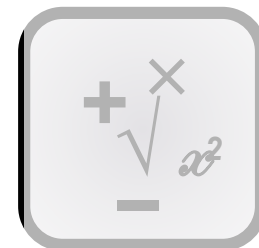


Rough estimated numbers

# Enough headroom?

# Unfortunately, hard tradeoffs

20mm

64-bit DP
15pJ

DRAM
Rd/Wr

15 pJ | 80 pJ | 2 **nJ**

256-bit
buses

Efficient
off-chip
link

500 pJ

15 pJ

256-bit access
8 kB SRAM

100 pJ

10nm

# Need **more headroom**
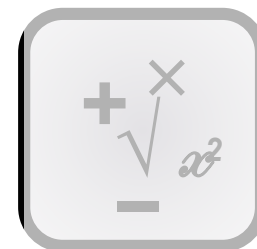 – *Minimize waste*

Do we care about single-unit performance?

Must all results be equally precise?

Must all results be correct?
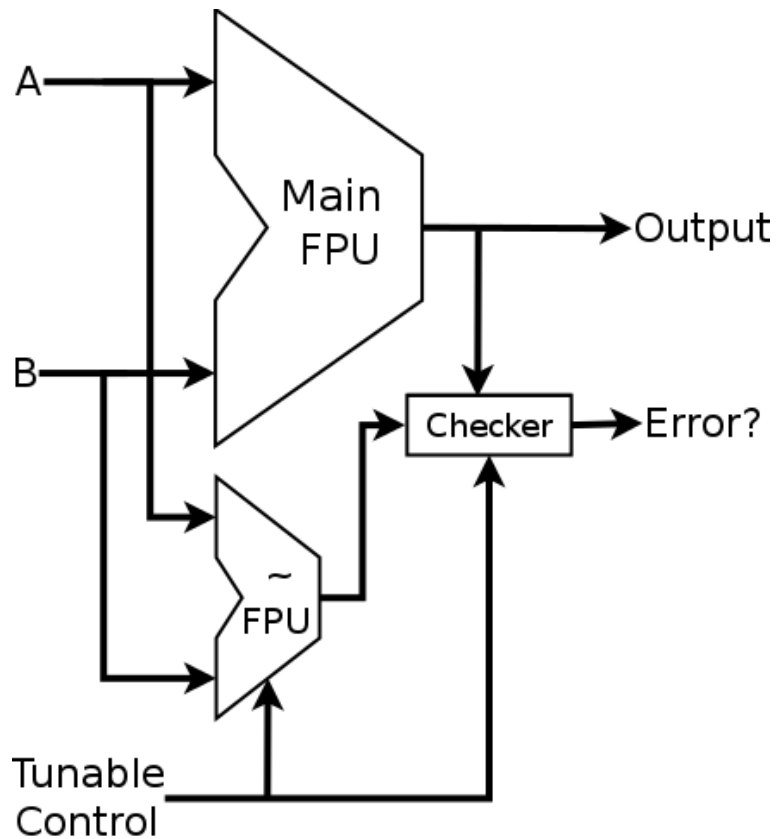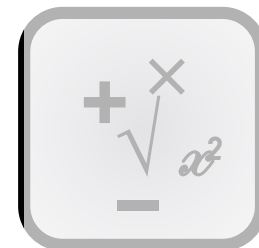
**Lunacy?**

# Relaxed reliability and precision

– Some lunacy
  (rare easy-to-detect errors + parallelism)

– **Lunatic fringe: bounded imprecision**

– Lunacy: live with real unpredictable errors



Legend: ■ Arith.   □ Headroom

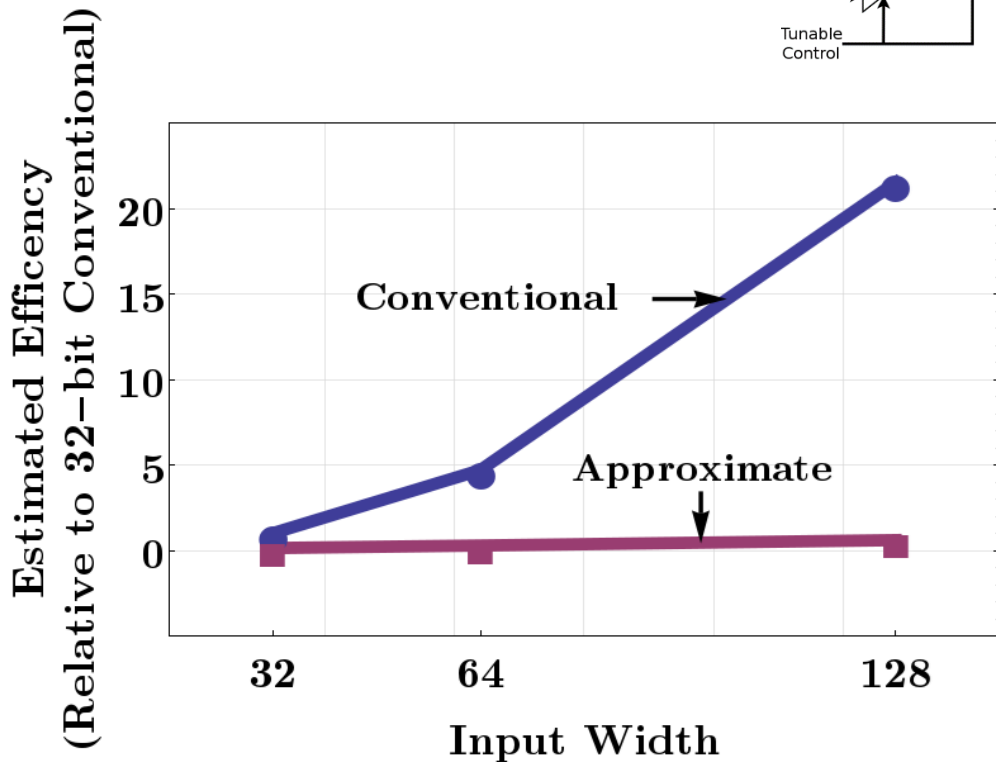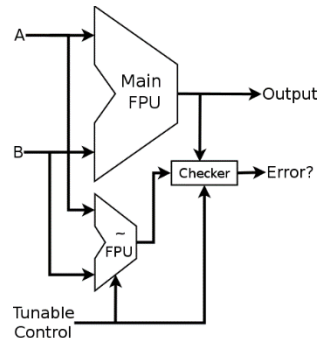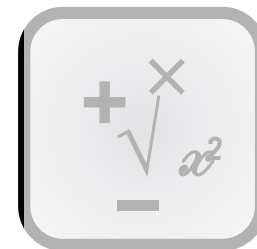| Category | Arith. | Headroom |
|---|---|---|
| Today | 40 | |
| Scaled | 20 | |
| Researchy | 15 | 5 |
| Some lunacy | 12 | 8 |
| Lunatic fringe | 8 | 12 |
| Lunacy | 2 | 18 |

Rough estimated numbers for illustration purposes

# **Bounded** Approximate Duplication

# **Bounded** Approximate Duplication
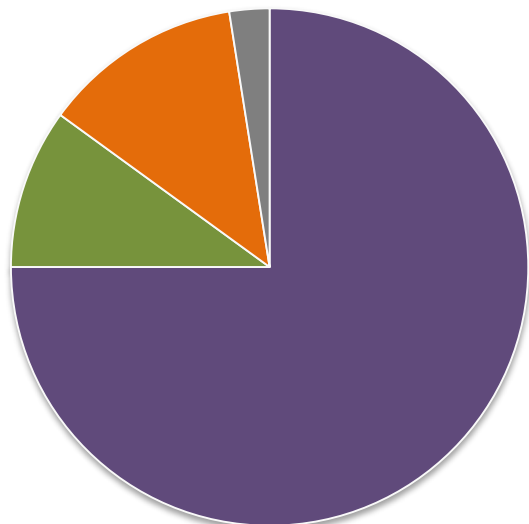
THE UNIVERSITY OF
TEXAS
— AT AUSTIN —

# Supercomputers are general
## → Dynamic **adaptivity** and tuning

– Some programmers crazier than others
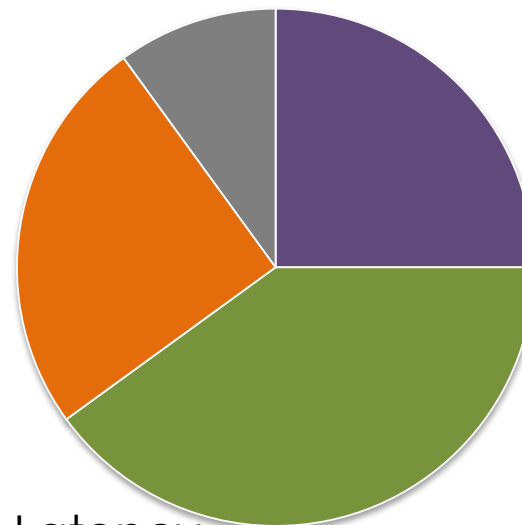
– Large effort in error-tolerant algorithms
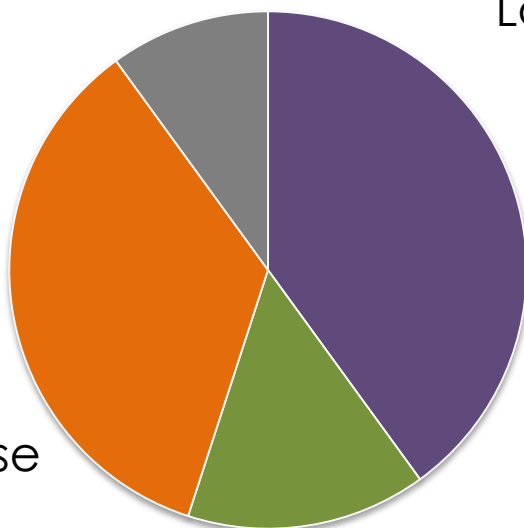
# **Proportionality** and overprovisioning



Dense

Latency

Sparse

Arithmetic
Control
Memory
Margin

Architecture goals:

– Balance possible and practical

– Enable generality

– Don't hurt the common case

It's all about the algorithm

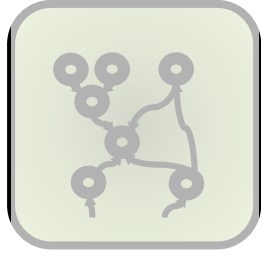THE UNIVERSITY OF
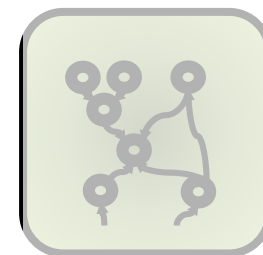**TEXAS**
—— AT AUSTIN ——

# Architecture goals:

– Balance possible and practical
– Enable generality
– Don't hurt the common case

# Architecture concepts so far:

– Proportionality
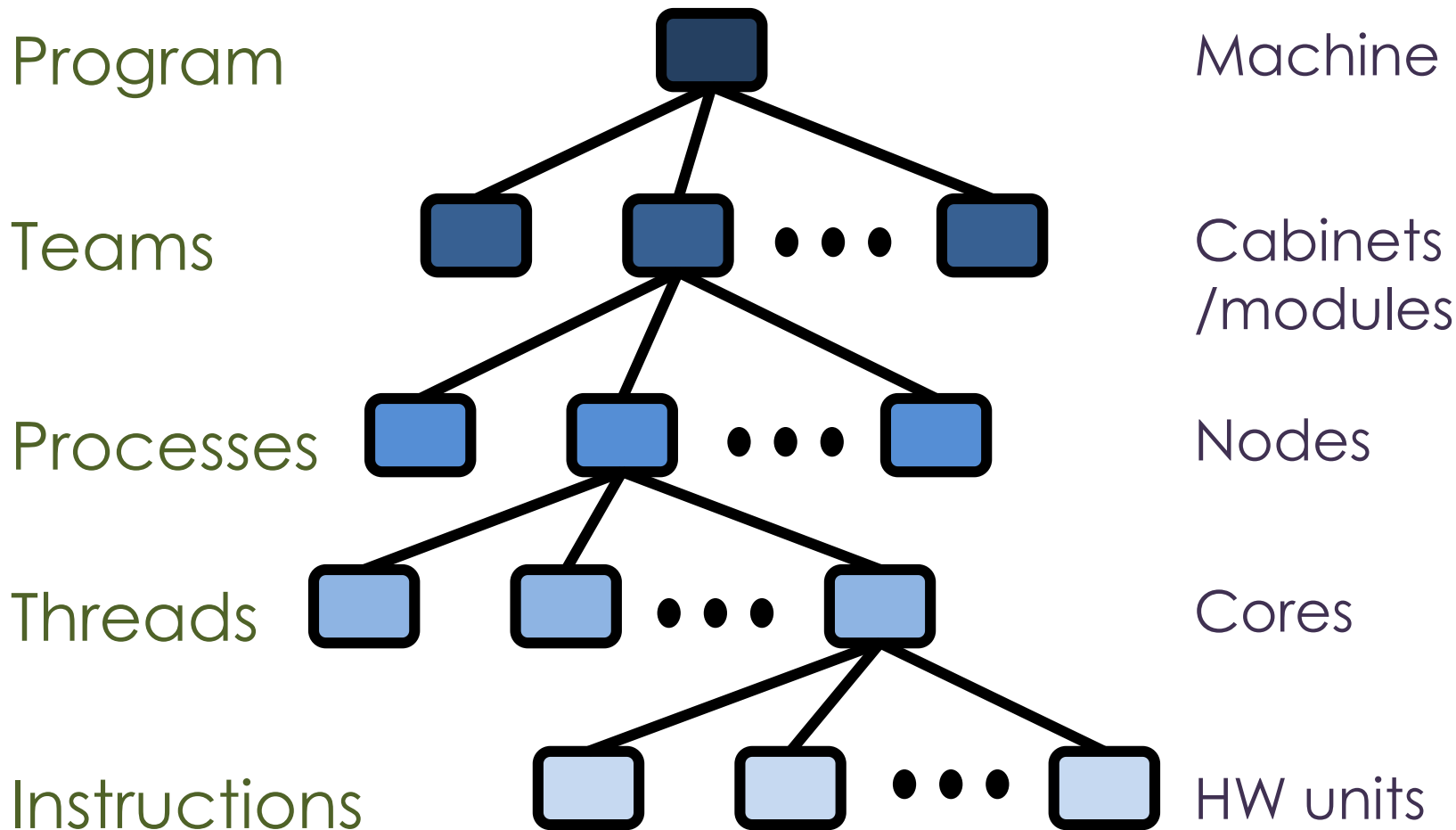  • Adaptivity
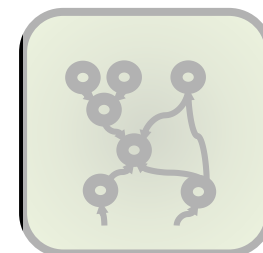  • SW-HW co-tuning
– Locality
– Parallelism

# How to control a **billion arithmetic units**?

# **Hierarchy** minimizes control cost

– Amortize control decisions

| Program | | Machine |
|---|---|---|
| Teams | | Cabinets /modules |
| Processes | | Nodes |
| Threads | | Cores |
| Instructions | | HW units |

# **Hierarchy** minimizes control cost

– Amortize control decisions
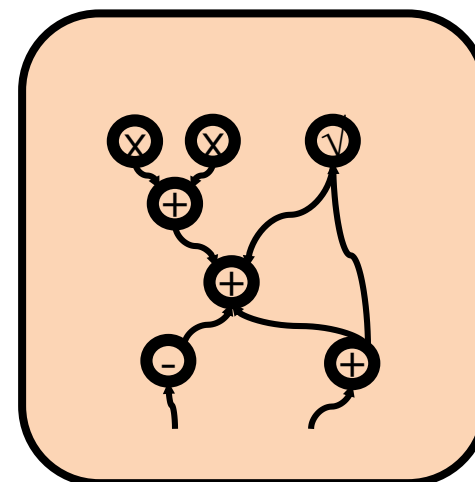


Program

Teams

Processes

**Threads**

**Instructions**

# What does a **HW instruction** do?

+
+
+
+
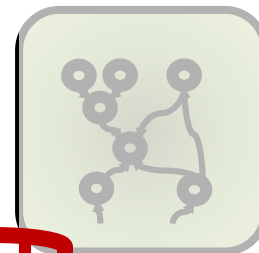x
x
x
x

**CPU**

+ + + +

x x x x

**GPU**

**Embedded**

Amortize/specialize more

THE UNIVERSITY OF
TEXAS
AT AUSTIN

# Single-thread or bulk-thread

latency
oriented
**CPU**

throughput
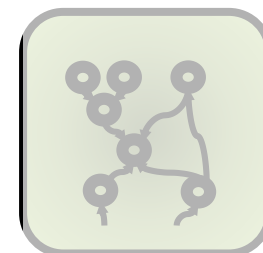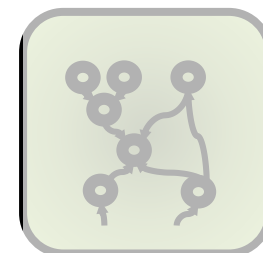oriented
**GPU**

# **Heterogeneity** is a necessity
### → **disciplined specialization**

Must balance generality and control cost

Over-specialization **stifles innovation**
– **Algorithms more important** than hardware

# Heterogeneity is **lunacy**

- Programing and tuning extremely tricky
- Diverse and rapidly evolving algorithms

# Disciplined heterogeneity
## **Scarcity of choice**

– **Throughput** oriented

– **Latency** oriented

– **Disciplined reconfigurable** accelerators

# Heterogeneous computing already here
 – "GPU" for throughput
 – CPU for latency
 – Abstracted FPGA accelerators

# Heterogeneous computing already here

## Titan node (Cray XK7)

Copyrighted
picture of
an XK7 node

Copyrighted
picture of
Kepler die

Copyrighted
picture of
Opteron die

## NVIDIA GPU + AMD CPU

## Stampede node

Copyrighted
picture of
Xeon Phi card

Copyrighted
picture of
Xeon Phi die

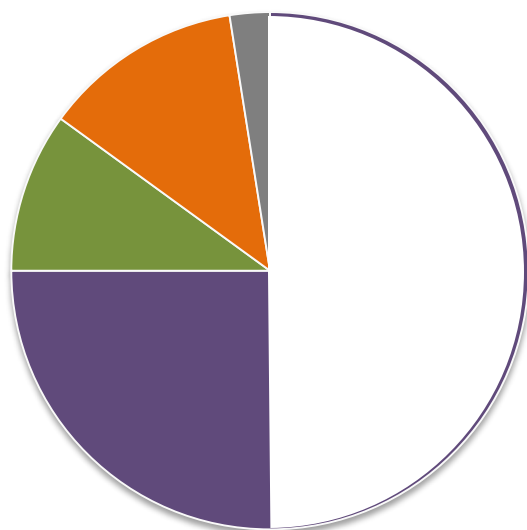Copyrighted
picture of
Stampede
node

Copyrighted
picture of
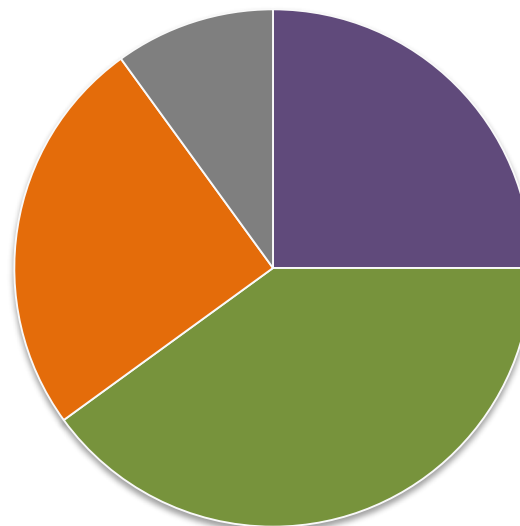Sany Bridge die

## Intel Xeon CPU + Xeon Phi

# **Choose** most **efficient** core
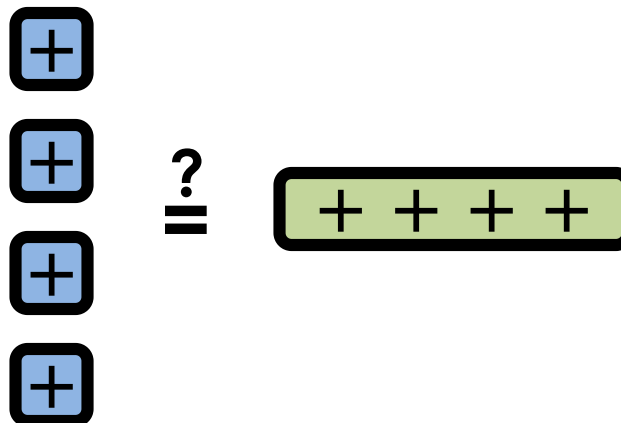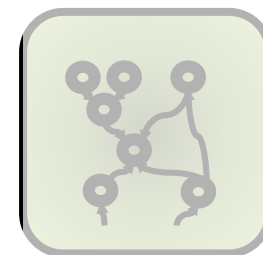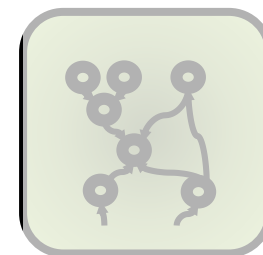
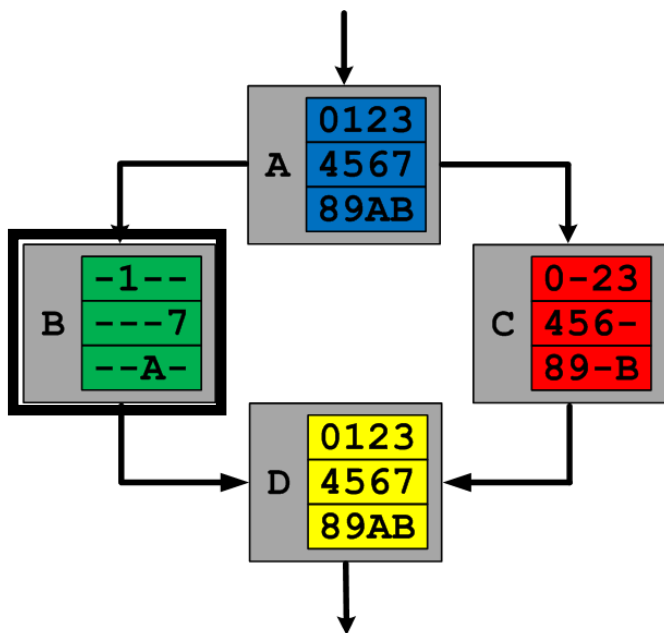# **Choose** most **efficient** core



- Arithmetic
- Control
- Memory
- Margin

# Generalize the throughput core

$$\begin{matrix} + \\ + \\ + \\ + \end{matrix} \overset{?}{=} \boxed{+ \ + \ + \ +}$$

# Generalize the throughput core



(a) Example control flow graph

(b) Execution flow without compaction

(c) Execution flow with compaction
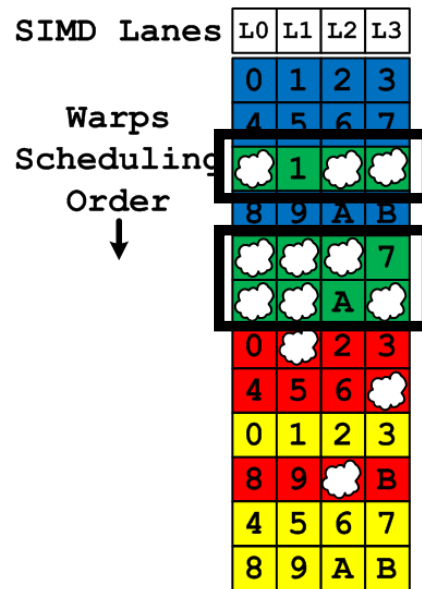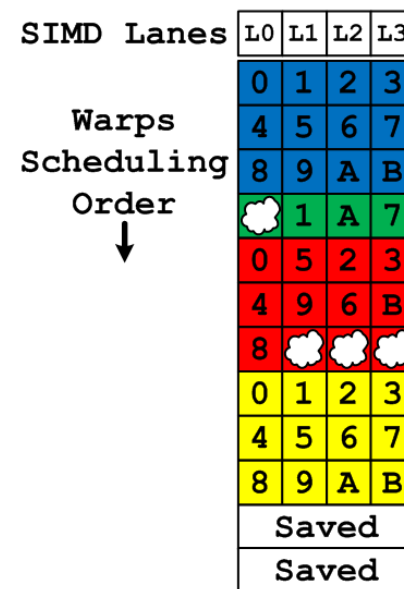
# Generalize the throughput core



(a) Example control flow graph

(b) Execution flow without compaction

(c) Execution flow with compaction

# Synchronization may impair execution

– Predict when necessary → robust optimization
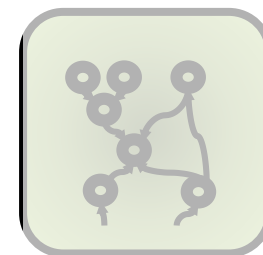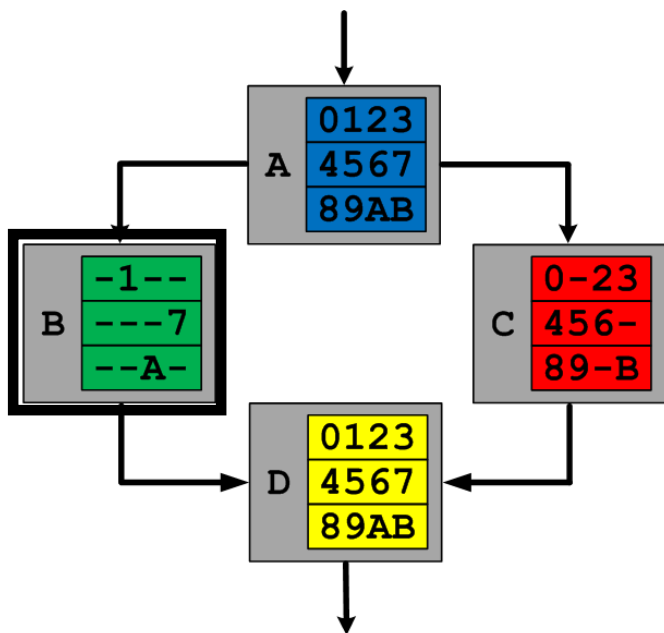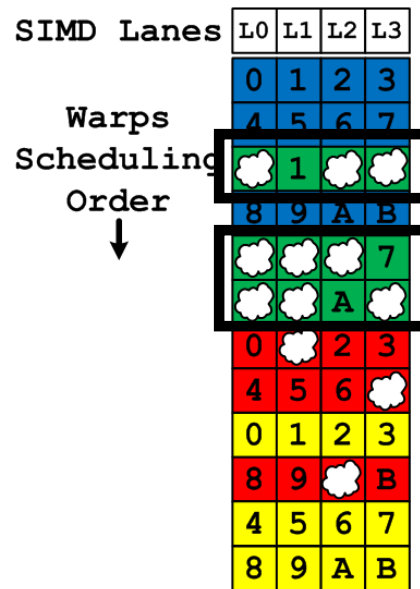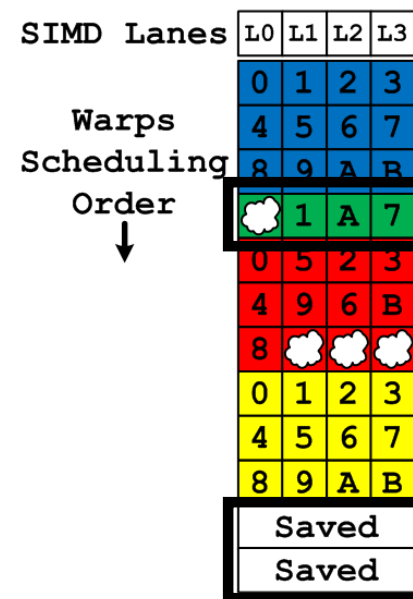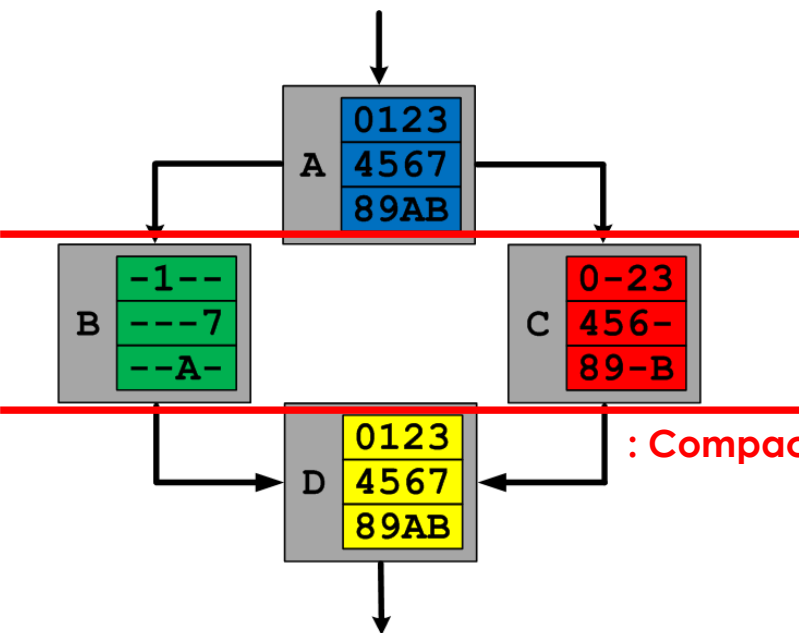– **Adaptive compaction**



: Compaction barrier

(a) Example control flow graph

(b) Execution flow without compaction

(c) Execution flow with compaction

# Architecture goals:
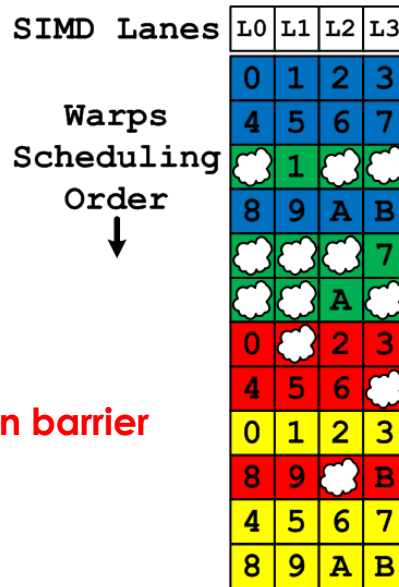
– Balance possible and practical
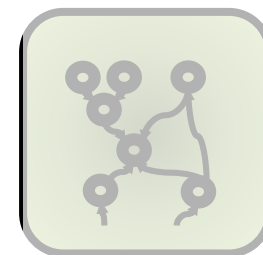
– Enable generality

– Don't hurt the common case

# Architecture so far:

– Proportionality
  - Adaptivity
  - Heterogeneity
  - SW-HW co-tuning
– Locality
– Parallelism
– Hierarchy

# Heterogeneity is **lunacy**
– Programing and tuning extremely tricky
– Diverse and rapidly evolving algorithms

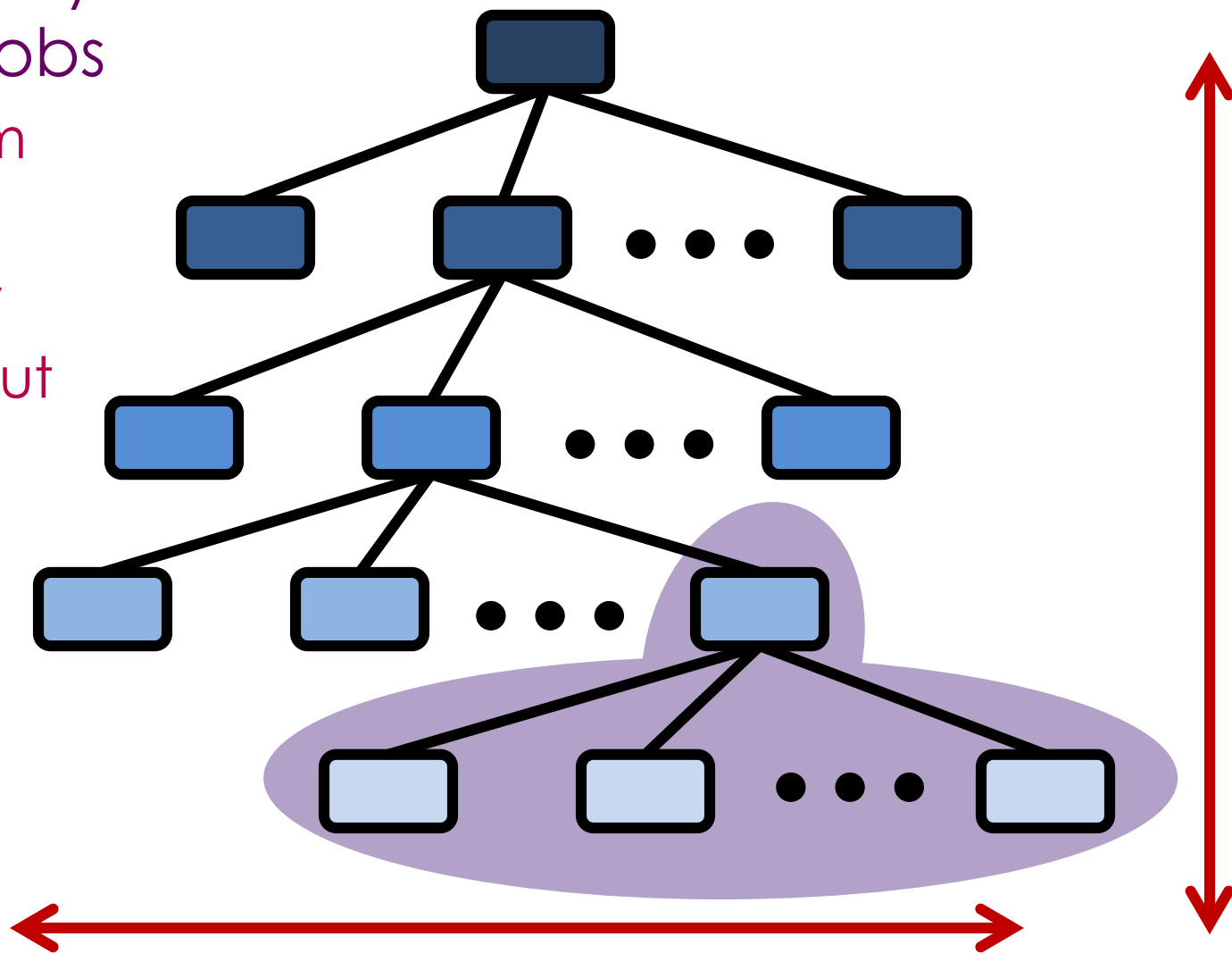# How can we **program** these machines

# **Hierarchical** languages restore **some sanity**

– **Abstract** key
  tuning knobs
  - Parallelism
  - Locality
  - Hierarchy
  - Throughput

# **Hierarchical languages** restore some sanity

- CUDA and OpenCL
- Sequoia (Stanford)
- Habanero (Rice)
- Phalanx (NVIDIA)
- Legion (Stanford)
- Working into X10, Chapel
- ...

Domain specific languages even better

# A counter example:
# DE Shaw Research **Anton**

Copyrighted
picture of
Anton package
(Google
image search)

Copyrighted
figure demonstrating
molecular dynamics
(Google
image search)

Copyrighted
figure of
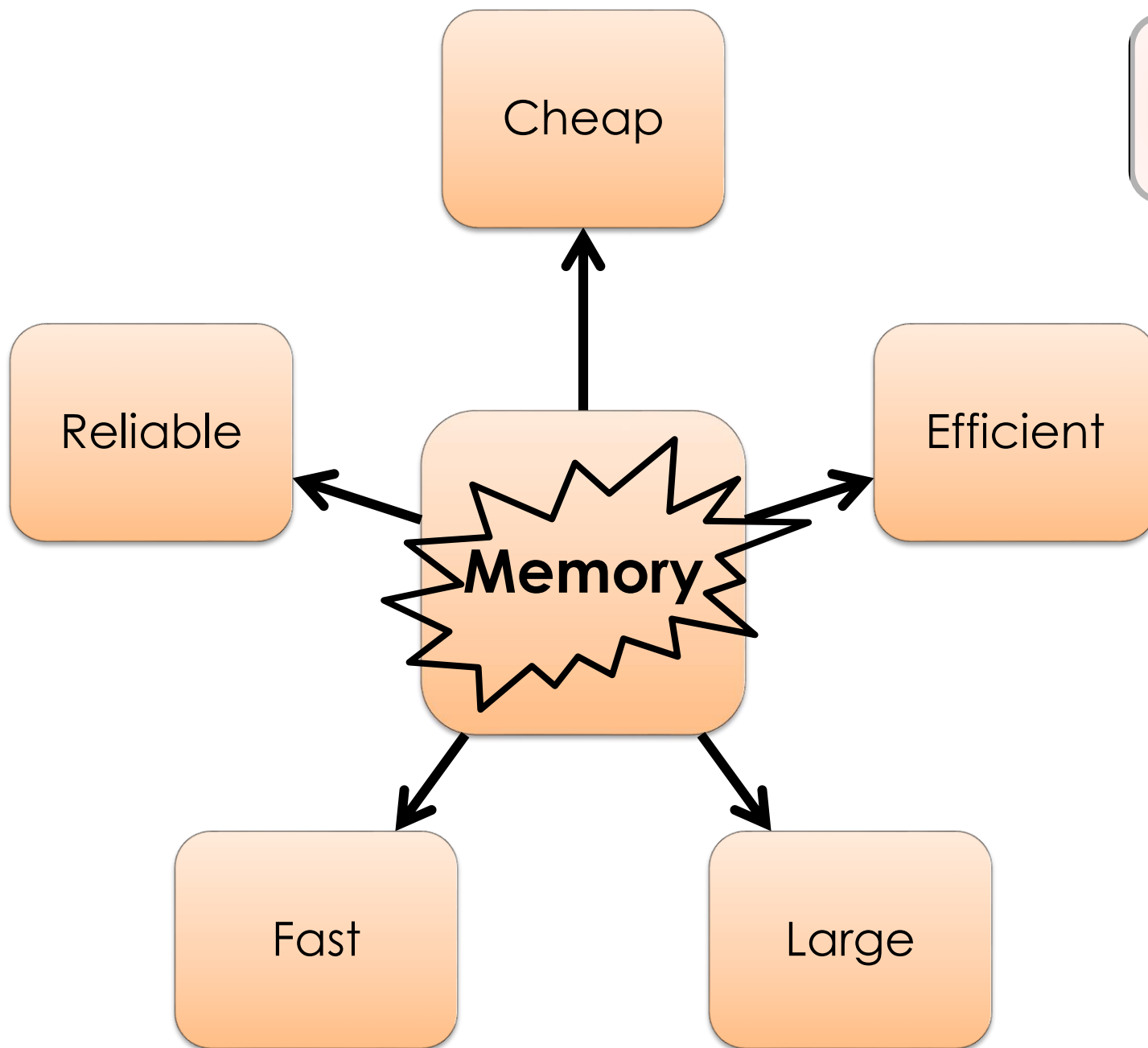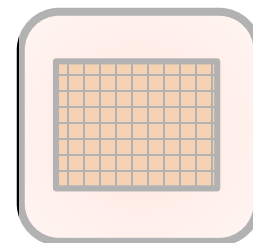Anton system
(Google
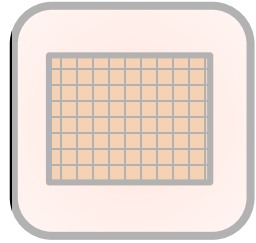image search)

# Another counter example?

## **Microsoft Catapult**

– Disciplined reconfigurability for the cloud

– Distributed FPGA accelerator

• Within form-factor and cost constraints

– Architected interfaces and compiler

• Compiler for software, not (just) hardeare

Copyrighted
figure of
datacenter
(Google image search)

Copyrighted
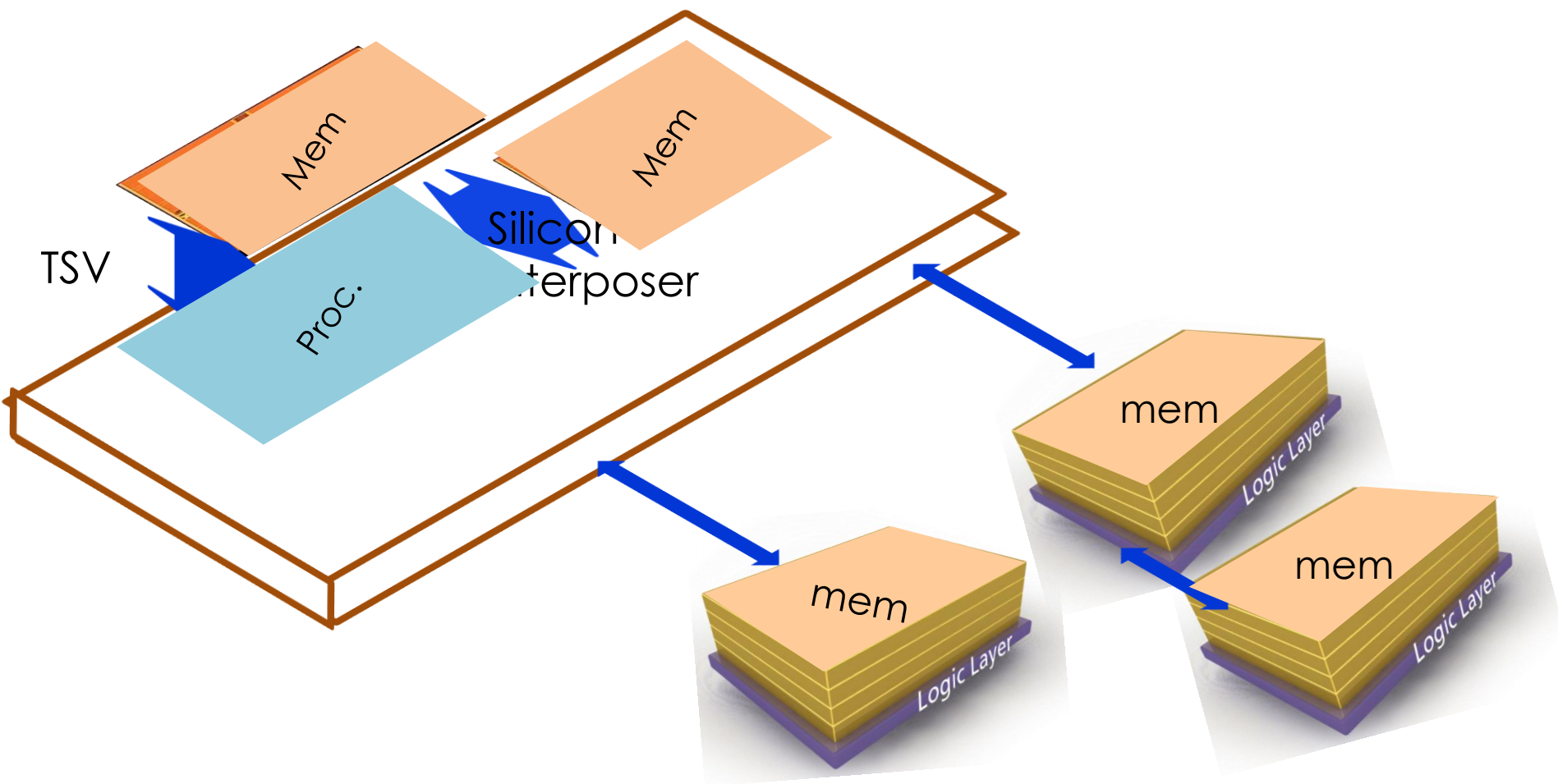figure of
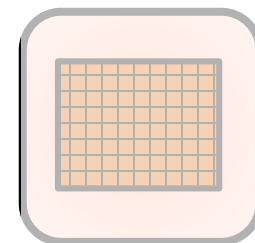Catapult diagram
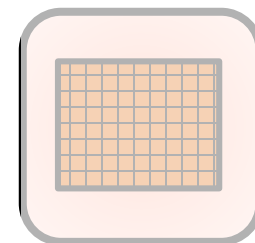(from ISCA 2014 paper)

THE UNIVERSITY OF
**TEXAS**
— AT AUSTIN —

# Fast + large →
## **hierarchy**

# Fast + large → **hierarchical**



TSV

Mem

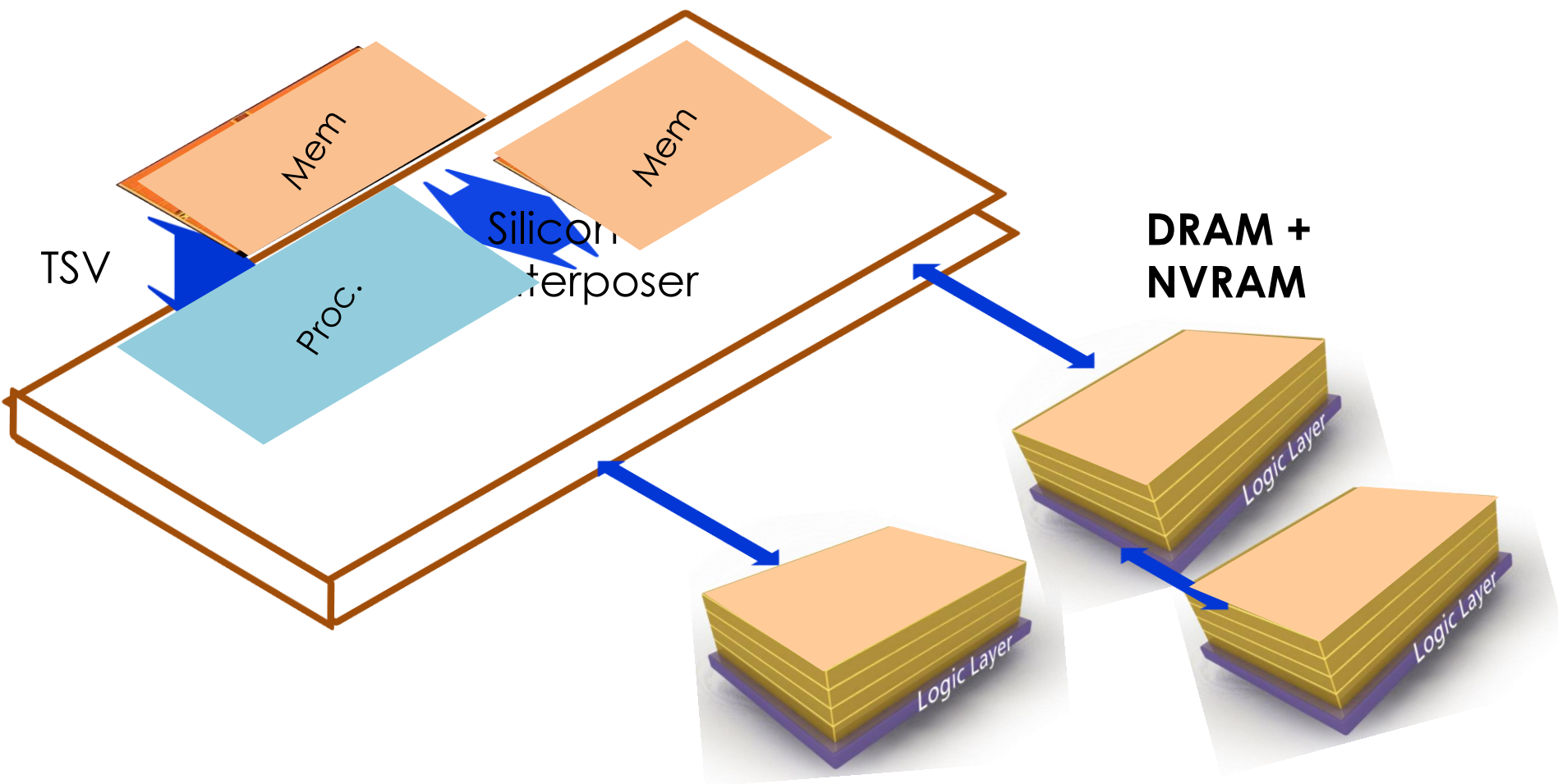Mem

Proc.

Silicon Interposer

mem

mem

mem

Logic Layer

# Large + efficient →
# **heterogeneity**

# Large + efficient → **heterogeneous**

TSV

Mem

Mem

Proc.

Silicon Interposer

**DRAM + NVRAM**

Logic Layer

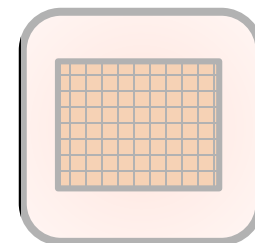Logic Layer

Logic Layer

# Heterogeneous + hierarchical
## → **big mess**

# Research just starting

– New mechanisms needed

– Very interesting tradeoffs with reliability
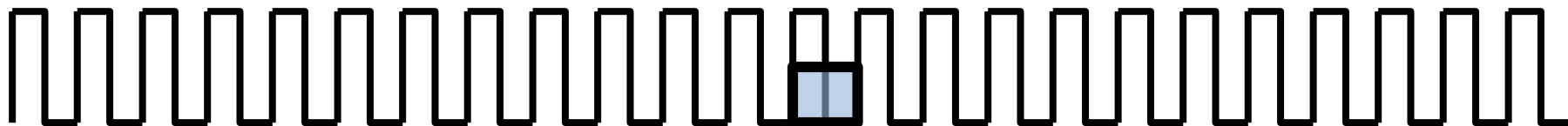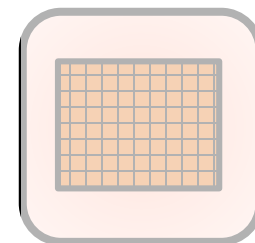
# Fast + efficient
## → **control hierarchy**
## **+ parallelism**

# Fast + efficient
### → **hierarchy + parallelism**
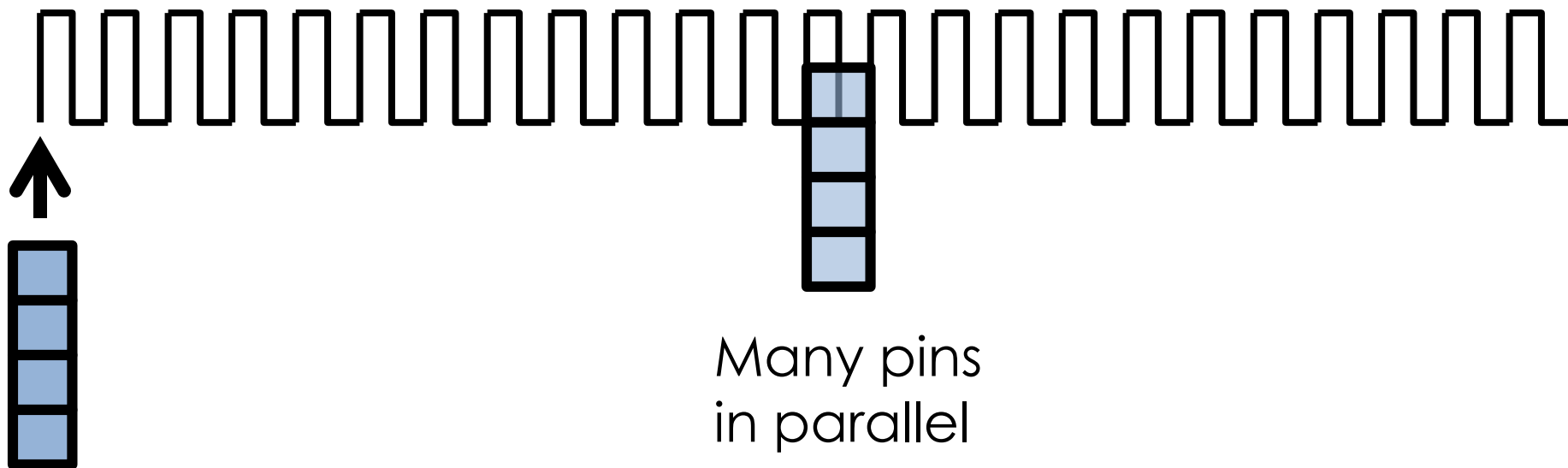
Access cell

# Fast + efficient
## → **hierarchy + parallelism**

Many pins
in parallel

Access many
cells in parallel

# Fast + efficient
### → **hierarchy + parallelism**

Many pins
in parallel over
multiple cycles

Access even more
cells in parallel
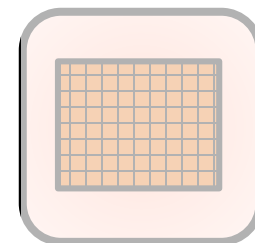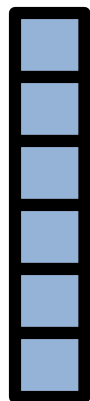
# Fast + efficient
## → **hierarchy + parallelism**

Many pins
in parallel over
multiple cycles
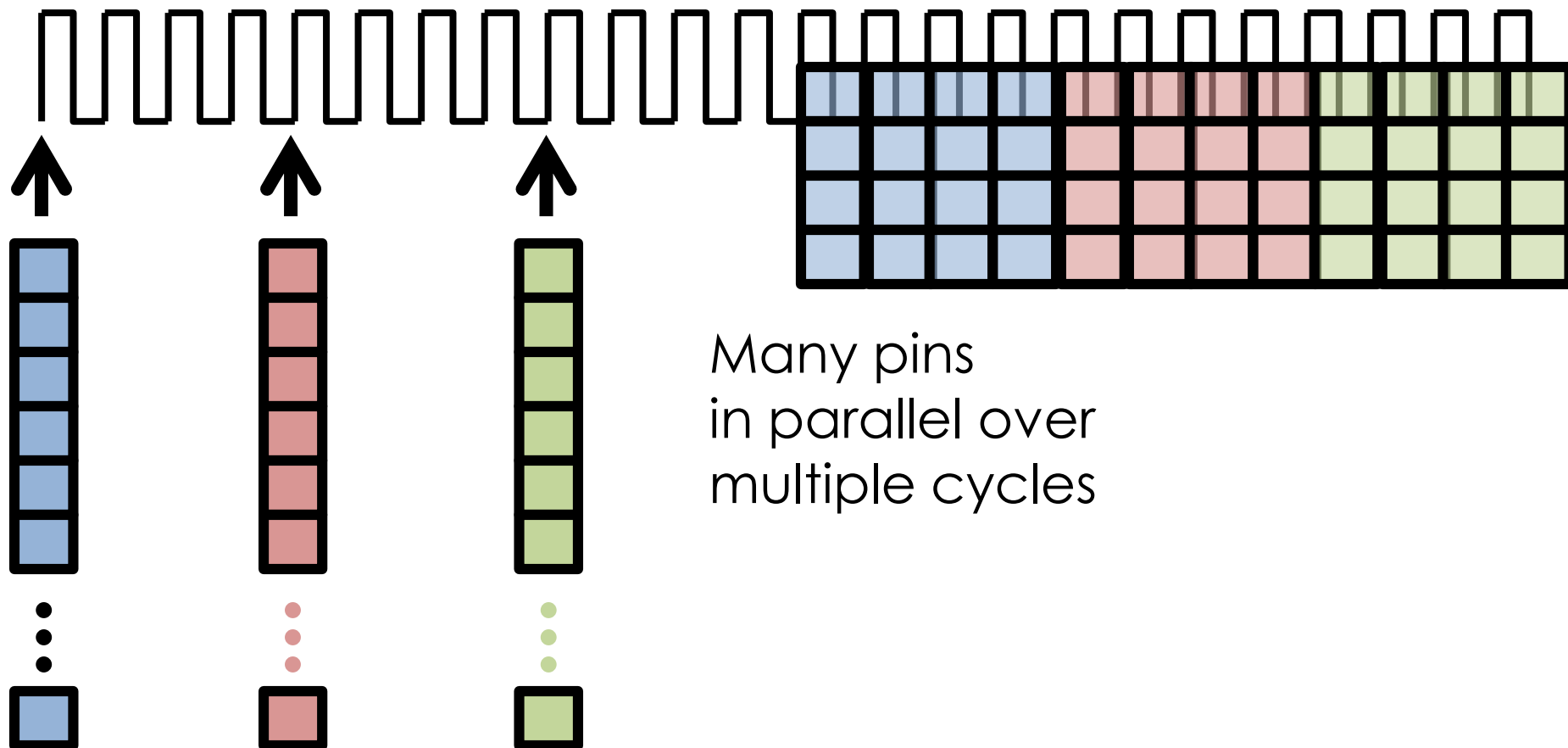
Access even more
cells in parallel

# Hierarchy + parallelism
## → **potential waste**

# Is fine-grained access better?

– Wasteful when all data is needed

CG  C    Data    ECC

FG  C0 C1 C2 C3 C4 C5 C6 C7   D0 E0   D1 E1   D2 E2   D3 E3   D4 E4   D5 E5   D6 E6   D7 E7

CG  C    Data    ECC

FG  C0   D0 E0

# Dynamically **adapt** granularity → best of both worlds

**Locality Predictor**

**Sector cache**
**(8 8B sectors per line)**

**Co-scheduling CG/FG w/ split CG**

**Sub-Ranked DRAM w/ 2X ABUS**
**Support both CG and FG**

SLP

Processor Core

$D    $I

L2

Core 0    Core 1    ...    Core N-1

Last Level Cache

MC

DRAM

# Dynamically **adapt** granularity

Cost + ...        → **poor reliability**

Efficiency + ... → **poor reliability**

New + ...        → **poor reliability**

# **Compensate** with error protection
 – ECC

# Compensate with **proportional protection**

# Adaptive reliability

– Adapt the **mechanism**

– Adapt the **error rate**

- precision (stochastic precision)

# Adaptive reliability

– By **type**
  - Application guided
– By **location**
  - Machine-state guided

# Example: **Virtualized ECC**
– Adapt the mechanism

THE UNIVERSITY OF
TEXAS
AT AUSTIN

**Virtual Address space**

**Physical Memory**

**Protection Level**

**Low**

Virtual page – x

**Virtual Page to Physical Frame mapping**

Page frame – x

Page frame – y

**Medium**

Virtual page – y

**High**

Virtual page – z

Page frame – z

**Data**      **ECC**

**Virtual Address space**

**Physical Memory**

**Protection Level**

**Low** — Virtual page – x → Page frame – x

**Virtual Page to Physical Frame mapping**

**Medium** — Virtual page – y → Page frame – y

**High** — Virtual page – z → Page frame – z

**ECC**

ECC page – y

ECC page – z

**Physical Frame to ECC Page mapping**

**Stronger ECC**

**Data**

THE UNIVERSITY OF
**TEXAS**
AT AUSTIN

**Virtual Address space**

**Physical Memory**

**Protection Level**

**Low**

Virtual page – x

Page frame – x

**Virtual Page to Physical Frame mapping**

**Medium**

Virtual page – y

Page frame – y

**High**

Virtual page – z

Page frame – z

**T2EC for Chipkill**

ECC page – y

ECC page – z

**Physical Frame to ECC Page mapping**

**T2EC for Double Chipkill**

**Data**     **T1EC**

**Write** **Read**

T1EC
Decoding

T1EC/T2EC
Encoding

**Two-Tiered ECC** **T1EC** **T2EC**

**Data**

# Caching work great



**Normalized Execution Time**

# Bonus: **flexibility** of ECC



**Normalized Energy-Delay Product**

# Can we do even better?

– Hide second-tier

# FrugalECC = Compression + VECC

# Adapt resilience scheme or adapt storage?



Figure 3: Frugal ECC layout for chipkill-correct (64-bit redundancy).

# Adapt compression scheme too! Coverage-oriented-Compression (CoC)



←————(443=512-64-5)-bit compressed data————→

Compression flag (1 to 4-bit)

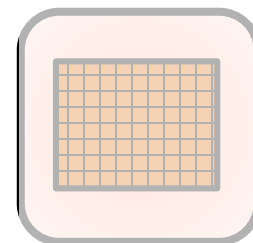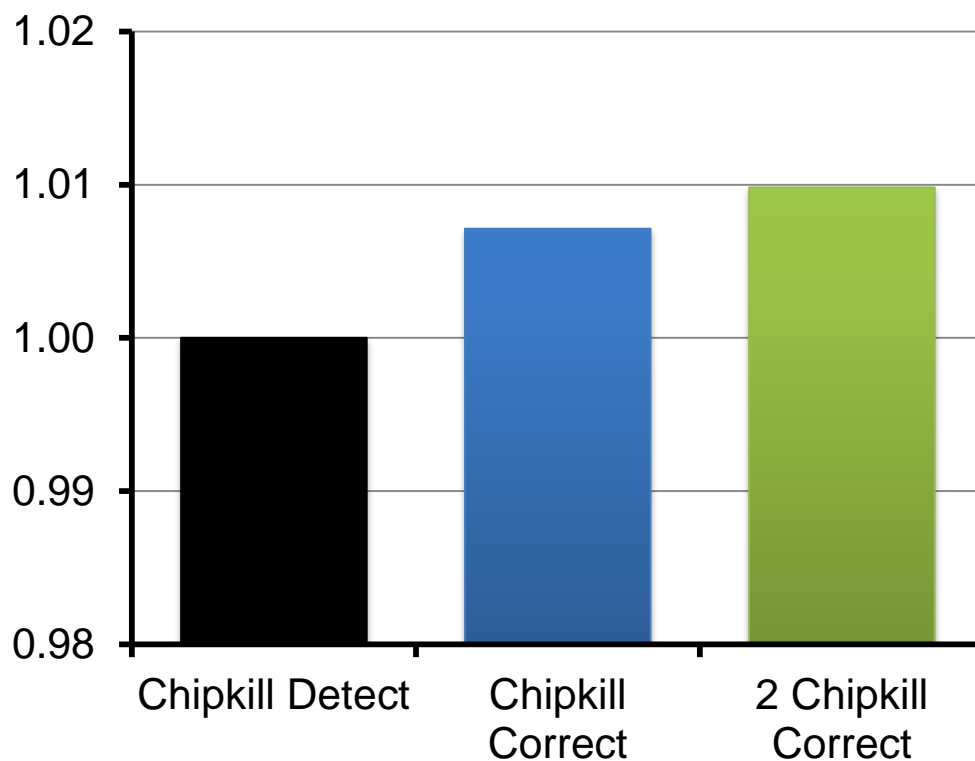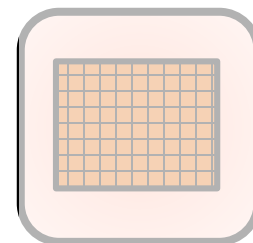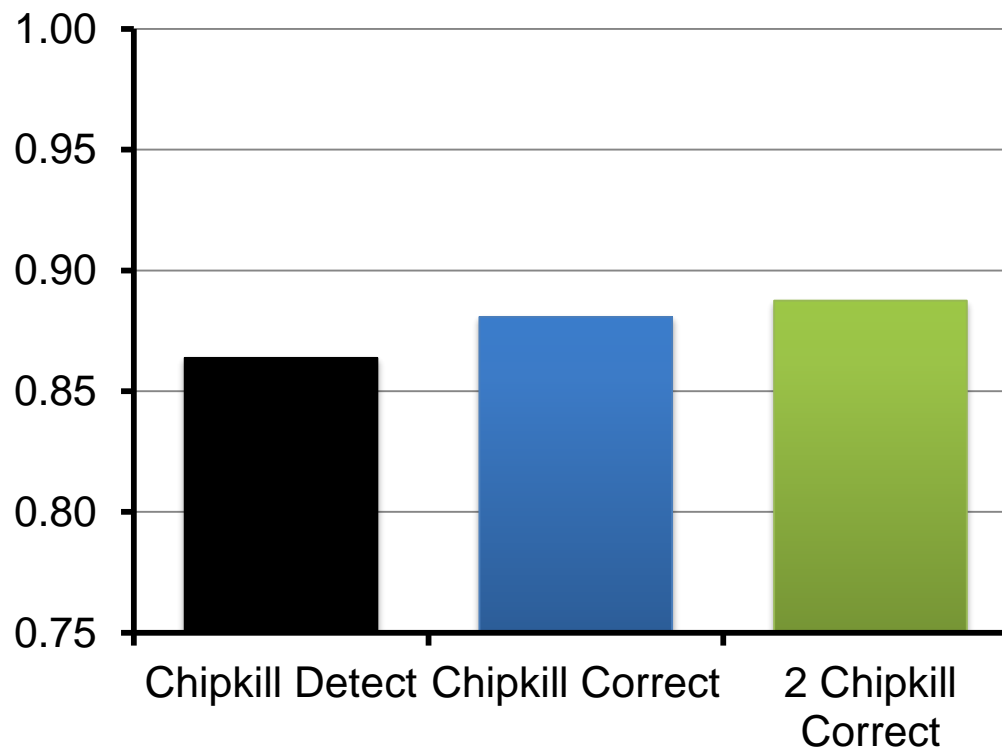| 1100 | 64b base (quadword) | {54b delta} x 4 + {53b delta} x 3 |
| 1101 | 32b base | {28b delta} x 2 + {27b delta} x 13 |
| 1110 | 16b B | {14b delta} x 20 + {13b delta} x 11 |
| 1111 | 8b B | {7b delta} x 53 + {6b delta} x 10 |
| 0 | 64b base (double) | {54b delta} x 7 |
| 101 | 32b base (flt) | {28b delta} x 3 + {27b delta} x 12 |
| 100 | ID | Comp.data | ID | Comp.data | ID | Comp.data |

Compressed words (up to 440b)

Experiments promising:

Match or exceed reliability

Improve power

Minimal impact on performance

THE UNIVERSITY OF
TEXAS
— AT AUSTIN —

# Architecture goals:

– Balance possible and practical
– Enable generality
– Don't hurt the common case

# Architecture so far:

– Proportionality
  • Adaptivity
  • SW-HW co-tuning
  • Heterogeneity
– Locality
– Parallelism
– Hierarchy

# The **reliability** challenge

# Detect → Contain → Repair → Recover

# Today's techniques not good enough
– Hardware support expensive
– Checkpoint-restart doesn't scale

# Exascale resilience **must be**:

– Proportional

– Parallel

– Adaptive

– Hierarchical

– Analyzable

# Containment domains

– Embed resilience within application
– Match system hierarchy and characteristics
– Analyzable abstraction consistent across layers
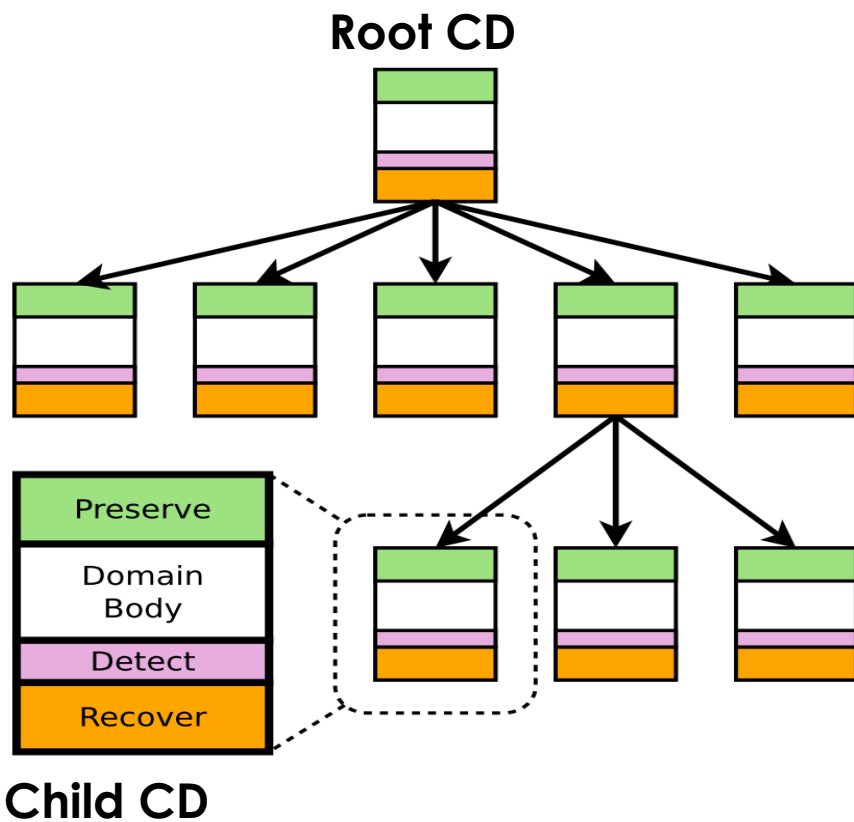


**Root CD**

**Child CD**

Preserve
Domain Body
Detect
Recover

CDs don't communicate erroneous data

CDs have means to recover

Phalanx C++ program

CD Annotations
resiliency model

efficiency-oriented
programming model

CD API
resiliency interface

Runtime Library Interface

CD control and persistence

Microkernel

Hardware Abstraction Layer

Error
Reporting
Architecture

ECC, status

Machine

101, 106

# Mapping example: SpMV



Matrix **M**          Vector **V**

```
void task<inner> SpMV( in M, in Vi,
out Ri){
    forall(…) reduce(…)
        SpMV(M[…],Vi[…],Ri[…]);
}


void task<leaf> SpMV(…){

  for r=0..N

    for c=rowS[r]..rowS[r+1] {
        resi[r]+=data[c]*Vi[cIdx[c]];

        prevC=c;

    }

}
```

THE UNIVERSITY OF
TEXAS
— AT AUSTIN —

# Mapping example: SpMV

$$M_{00} \quad M_{01}$$

$$M_{10} \quad M_{11}$$

Matrix **M**

$$V_0$$

$$V_1$$

Vector **V**

```
void task<inner> SpMV( in M, in V_i,
out R_i){
    forall(…) reduce(…)
      SpMV(M[…],V_i[…],R_i[…]);
}


void task<leaf> SpMV(…){

  for r=0..N

    for c=rowS[r]..rowS[r+1] {
       res_i[r]+=data[c]*V_i[cIdx[c]];

      prevC=c;

    }

}
```

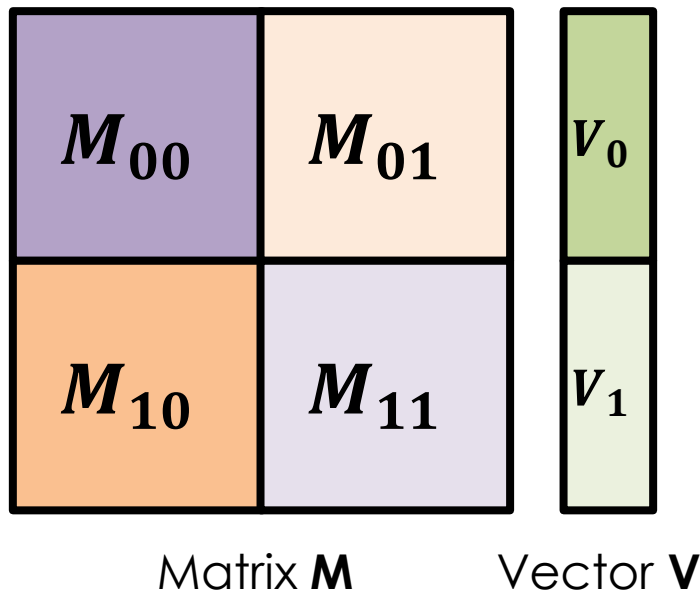# Mapping example: SpMV



```
void task<inner> SpMV( in M, in V_i,
out R_i){
    forall(…) reduce(…)
        SpMV(M[…],V_i[…],R_i[…]);
}


void task<leaf> SpMV(…){
  for r=0..N
    for c=rowS[r]..rowS[r+1] {
        res_i[r]+=data[c]*V_i[cIdx[c]];
        prevC=c;
    }
}
```
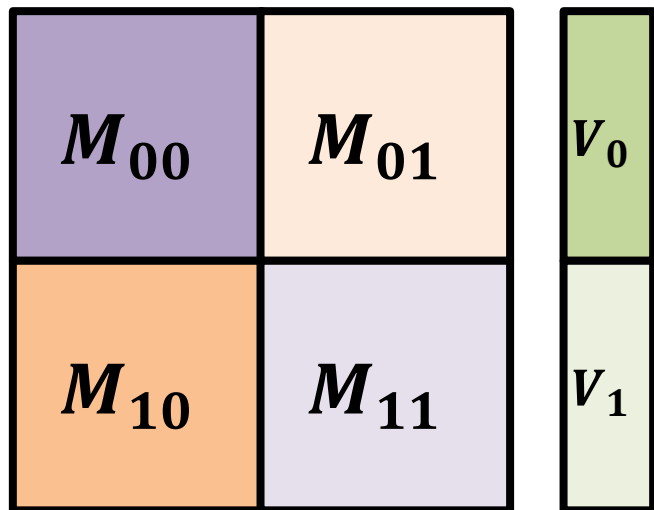
Matrix **M**    Vector **V**

**Distributed to 4 nodes**

$M_{00}$ $V_0$     $M_{10}$ $V_0$     $M_{01}$ $V_1$     $M_{11}$ $V_1$
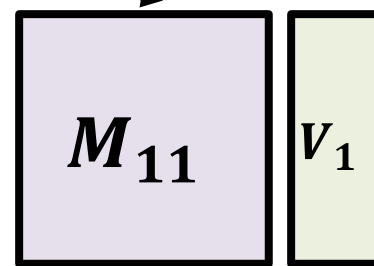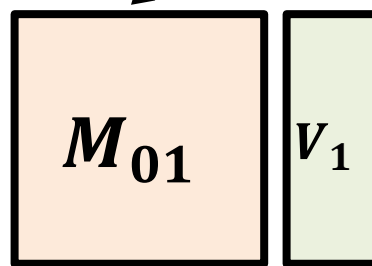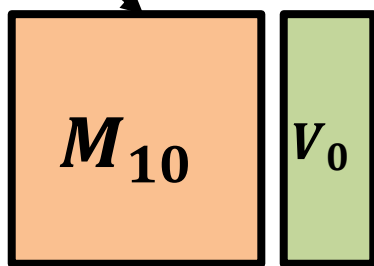
# Mapping example: SpMV
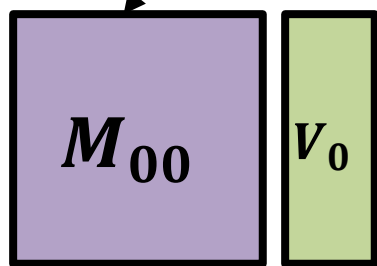


Matrix **M**    Vector **V**

```
void task<inner> SpMV( in M, in Vi,
out Ri){
    forall(…) reduce(…)
      SpMV(M[…],Vi[…],Ri[…]);
}


void task<leaf> SpMV(…){
  for r=0..N
    for c=rowS[r]..rowS[r+1] {
      resi[r]+=data[c]*Vi[cIdx[c]];
      prevC=c;
    }
}
```
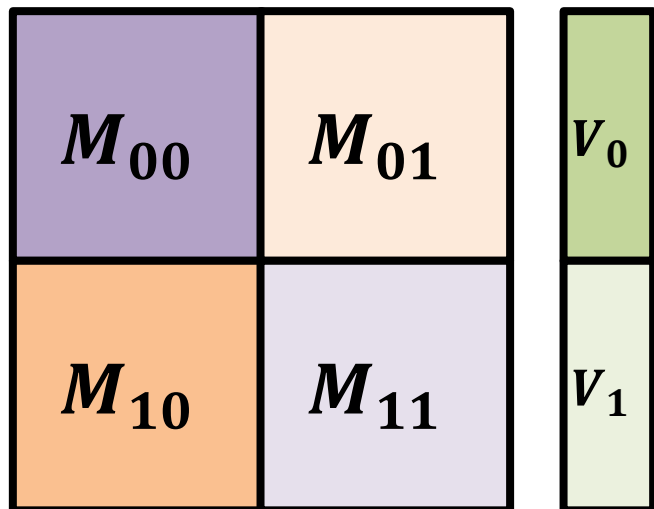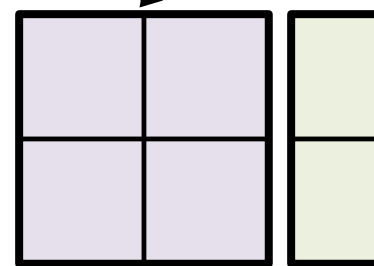
**Distributed to 4 nodes**
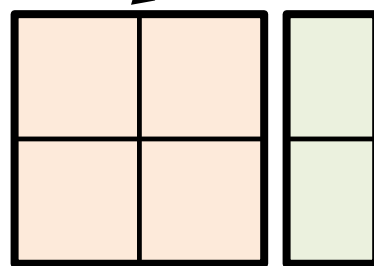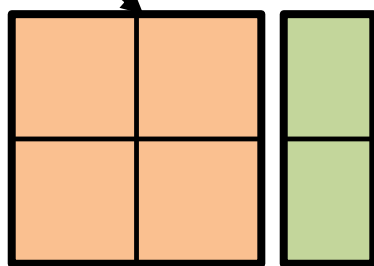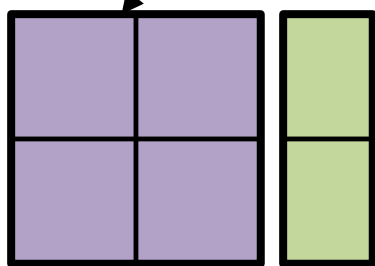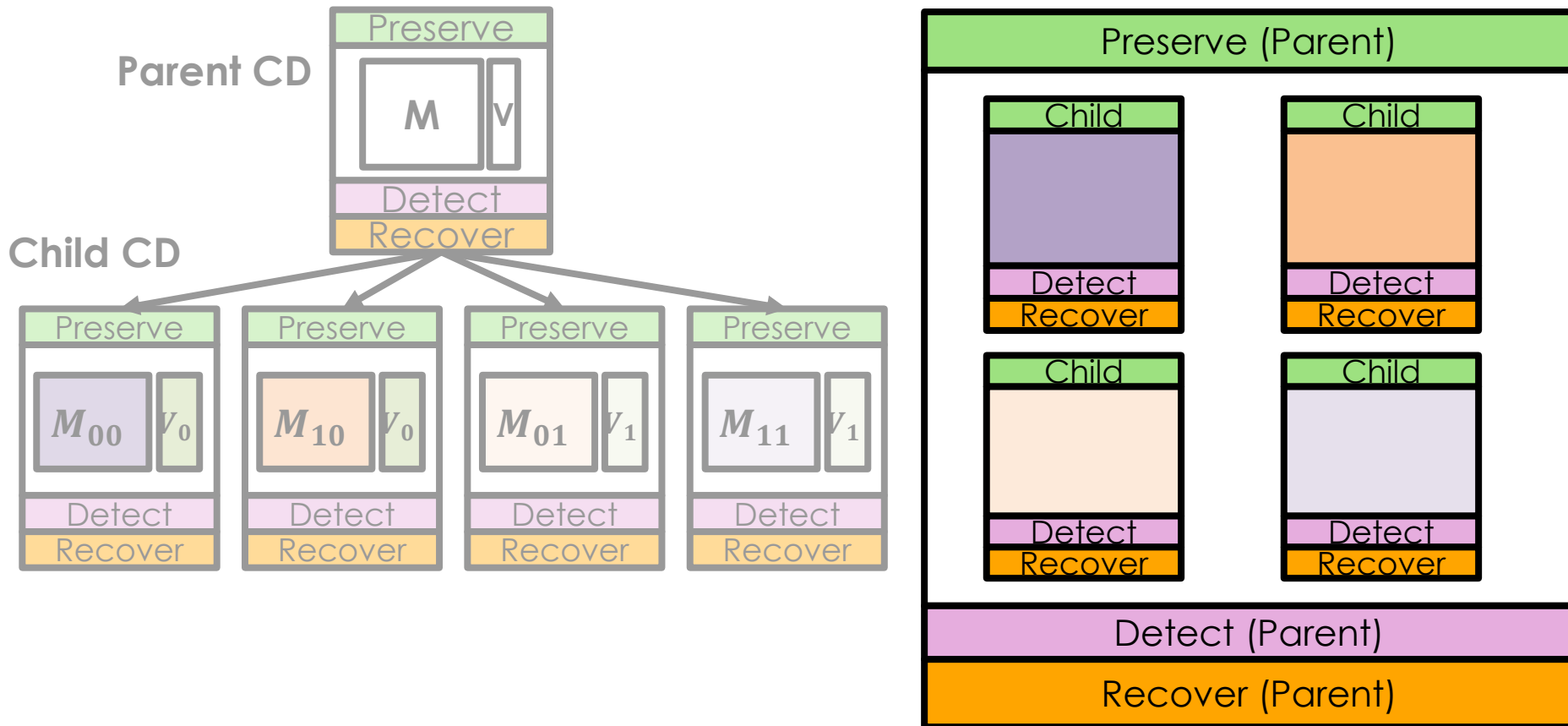
# Mapping example: SpMV

```
void task<inner> SpMV(in M, in Vi, out Ri) {
  cd = begin(parentCD);
  preserve_via_copy(cd, matrix, …);
  forall(…) reduce(…)
    SpMV(M[…],Vi[…],Ri[…]);
  complete(cd);
}

void task<leaf> SpMV(…) {
  cd = begin(parentCD);
  preserve_via_copy(cd, matrix, …);
  preserve_via_parent(cd, veci, …);
  for r=0..N
   for c=rowS[r]..rowS[r+1] {
     resi[r]+=data[c]*Vi[cIdx[c]];
     check {fault<fail>(c > prevC);}
     prevC=c;
   }
  complete(cd);
}
```

| API |
| --- |
| begin |
| configure |
| preserve |
| check |
| complete |

# SpMV preservation tuning

$M_{00}$ $M_{01}$ $M_{10}$ $M_{11}$ $V_0$ $V_1$

```
void task<leaf> SpMV(…) {
  cd = begin(parentCD);
  preserve_via_copy(cd, matrix, …);
  preserve_via_parent(cd, vec_i, …);
  for r=0..N
    for c=rowS[r]..rowS[r+1] {
      res_i[r]+=data[c]*V_i[cIdx[c]];
      check {fault<fail>(c > prevC);}
      prevC=c;
    }
  complete(cd);
}
```
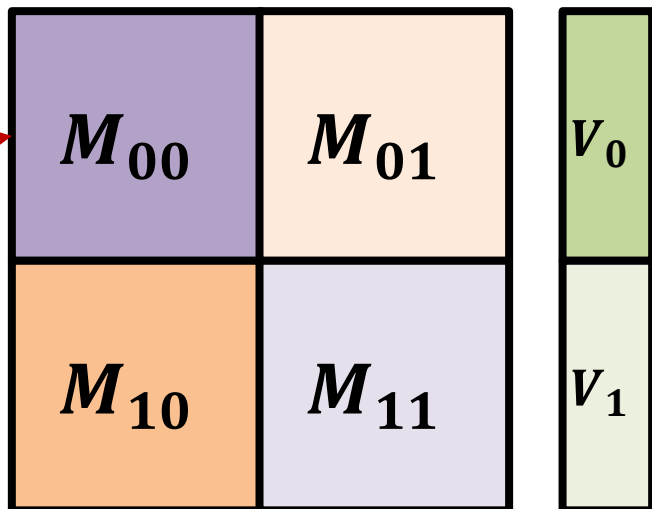
**Hierarchy** Matrix **M**        Vector **V**

$M_{00}$ $V_0$    $M_{10}$ $V_0$    $M_{01}$ $V_1$    $M_{11}$ $V_1$

**Natural redundancy**

# Concise abstraction for complex behavior
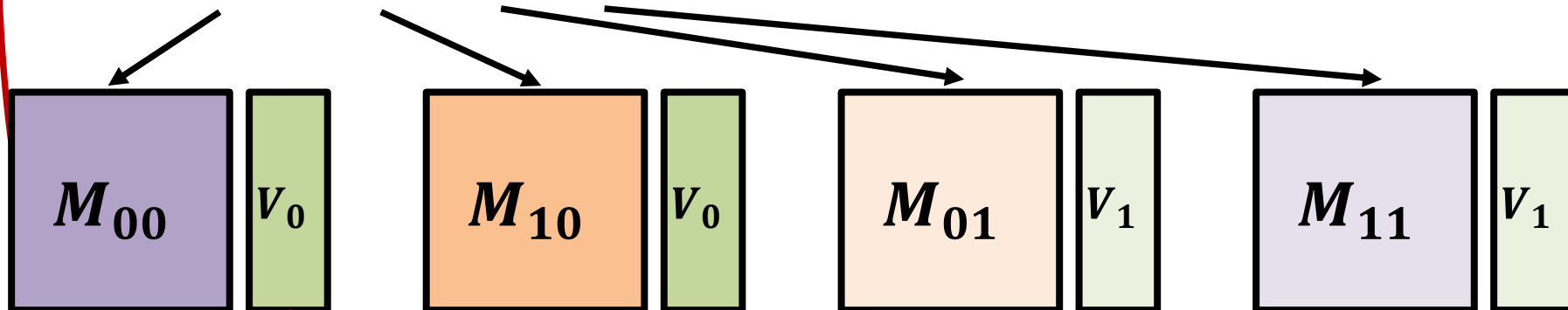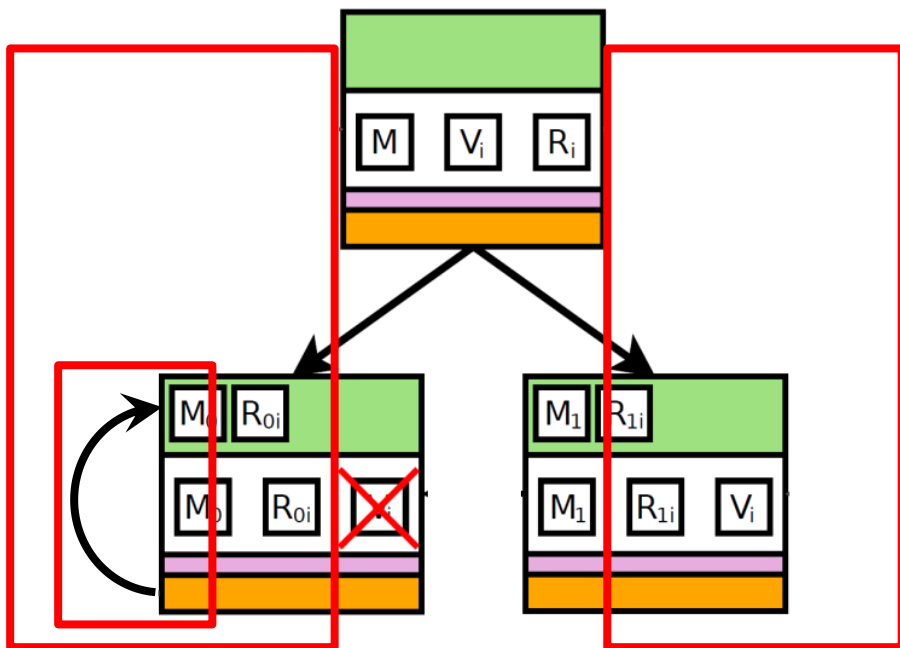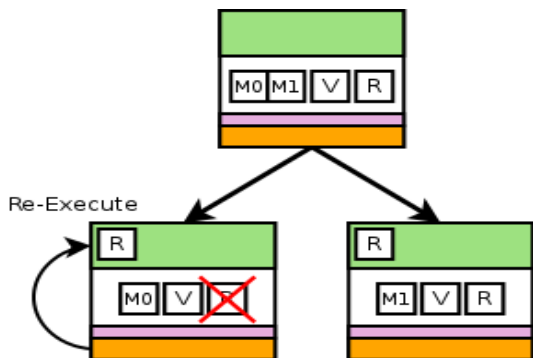


```
void task<leaf> SpMV(…) {
  cd = begin(parentCD);
  preserve_via_copy(cd, matrix, …);
  preserve_via_parent(cd, vec_i, …);
  for r=0..N
    for c=rowS[r]..rowS[r+1] {
      res_i[r]+=data[c]*V_i[cIdx[c]];
      check {fault<fail>(c > prevC);}
      prevC=c;
    }
  complete(cd);}
```
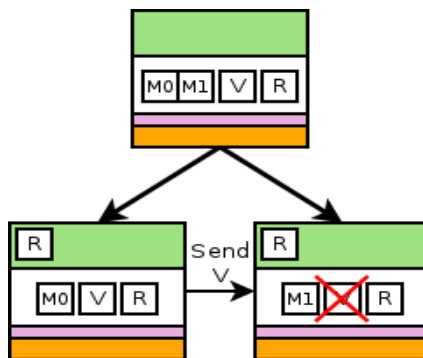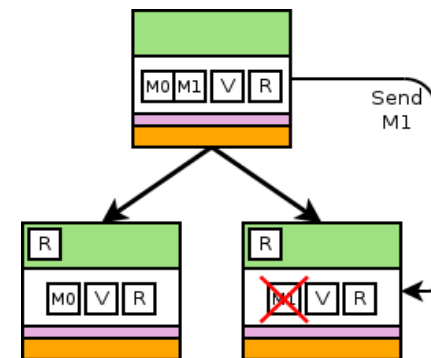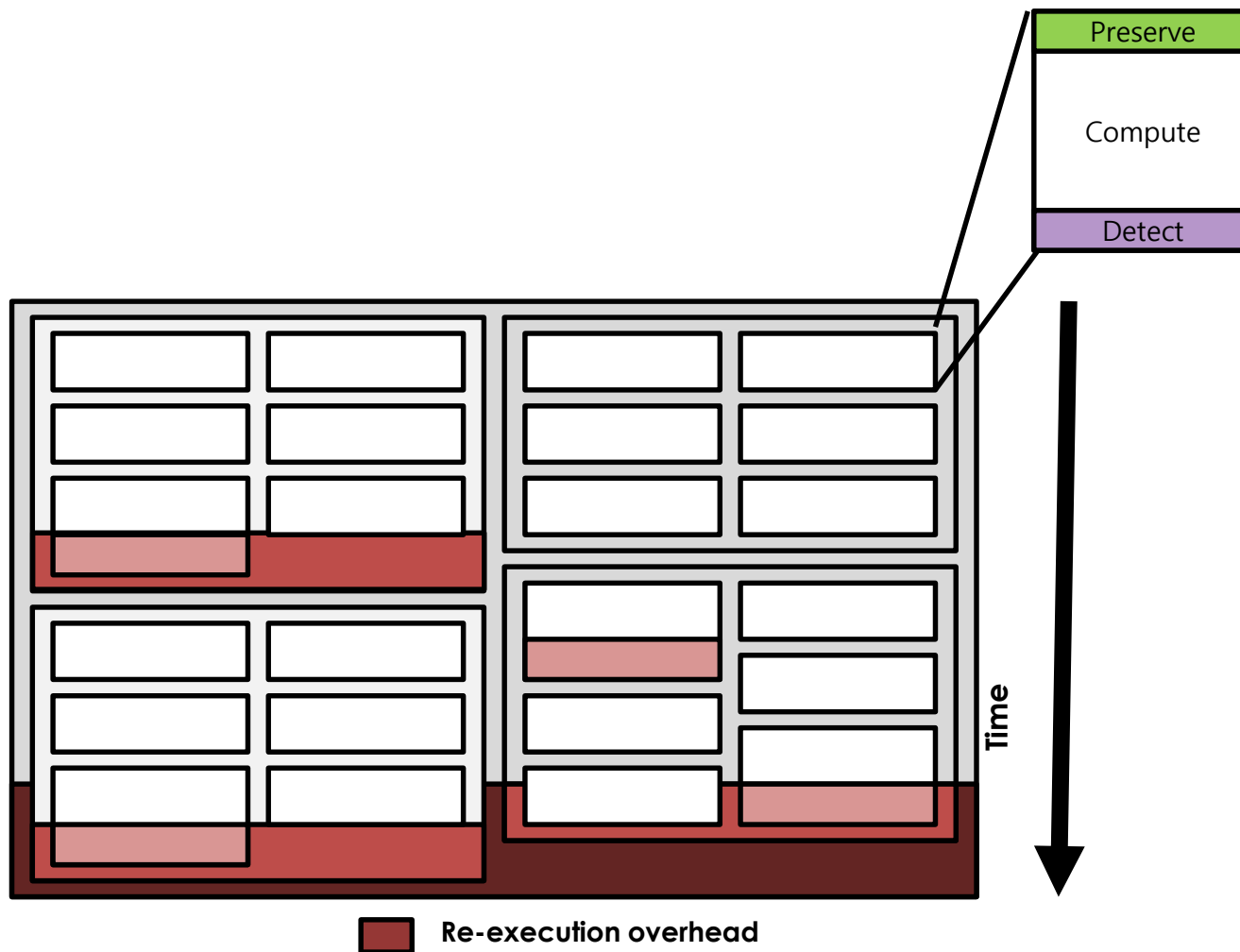
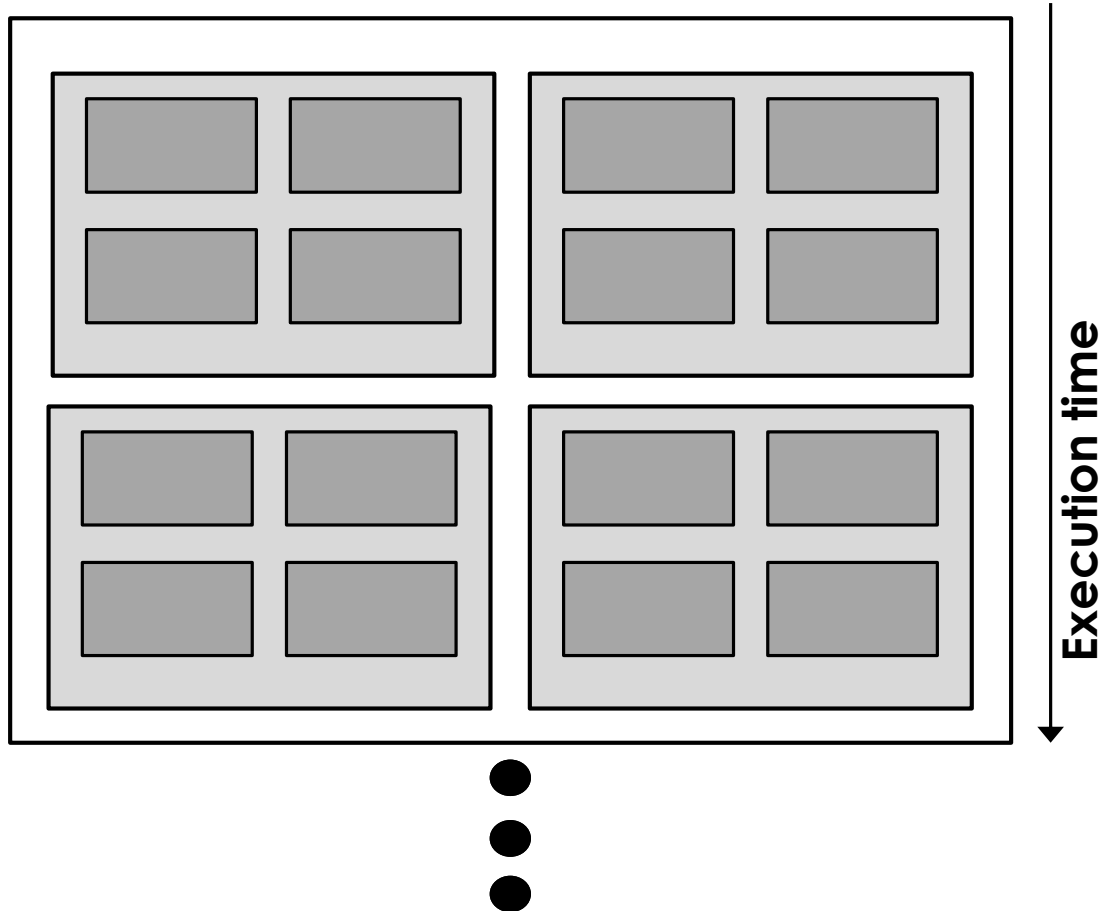**Local copy or regen**       **Sibling**       **Parent (unchanged)**

# Recovery:
## concurrent, hierarchical, adaptive



Re-execution overhead

# Analyzable

– Application model (tree)

– Overhead model

– Fault model

**Execution time**

# Analyzable

$$q[0, m, n] = (1 - p_c)^{m \cdot n}$$

- Probability that all child CDs experienced no failures

$$q[x, m, n] = \left( \sum_{i=0}^{x} \left( \binom{i+m-1}{i} p_c^i (1 - p_c)^m \right) \right)^n$$

- Probability that all child CDs experienced at most x failures in the m serial CDs
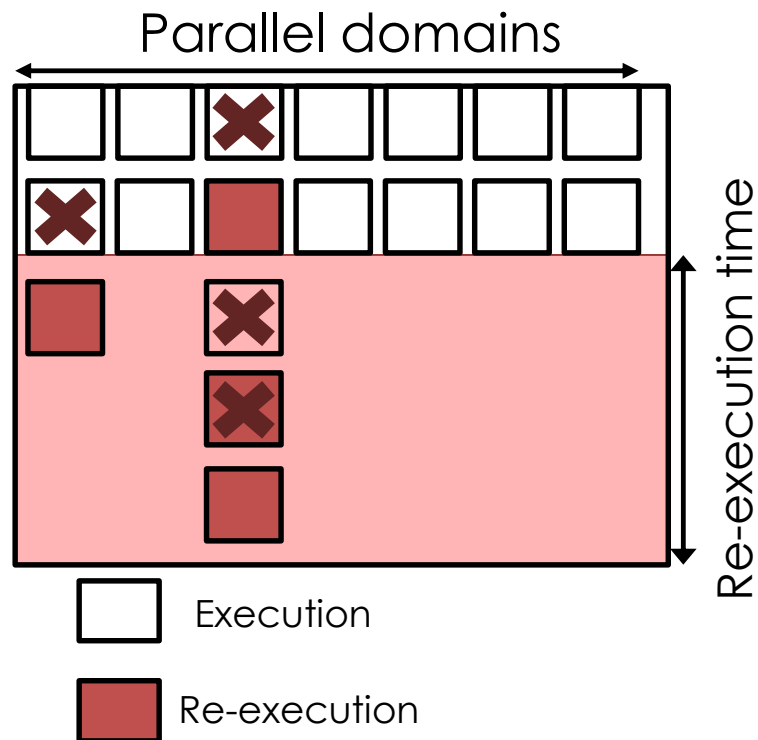
$$d[0, m, n] = q[0, m, n]$$

- Probability that all child CDs experienced exactly 0 failures

$$d[x, m, n] = q[x, m, n] - q[x - 1, m, n]$$

- Probability that sibling with the most failure experiences exactly x failures

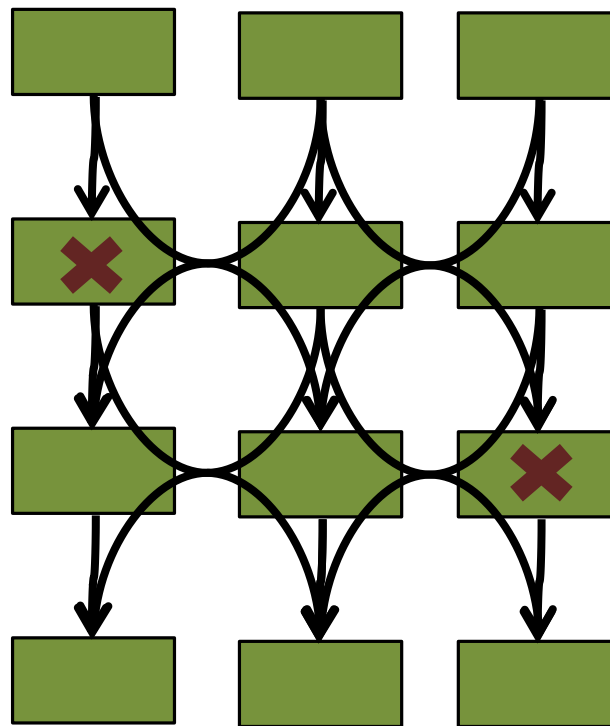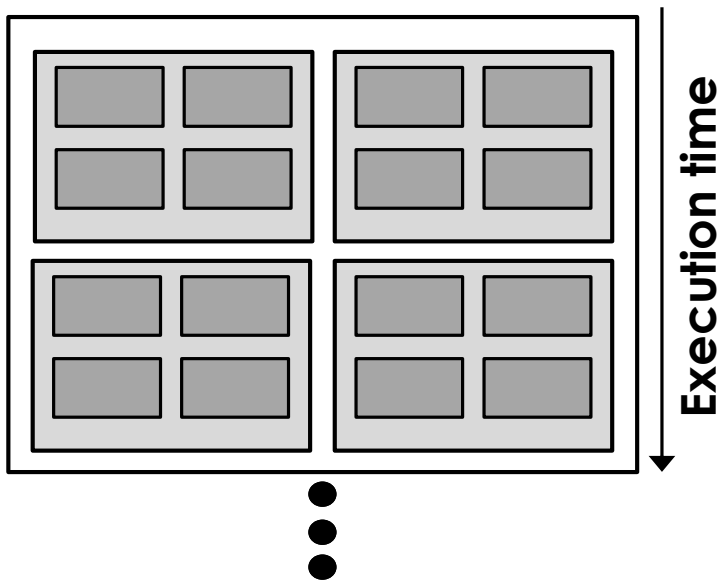$$T_{parent}[x, m, n] = \sum_{i=0}^{\infty} (i + m) T_c d[i, m, n]$$

Parallel domains

Re-execution time

Execution

Re-execution

Heterogeneous CDs
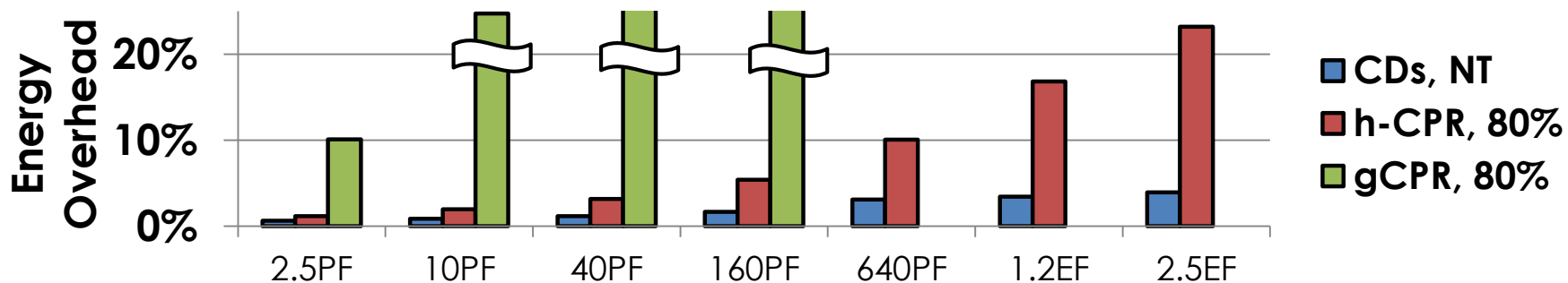
Asymmetric preservation and restoration

# Analyzable?

**Execution time**

# CDs are scalable

# CDs are proportional and efficient

Architecture goals:

– Balance possible and practical

– Enable generality

– Don't hurt the common case

It's all about the algorithm

# Architecture goals:
– Balance possible and practical
– Enable generality
– Don't hurt the common case

## Architecture so far:
– Proportionality
  • Adaptivity
  • SW-HW co-tuning **+ analyzability**
  • Heterogeneity
– Locality
– Parallelism
– Hierarchy

**Arithmetic**
**Control**
**Memory**
**Reliability**
**Programming**

# Exascale computers are **lunacy**

# Exascale computers are lunacy
## **science**

Harbingers of things to come

Discovery awaits

Enable and promote **open innovation**

Math +
   Systems +
      Engineering +
         Technology