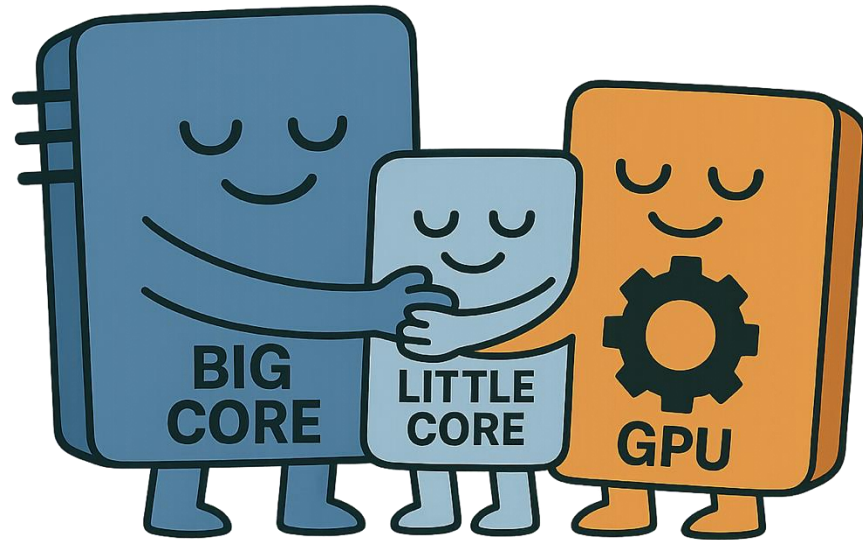


# BetterTogether



## An Interference-Aware Framework for Fine-grained Software Pipelining on Heterogeneous SoCs

*Yanwen Xu, Rithik Sharma, Zheyuan Chen,  
Shaan Mistry, Tyler Sorensen*

UC SANTA CRUZ

**BE** Baskin  
Engineering

 Microsoft

# Motivation: Accelerating Computation at **Edge**

- Lower latency
- Energy efficiency
- Privacy benefits



*NVIDIA Jetson Thor*

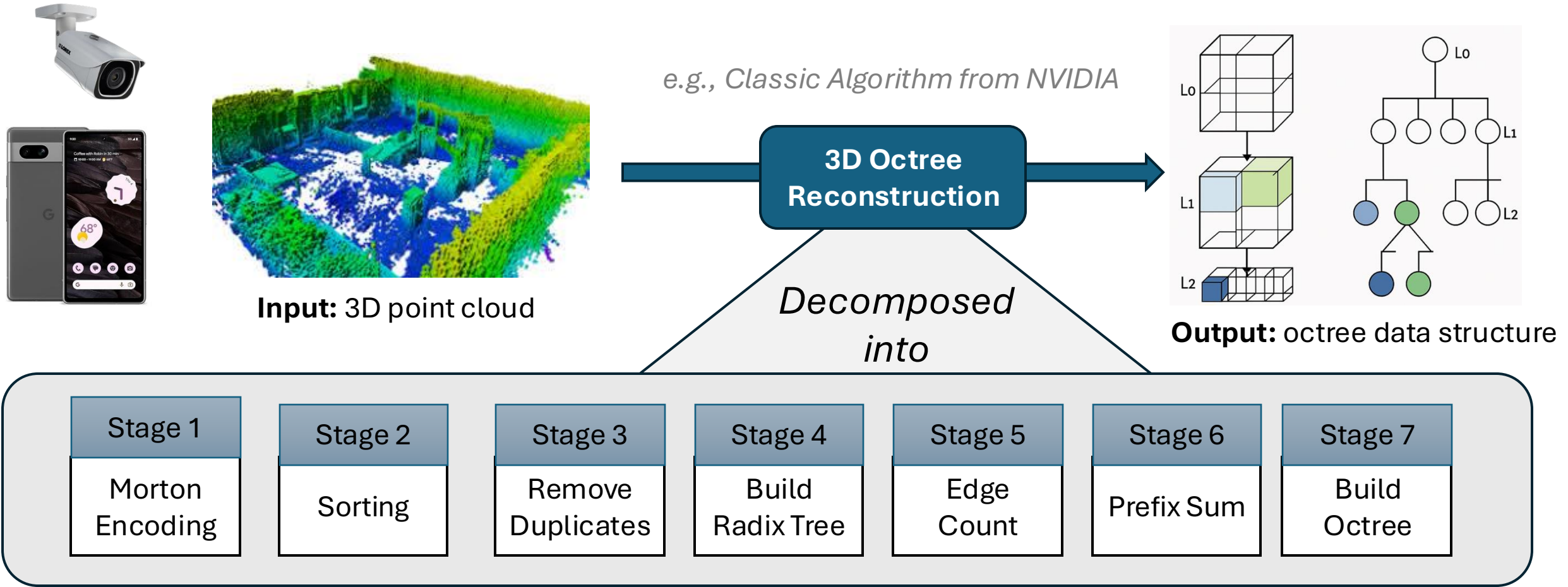


*Google Pixel*



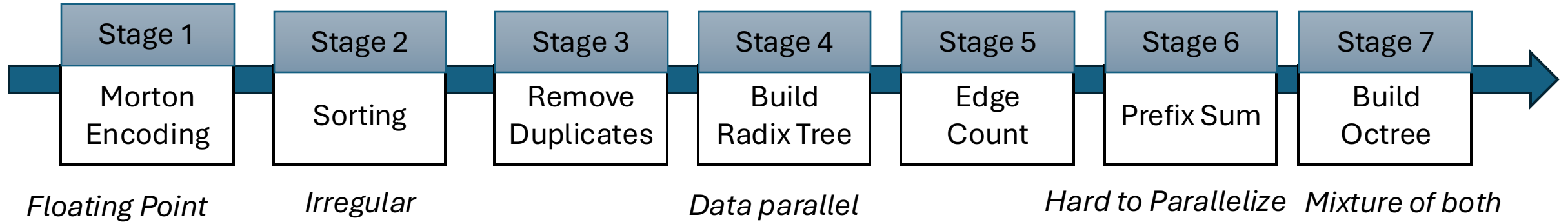
*NVIDIA Jetson Orin Nano*

# Motivation: Accelerating Computation at Edge

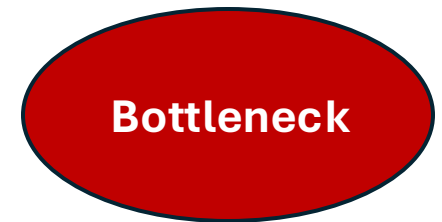
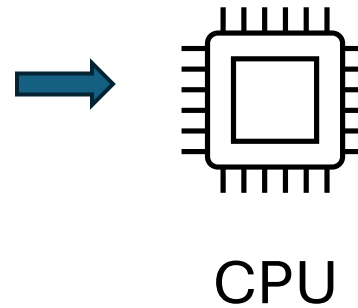


Stages depends on data from previous stages

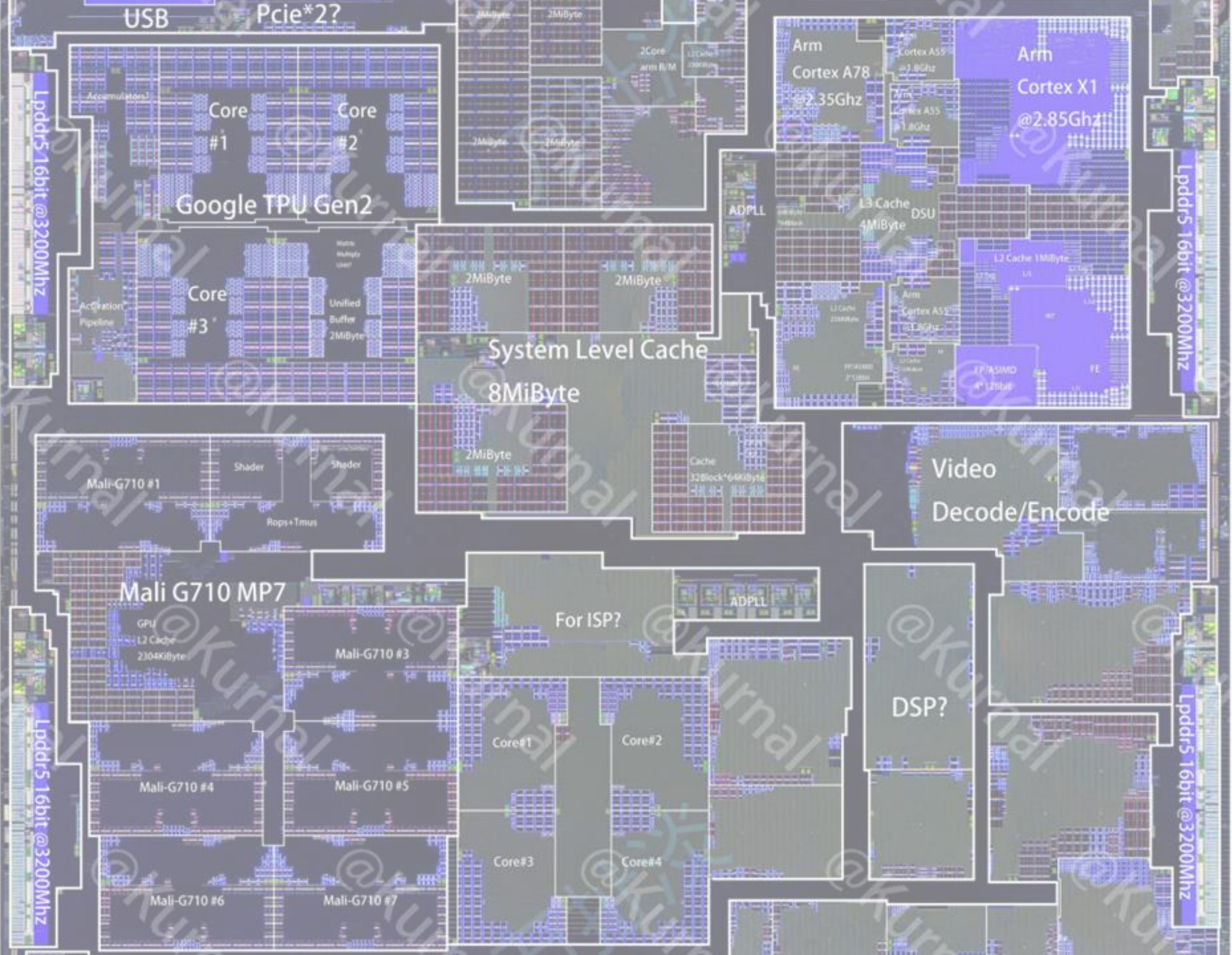
# Motivation: Accelerating Computation at Edge



```
// -----  
// Encode  
// -----  
[[nodiscard]] constexpr uint32_t morton3D_SplitBy3bits(const uint32_t a) {  
    auto x = static_cast<uint32_t>(a) & 0x0000003ff;  
    x = (x | x << 16) & 0x300000ff;  
    x = (x | x << 8) & 0x0300f00f;  
    x = (x | x << 4) & 0x30c30c3;  
    x = (x | x << 2) & 0x9249249;  
    return x;  
}  
[[nodiscard]] constexpr uint32_t m3D_e_magicbits(const uint32_t x,  
    const uint32_t y,  
    const uint32_t z) {  
    return morton3D_SplitBy3bits(x) | (morton3D_SplitBy3bits(y) << 1) |  
        (morton3D_SplitBy3bits(z) << 2);  
}  
[[nodiscard]] constexpr uint32_t xyz_to_morton32(const glm::vec4 &xyz,  
    const float min_coord,  
    const float range) {  
    constexpr auto bit_scale = 1024;  
    const auto i = static_cast<uint32_t>((xyz.x - min_coord) / range * bit_scale);  
    const auto j = static_cast<uint32_t>((xyz.y - min_coord) / range * bit_scale);  
    const auto k = static_cast<uint32_t>((xyz.z - min_coord) / range * bit_scale);  
    return m3D_e_magicbits(i, j, k);  
}
```







Google Tensor G2





# Mali-G710 MP7

## 2 Cortex-X1 (P-cores)

## 2 Cortex-A78 (medium)

## 4 Cortex-A55 (E-cores)

## Modern SoCs integrate **diverse** Processing Units (PU)

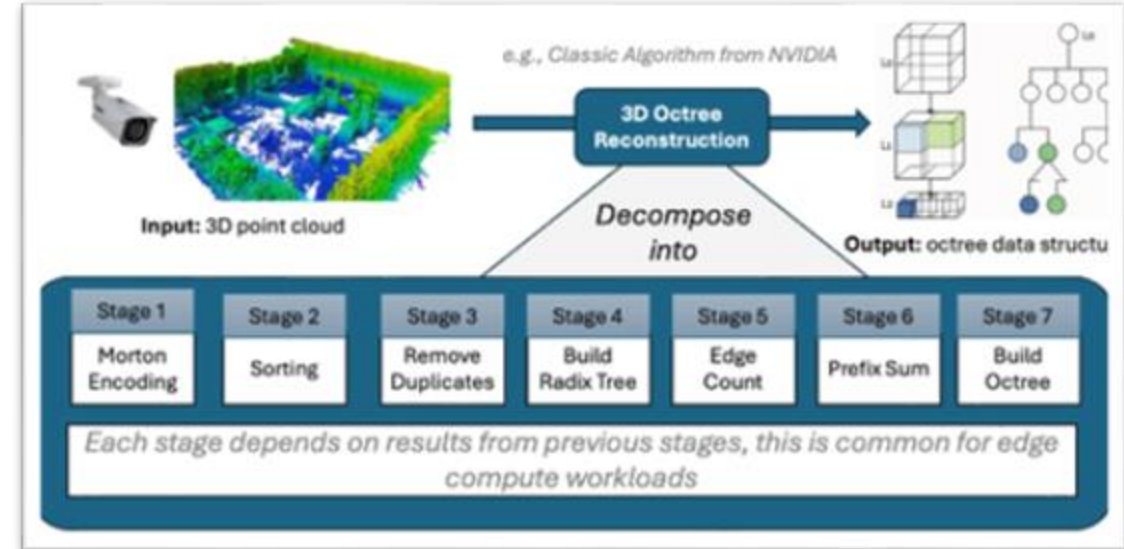
# PU Profiling

Mali-G710 MP7  
2 Cortex-X1  
2 Cortex-A78  
4 Cortex-A55



*All Available PUs on Pixel*

- Ran each stage on each available PUs



*All Stages*

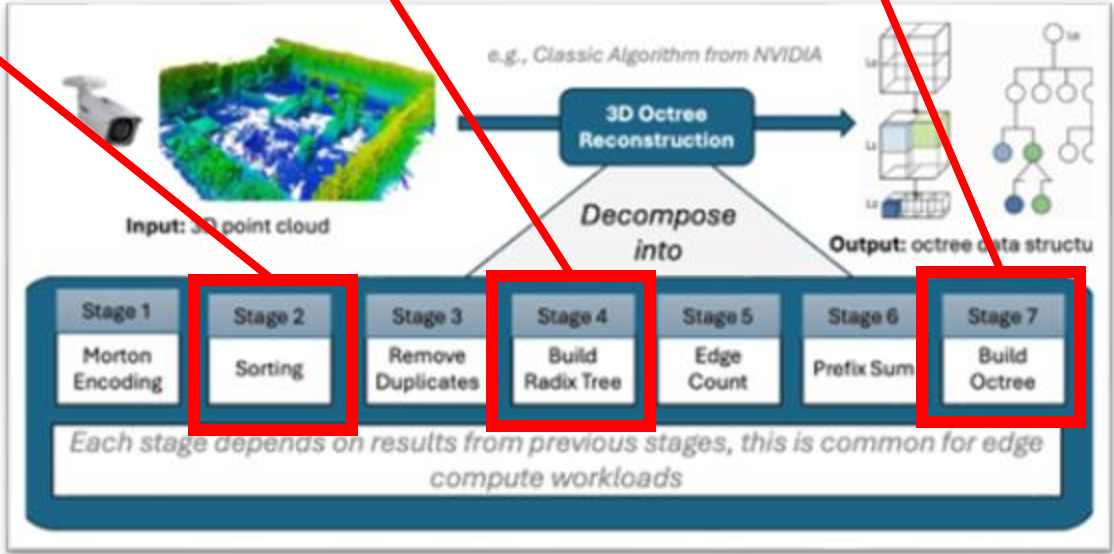
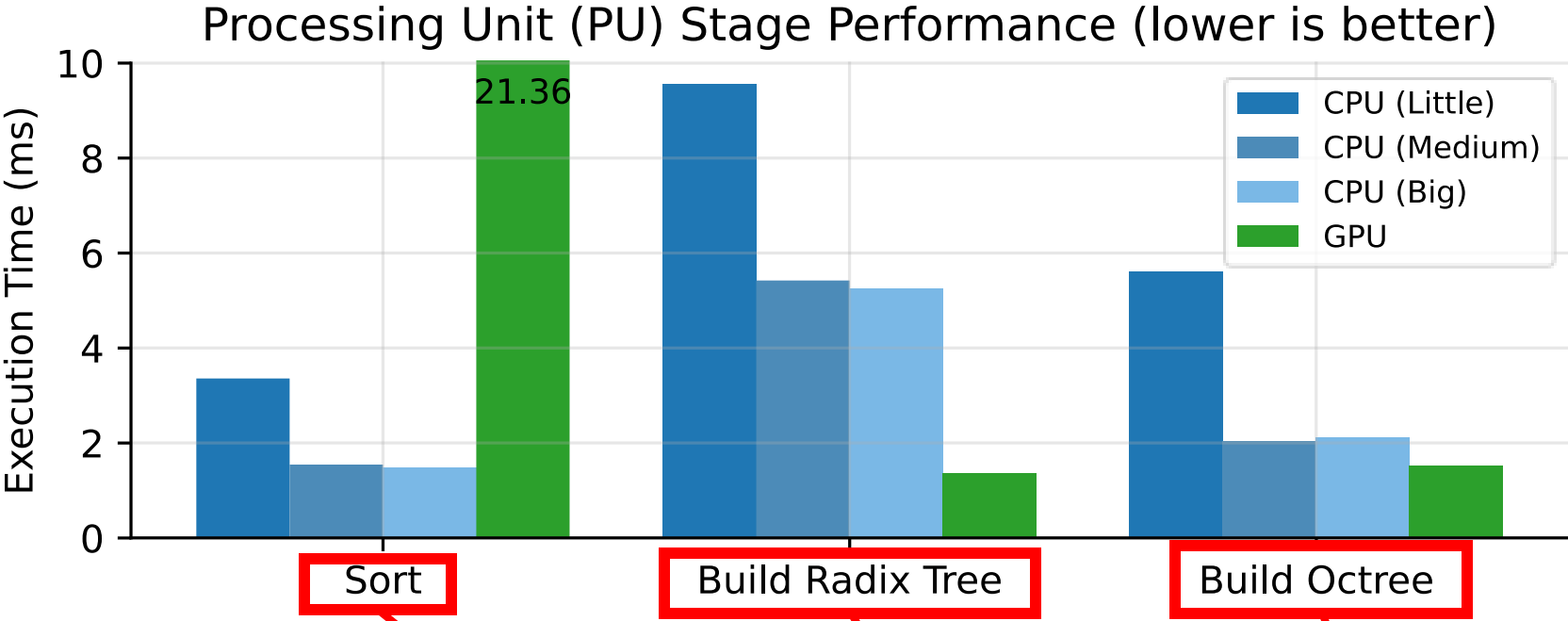
# PU Profiling

Mali-G710 MP7  
2 Cortex-X1  
2 Cortex-A78  
4 Cortex-A55



All Available PUs on Pixel

- Ran each stage on each available PUs



All Stages



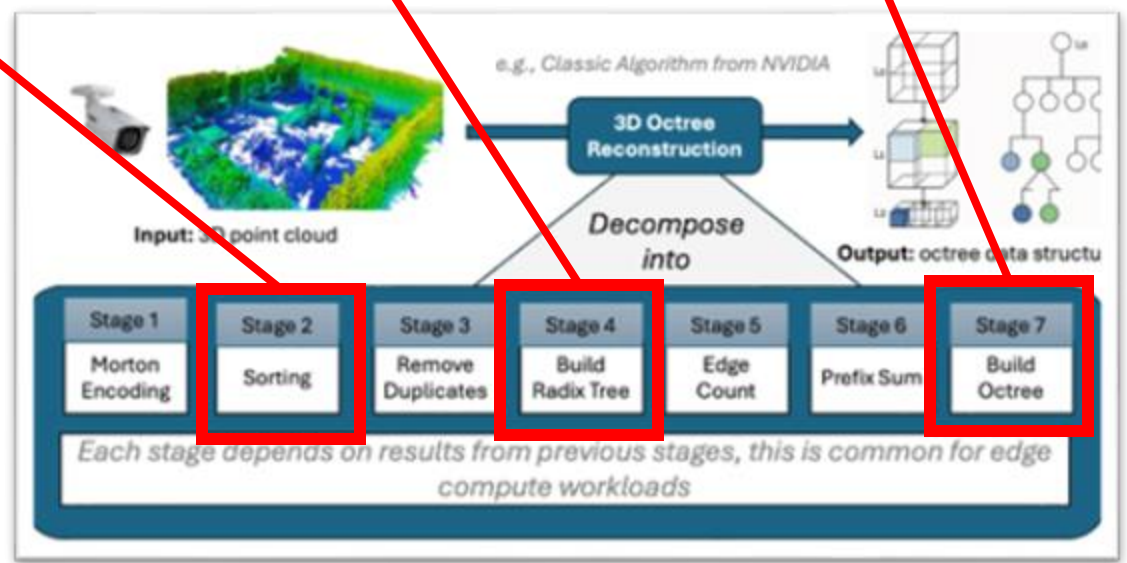
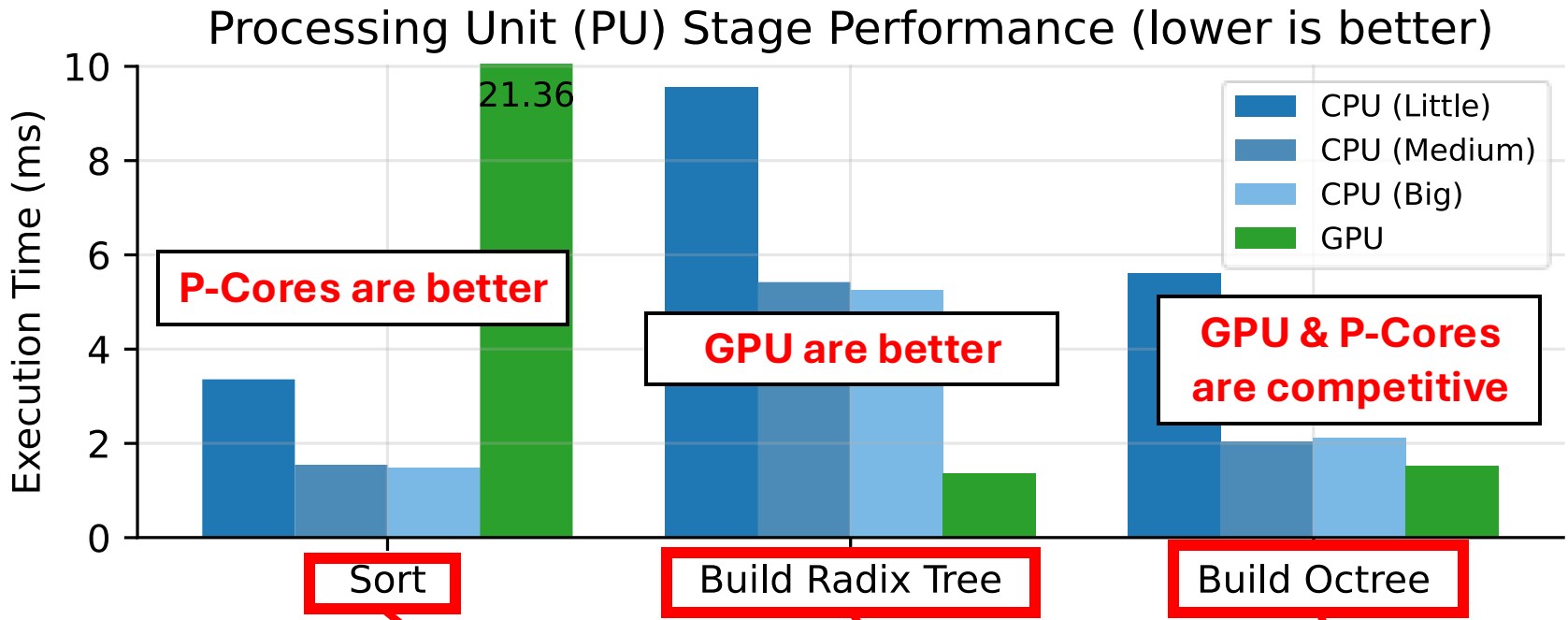
# PU Profiling

Mali-G710 MP7  
2 Cortex-X1  
2 Cortex-A78  
4 Cortex-A55



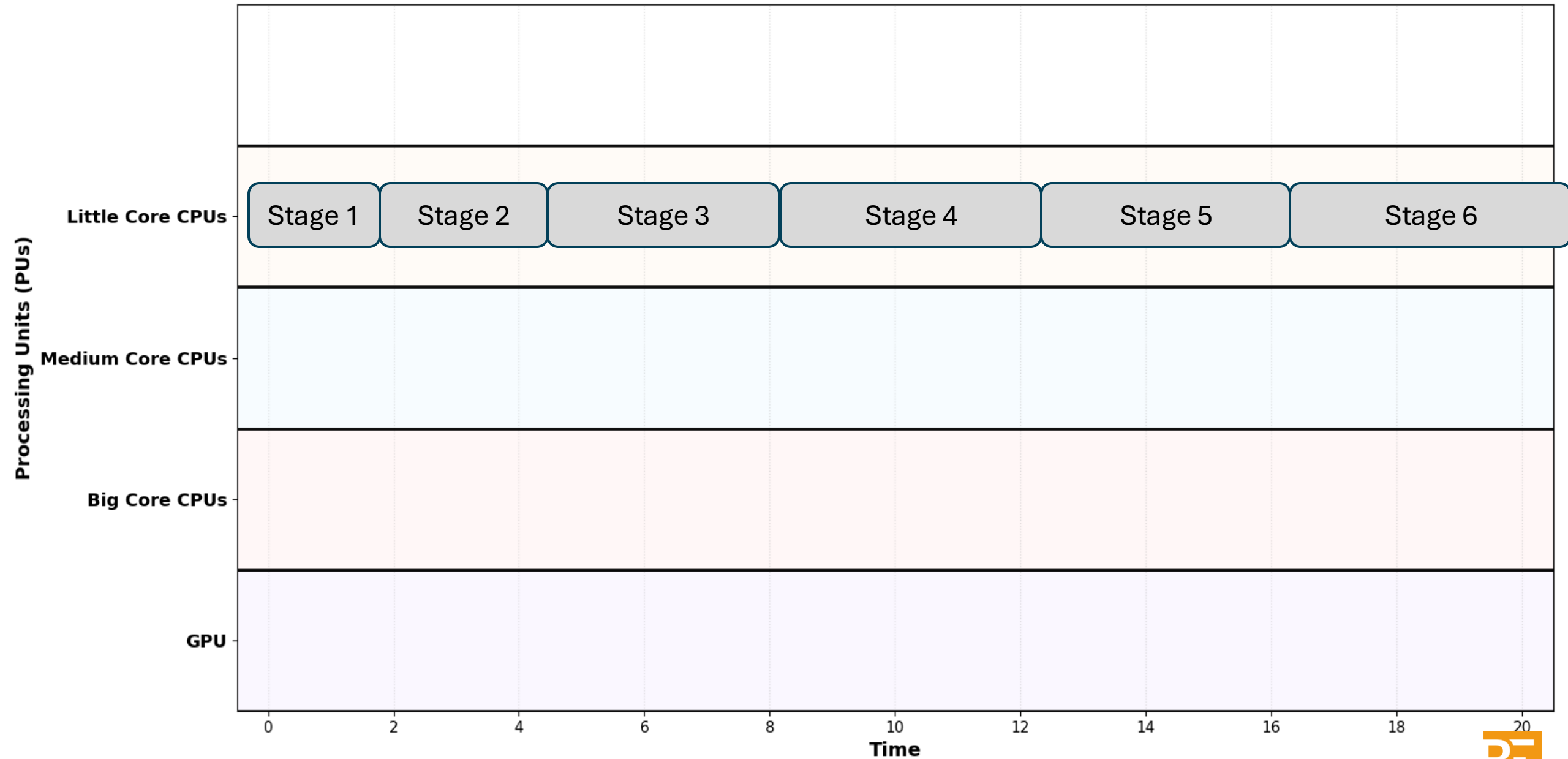
All Available PUs on Pixel

- Ran each stage on each available PUs
- Found optimal **Stage->PU** mapping

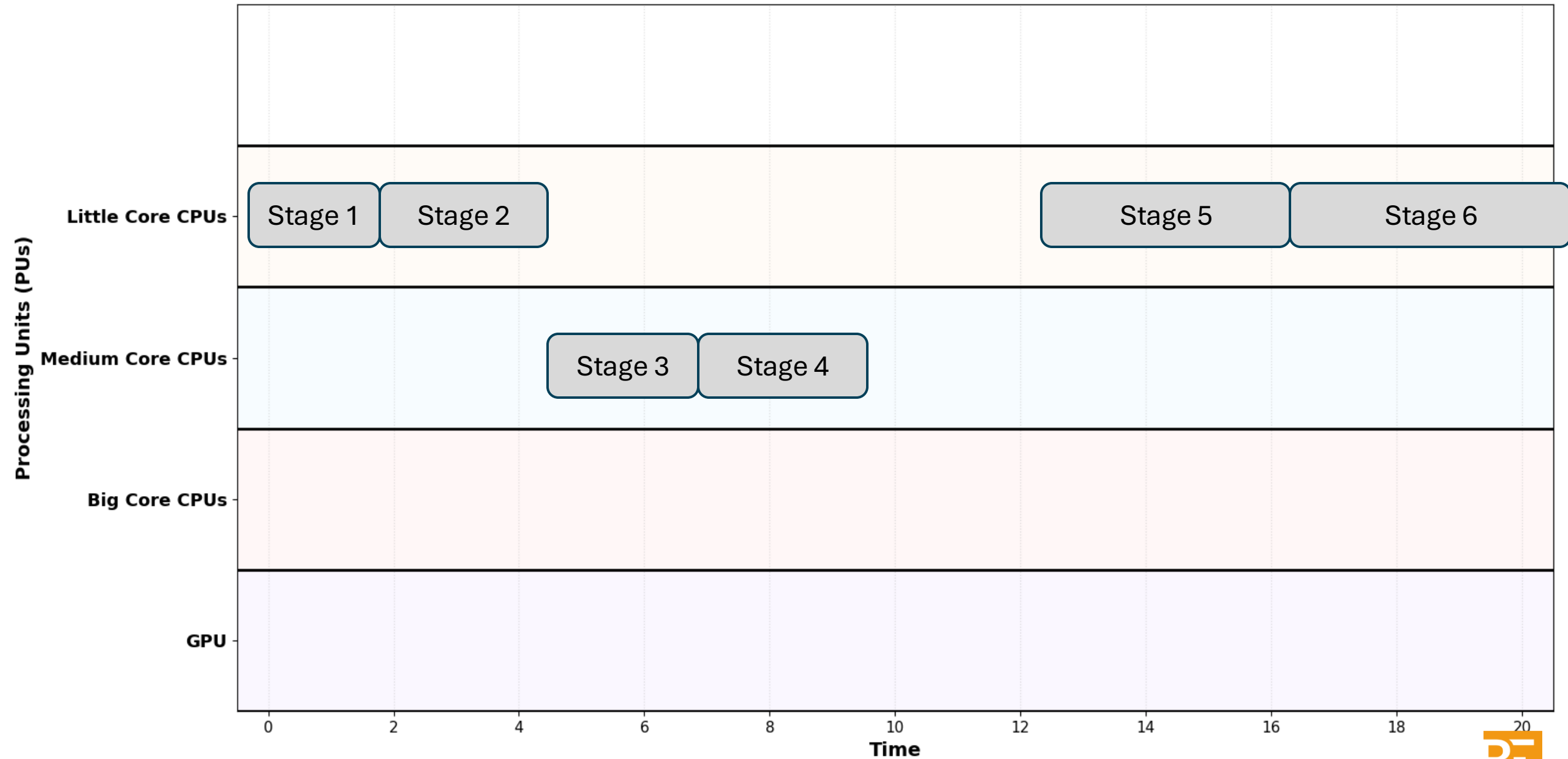


All Stages

# Efficient Pipelined Scheduling

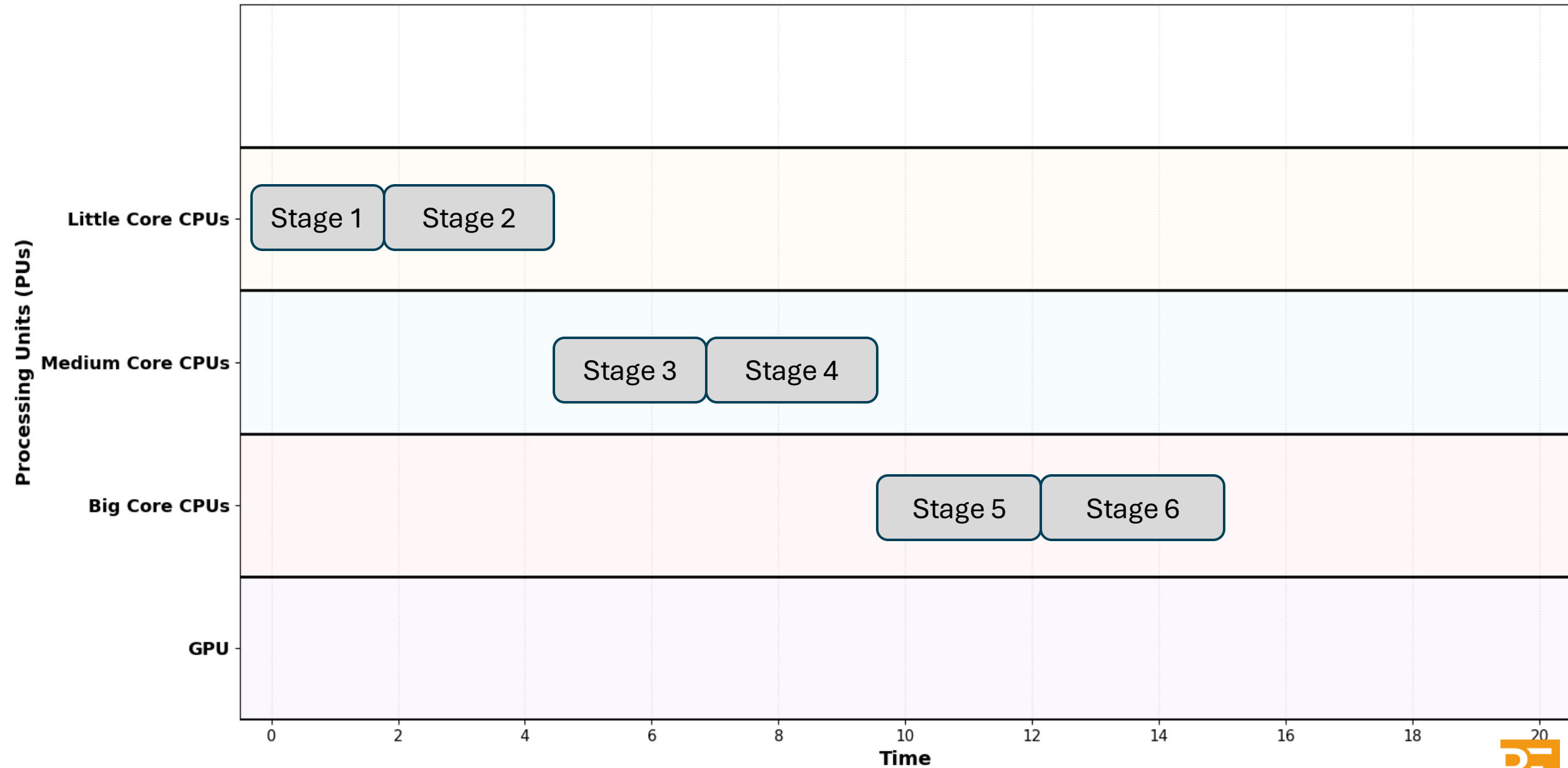


# Efficient Pipelined Scheduling

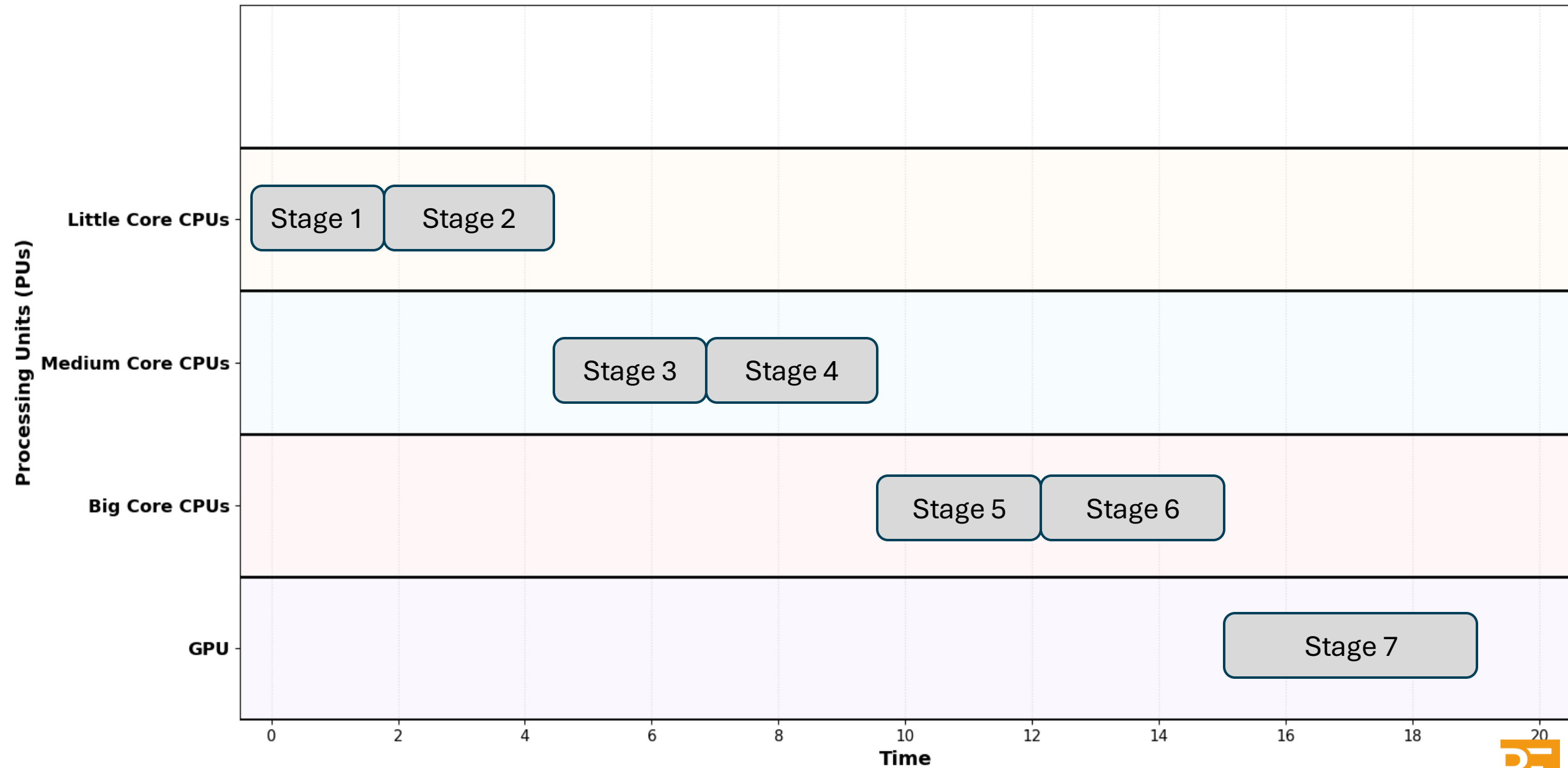




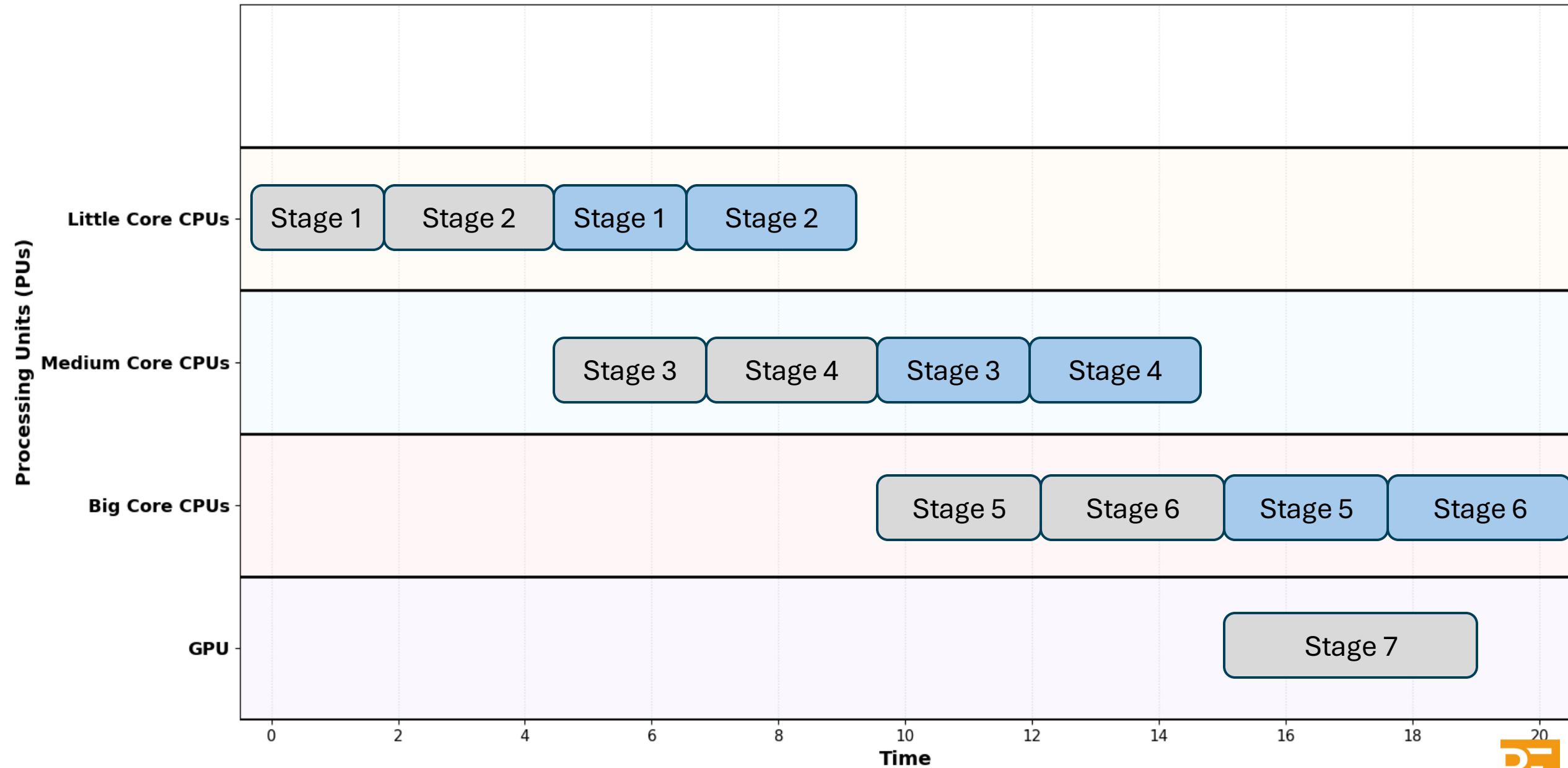
# Efficient Pipelined Scheduling



# Efficient Pipelined Scheduling

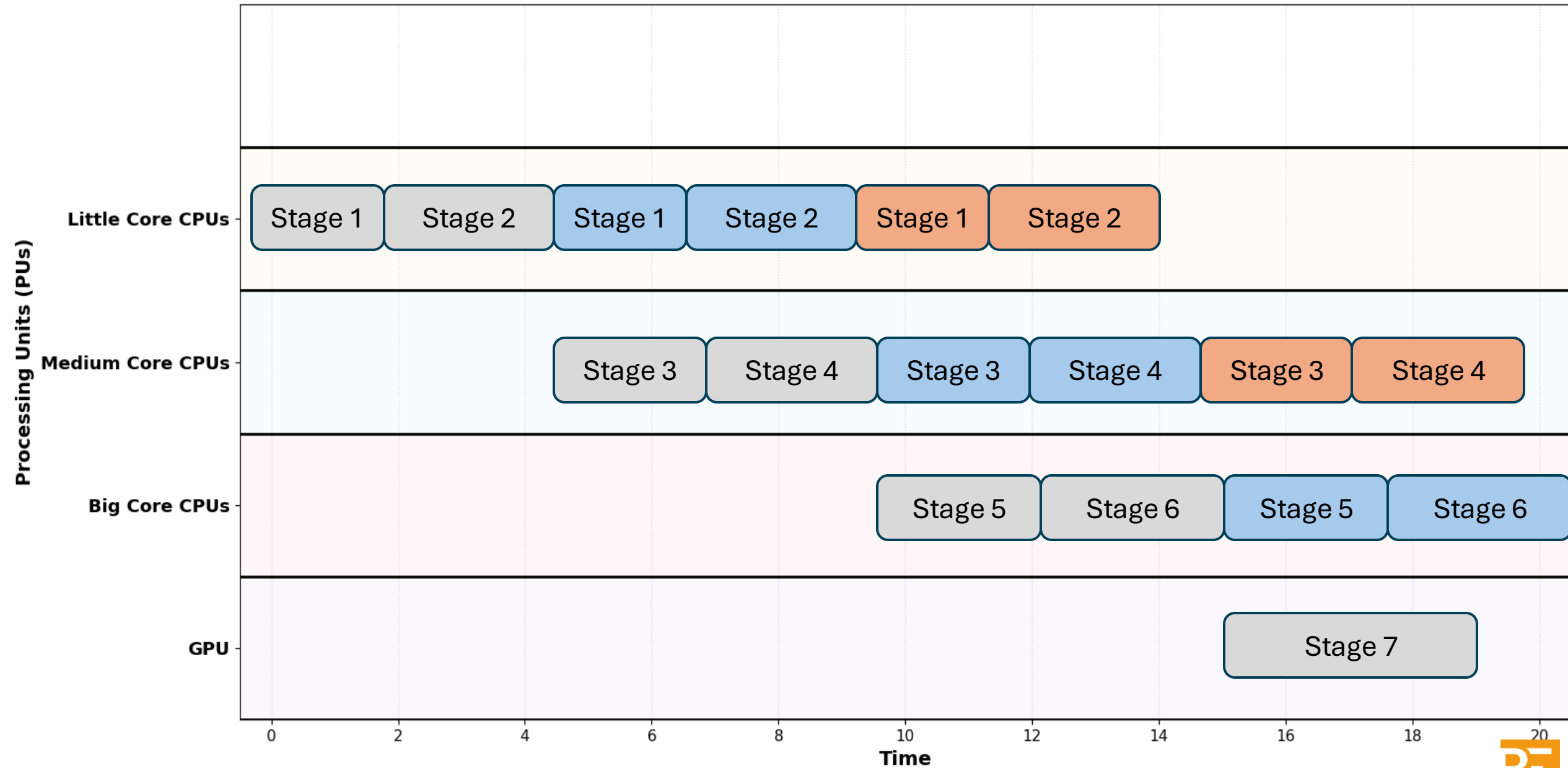


# Efficient Pipelined Scheduling

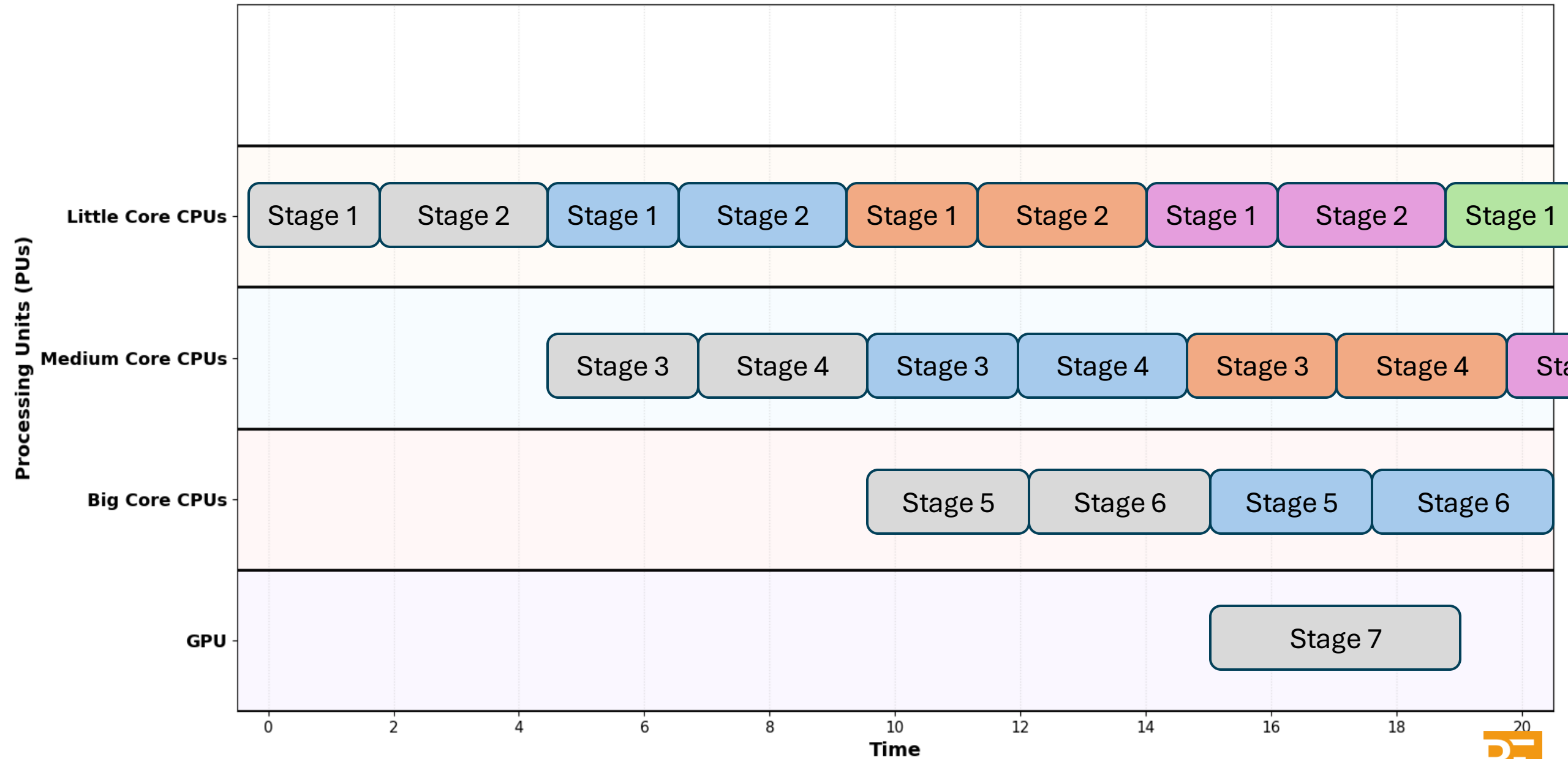




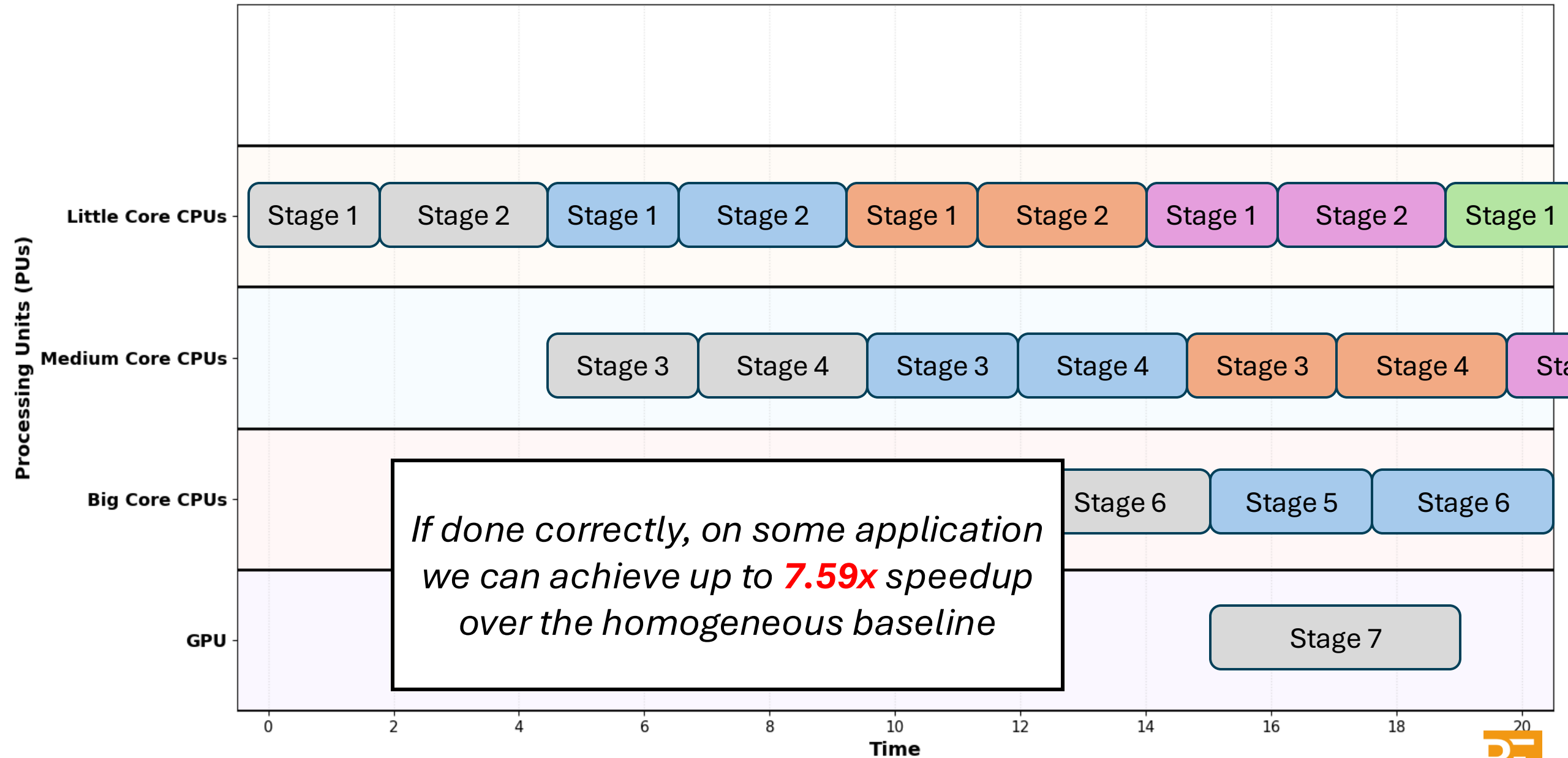
# Efficient Pipelined Scheduling



# Efficient Pipelined Scheduling

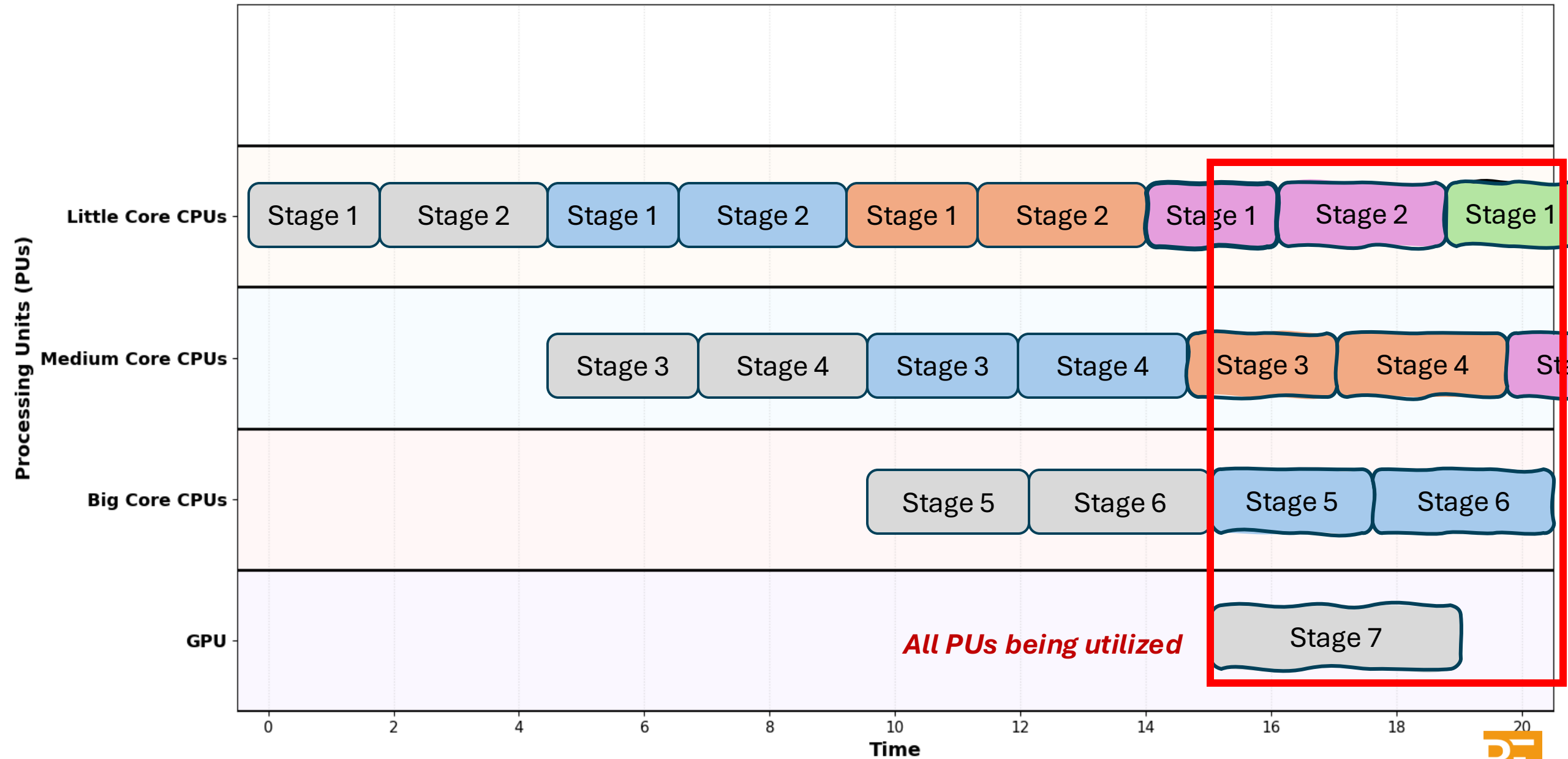


# Efficient Pipelined Scheduling

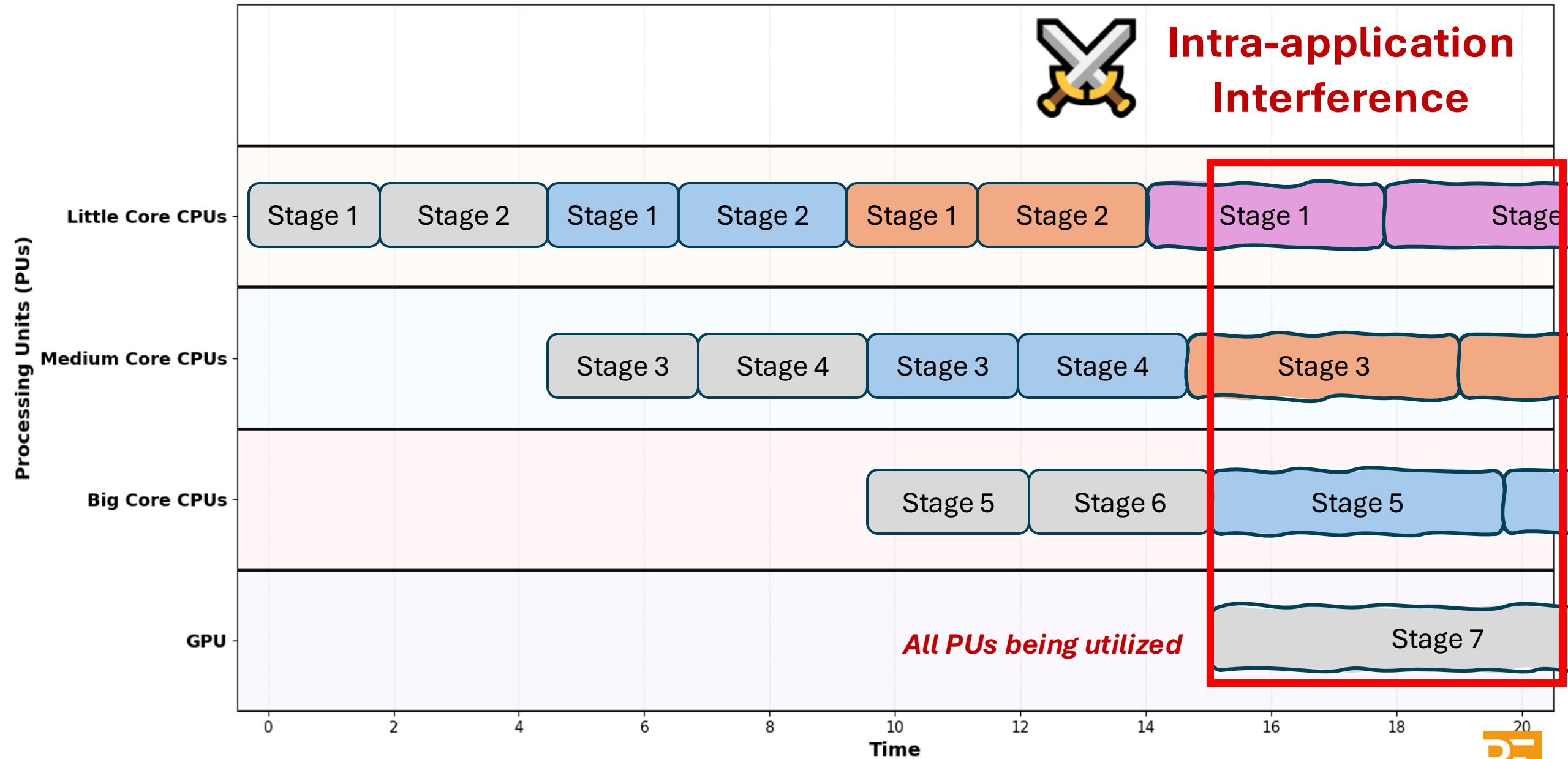




# Efficient Pipelined Scheduling



# Efficient Pipelined Scheduling

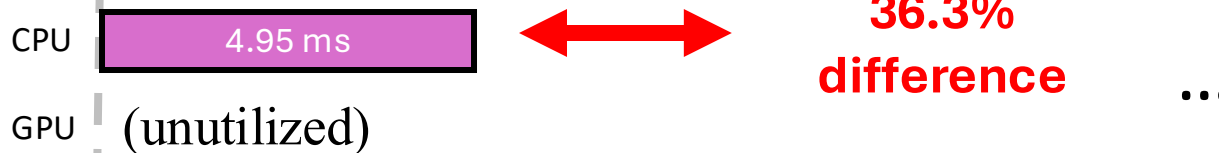


- Mobile systems are prone to **intra-application interferences**

- E.g., we think **4.95** ms, but real measurement was **7.77** ms, **~57% slower**

#### Execution Timeline

(a) Expected



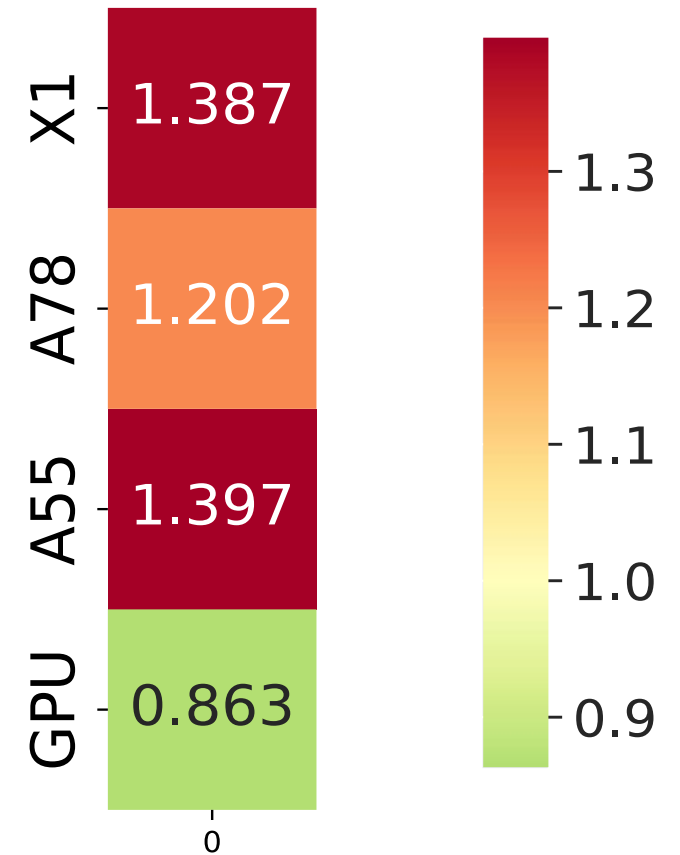
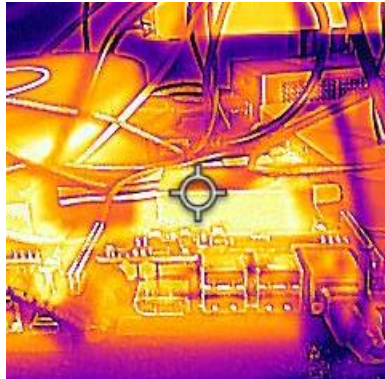
(b) Measured





# Challenge I: Interference

- When PUs fully utilized
- **Slowdowns** and **Speedups\*** due to
  - *Resource contention*
  - *Dynamic voltage and frequency scaling (DVFS)*
  - *Thermal throttling*
  - *Power management*
  - ...

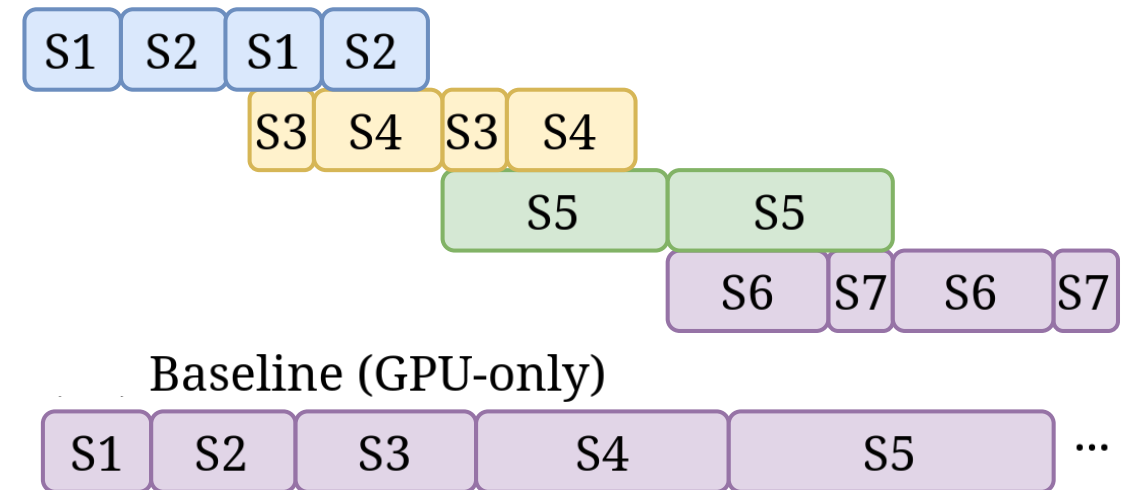
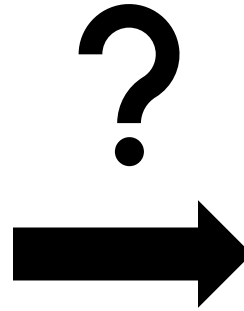
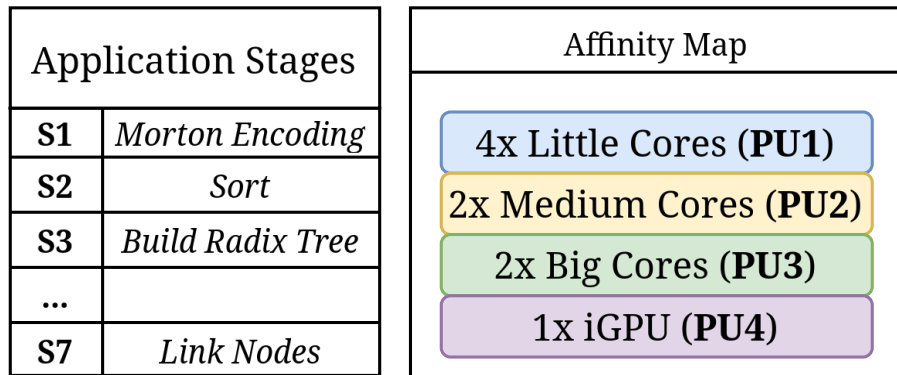


**Red = Slower down**  
**Green = Speedup**

\*We consulted with engineers from a major mobile vendor, whose insights were consistent with our observations<sup>21</sup>

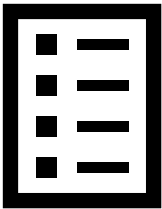
# Finding the Optimal Schedule is Hard

- **Schedule** = mapping from program **stages** to appropriate **PU**

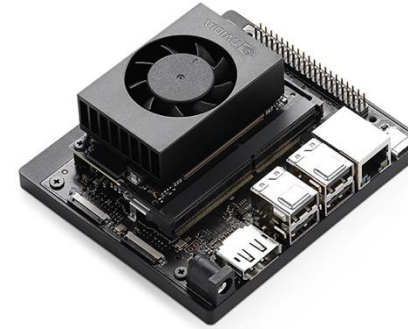


# Challenge II: Portability

- Large design exploration space
  - e.g., 9 stage AlexNet  $5^9 \approx 1.9 \text{ M}$  potential schedules
  - **~37 years** for Google Pixel 7a
- **Schedules** are not portable
  - Optimal schedule on Pixel does not work on NVIDIA Jetson



The scheduling framework need to be **portable** and **flexible**, and **suitable** for rapid development



NVIDIA Jetson

8 GPU-SMs  
6 E-cores

Apple A18 SoC

5 GPU-cores  
8 NPU-cores  
2 P-cores  
4 E-cores

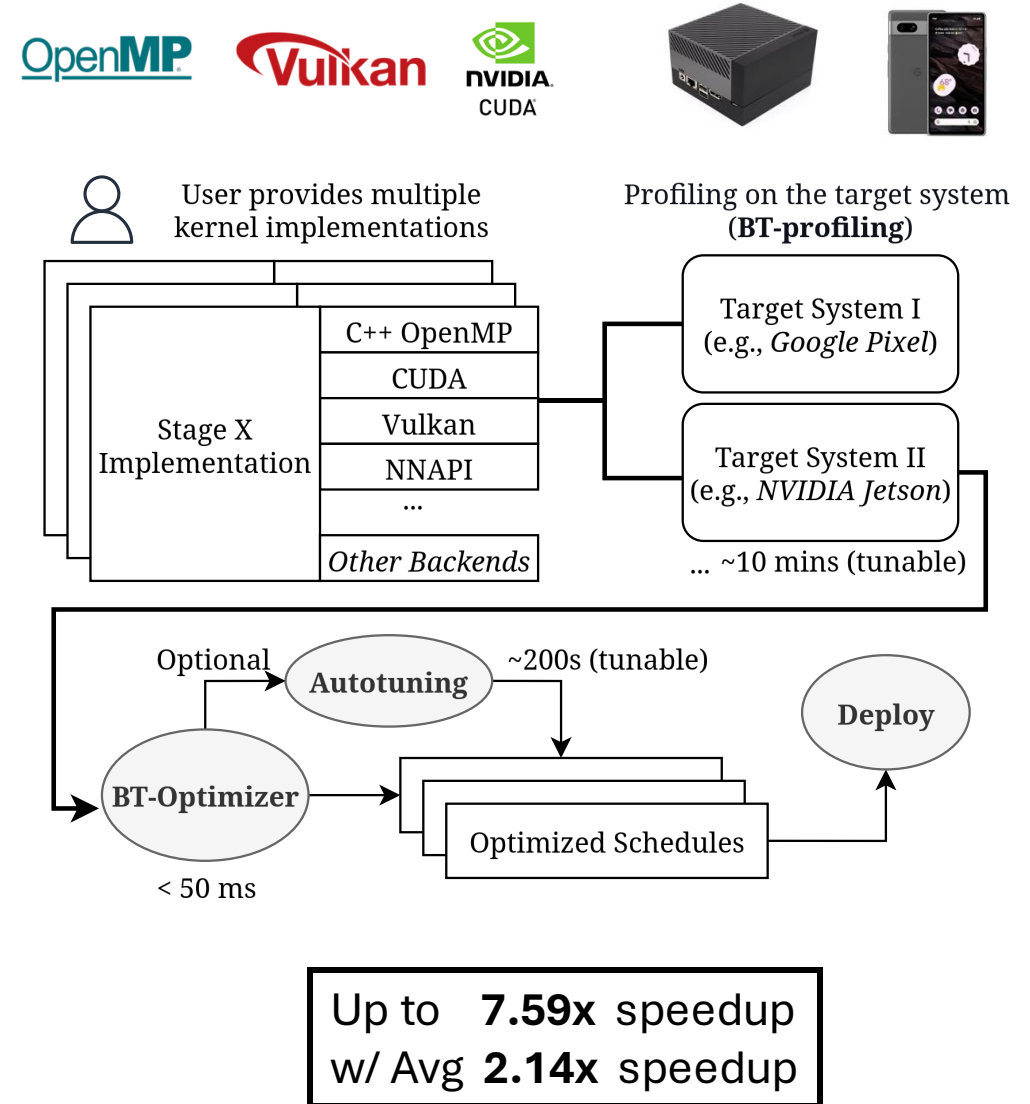
Google Pixel 7a

7 GPU-cores  
1 TPU  
2 P-cores  
2 M-cores  
4 E-cores



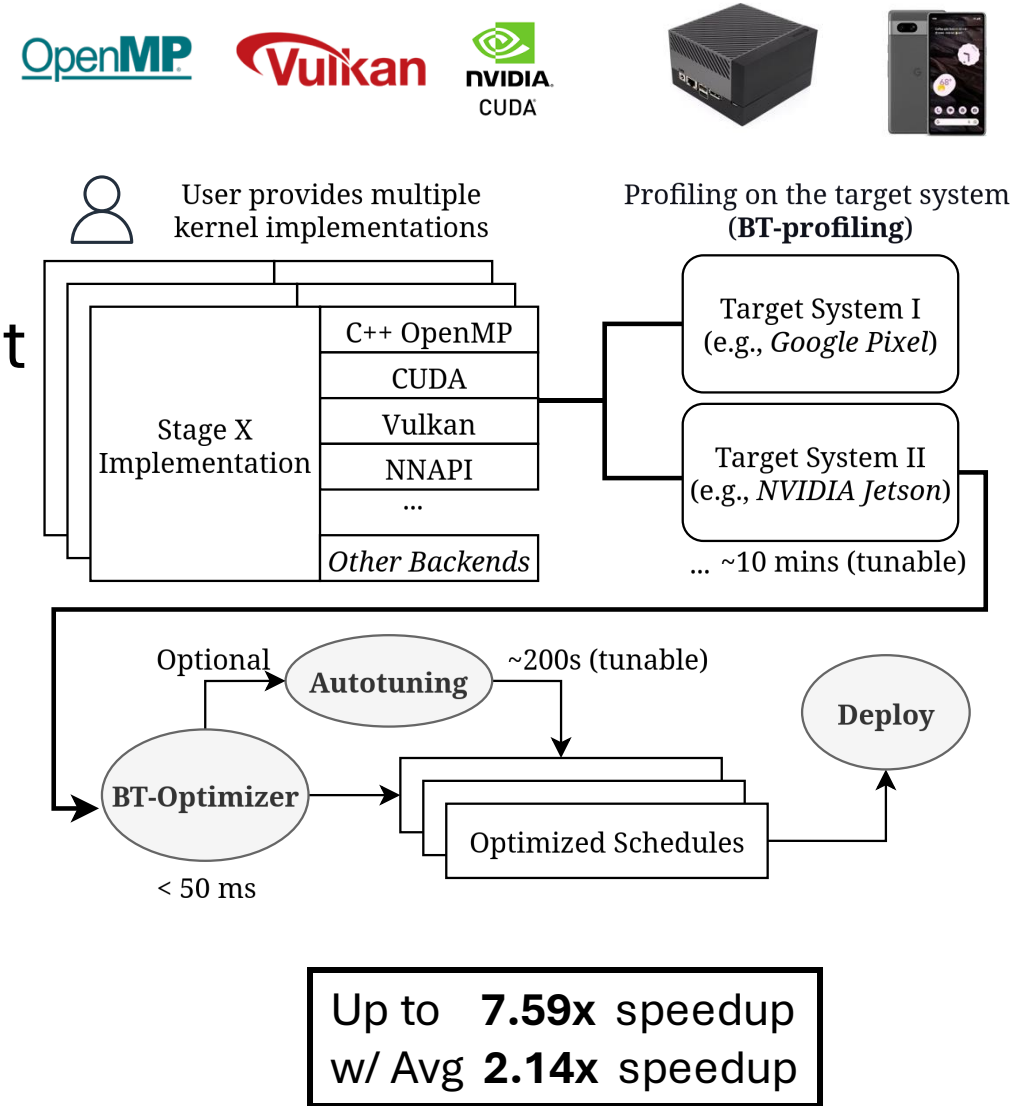
# We present *BetterTogether*

- A performance modeling approach that accounts for intra-application interference



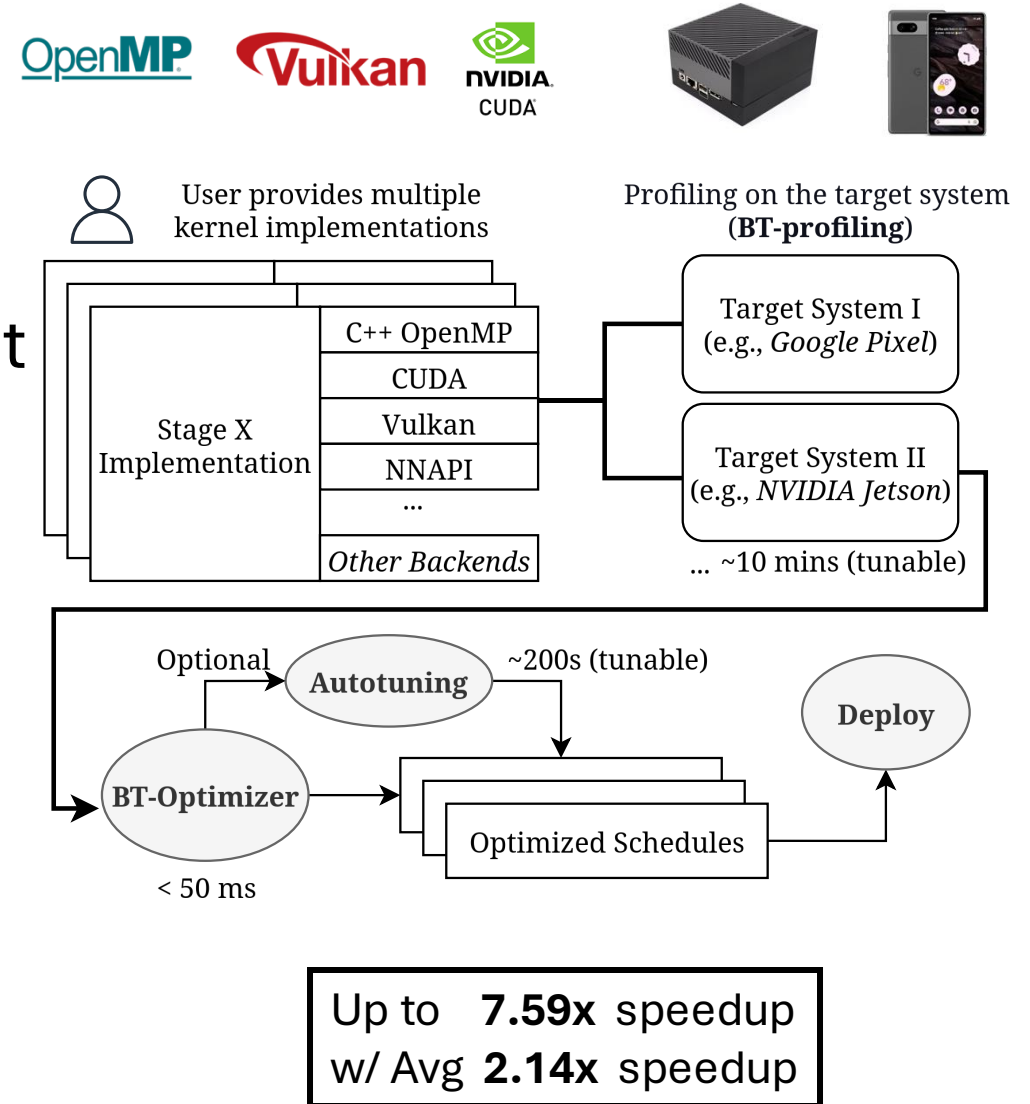
# We present *BetterTogether*

- A performance modeling approach that accounts for intra-application interference
- An end-to-end static pipeline generator that works across a variety of devices



# We present *BetterTogether*

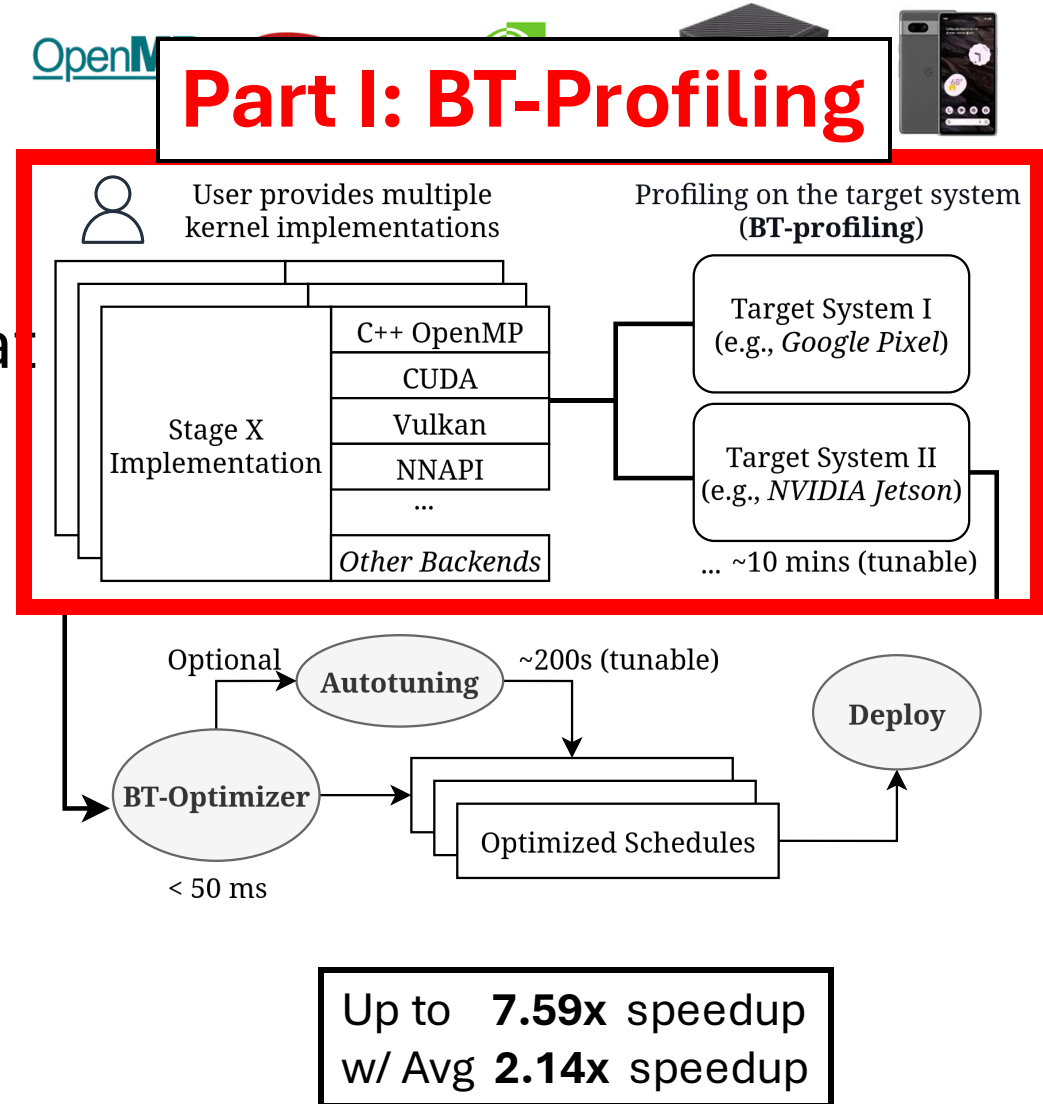
- A performance modeling approach that accounts for intra-application interference
- An end-to-end static pipeline generator that works across a variety of devices
- Consists of 3 major components





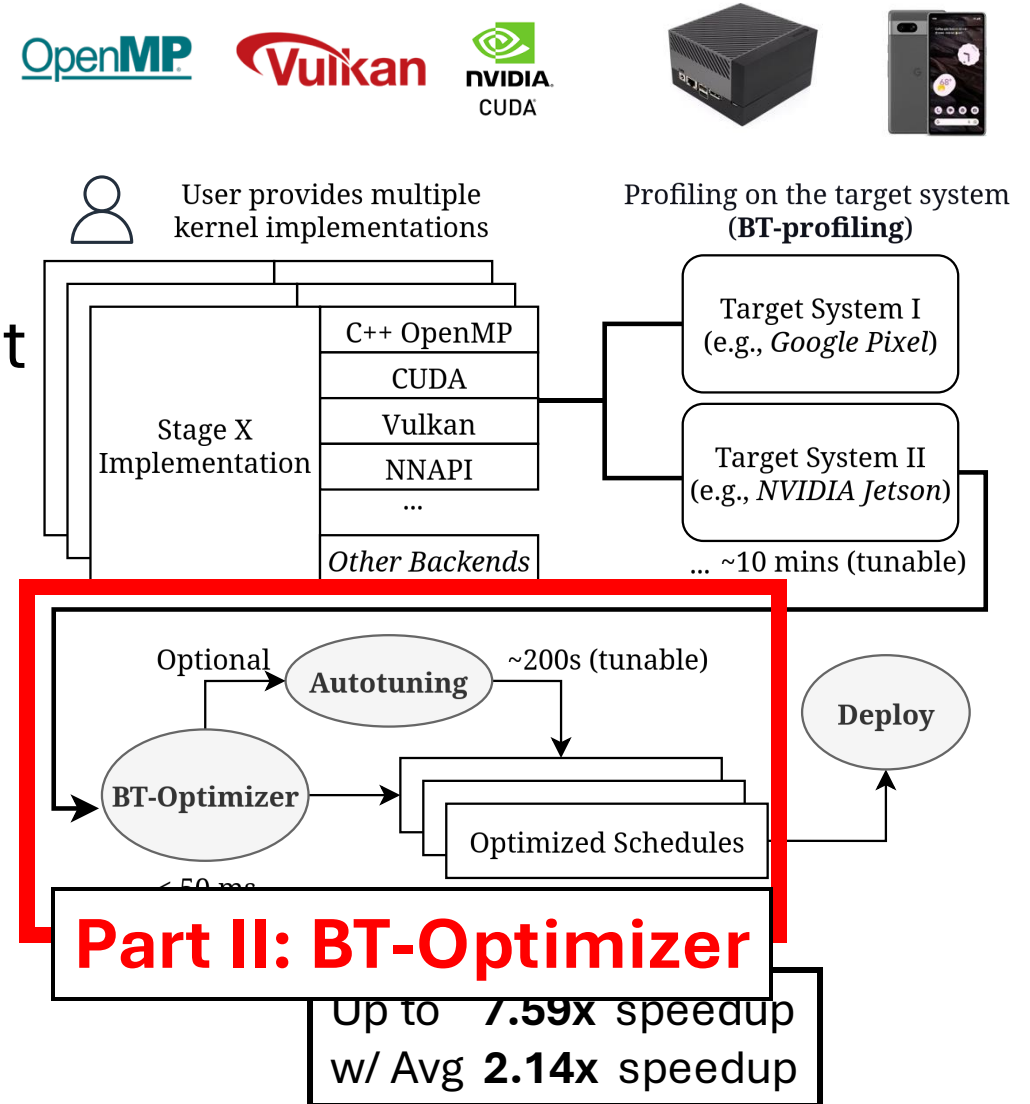
# We present ***BetterTogether***

- A performance modeling approach that accounts for intra-application interference
- An end-to-end static pipeline generator that works across a variety of devices
- Consists of 3 major components
  - *BetterTogether* **Profiling**
    - Attack **Challenge I (Interference)**



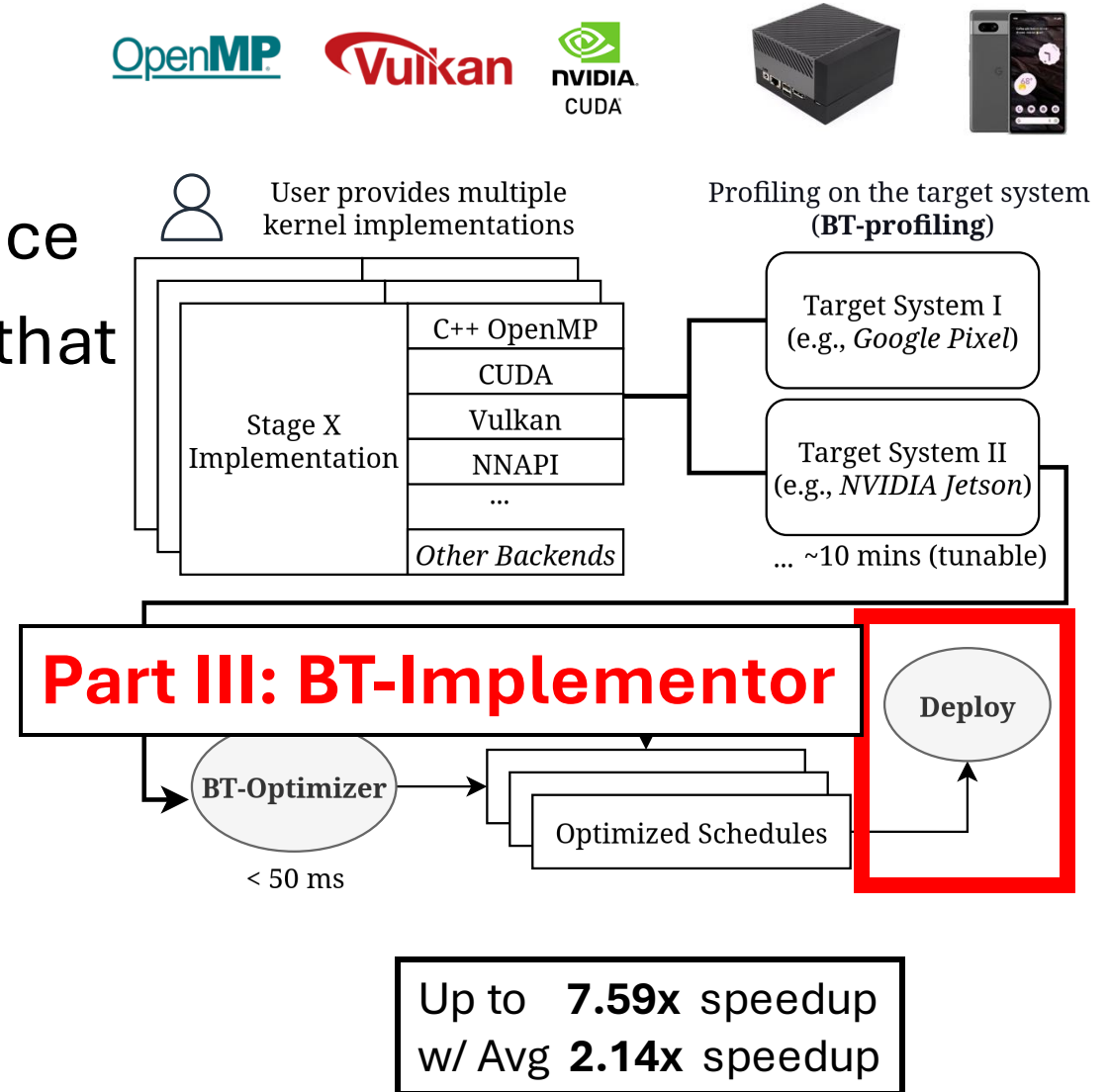
# We present **BetterTogether**

- A performance modeling approach that accounts for intra-application interference
- An end-to-end static pipeline generator that works across a variety of devices
- Consists of 3 major components
  - **BetterTogether Profiling**
    - Attack **Challenge I (Interference)**
  - **BetterTogether Optimizer**
    - Attack **Challenge II (Portability)**

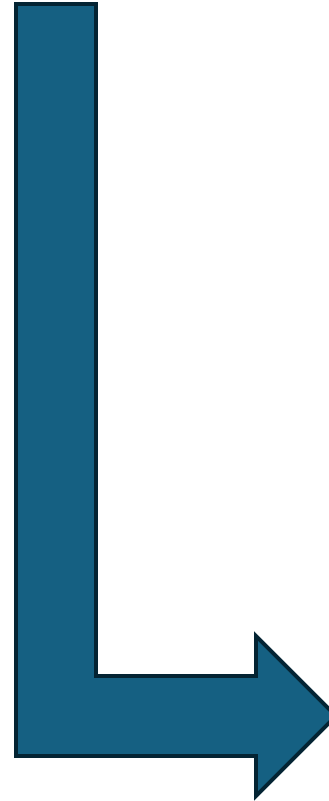
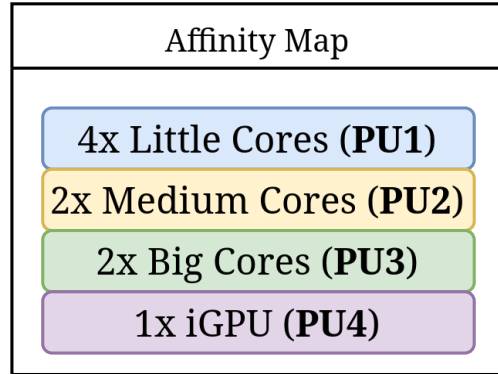


# We present *BetterTogether*

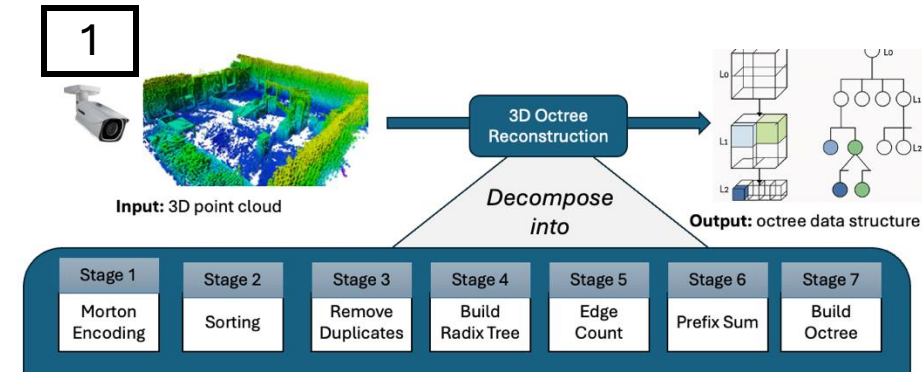
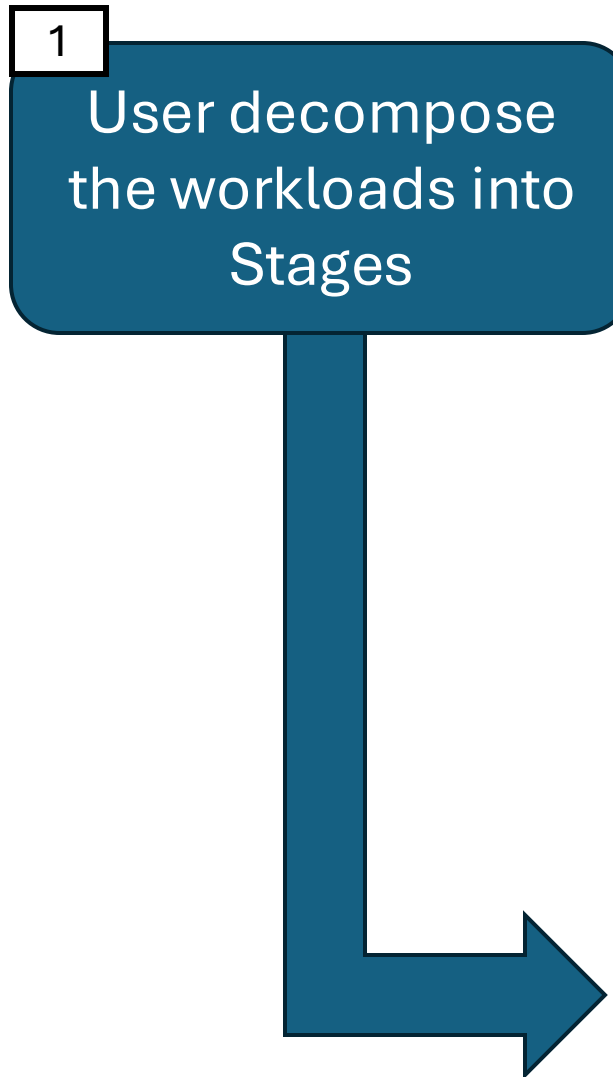
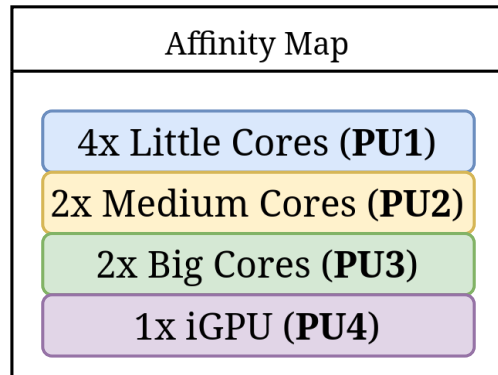
- A performance modeling approach that accounts for intra-application interference
- An end-to-end static pipeline generator that works across a variety of devices
- Consists of 3 major components
  - *BetterTogether* **Profiling**
    - Attack **Challenge I (Interference)**
  - *BetterTogether* **Optimizer**
    - Attack **Challenge II (Portability)**
  - *BetterTogether* **Implementor**
    - Efficient static heterogeneous pipeline execution



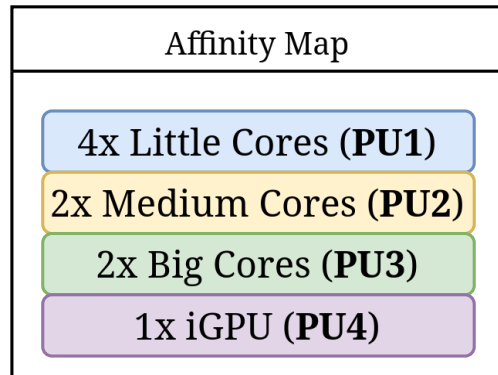
# This Work: *BetterTogether* Overview



# This Work: *BetterTogether* Overview



# This Work: *BetterTogether* Overview



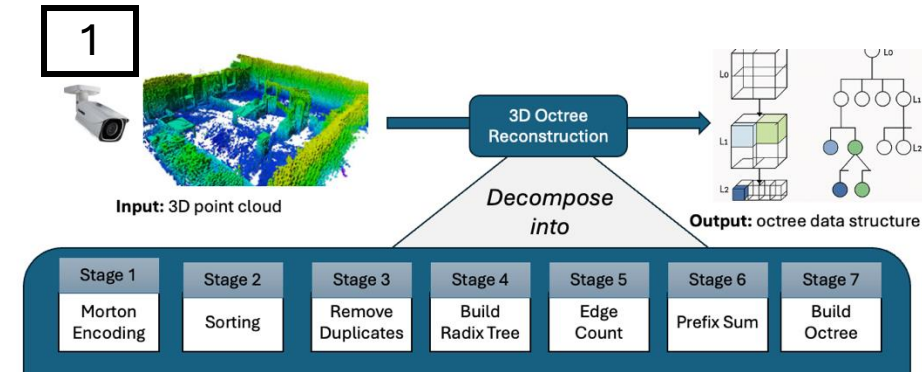
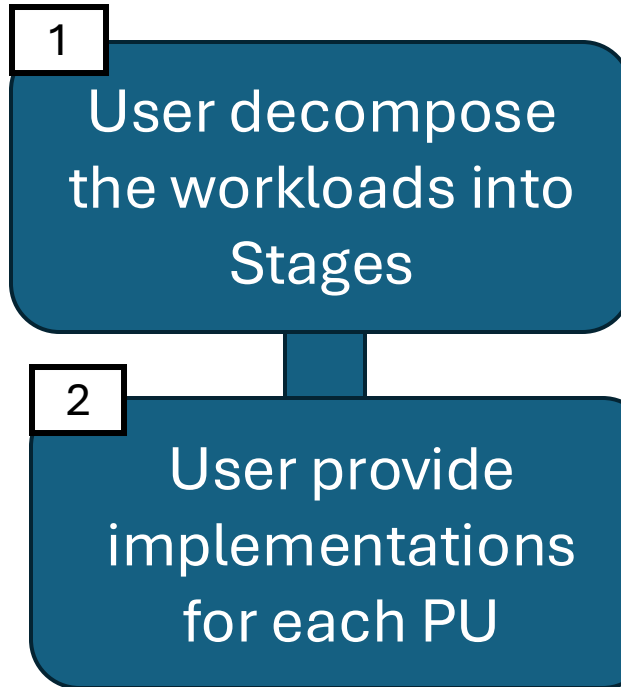
2

```
1 void run_stage_1_cpu(in, out, N) {  
2   #pragma omp parallel for  
3   for (int i = 0; i < N; ++i)  
4     out[i] = morton32(in[i]);  
5 }
```

CPU Code (e.g., OpenMP)

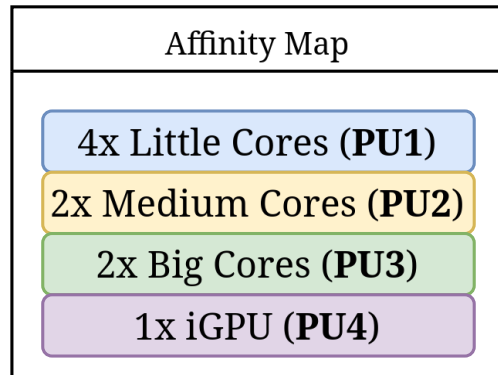
```
1 __global__ void run_stage_1_gpu(in, out, N) {  
2   int idx = threadIdx.x + blockDim.x * blockIdx.x;  
3   int stride = blockDim.x * gridDim.x;  
4   for (int i = idx; i < N; i += stride)  
5     out[i] = morton32(in[i]);  
6 }
```

GPU Code (e.g., CUDA)



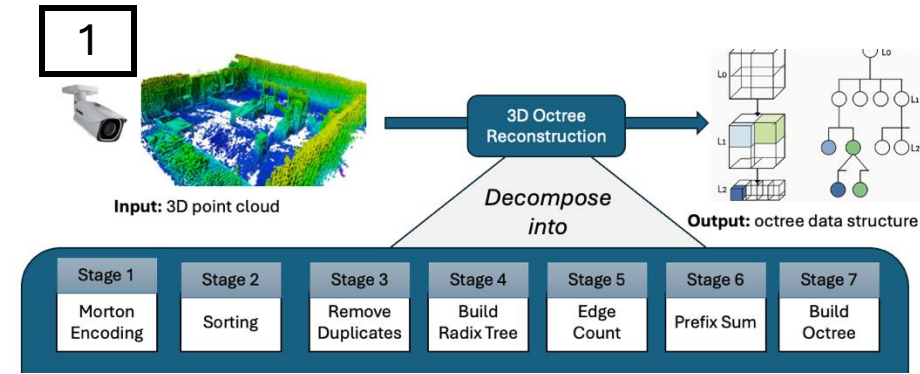


# This Work: *BetterTogether* Overview



1 User decompose the workloads into Stages

2 User provide implementations for each PU



2

```
1 void run_stage_1_cpu(in, out, N) {  
2     #pragma omp parallel for  
3     for (int i = 0; i < N; ++i)  
4         out[i] = morton32(in[i]);  
5 }
```

CPU Code (e.g., OpenMP)

```
1 __global__ void run_stage_1_gpu(in, out, N) {  
2     int idx = threadIdx.x + blockDim.x * blockIdx.x;  
3     int stride = blockDim.x * gridDim.x;  
4     for (int i = idx; i < N; i += stride)  
5         out[i] = morton32(in[i]);  
6 }
```

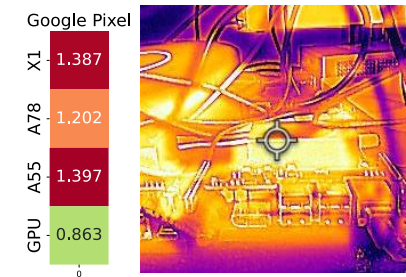
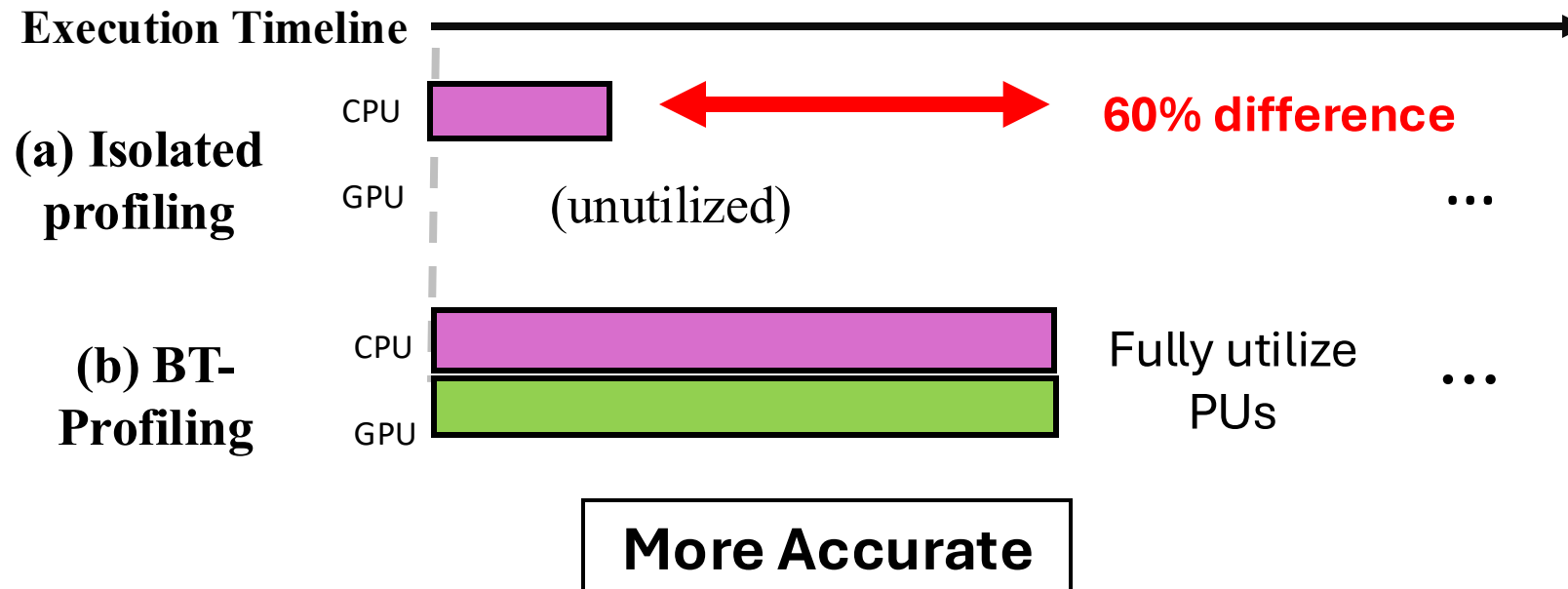
GPU Code (e.g., CUDA)

Interference-aware Profiling Table				
	S1	S2	..	S7
PU0	2.6	3.3		5.8
PU1	0.8	1.5		1.9
PU2	0.6	1.4		2.1
PU3	0.8	9.0		1.5

Interference aware  
**BT-Profiling**

# BT-Profiling - *Interference aware profiling*

- While profiling each {PU × Stage} pair:
  - Concurrently execute similar stages on other PUs
  - Simulate **whole-application** execution to capture resource contention



Overcome

# BT-Optimizer

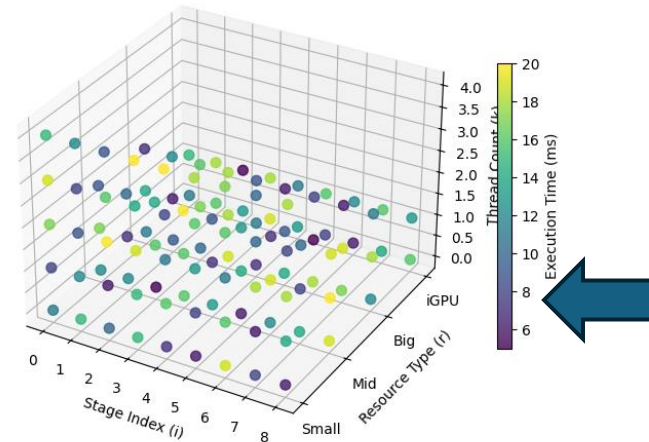
- We express our optimization problem as linear constraints

## d) Notation and Decision Variables:

$N$	Total number of pipeline stages
$N_i$	The pipeline stage $i$ , with $i \in \mathcal{N} = \{0, \dots, N-1\}$
$\mathcal{C}$	PU classes: $\mathcal{C} = \{c_1, \dots, c_M\}$
$t_{i,c}$	Profiled latency of stage $i$ on PU $c$
$x_{i,c}$	Decision variable: $x_{i,c} \in \{0, 1\}$ ; $x_{i,c} = 1 \Leftrightarrow$ stage $i$ runs on PU $c$

- We propose a **three-step** optimization approach

3D Visualization of  $T_{\{i,r,k\}}$  (Execution Times)



**SMT  
Solver**

*Raw BT  
Profiling  
results*

584	9b834f1b	Tree	OMP	4	medium	2	1.52
585	9b834f1b	Tree	OMP	5	little	1	4.35
586	9b834f1b	Tree	OMP	5	little	2	4.37
587	9b834f1b	Tree	OMP	5	little	3	4.73
588	9b834f1b	Tree	OMP	5	medium	1	1.52
589	9b834f1b	Tree	OMP	5	medium	2	1.52
590	9b834f1b	Tree	OMP	6	little	1	0.825
591	9b834f1b	Tree	OMP	6	little	2	0.838
592	9b834f1b	Tree	OMP	6	little	3	0.975
593	9b834f1b	Tree	OMP	6	medium	1	0.212
594	9b834f1b	Tree	OMP	6	medium	2	0.281
595	9b834f1b	Tree	OMP	7	little	1	23.6
596	9b834f1b	Tree	OMP	7	little	2	11.7
597	9b834f1b	Tree	OMP	7	little	3	9.13
598	9b834f1b	Tree	OMP	7	medium	1	2.86
599	9b834f1b	Tree	OMP	7	medium	2	1.43
600	9b834f1b	Tree	OMP	7	medium	3	1.43
601	9b834f1b	CifarDense	VK	0	None	None	31.8
602	9b834f1b	CifarDense	VK	1	None	None	1.08
603	9b834f1b	CifarDense	VK	2	None	None	0.886
604	9b834f1b	CifarDense	VK	3	None	None	4.12
605	9b834f1b	CifarDense	VK	4	None	None	0.837
606	9b834f1b	CifarDense	VK	5	None	None	5.34
607	9b834f1b	CifarDense	VK	6	None	None	7.98
608	9b834f1b	CifarDense	VK	7	None	None	5.64
609	9b834f1b	CifarDense	VK	8	None	None	0.918
610	9b834f1b	CifarDense	VK	9	None	None	3.68
611	9b834f1b	CifarSparse	VK	0	None	None	19.5
612	9b834f1b	CifarSparse	VK	1	None	None	2.44
613	9b834f1b	CifarSparse	VK	2	None	None	0.984
614	9b834f1b	CifarSparse	VK	3	None	None	1.39
615	9b834f1b	CifarSparse	VK	4	None	None	0.895
616	9b834f1b	CifarSparse	VK	5	None	None	1.84
617	9b834f1b	CifarSparse	VK	6	None	None	1.88
618	9b834f1b	CifarSparse	VK	7	None	None	1.84
619	9b834f1b	CifarSparse	VK	8	None	None	0.856
620	9b834f1b	CifarSparse	VK	9	None	None	0.853



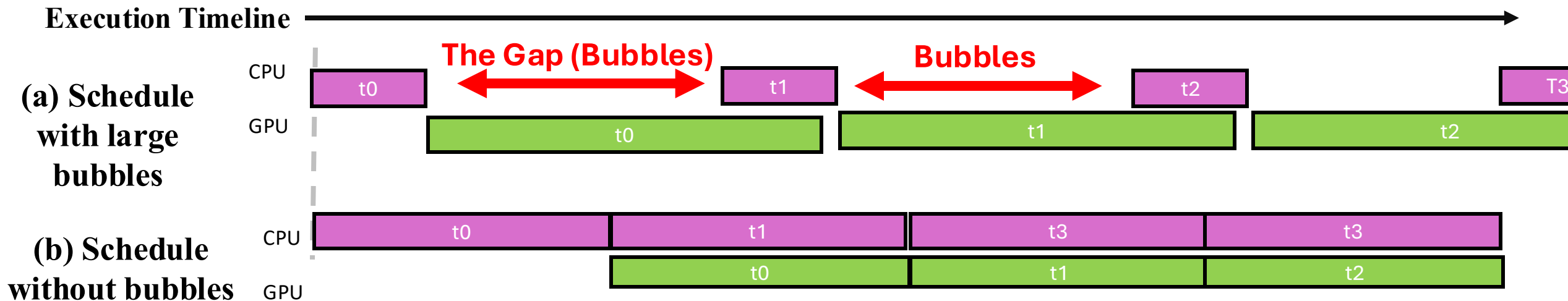
Interference-aware Profiling Table				
	S1	S2	..	S7
PU0	2.6	3.3		5.8
PU1	0.8	1.5		1.9
PU2	0.6	1.4		2.1
PU3	0.8	9.0		1.5

*BT Profiling Table*

# BT-Optimizer **Step 1** - Minimizing Pipeline Bubbles

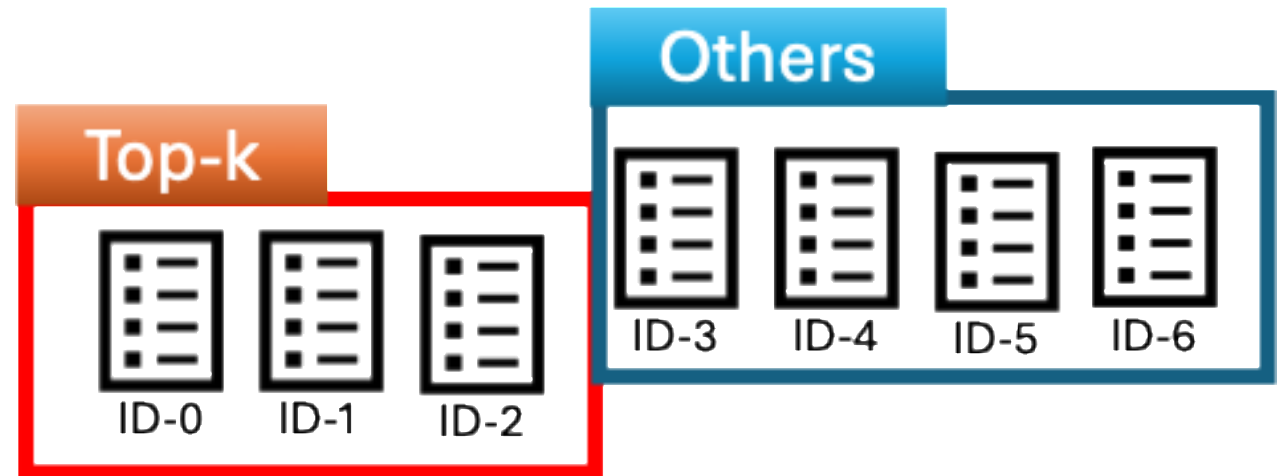


- Improve utilization by reducing idle gaps (*pipeline bubbles*) across PUs
- By **reducing bubbles**, we **keep all PUs busy** and improve pipeline throughput.

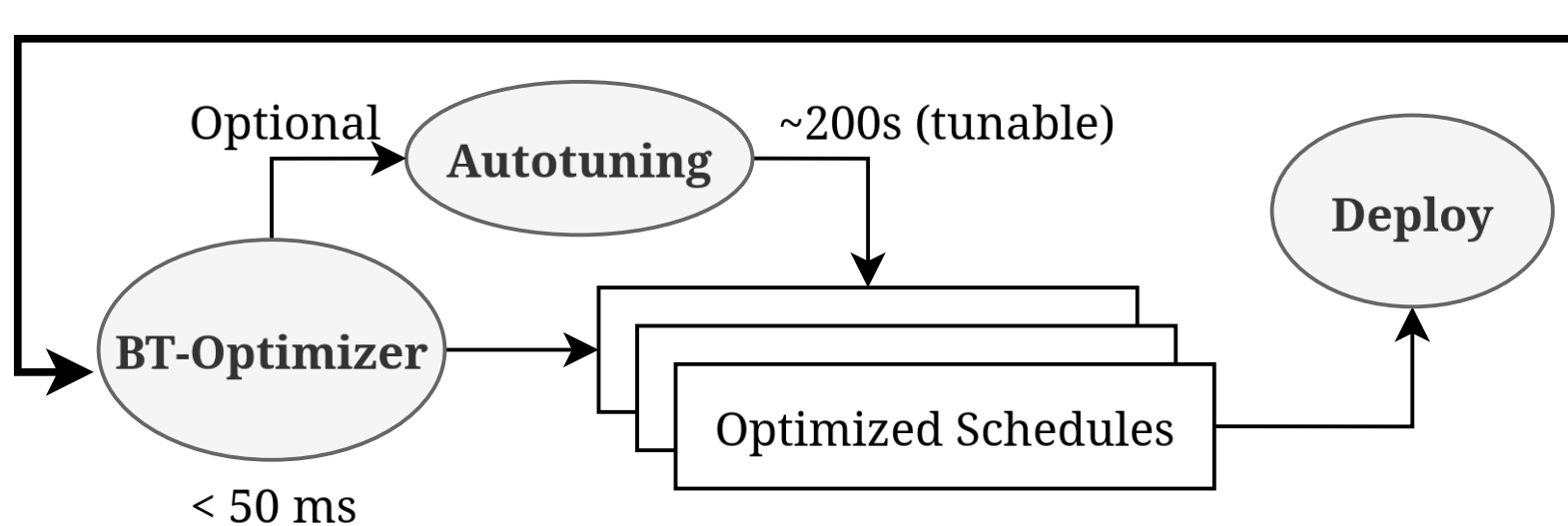


# BT-Optimizer **Step 2: Optimizing Latency**

- High utilization  $\neq$  **low latency**
- Generate K schedules (e.g.,  $k = 20$ ),
  - each with a different assignment of stages to PUs.
- Minimize latency

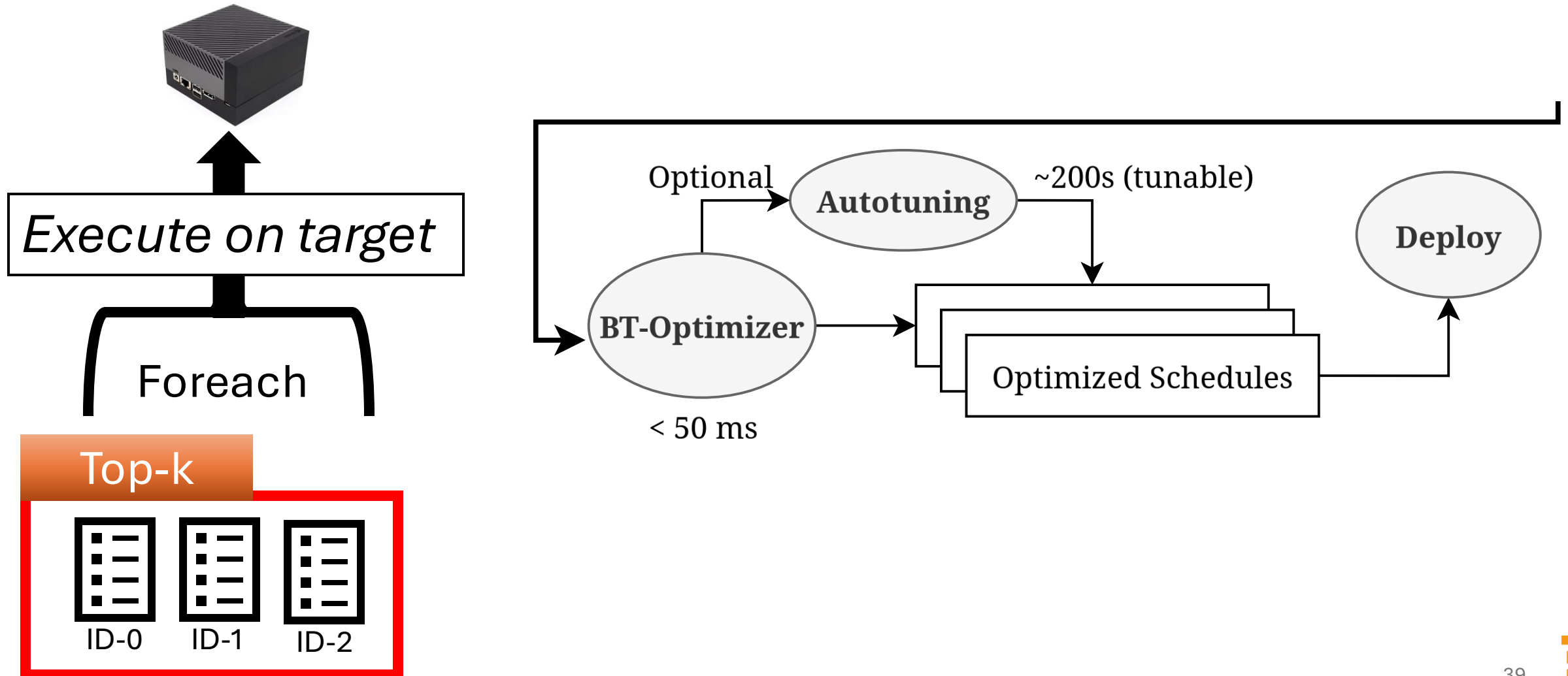


# BT-Optimizer Optional Step 3: Autotuning





# BT-Optimizer Optional Step 3: Autotuning

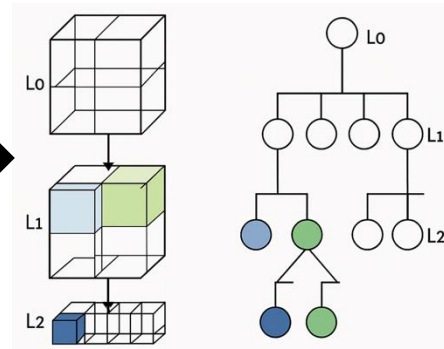
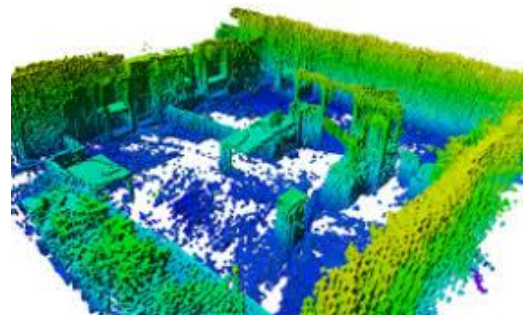


# Workloads Evaluated

- Three edge compute vision tasks

1. AlexNet-dense
2. AlexNet-sparse
3. Octree Construction

- Common in resource-constrained environments



- *AlexNet-dense*

- Dense linear algebra
- CIFAR-10 dataset

- *AlexNet-sparse*

- Sparse linear algebra
- Pruned w/ CONDENSEA\*
- Stored in CSR

- *Octree*

- Tree traversals
- Irregular memory access
- Sorting
- Prefix Sum

\* <https://github.com/NVlabs/condensa>

# Platforms Evaluated



Less Powerful GPUs  
Mix CPUs

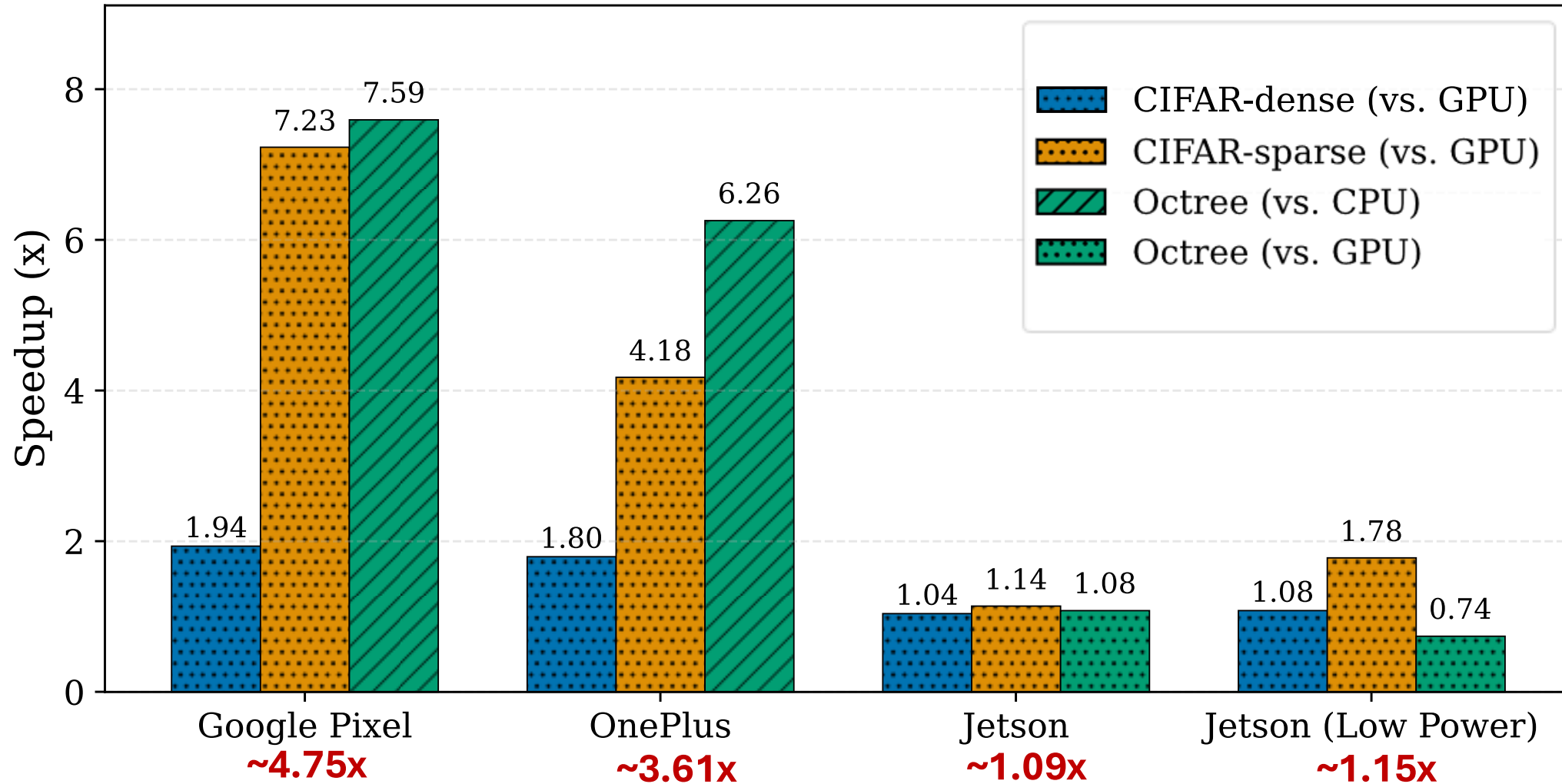
Powerful GPUs  
Little ARM CPUs



Platform	Backend	CPU	CPU Frequency	Integrated GPU
Google Pixel 7A	Vulkan	2x Cortex-X1 2x Cortex-A78 4x Cortex-A55	2.85 GHz 2.35 GHz 1.80 GHz	Mali-G710 MP7
OnePlus 11	Vulkan	1x Coretex-X3 2x Coretex-A715 2x Cortex-A710 3x Cortex-A510	3.2 GHz 2.8 GHz 2.8 GHz 2.0 GHz	Adreno 740
NVIDIA Jetson Orin Nano	CUDA Vulkan	6x Cortex-A78AE	1.7 GHz	Ampere GPU
*NVIDIA Jetson Orin Nano (low-power mode)	CUDA Vulkan	4x Cortex-A78AE	~0.85 GHz	Ampere GPU

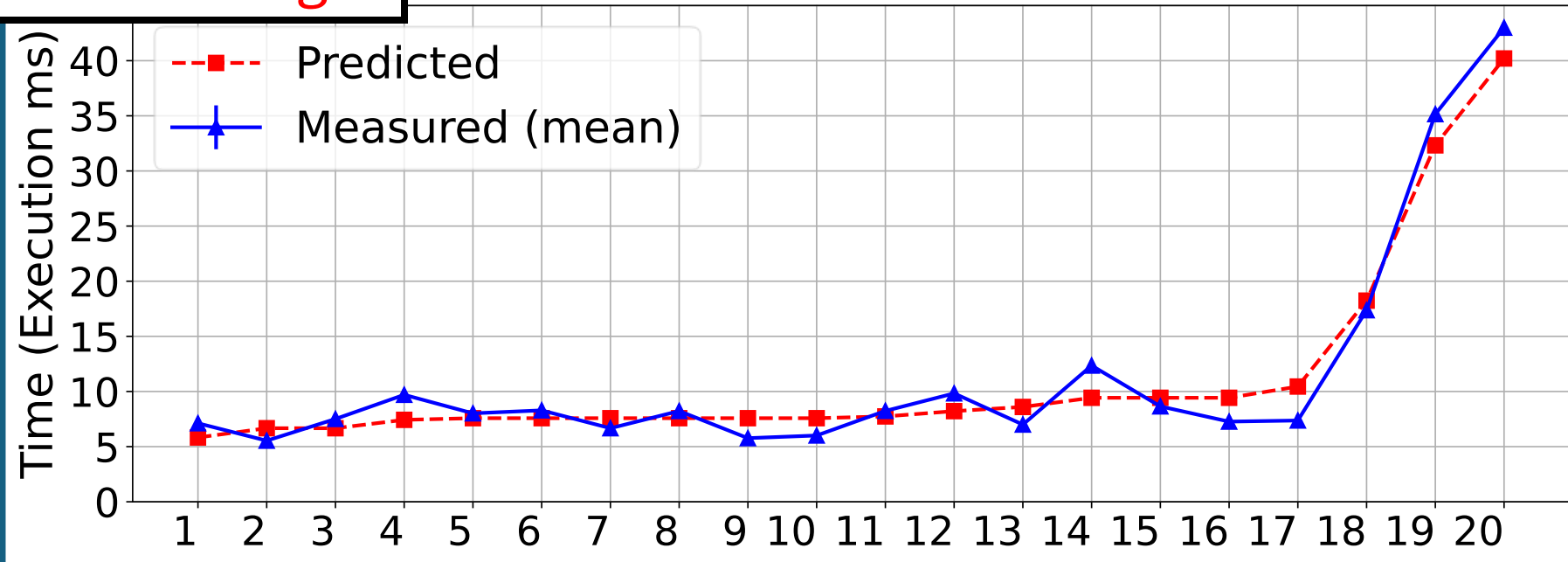
\* In low-power mode, 2 cores are shutdown, and overall CPU frequency is reduced by half

# Results Overview

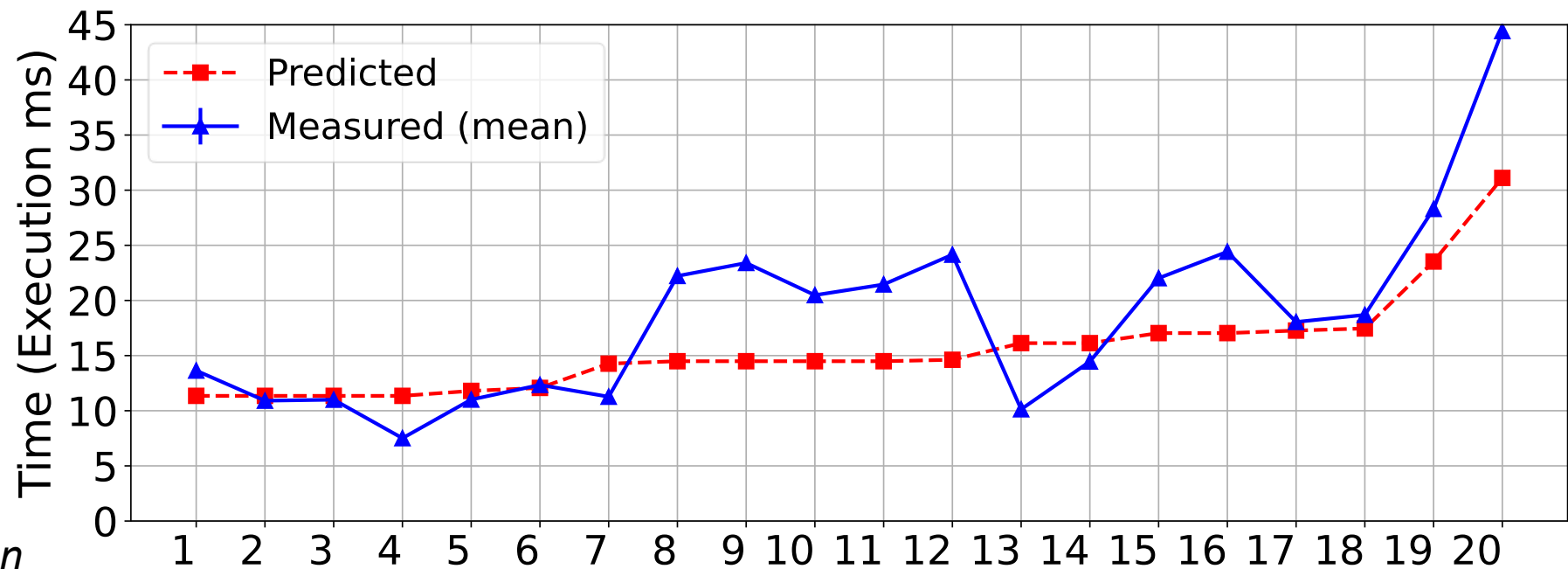


- Geomean speedup of **2.14x** across all workloads, w/ peak of **7.59x**
- In 1 case, slowdown

## BT-Profiling



*BetterTogether*  
produces predictions  
that closely match  
measured execution  
time

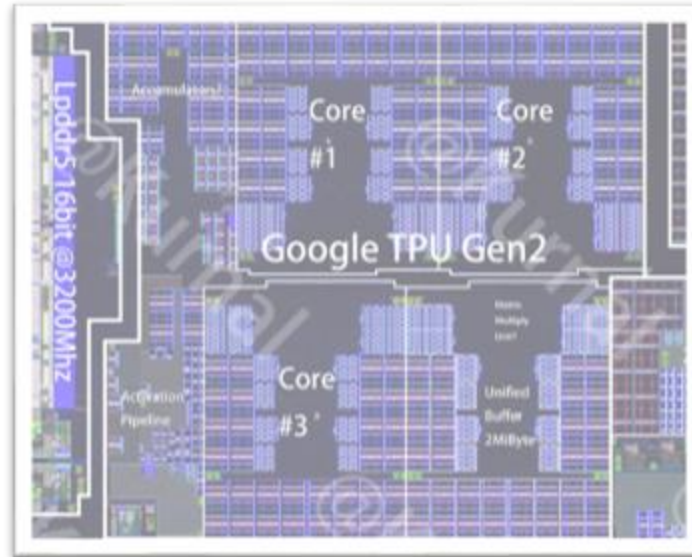


Without BT show  
discrepancies

*predicted and measured execution*

# We have additional insights in the paper

- Preliminary results on
  - Using Google's **EdgeTPU**
  - Implemented using *nnapi*
  - Showing a **1.25x** speedup for AlexNet-dense application on top of existing results
  - Showcasing the flexibility and extensibility of *BetterTogether*





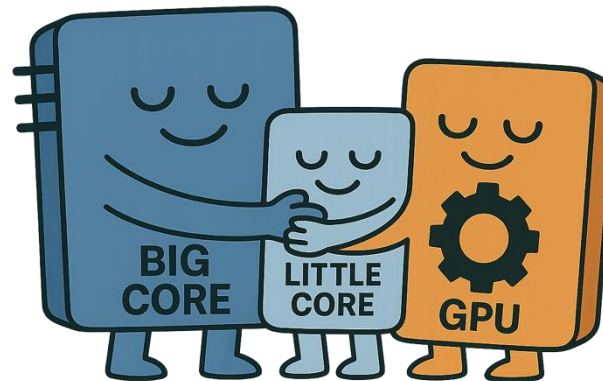
# Conclusion

- We propose **BetterTogether**
  - An **interference-aware profiling** method that produce accurate profiling tables by accounting for *intra-application interference*
  - an end-to-end static pipeline generator for edge SoCs
- Using *BetterTogether*, we implemented **3** class of applications
- Evaluated across **4** diverse devices, achieving up to **7.59x** (geo. **2.14x**)



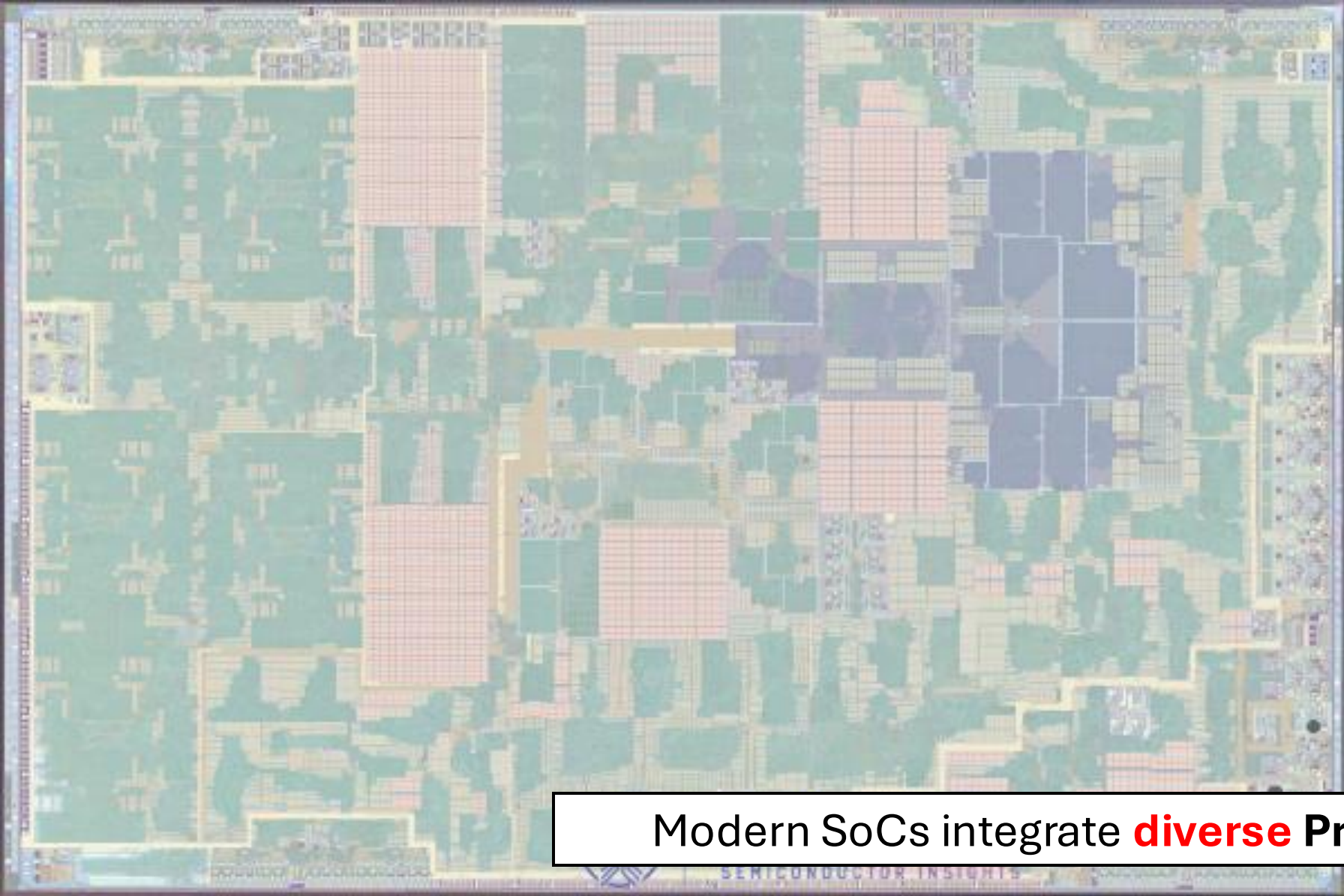
## Team

Yanwen Xu, Rithik Sharma, Zheyuan Chen  
Shaan Mistry, Tyler Sorensen  
{[yxu83](#), [riksharm](#), [zchen406](#), [sdmistry](#),  
[tyler.sorensen](#)}@ucsc.edu



**Open-Source Repo**  
[github.com/ucsc-redwood/better-together](https://github.com/ucsc-redwood/better-together)

# Backup slides

A detailed die shot of the Apple A18 SoC, showing a complex arrangement of various processing units in different colors (green, red, blue, yellow) and patterns. The chip is rectangular with a dense grid of components.

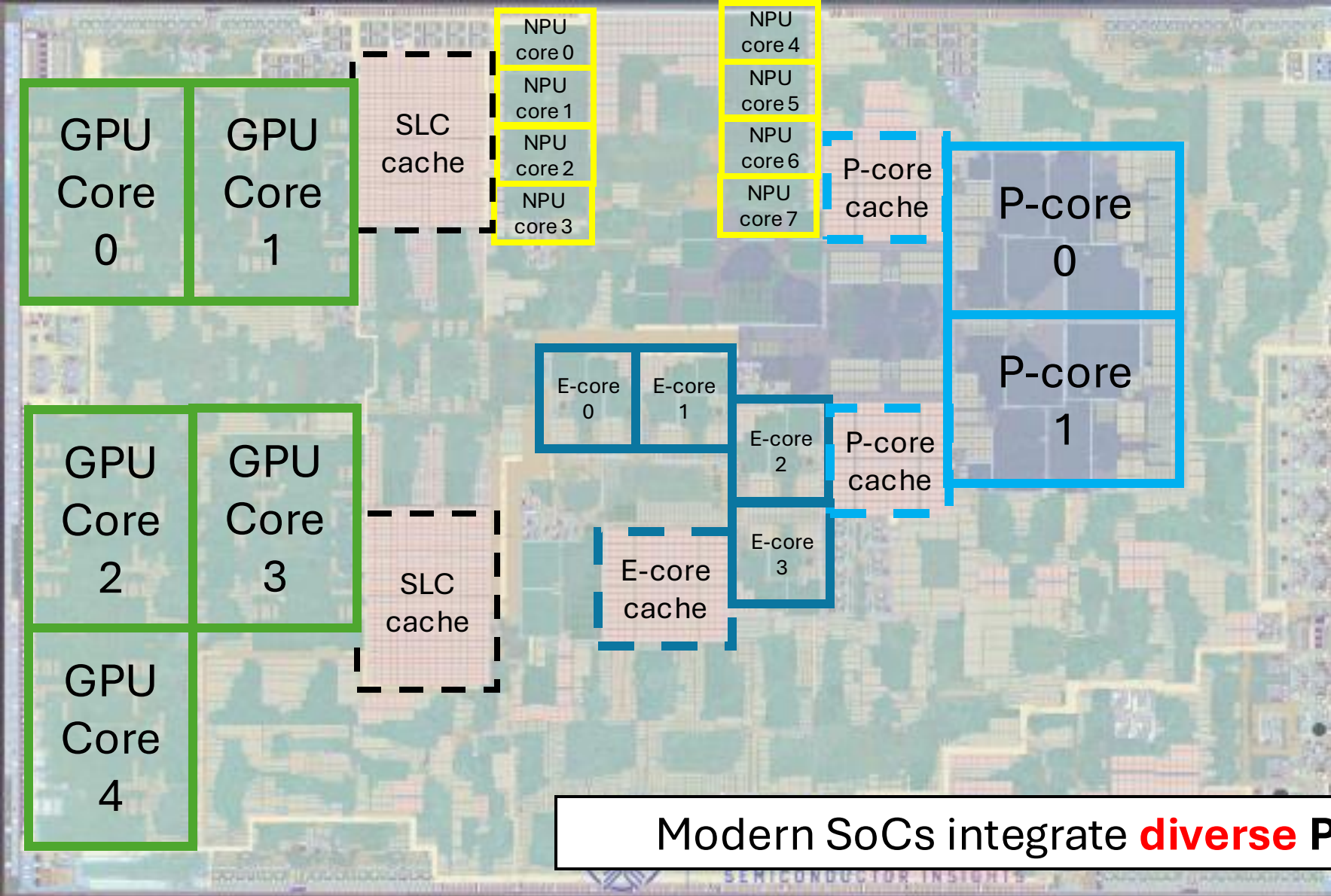
Apple A18 SoC  
90mm<sup>2</sup> @ TSMC N3E

- *big.LITTLE CPU*
- *Integrated GPU*
- *Domain Specific Accelerators (DSA)*
- ...

Modern SoCs integrate **diverse Processing Units (PU)**

Image by ChipWise: <https://chipwise.tech/our-portfolio/apple-a18-a18-pro-die-shot/>





Apple A18 SoC  
90mm<sup>2</sup> @ TSMC N3E

5 GPU-cores  
8 NPU-cores  
2 P-cores  
4 E-cores

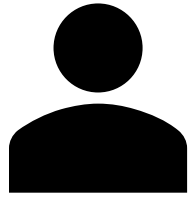
- *big.LITTLE CPU*
- *Integrated GPU*
- *Domain Specific Accelerators (DSA)*
- ...

Modern SoCs integrate **diverse Processing Units (PU)**

Image by ChipWise: <https://chipwise.tech/our-portfolio/apple-a18-a18-pro-die-shot/>

*We want to utilize all system resources (i.e., PUs)*

# Profiling-guided approach (isolated benchmarking)



Developer

- Given  $N$  stages
- Given  $M$  types of PU



Benchmarking each  $\{PU \times stage\}$  pair on the target system



*Get a  $N \times M$  profiling table*

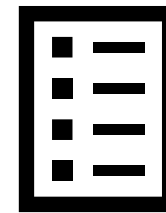
	$P_0$	...	$P_m$
$S_0$	1.2		2.5
...		...	
$S_N$	6.1		4.5



Optimize

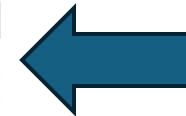


**Optimal Schedule**



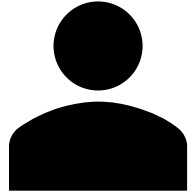
*implementation*

Deploy



*Expect  $X$  speedup*

# Profiling-guided approach (isolated benchmarking)



- Given  $N$  stages
- Given  $M$  types of PU

Developer

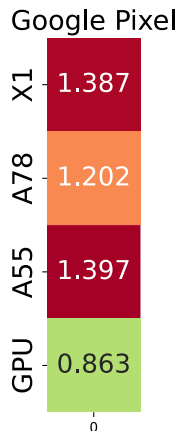


Benchmarking each  $\{PU \times stage\}$  pair on the target system



Get a  $N \times M$  profiling table

**Inaccurate profiling**



	$P_0$	...	$P_m$
$S_0$	1.2		2.5
...		...	
$S_N$	6.1		4.5



Optimize

**Inaccuracy propagate**



**Suboptimal Schedule**

Deploy



Expect ~~X~~ speedup

**e.g., Only X/2 speedup**

implementation



# BT Implementor

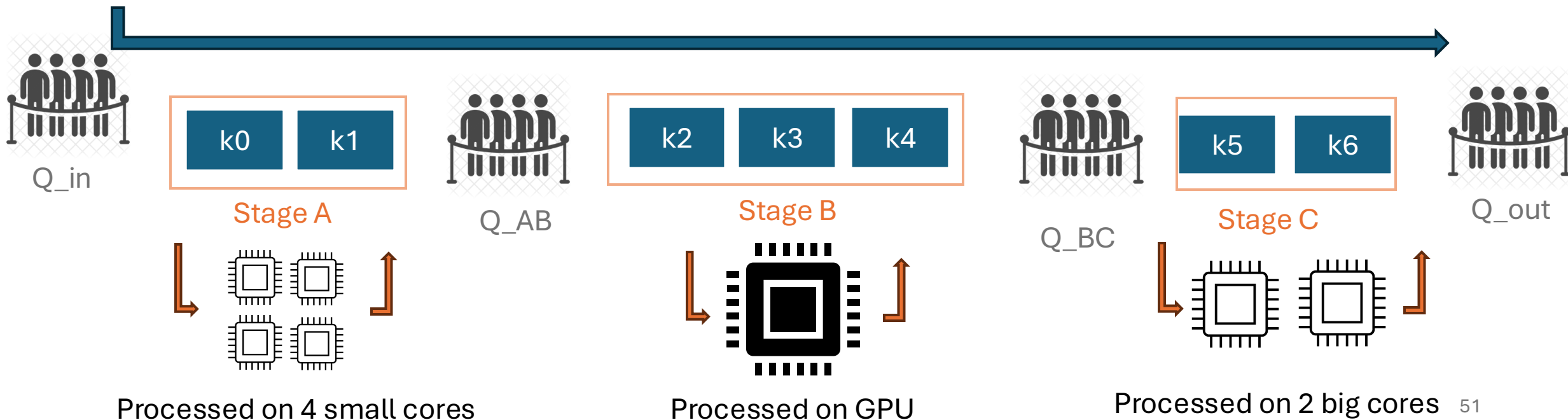
- We define Task
  - light weight, pointing to CPU/iGPU shared memory

```
struct Task {  
    uint32_t uid;  
    float* u_input;  
    float* u_output;  
    size_t n;  
};
```

- Using concurrent Queue to pass Tasks around stages.

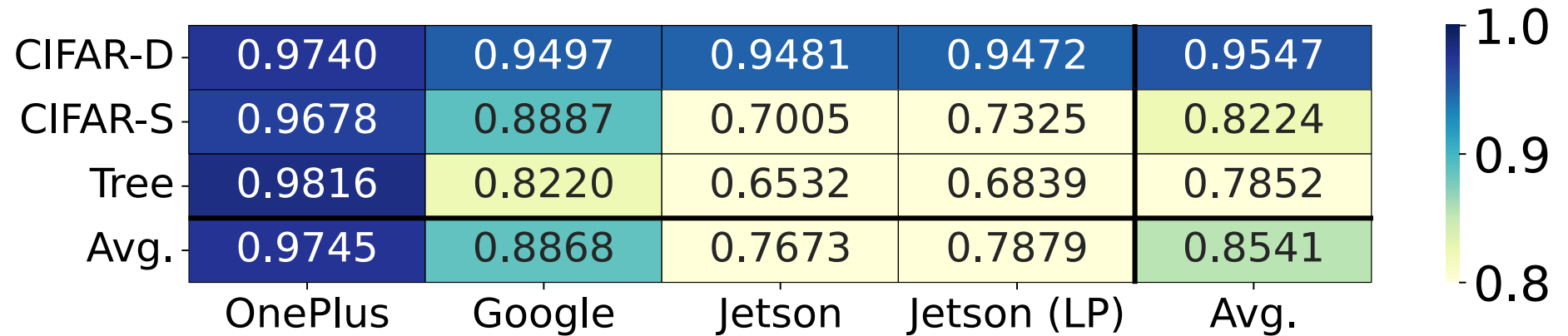
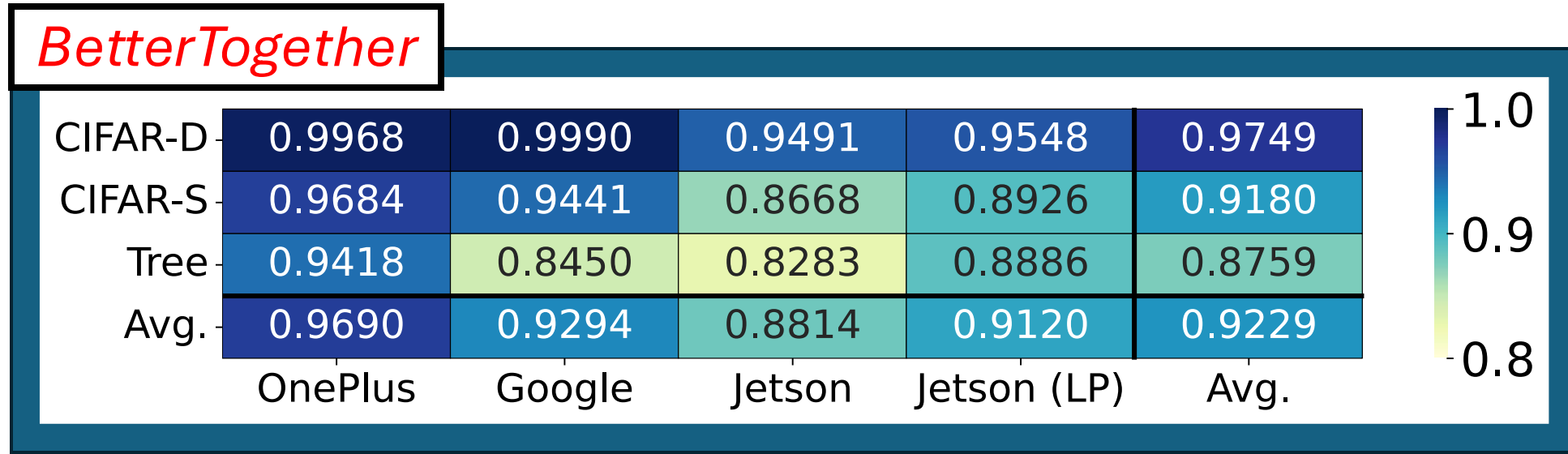
Each chunk process the incoming tasks in respective type for Cores

CPU use thread affinity





- BetterTogether yields higher correlations



Correlation (0.0–1.0) between predicted and actual times across all applications and platforms. **Higher is better.**