# Redwood: Flexible and Portable Heterogeneous Tree Traversal Workloads

**Yanwen Xu**

Ang Li

Tyler Sorensen

University of California, Santa Cruz

Princeton University

University of California, Santa Cruz

# Highlights

Many applications in edge computing can utilize tree data structures to accelerate their workloads
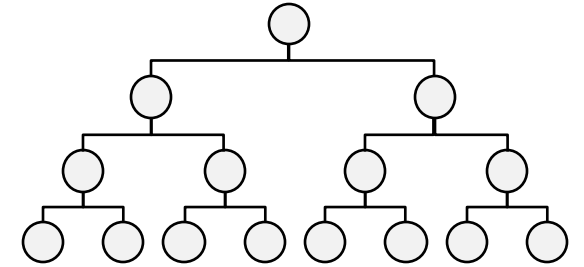
**Developed** a framework, ***Redwood***, for writing a special class of tree algorithms, which we call *traverse-compute*

**Evaluated** 5 shared memory heterogeneous systems using ***Grove***, a suite of 9 heterogeneous traverse-compute workloads

**Achieved** speedups up to
- 8.12x on commodity systems
- 13.53x on academic system

**Insight:** Traverse-compute workload has natural heterogeneous decompositions on modern shared memory system-on-chips

# Motivation: Accelerating Computations at Edge

Edge computing are getting popular …

But they has **constraints**
- ◦ *e.g., energy or latency requirement*

Application of edge computing
- ◦ *Surveillance cameras*
- ◦ *Autonomous vehicles*
- ◦ *Mobile gaming*

# Motivation: Accelerating Computations at Edge

Edge computing are getting popular …

But they has **constraints**
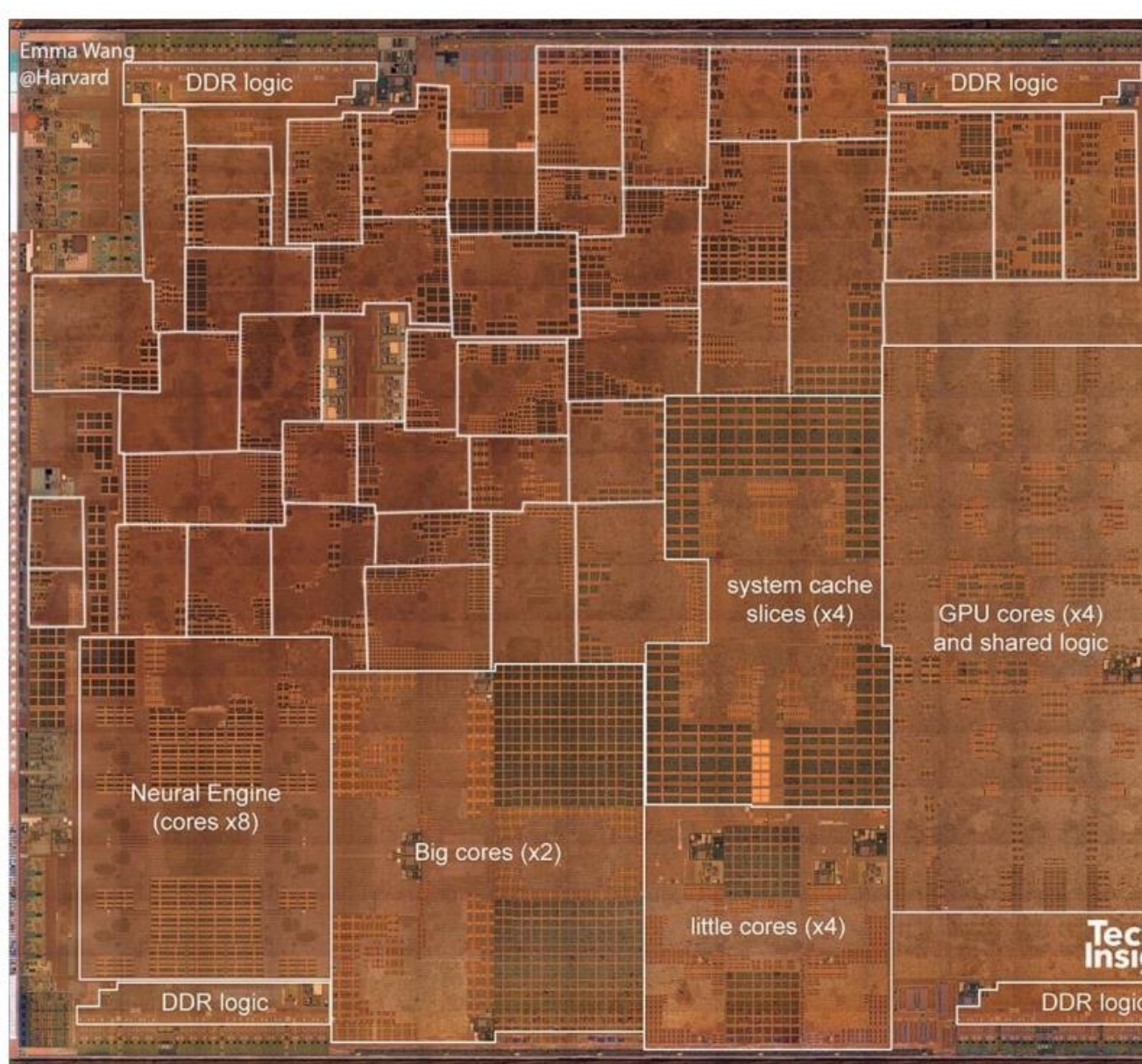- *e.g., energy or latency requirement*

Application of edge computing
- *Surveillance cameras*
- *Autonomous vehicles*
- *Mobile gaming*

**Modern edge devices are becoming increasingly heterogeneous**
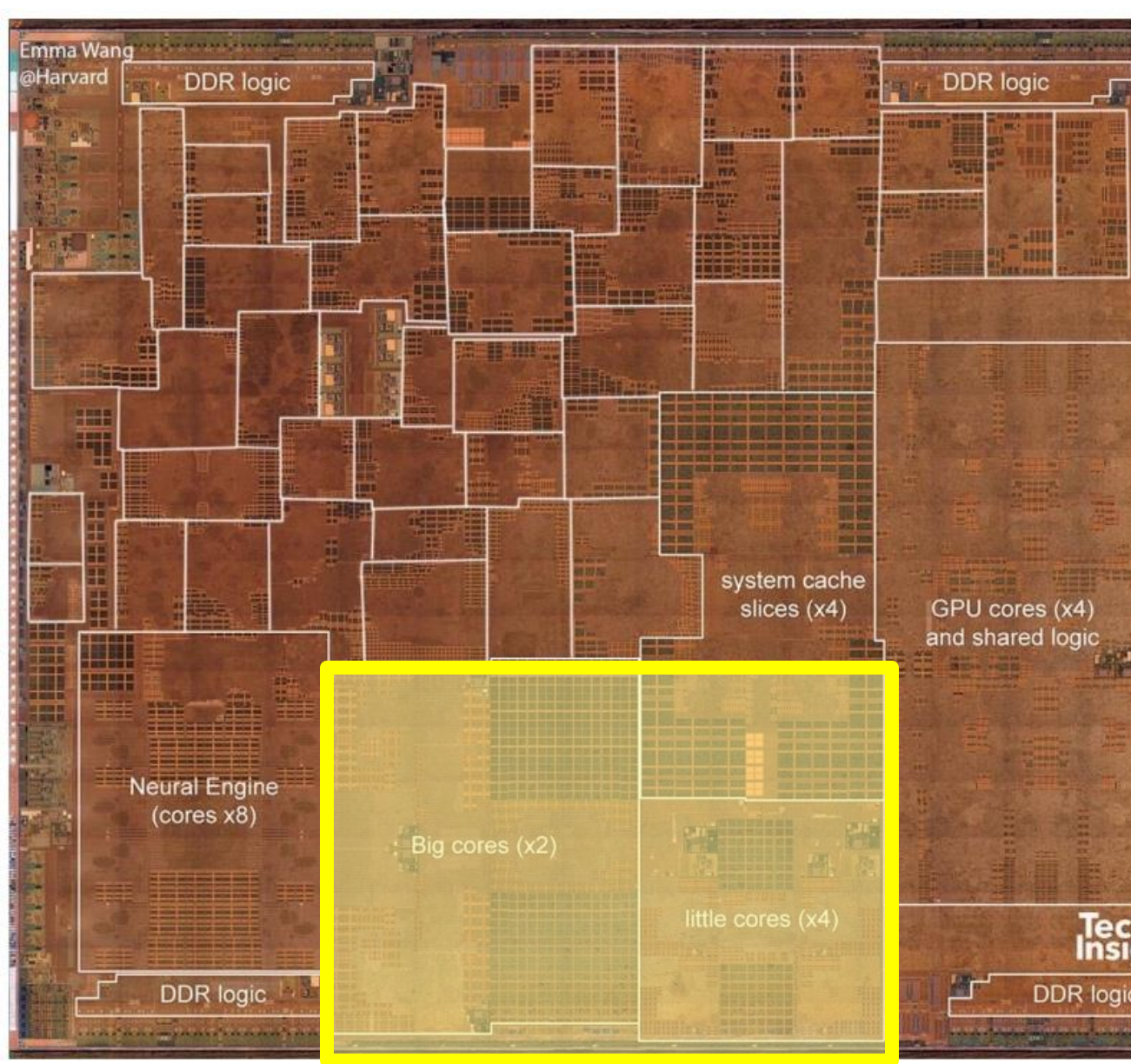- w/ specialized *Processing Units* (PUs)

**We need to efficiently utilize these available system resources**
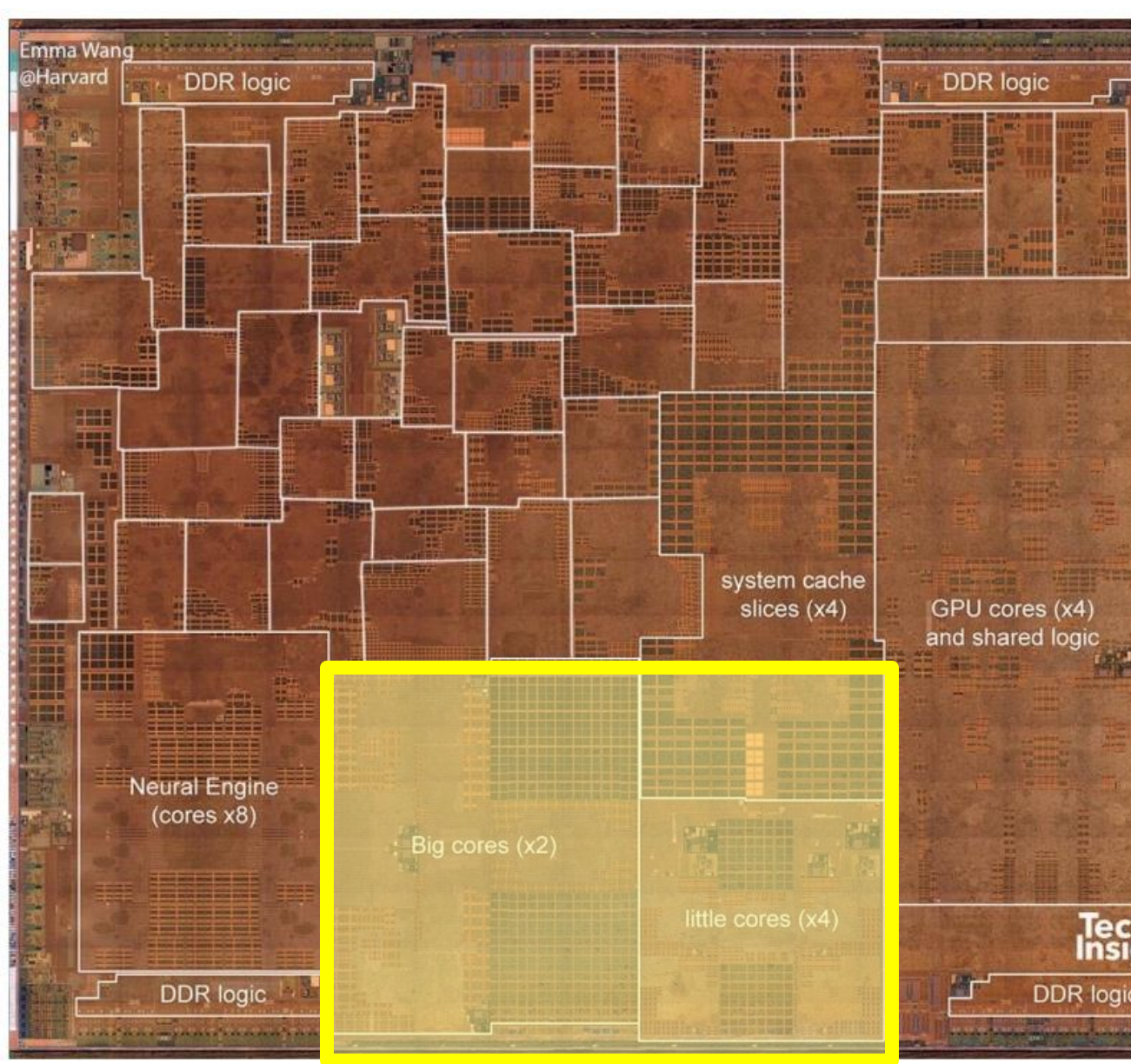
# What do we mean by resources?

*From David Brooks lab at Harvard:*
*https://vlsiarch.eecs.harvard.edu/research/accelerators/die-photo-analysis*

5

From David Brooks lab at Harvard:
https://vlsiarch.eecs.harvard.edu/research/accelerators/die-photo-analysis

# What do we mean by resources?

- E.g., less than 20% of the die area of an iPhone contains the CPU

- The rest contains specialize *Programmable Accelerating PUs* (PAPU)
  - e.g., integrated GPUs, FPGAs
  - Interconnected to a shared memory hierarchy

- *Shared Memory Heterogeneous System* (SMHS) enables efficient communication between PUs

# What do we mean by resources?



- E.g., less than 20% of the die area of an iPhone contains the CPU

- The rest contains specialize *Programmable Accelerating PUs* (PAPU)
  - e.g., integrated GPUs, FPGAs
  - Interconnected to a shared memory hierarchy

- *Shared Memory Heterogeneous System* (SMHS) enables efficient communication between PUs

*From David Brooks lab at Harvard:*
*https://vlsiarch.eecs.harvard.edu/research/accelerators/die-pho*

**How does workloads utilize each PU?**

# Processing Units (PU) Characteristics

| CPU | Programmable Accelerating PUs (PAPU) | |
|---|---|---|
| | **GPU** | **FPGA** |
| Features: High-performance cores, reorder buffer, load store queue, … | *Features:* SIMT (*Single Instruction, Multiple Threads*) execution, coalesced memory access | Features: Specialized tasks, Pipeline parallelism |
| + **Latency optimized**<br><br>- **Limited throughput** | <br><br>+ **Throughput optimized**<br>- **Warp Divergence** | + **Close to ASIC performance**<br>- **Orders-of-magnitude harder to program** |

Good for **irregular** programs

Good for accelerating **compute-intense** programs

# Trees on the edge



*Root Node*

R

- Edge applications process a large amount of data
- They can utilize tree structures and traversals to perform edge tasks
  - E.g., *in classifications and security*



Input data



Spatial Partition

# Trees on the edge



*Root Node*

R

- Edge applications process a large amount of data
- They can utilize tree structures and traversals to perform edge tasks
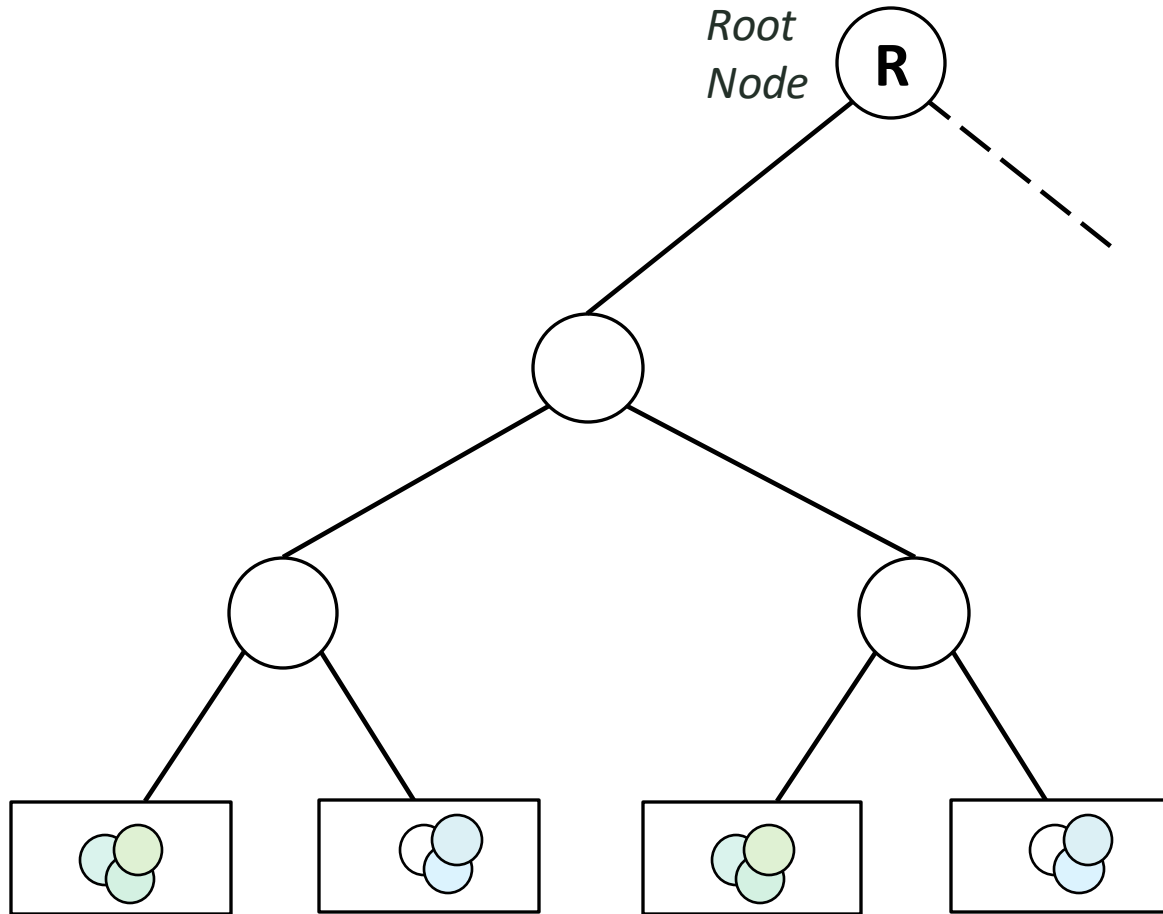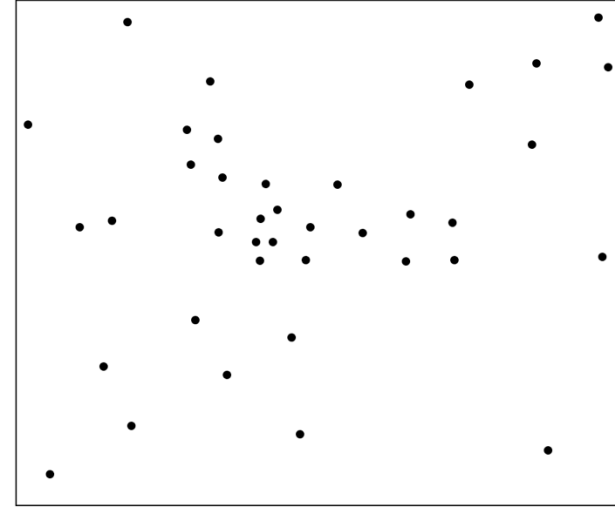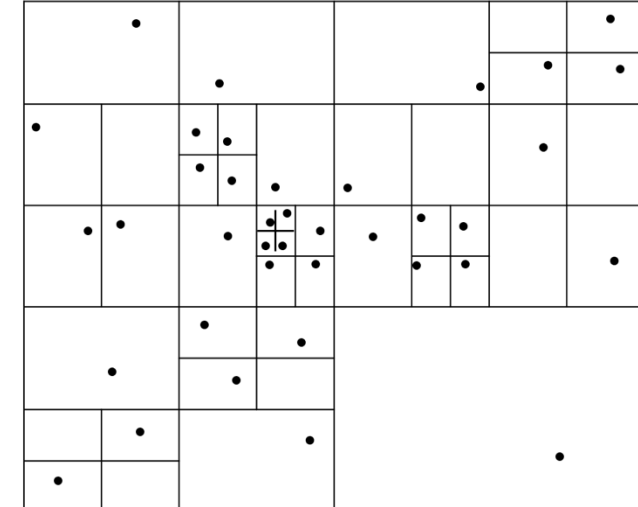  - E.g., *in classifications and security*
- The dataset are organized into a hierarchical tree structure, allowing data to be **efficiently** searched from $O(n)$ to $O(\log n)$
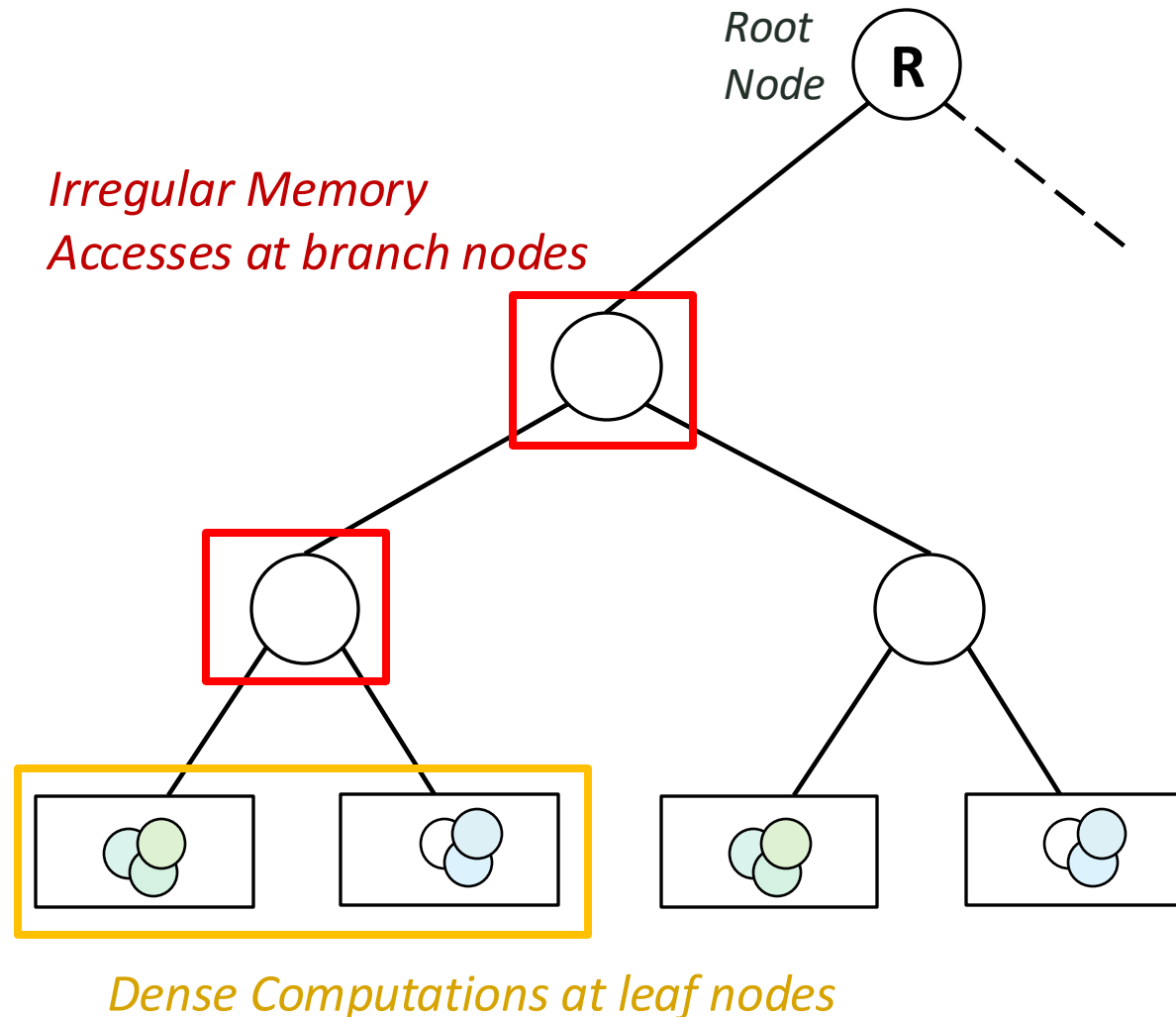


Input data



Spatial Partition

# Traverse-Compute Workloads



*Root Node*

*Irregular Memory Accesses at branch nodes*

*Dense Computations at leaf nodes*

- Repeatedly traversing a sparse tree structure
- Each traversal consists of
  - Indirect memory loads at branch nodes (Red box)
  - Dense data to be processed at leaf nodes visited (**Orange box**)
    - Computing pairwise interactions (e.g., Euclidean distance)
    - Reductions (e.g., sum, min)
- Example workloads:
  - Barnes-hut Algorithm (octree)
  - Nearest Neighbor Search (k-dimensional tree)
  - Ray Tracing (bounding volume hierarchy)

# Decomposing Traverse Compute Workloads

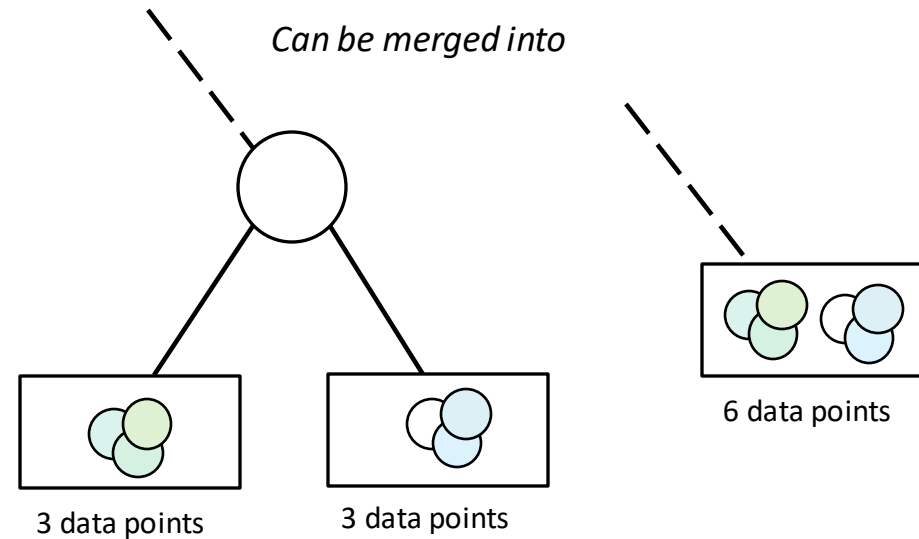**CPUs** are good at handling dynamic control flows and tolerating indirect memory loads

**PAPUs** are good at accelerating dense, compute-intense operations

**A natural Heterogeneous Approach is to**

CPU traverses the tree

Computation at leaf nodes is offloaded to accelerator

# Accelerating Traverse-compute workloads on SMHSs

- The tree can be *parameterized* by how many data points exist on the leaf nodes.

*Can be merged into*

6 data points

3 data points          3 data points

**Larger node sizes tradeoffs**

+ Larger chunk of contiguous data

+ Less irregular accesses in tree traversals

- Potentially unneeded computation

**Heterogeneous Approach**

CPU traverses the tree

Computation is offloaded to accelerator

# Accelerating Traverse-compute workloads on SMHSs

- The tree can be *parameterized* by how many data points exist on the leaf nodes.

*Can be merged into*

6 data points

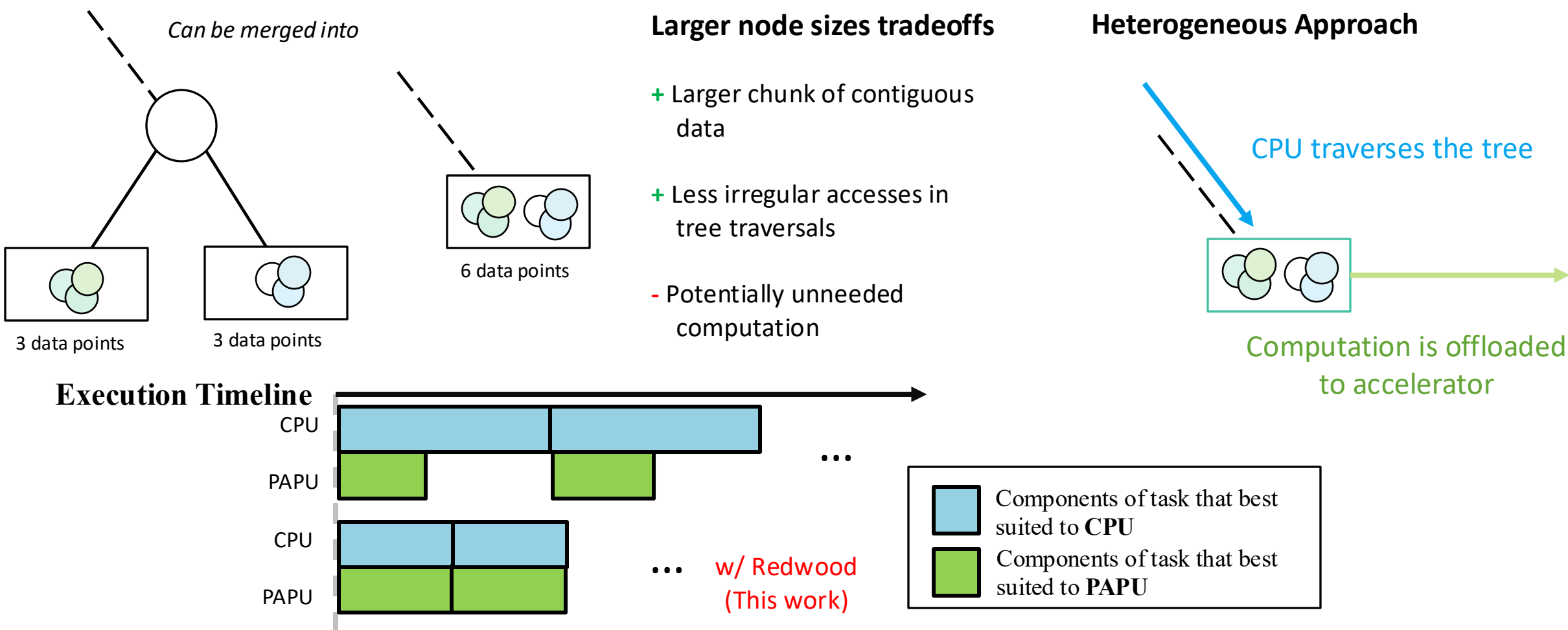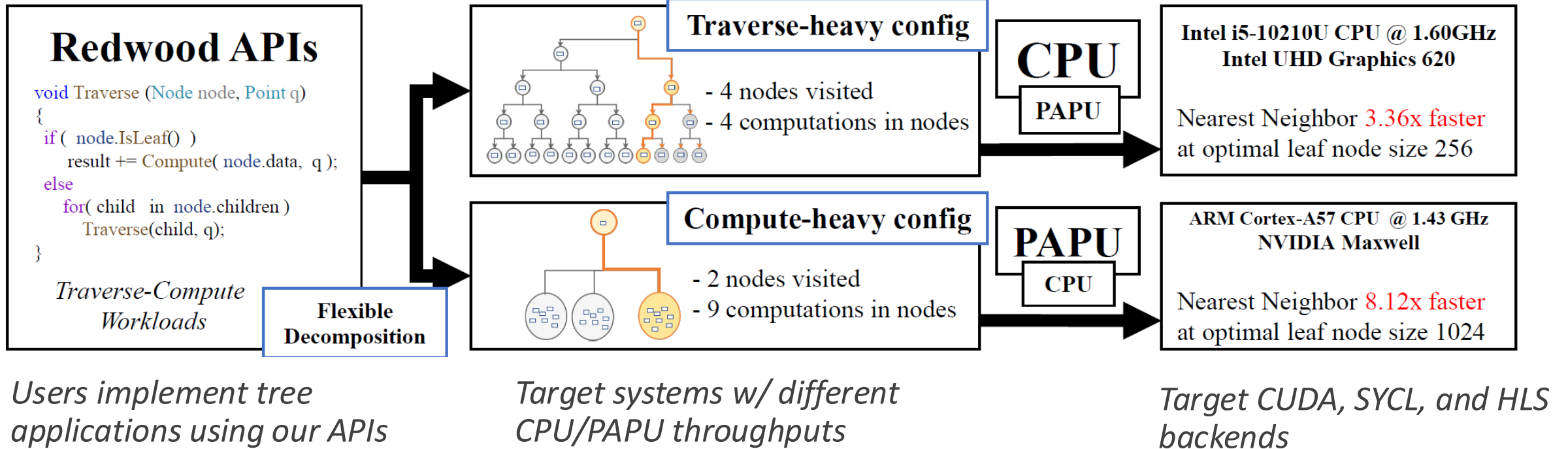3 data points        3 data points

## Larger node sizes tradeoffs

**+** Larger chunk of contiguous data

**+** Less irregular accesses in tree traversals

**-** Potentially unneeded computation

## Heterogeneous Approach

CPU traverses the tree

Computation is offloaded to accelerator

**Execution Timeline**

CPU

PAPU

...

CPU

PAPU

...  w/ Redwood (This work)

| | Components of task that best suited to **CPU** |
| --- | --- |
| | Components of task that best suited to **PAPU** |

**Flexible & Specialize Heterogeneous Decomposition**

# This Work: Redwood Overview



**Redwood APIs**

```
void Traverse (Node node, Point q)
{
  if ( node.IsLeaf() )
    result += Compute( node.data, q );
  else
    for( child  in  node.children )
      Traverse(child, q);
}
```

*Traverse-Compute Workloads*

**Flexible Decomposition**

**Traverse-heavy config**

- 4 nodes visited
- 4 computations in nodes

**Compute-heavy config**

- 2 nodes visited
- 9 computations in nodes

**CPU**
PAPU

**PAPU**
CPU

Intel i5-10210U CPU @ 1.60GHz
Intel UHD Graphics 620

Nearest Neighbor 3.36x faster
at optimal leaf node size 256

ARM Cortex-A57 CPU @ 1.43 GHz
NVIDIA Maxwell

Nearest Neighbor 8.12x faster
at optimal leaf node size 1024

*Users implement tree applications using our APIs*

*Target systems w/ different CPU/PAPU throughputs*

*Target CUDA, SYCL, and HLS backends*

Tyler 98.3%

KNN based Facial recognition

*Intel SoCs*

*Nvidia SoCs*

# Redwood: APIs and Data Structures

CPU Sequential Code (NN)

```python
tree = KDTree()
min_dist = 99999.9999
def traverse(node, q):
    if is_leaf(node):
```
```python
        # Reduce Leaf Node
        for i in range(node.leaf_size):
            kernel_func(q, node.data[i])
```
```python
    else:
        dist = compute_dist(q, node.data[0])
        min_dist = min(min_dist, dist)
        traverse(node.leaf_child)
        if check_other_side(dist):
            traverse(node.right_child)
```
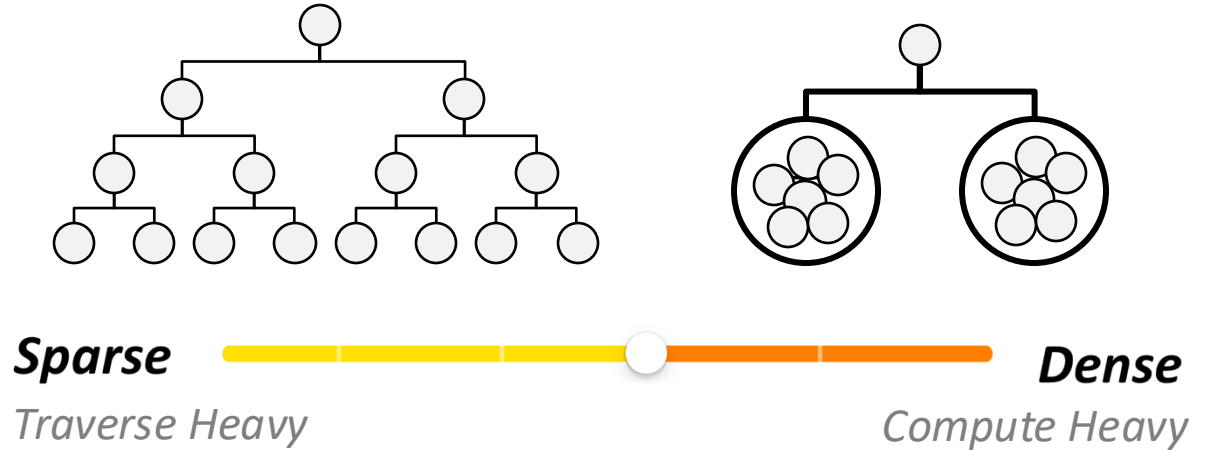
*Implemented using:*

w/ Redwood API

```python
tree = KDTree(leaf_size=32)

redwood_set_query(q)

def traverse(node, q):
```
```python
    if is_leaf(node):
        redwood_compute_leaf(node.data())
```
```python
    else:
        dist = compute_dist(q, node.data[0])
        min_dist = min(min_dist, dist)
        traverse(node.leaf_child)
        if check_other_side(dist):
            traverse(node.right_child)
```

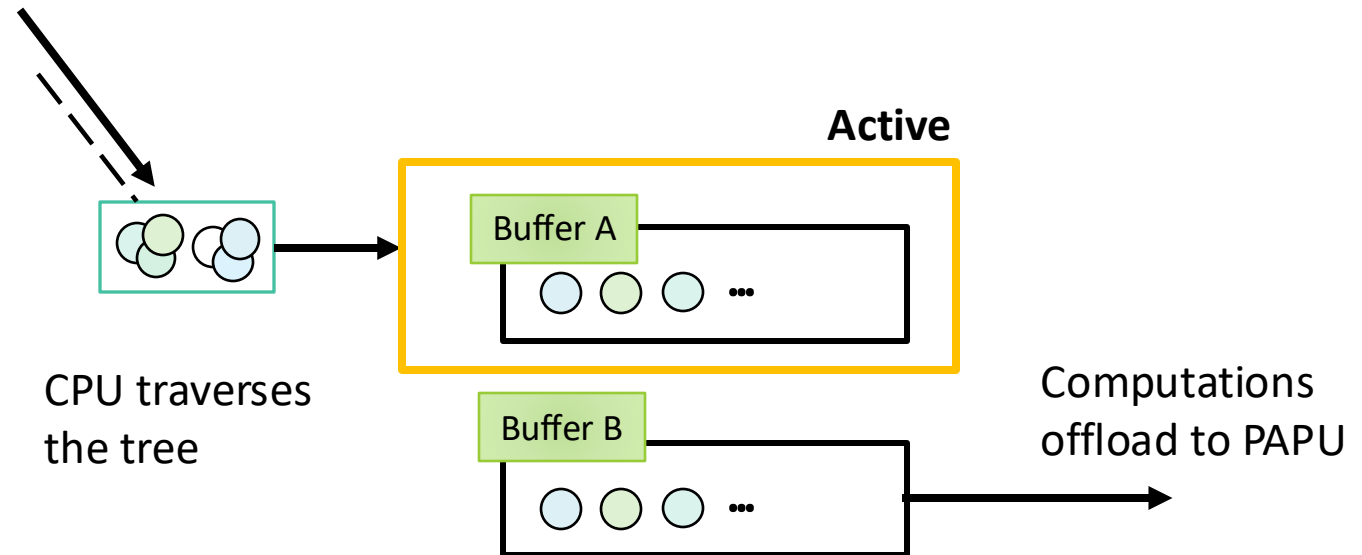# Redwood Heterogenous Optimizations

## Flexible Leaf Size Configuration

- Adapt to various heterogeneous systems with different relative throughput between the CPU and the PAPU

**Sparse** — *Traverse Heavy*

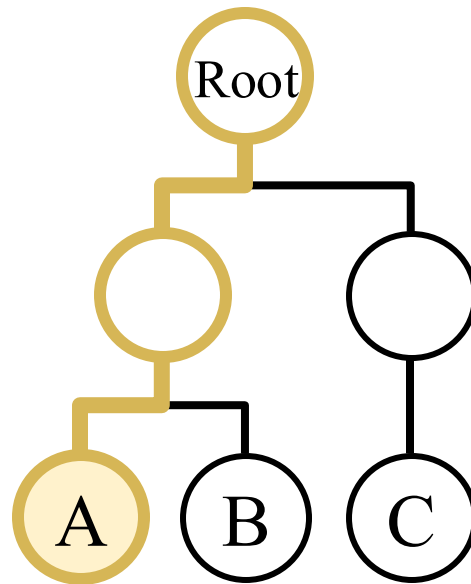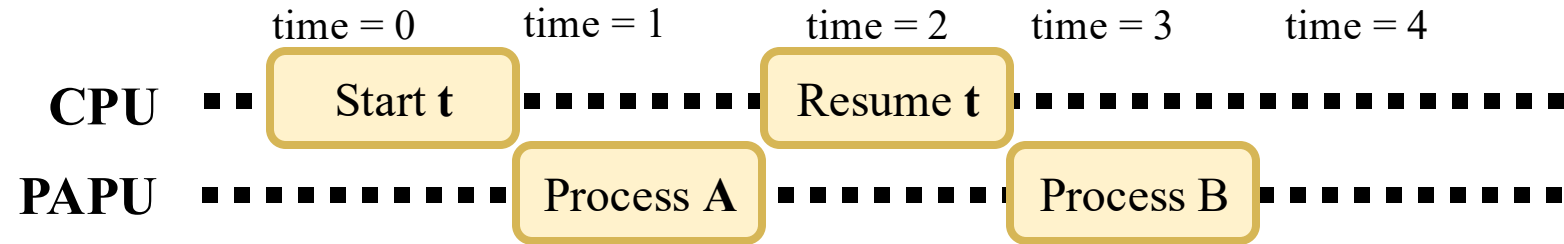**Dense** — *Compute Heavy*

## Automatic Batching & Ping-pong Buffering

- Fuse multiple small computations to a larger GPU kernel

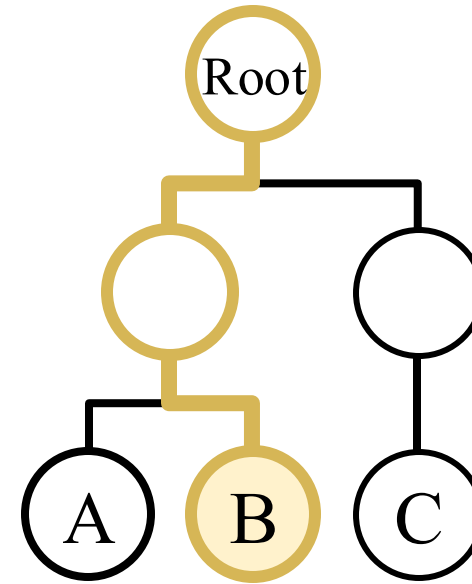- Overcome overheads related to GPU kernel launching & synchronization

CPU traverses the tree

**Active**

Buffer A

Buffer B

Computations offload to PAPU

# Redwood Heterogenous Optimizations

**Traverser Runtime**

- Allow a single CPU thread to execute many traversals concurrently to avoid stalling when a traversal depends on a PAPU accelerated value
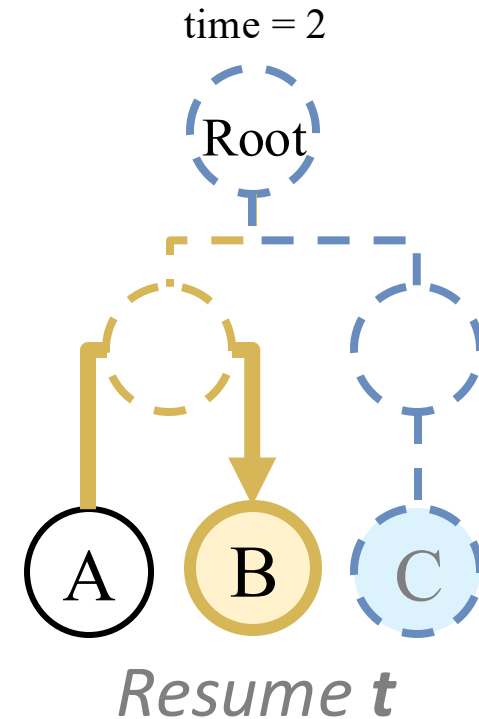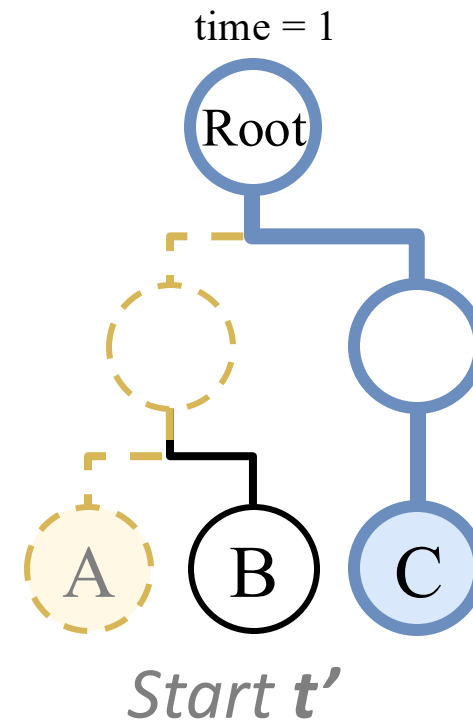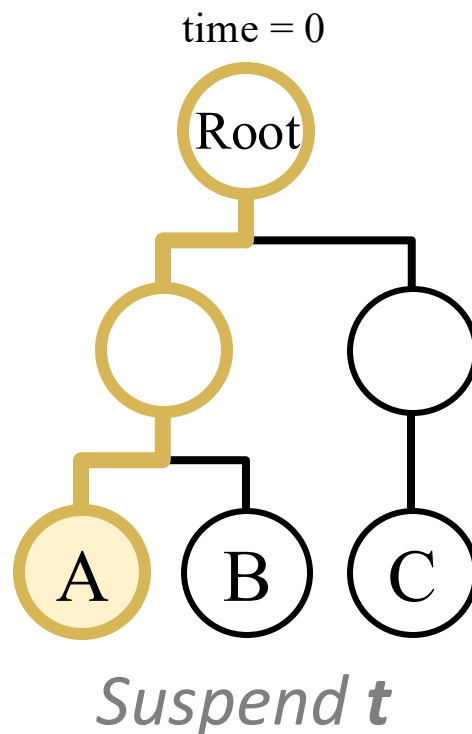


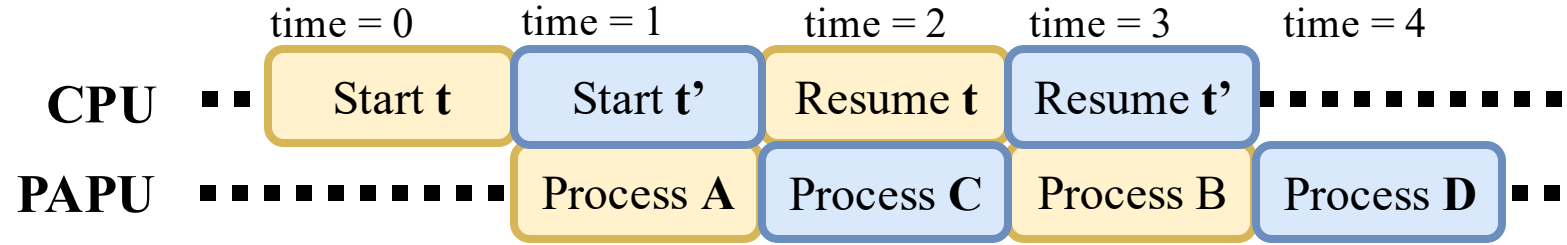*Start t*                    *continue when result ready*

# Redwood Heterogenous Optimizations

## **Traverser Runtime**

- Allow a single CPU thread to execute many traversals concurrently to avoid stalling when a traversal depends on a PAPU accelerated value

- *Light weight Coroutine*
  - **Suspend**
  - **Resume**

# **Grove** Benchmark Suite for SMHS

Grove contains **9** traverse-compute workloads

Can be found in many applications
- *Facial recognition*
- *Anomaly detection*
- *Outlier detection*
- *Particle simulation*

Tree Structures
- Octree/quadtree
- k-d tree

Three Algorithms
- Barnes Hut
- Nearest Neighbor
- k Nearest Neighbor

Computation Patterns
- Aggregation *(sum)*
- Reduction *(e.g., min)*
- Sorting

Various Distance Metrics
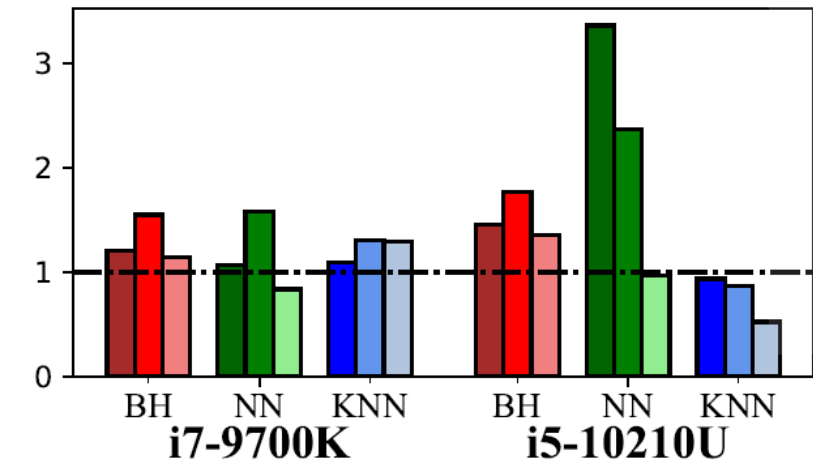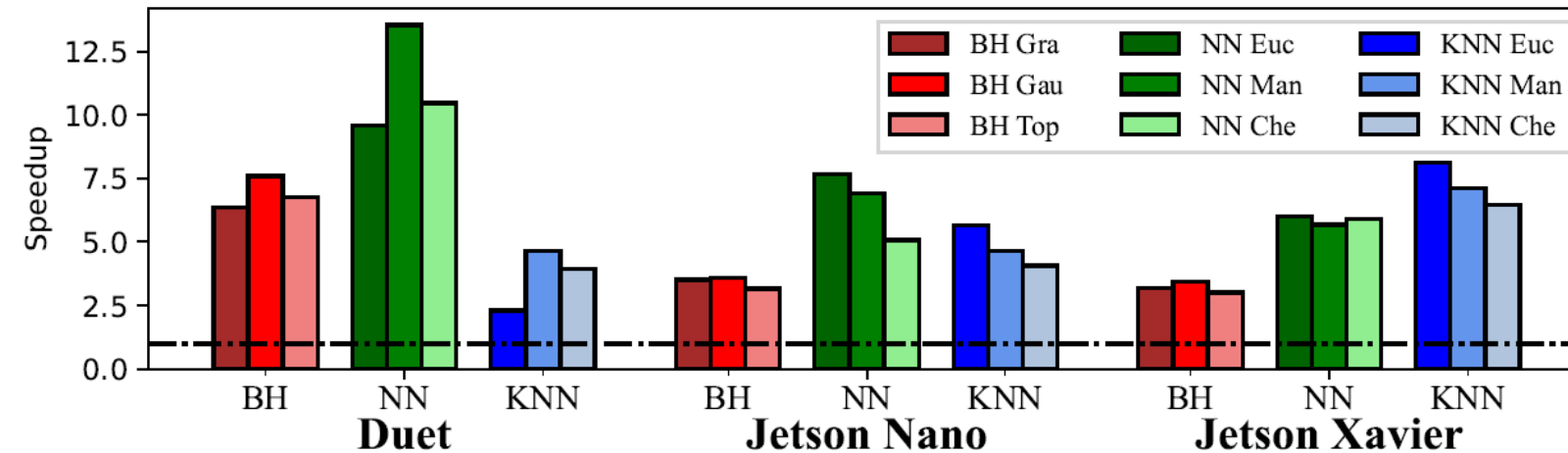- Euclidean
- Manhattan
- Chebyshev

# Platforms Evaluated

| Platform | Backend | CPU | CPU Frequency | Accelerator | Accelerator frequency |
|----------|---------|-----|---------------|-------------|----------------------|
| **NVIDIA Jetson Nano** | CUDA | ARM Cortex-A57 | 1.43 GHz | 128-core Maxwell | 921 Mhz |
| **NVIDIA Jetson Xavier NX** | CUDA | Carmel ARMv8.2 | 1.5 GHz | 384-core Volta | 1.21 GHz |
| **Intel i7-9700K** | SYCL | i7-9700K | 3.60 GHz/ 4.20 GHz | Intel UHD Graphics, 24 EUs | 350 MHz/ 1.20 GHz |
| **Intel i5-10210U** | SYCL | i5-10210U | 1.60 GHz/ 4.90 GHz | Intel UHD Graphics, 24 EUs | 300 MHz/ 1.10 GHz |
| **Duet** [1] *(simulated in **gem5**)* | HLS | RISC-V TimingSimpleCPU | 1.5 GHz | Duet eFPGA | 333MHz |

Powerful NVIDIA GPUs
+
Little ARM CPUs

Little Intel GPUs
+
Powerful Intel CPUs

Tiny in-order CPUs
w/ Efficient eFPGA
integration

*[1] Li, Ang, August Ning, and David Wentzlaff. HPCA'23*

# Grove Results Overview



Speedups of the best heterogeneous configuration vs. the best homogeneous configuration of Grove.

# Grove Results Overview



## Highlights

**Duet**
highest 13.53x
geomean 6.43x

*With minimal kernel launching overhead*

**Xavier**
highest 8.12x
geomean 5.13x

*Powerful NVIDIA GPUs But has small CPUs*

**Nano**
highest 6.9x
geomean 4.71x

**Intels**
highest 3.36x
**even has slow downs**

*Intel OoO is already fast, but has small GPUs*

# Grove Results Overview



## Highlights

**Duet**
highest 13.53x
geomean 6.43x

**Xavier**
highest 8.12x
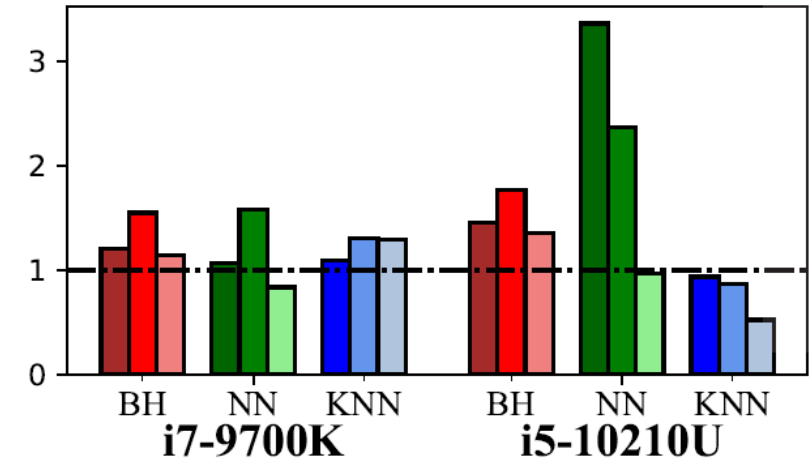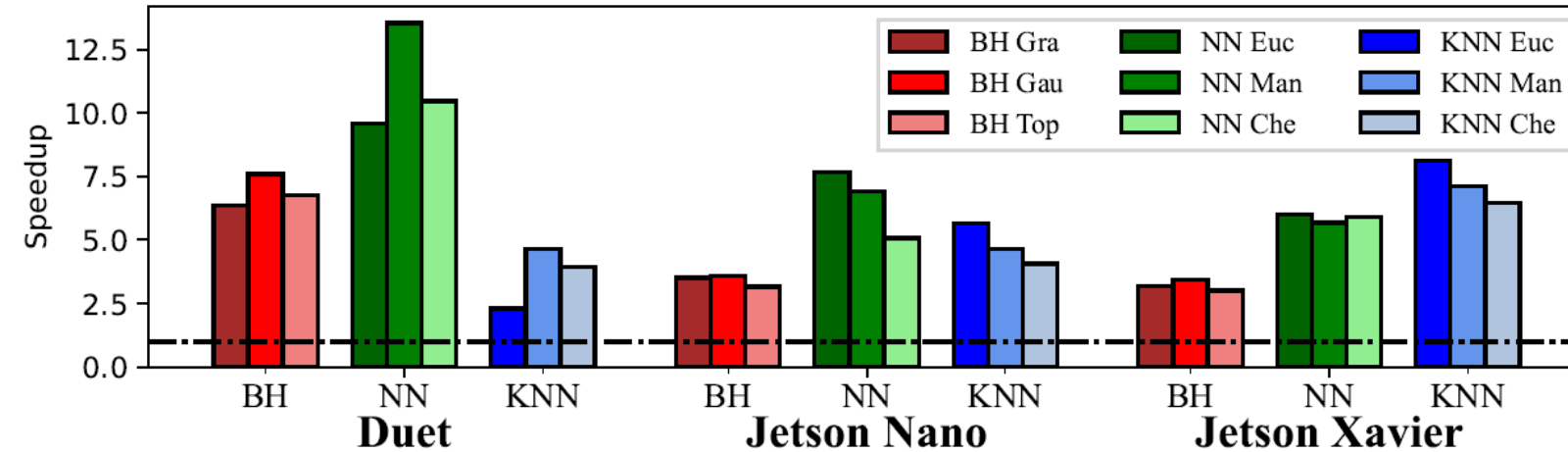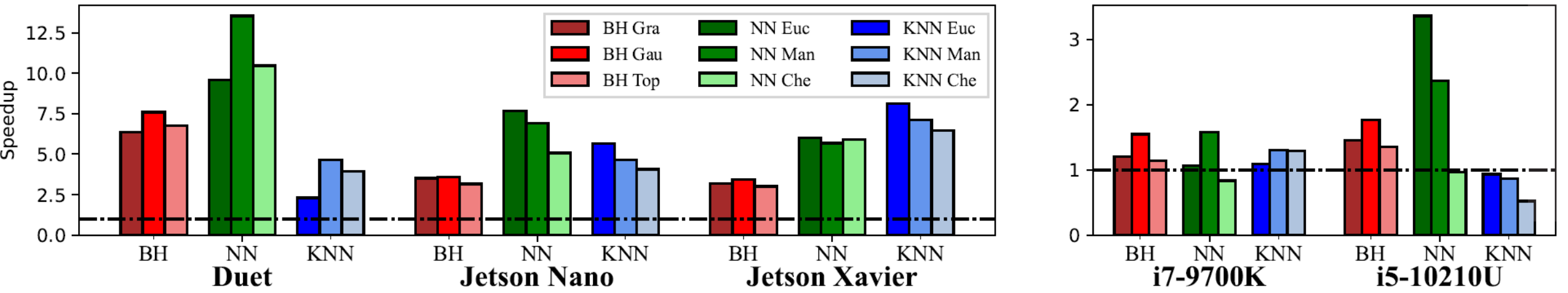geomean 5.13x

**Nano**
highest 6.9x
geomean 4.71x

**Intels**
highest 3.36x
**even has slow downs**

Original Duet Paper[1] reported **~3x** speedup in BH algorithm. With Redwood's flexible task decomposition, we achieved **~6x** speedup

*Intel OoO is already fast, but has small GPUs*

[1] Li, Ang, August Ning, and David Wentzlaff. HPCA'23

# Grove Results Overview



With Powerful Accelerators, more leaf node computations can be offloaded to GPU/FPGAs.

| | BH Gra | BH Gau | BH Top | NN Euc | NN Man | NN Che | KNN Euc | KNN Man | KNN Che | Average | Ratio |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Jetson Nano | 256/8 | 256/8 | 128/8 | 512/256 | 512/256 | 512/256 | 256/128 | 256/128 | 256/256 | 327/115 | 2.26 |
| Jetson Xavier | 128/16 | 128/8 | 128/8 | 1024/128 | 1024/256 | 512/64 | 512/64 | 512/64 | 512/64 | 498/75 | 6.67 |
| i7-9700k | 64/8 | 128/8 | 128/8 | 256/64 | 256/64 | 128/64 | 64/64 | 64/64 | 64/64 | 128/45 | 2.82 |
| i5-10210U | 64/8 | 64/8 | 64/8 | 256/64 | 256/64 | 256/64 | 32/64 | 32/64 | 32/64 | 117/45 | 2.59 |
| Duet | 512/4 | 512/2 | 512/4 | 512/32 | 512/32 | 256/16 | 128/16 | 128/32 | 128/32 | 355/19 | 18.82 |
| **Average** | 205/9 | 218/7 | 192/7 | 512/109 | 512/134 | 333/93 | 198/67 | 198/70 | 198/96 | 285/66 | 4.33 |

Optimal leaf node size for heterogeneous (left) vs homogeneous (right) for each workload and platform
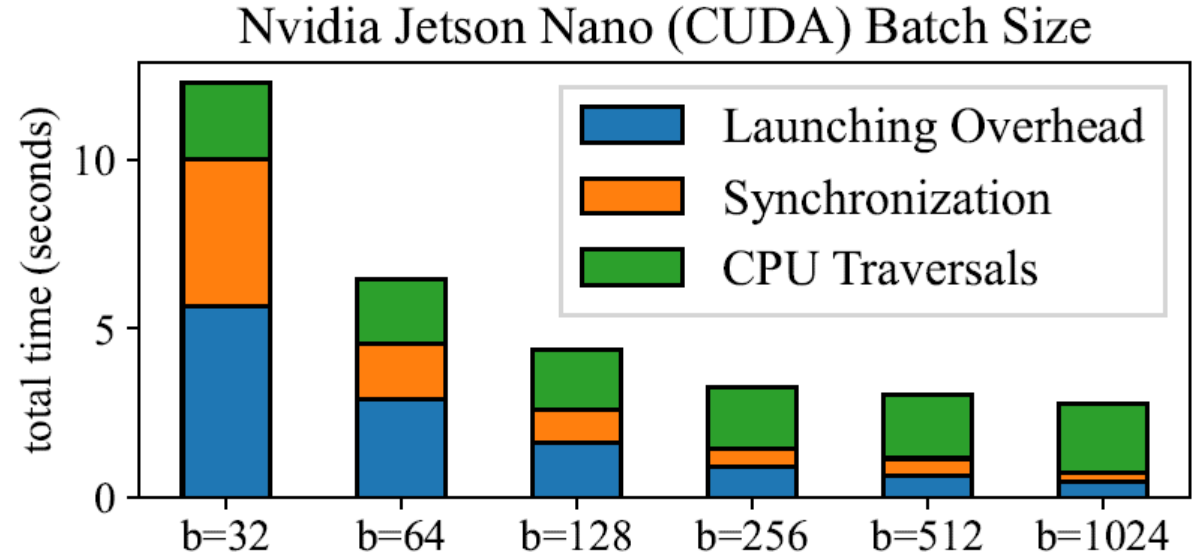
# We have several insights in the paper

**Kernel Submission Cost**

- Traverse-compute applications **frequently** invoke small kernels
- Useful works are shown in Green
- Orange/Blue are overheads
- Low-cost kernel submission is important for accelerating applications on edge devices

**Future work**

- We are exploring more efficient CPU-PAPU communication methods



Batching multiple GPU kernels into a single/larger kernel helps amortizing kernel launching overhead

# Conclusion

✓We identify a class of applications, *traverse-compute* applications, that are ideal for shared memory heterogeneous systems

✓We present <u>Redwood</u>: a framework for writing heterogeneous traverse-compute workloads

✓Using Redwood, we implemented <u>Grove</u>, a benchmark suite contains 9 traverse-compute applications

✓Finally evaluated 5 systems, achieving up to **13.53×** speedups (geomean of **3.01×**)



*UC Santa Cruz Redwood Grove*

**Team**
**Yanwen Xu** yxu83@ucsc.edu
Ang Li angl@princton.edu
Tyler Sorensen tyler.sorensen@ucsc.edu

**Open-Source Repo**
Redwood & Grove at
https://github.com/xuyanwen2012/redwood-rt