

Squeezing Operator Performance Potential for the Ascend Architecture

Yuhang Zhou¹, Zhibin Wang¹, Guyue Liu², Shipeng Li¹, Xi Lin¹, Zibo Wang¹, Yongzhong Wang³,
Fuchun Wei³, Jingyi Zhang³, Zhiheng Hu³, Yanlin Liu³, Chunsheng Li³, Ziyang Zhang³,
Yaoyuan Wang³, Bin Zhou⁴, Wanchun Dou¹, Guihai Chen¹, Chen Tian¹

¹ Nanjing University ² Peking University ³ Huawei Technologies Co., Ltd. ⁴ Shandong University



南京大學
NANJING UNIVERSITY





Outline



Introduction



System Design



Case Study



Evaluation



Conclusion





Outline



Introduction



System Design



Case Study



Evaluation



Conclusion



AI Domain-Specific Architecture (DSA)

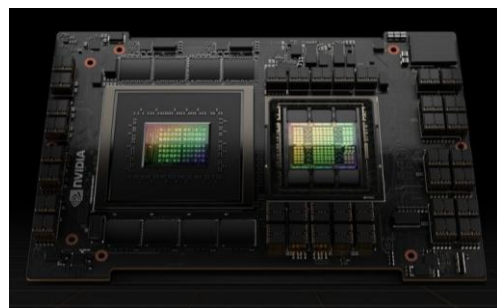


Deep learning models



Better arithmetic support

Domain-specific architecture



NVIDIA GPU



Google TPU

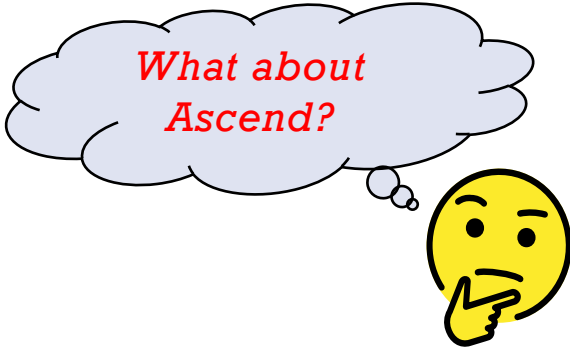
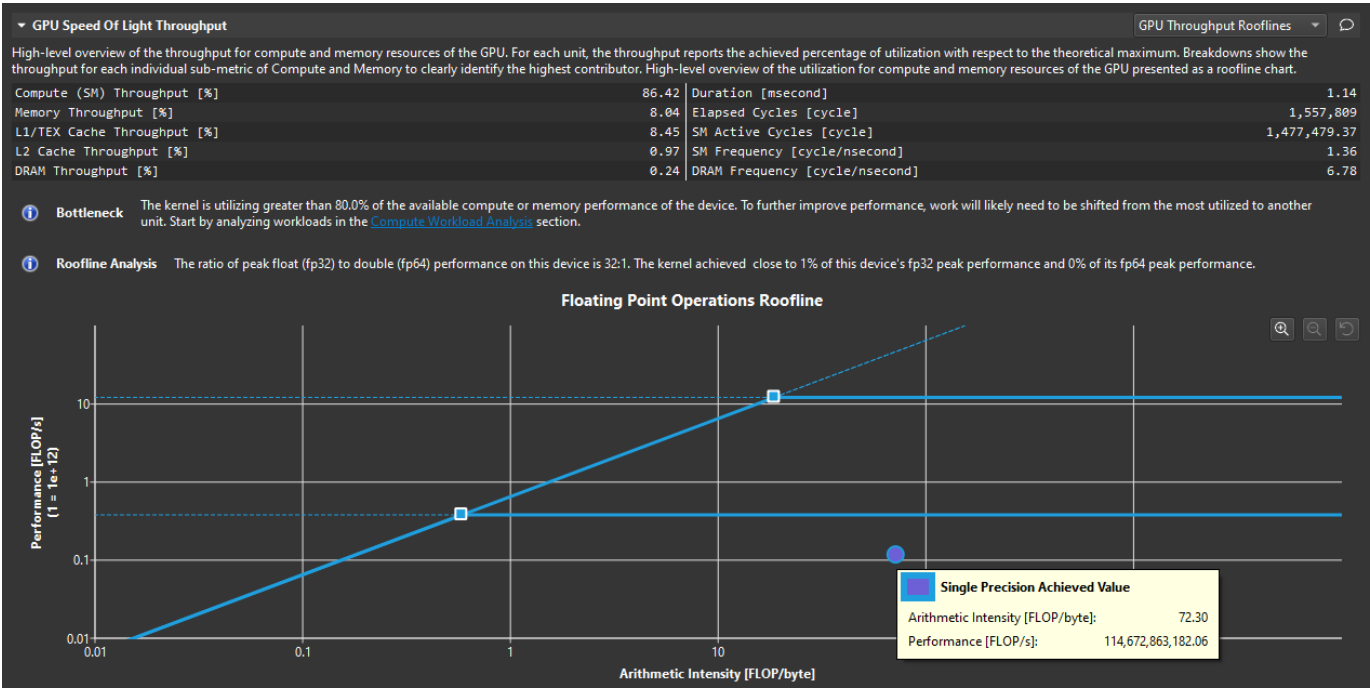
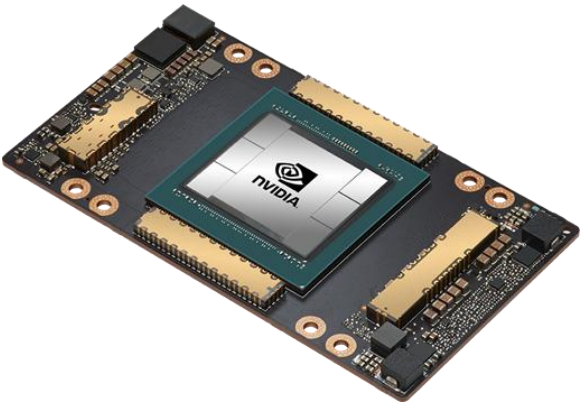


Huawei Ascend NPU

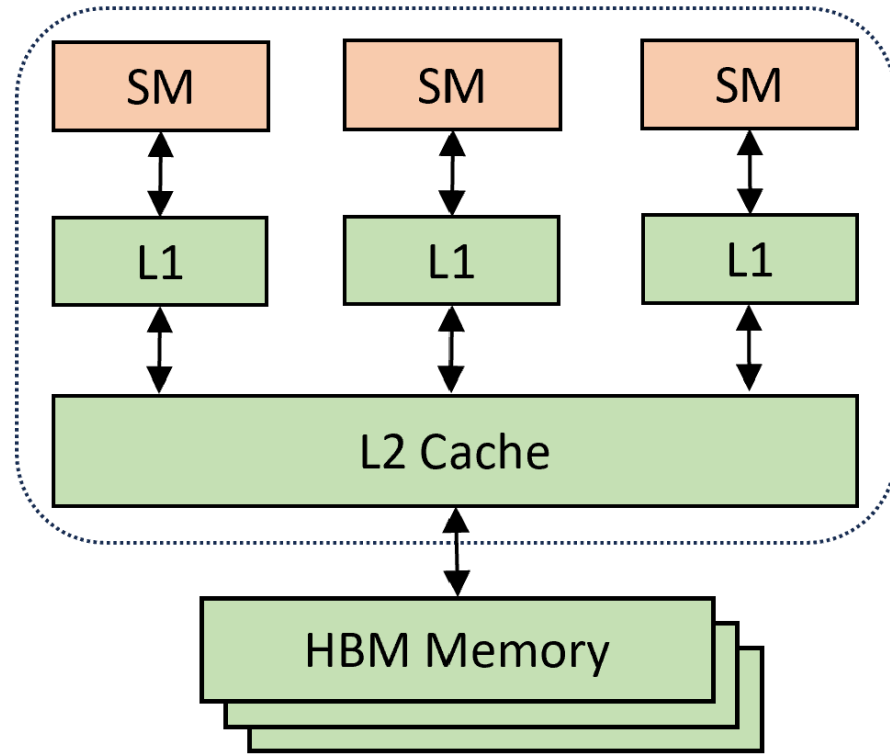


Cambricon MLU

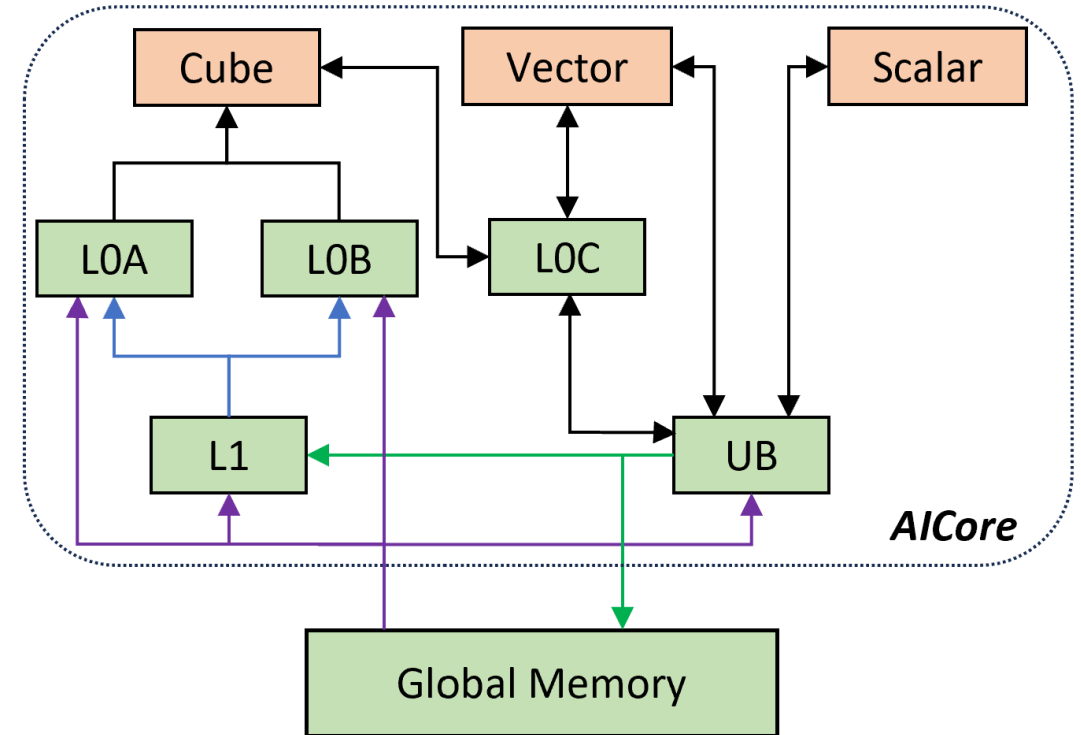
Operator Optimization Needs Profiling



Ascend Architecture



GPU (NVIDIA)

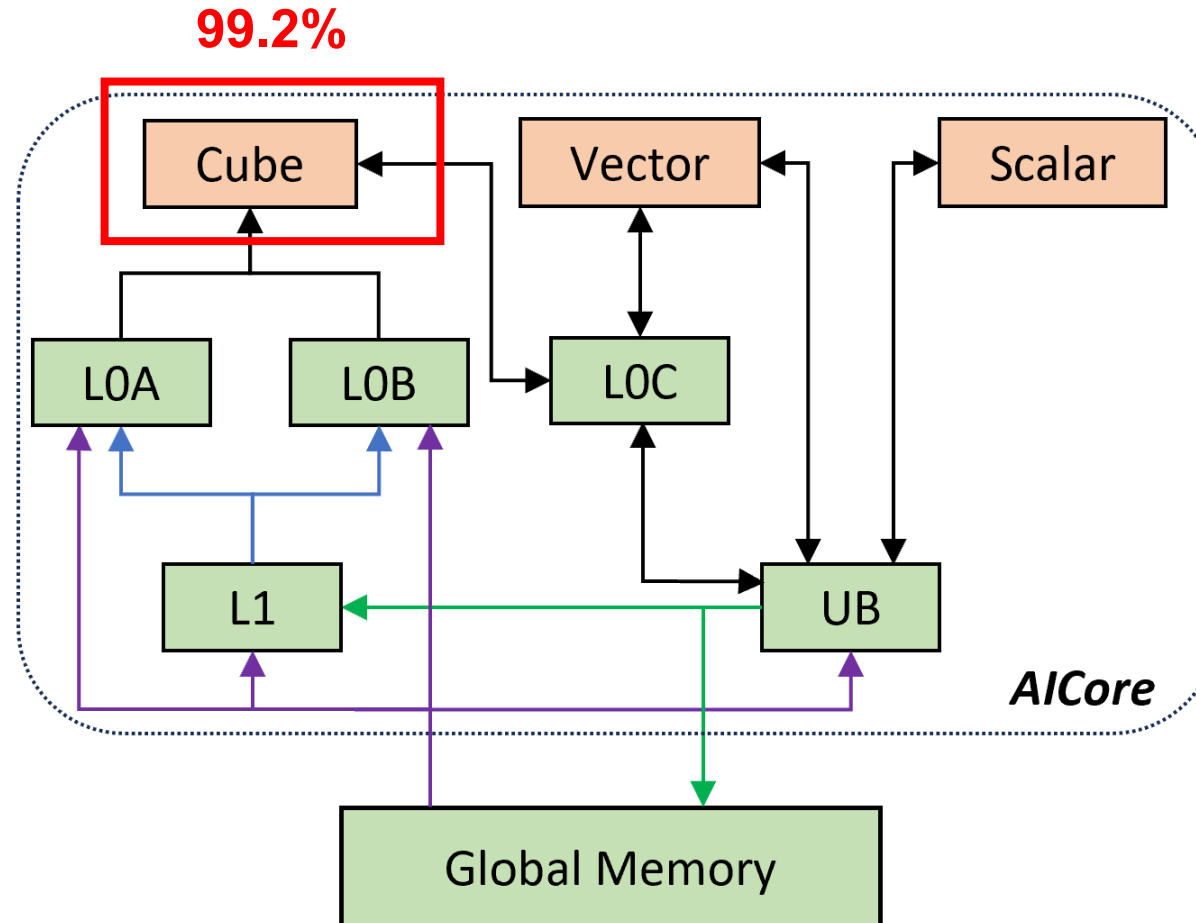
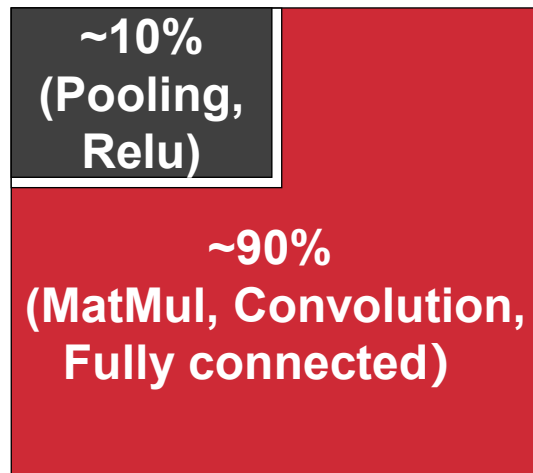


NPU (Ascend)

Compute Unit Memory Unit

Dedicated Compute Units

Transformer computation



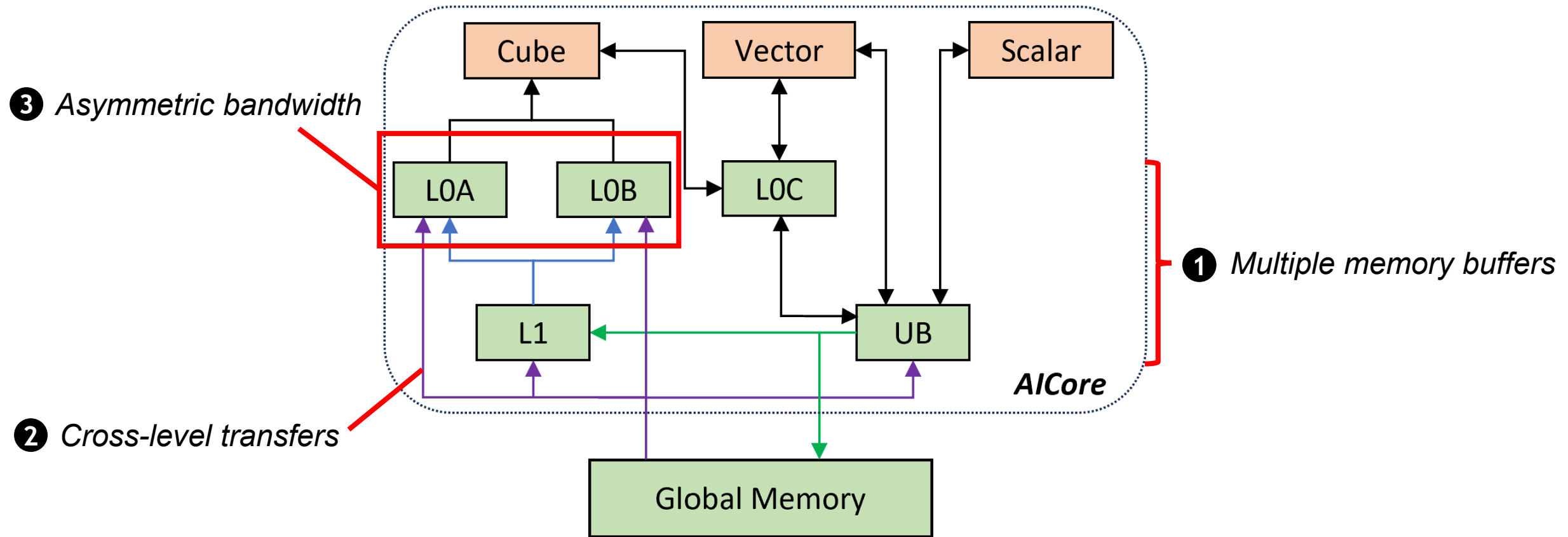
Computing precision

Cube: INT8/FP16

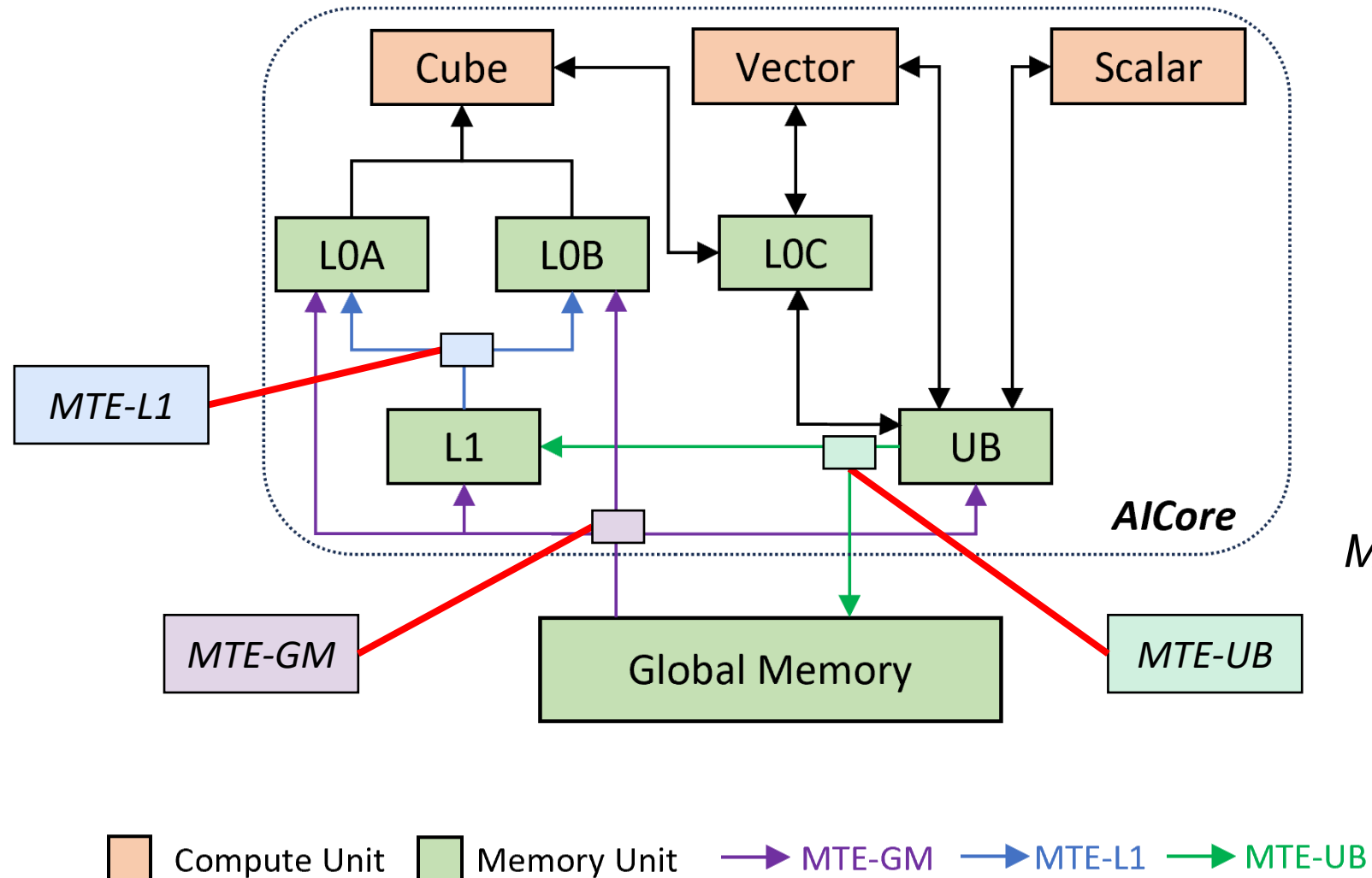
Vector: INT8/FP16/FP32

Scalar: INT32/FP32/...

Customized Memory Architecture

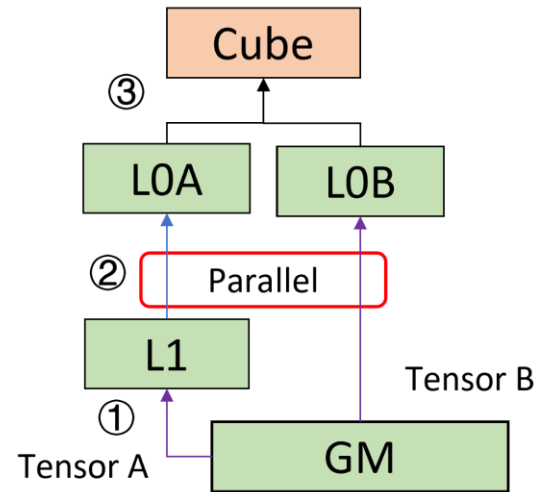
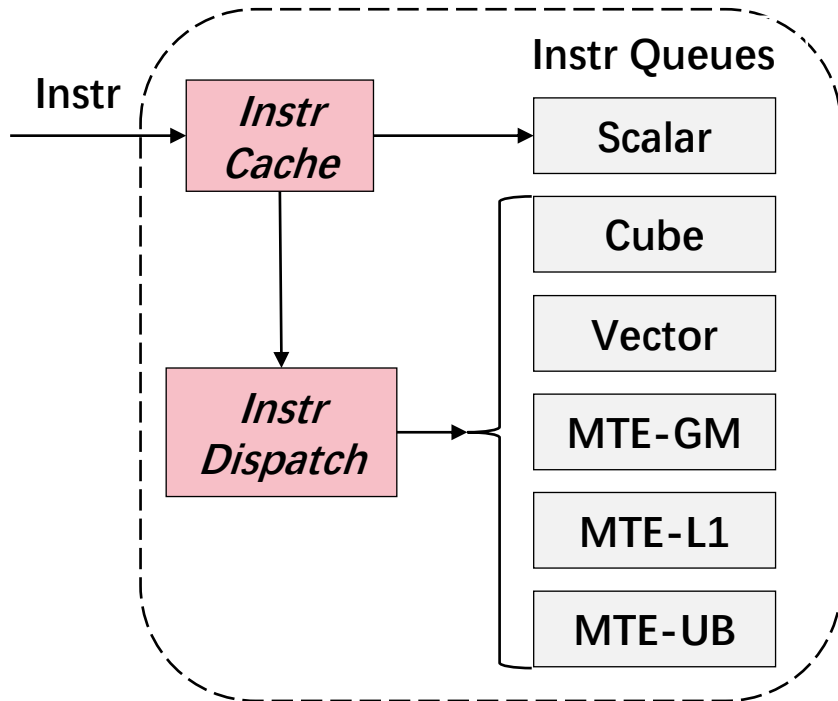


Efficient Transfer Control Units

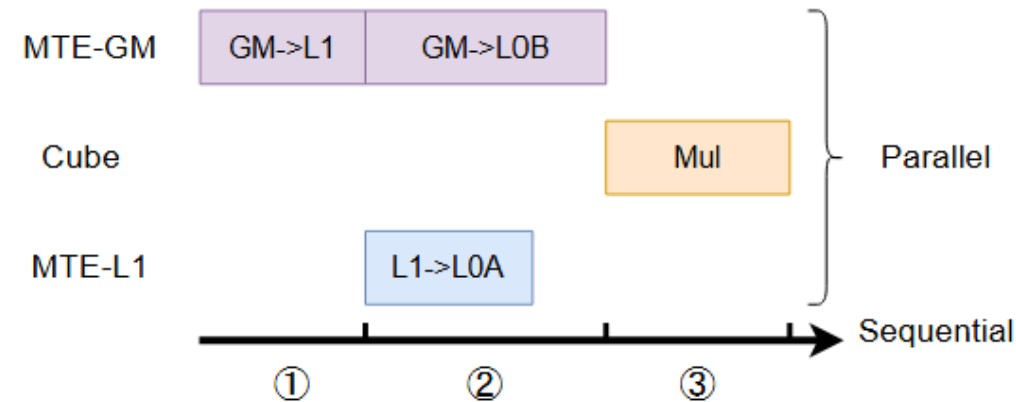


Memory transfer engine

Instruction Pipeline



Example of matrix multiplication $A \times B$



Inter parallelism, intra serialism.

Summary of Ascend Architecture

Pros

Ascend Architecture

Cons

Accurately identifying operator bottleneck is a challenging, but essential task!

Operational flexibility

Efficient transfer control and instruction pipeline

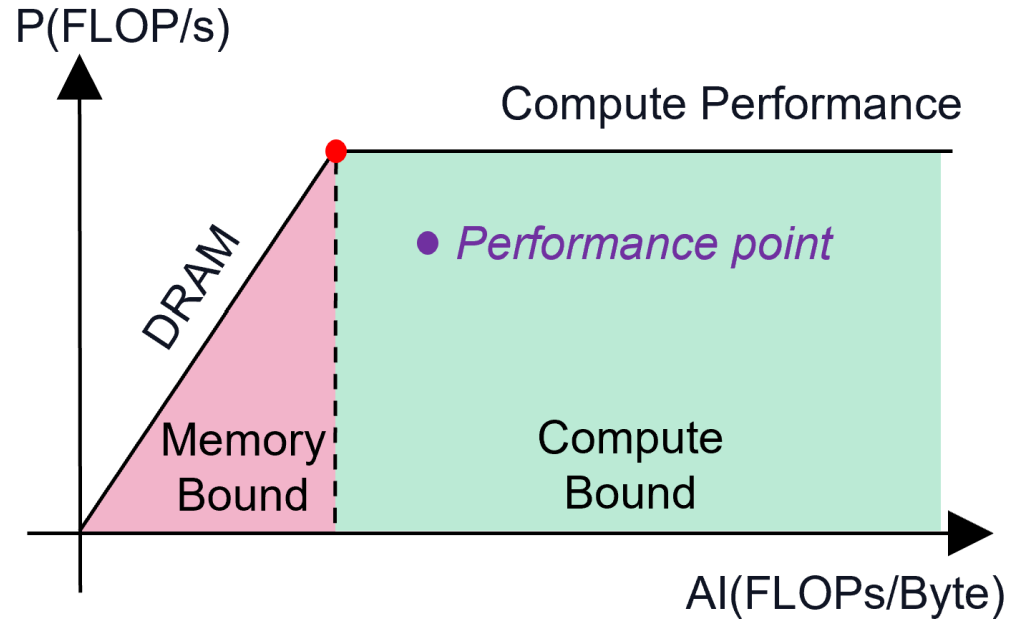
Inputting

ALU

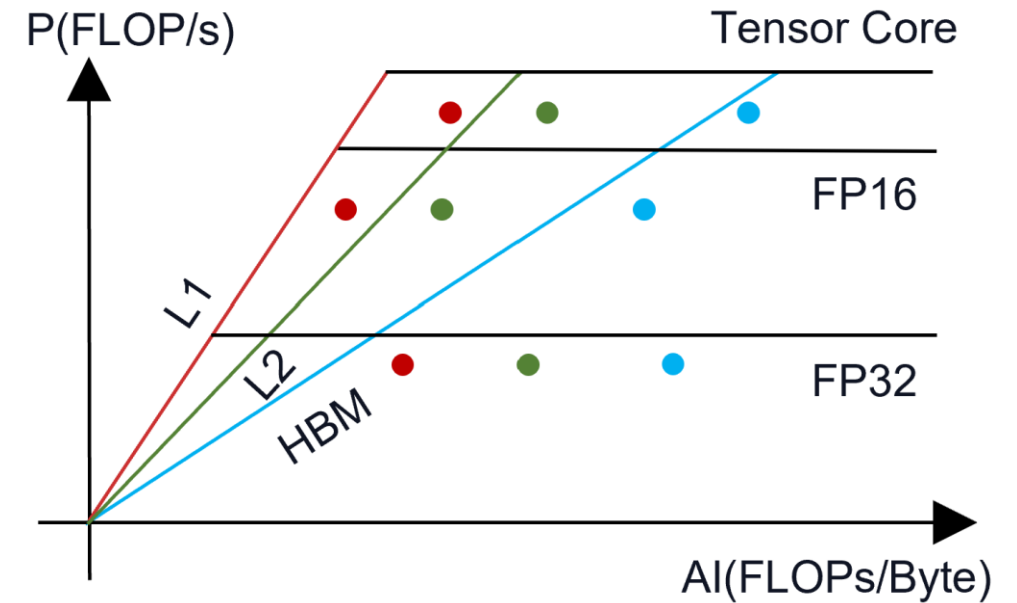
pipeline

...

Existing Operator Performance Analysis

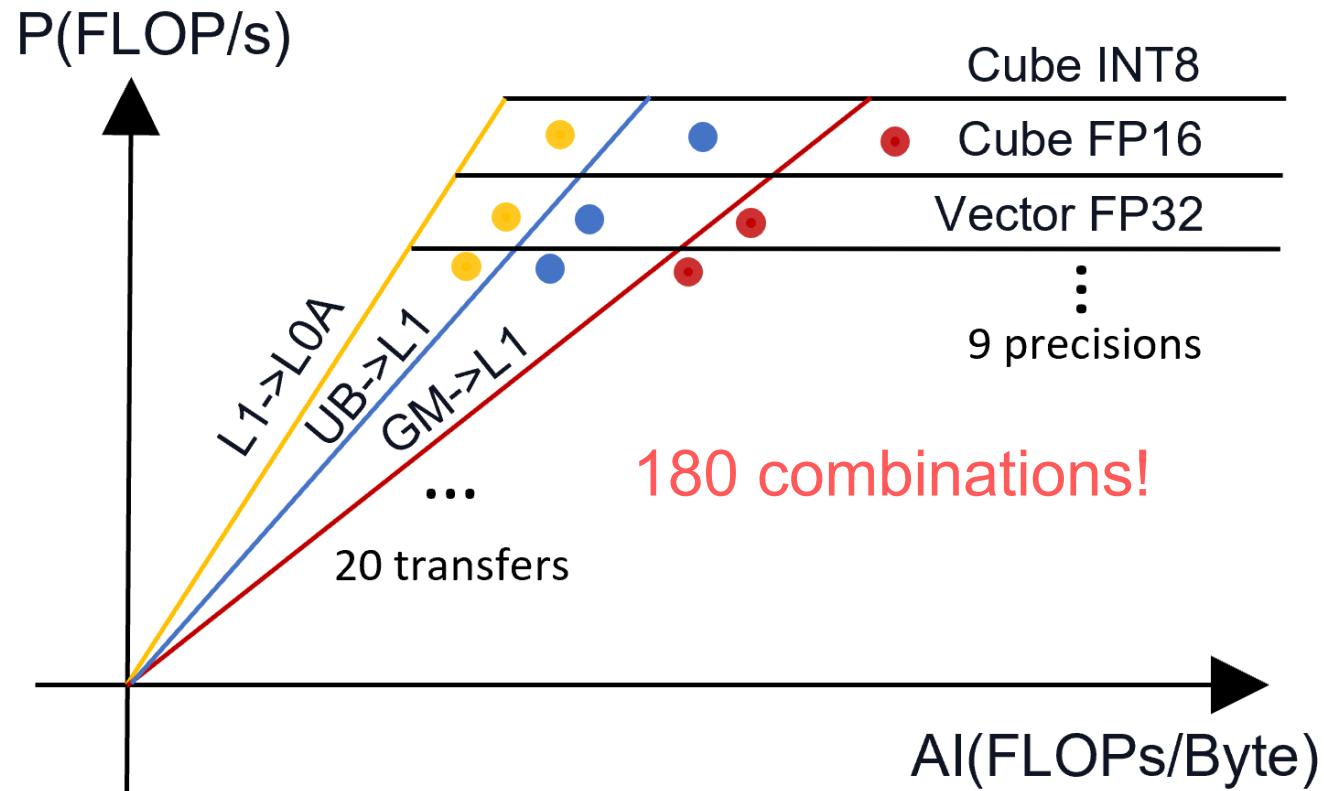


DRAM Roofline



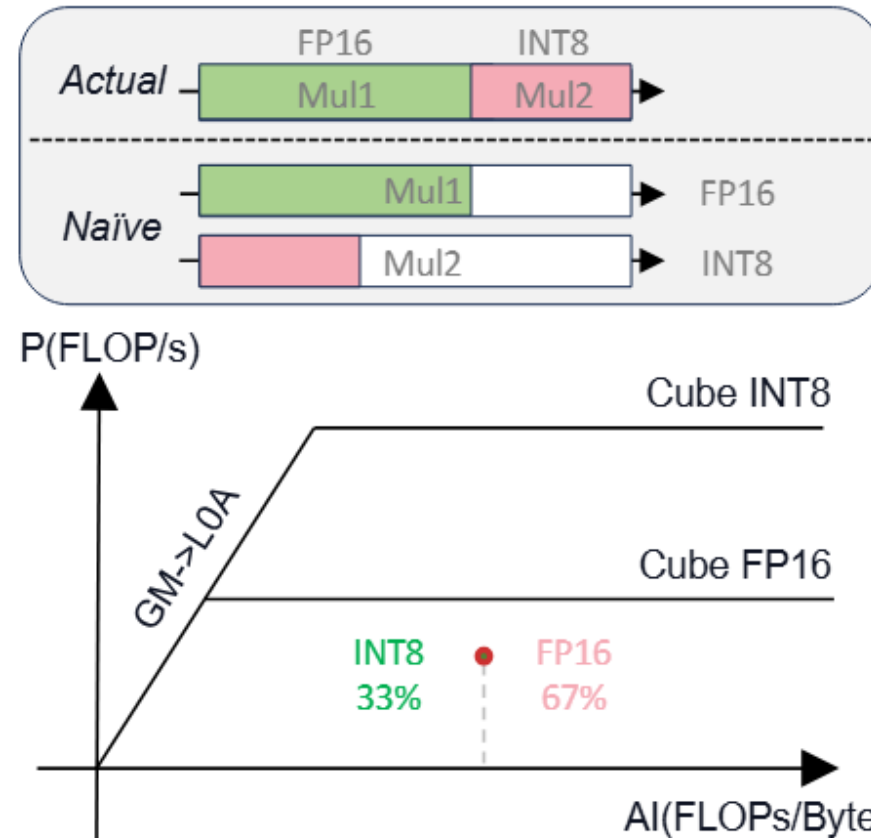
Hierarchical Roofline

Limitations of Performance Analysis



(i) Massive combinations between precisions and transfers

Limitations of Performance Analysis



Underutilization?



(ii) Incorrect analysis by ignoring the sequential execution

Our Goals

Problems

Complicated analysis

Incorrect analysis

Non-trivial optimization

Solutions

Simplified and accurate operator performance analysis

Effective guidance for operator optimization



Outline



Introduction



System Design



Case Study



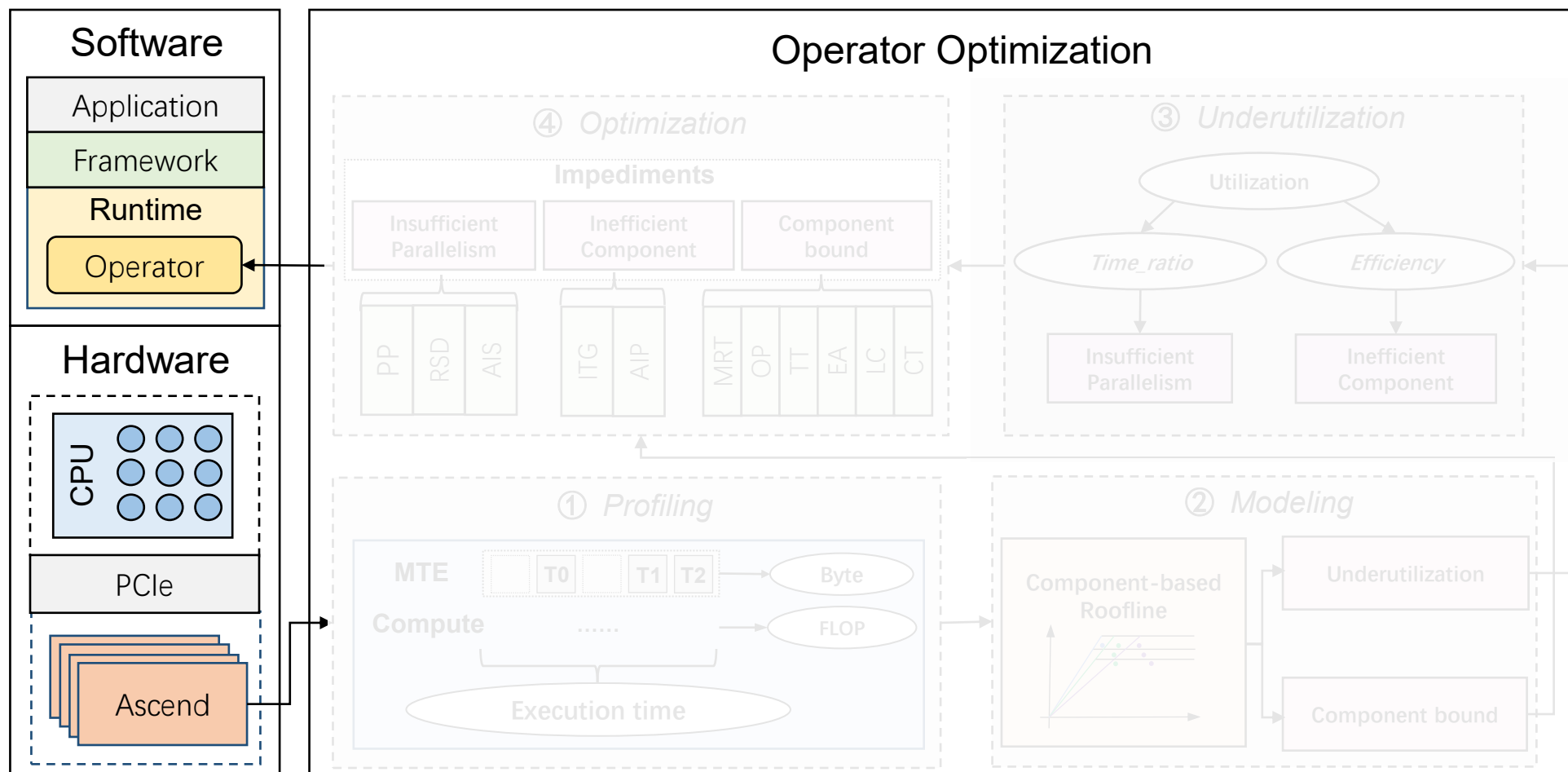
Evaluation



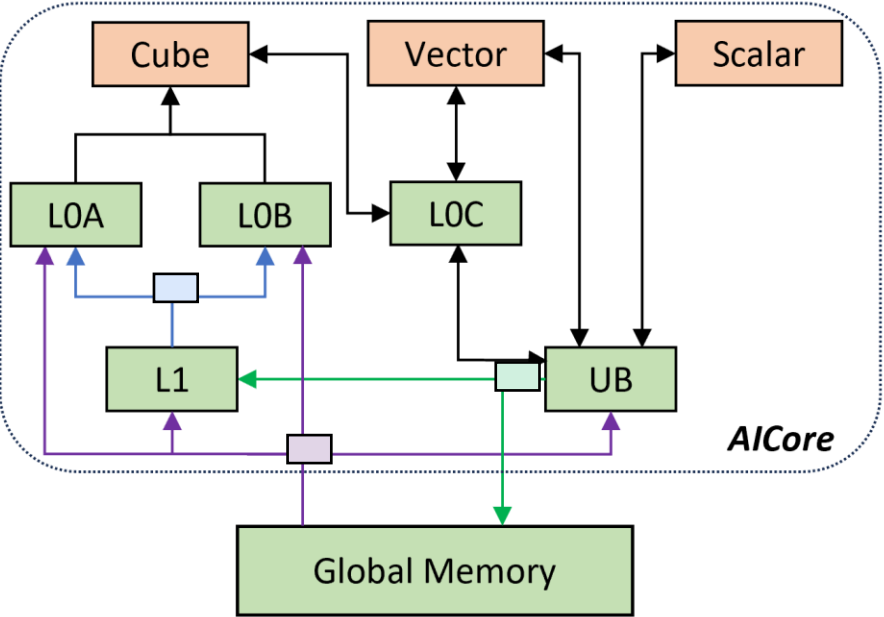
Conclusion



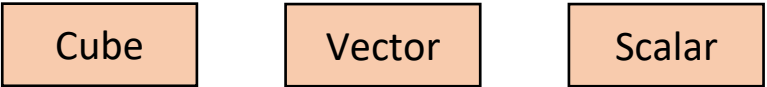
Overview



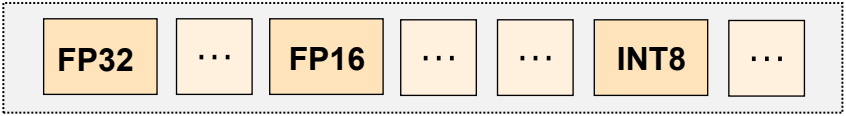
Profiling and Component Abstraction



Compute component

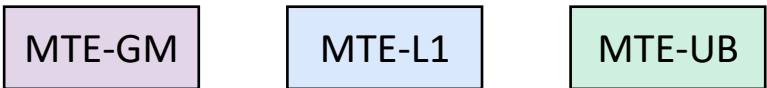


Compute Queue

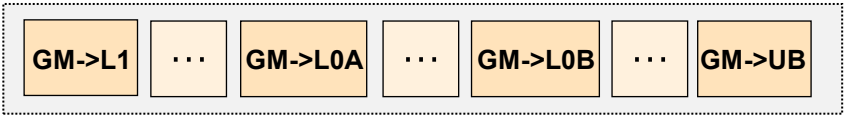


Operations for the precision

MTE component



MTE Queue

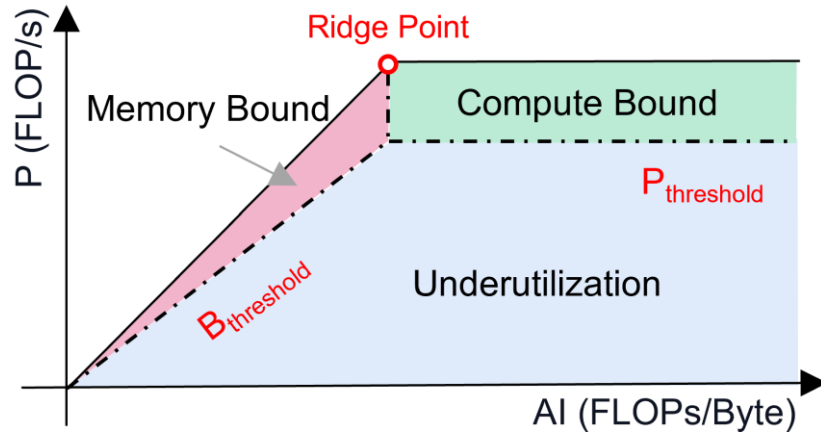


Bytes for the transfer



Actual execution time of the component

Component-based Roofline Model



Utilization of component can reflect the operator's bottleneck.

Operator-aware ideal performance

$$U_{cube} = \frac{A_{cube}}{I_{cube}} \begin{cases} A_{cube} = \frac{O_{cube}}{T_{total}} & \text{Profiling} \\ I_{cube} & \text{Different precisions?} \end{cases}$$

$$I_{cube} = \frac{\sum_{prec} O_{prec}}{\sum_{prec} \frac{O_{prec}}{P_{prec}}} \quad \begin{matrix} \text{Harmonic Mean} \\ \text{Arithmetic Power} \end{matrix}$$

Underutilization Analysis

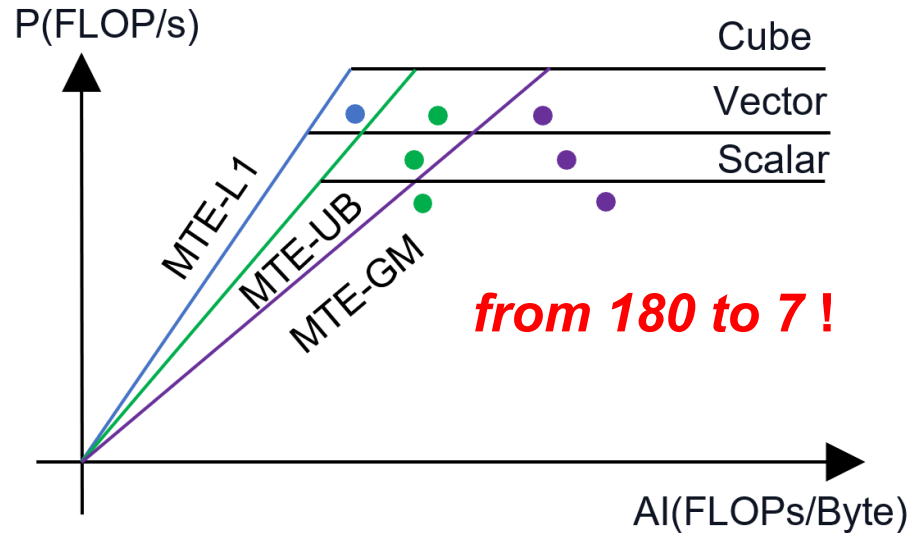
$$\textcircled{1} U_{cube} = \frac{A_{cube}}{I_{cube}} = \frac{O_{cube}}{\underbrace{T_{cube} \cdot I_{cube}}_{E_{cube}}} \cdot \frac{T_{cube}}{\underbrace{T_{total}}_{R_{cube}}}$$

$$\textcircled{2} E_{component} \leq \frac{R_{threshold}}{U_{threshold}} \quad \swarrow \quad \searrow \quad R_{component} < R_{threshold}$$

Inefficient Component **Insufficient Parallelism**

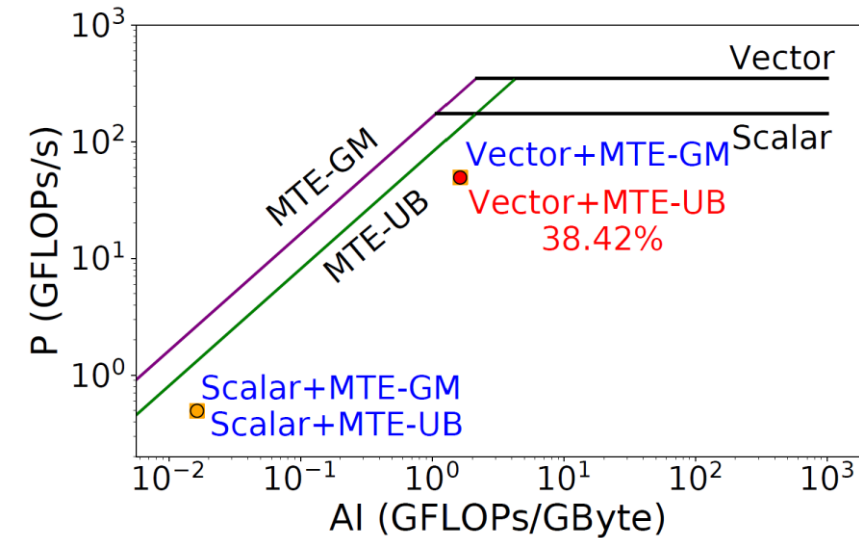
Pruning, Visualization and Analysis

Pruning results



- ✓ Component abstraction
- ✓ Remove irrelevant components
- ✓ Remove impossible combinations

Roofline Analysis of Add_ReLU Operator



$U_{component}$ of Vector+MTE-UB (38.42%):
Underutilization
 $R_{component}$ of MTE-GM (58.68%):
Insufficient parallelism



Outline



Introduction



System Design



Case Study



Evaluation

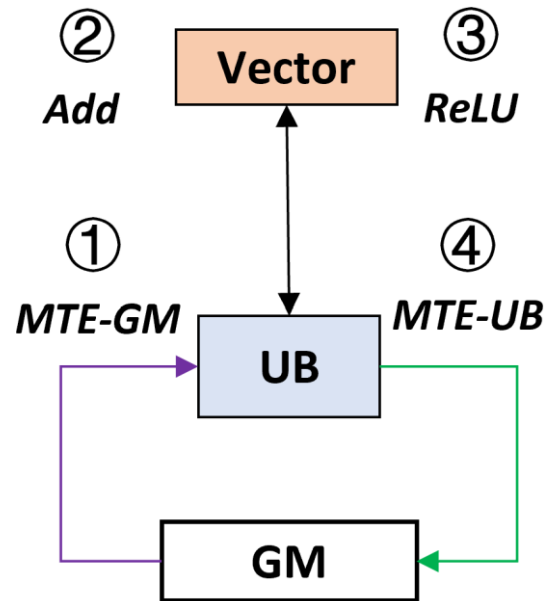


Conclusion

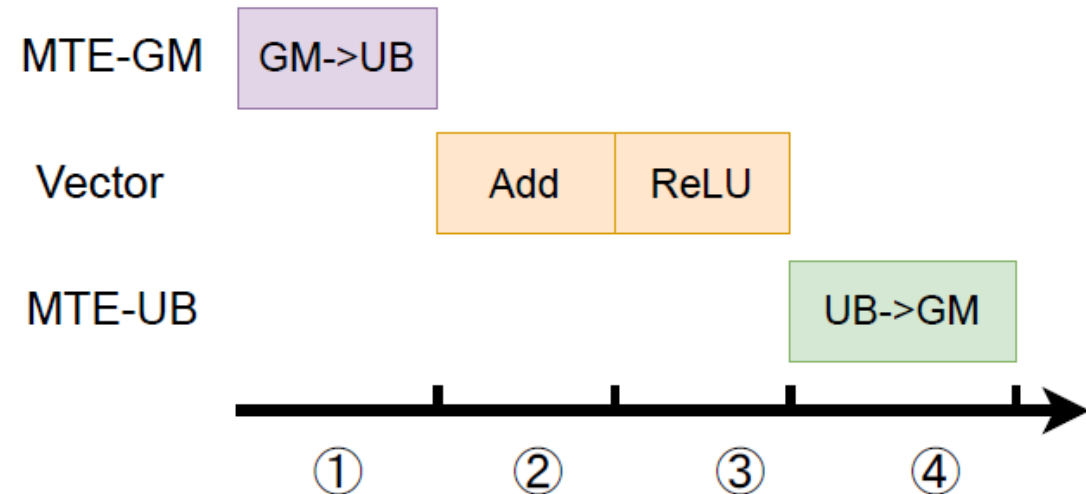


Case Study: Optimization of *Add_ReLU* Operator

$$\text{Add_ReLU}(x) = \text{ReLU}(x + c)$$



Data flow



Instruction timeline

Iteration 1: *Reducing spatial dependency*

Original Code

```

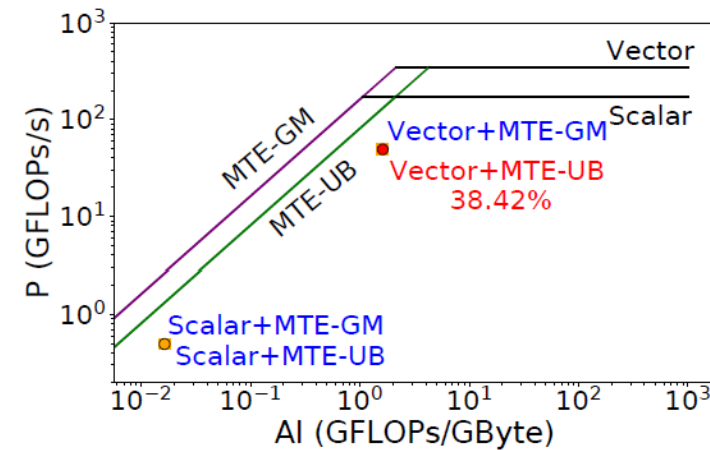
① ...
② ub_to_gm(gm_1, ub_1);
③ gm_to_ub(ub_1, gm_2);
④ ...
    
```



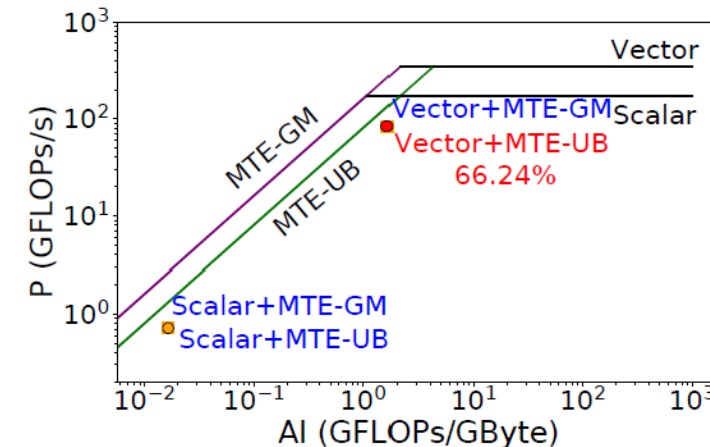
Optimized Code

```

① ...
② ub_to_gm(gm_1, ub_2);
③ gm_to_ub(ub_1, gm_2);
④ ...
    
```



**Insufficient parallelism
(38.42%)**



**MTE-UB bound
(66.24%)**

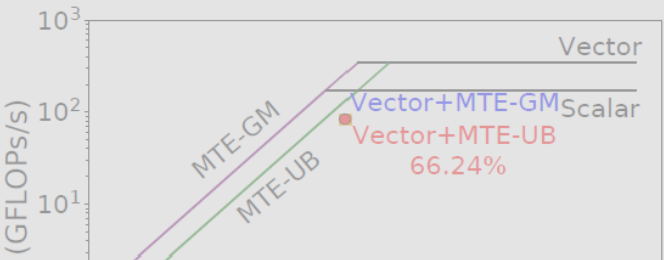
Iteration 2: *Minimizing redundant transfer*

Original Code

```

① for i = 1 to n do
②   gm_to_ub(ub_1, c);
③   ...

```



MTE-UB bound
(66.24%)

The single operator time reduced by **1.73×**.

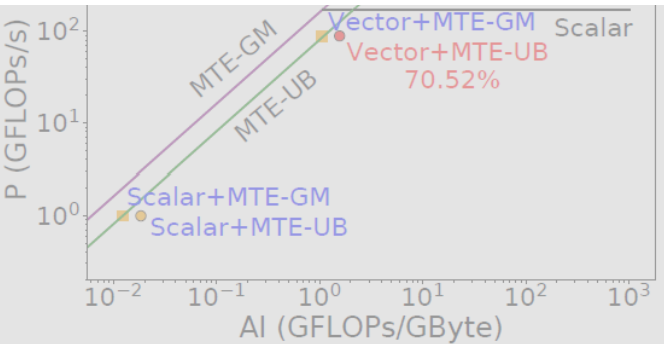
The *component_utilization* up by **32.1%**.

The total inference latency down by **244.261 μs**.

```

① gm_to_ub(ub_1, c);
② for i = 1 to n do
③   ...
④ end for

```



MTE-UB bound
(70.52%)

Optimization Experience

We summarize the common bottleneck causes and optimization strategies.

Bottleneck Cause	Compute Bound	MTE Bound	Insufficient Parallelism	Inefficient Compute	Inefficient MTE		
Strategy	Operator	Bottleneck Cause and Optimization Strategy					Speedup
		Compute Bound	MTE Bound	Insufficient Parallelism	Inefficient MTE	Inefficient Compute	
	Add_ReLU		MRT	RSD			1.72
	Depthwise		MRT	AIS,RUS,PP	ITG		1.26
	AvgPool					AIP	4.31
	Mul			RSD			1.34
	Conv2D		MRT	RSD			2.65
	FullyConnection				ITG		1.22
	MatMul		OF				1.10
	GeLU	EA					1.06

In MobileNetV3 inference, Our operator optimizations perform well with speedups of 1.06-4.31×.

More cases can be found in the paper.



Outline



Introduction



System Design



Case Study



Evaluation



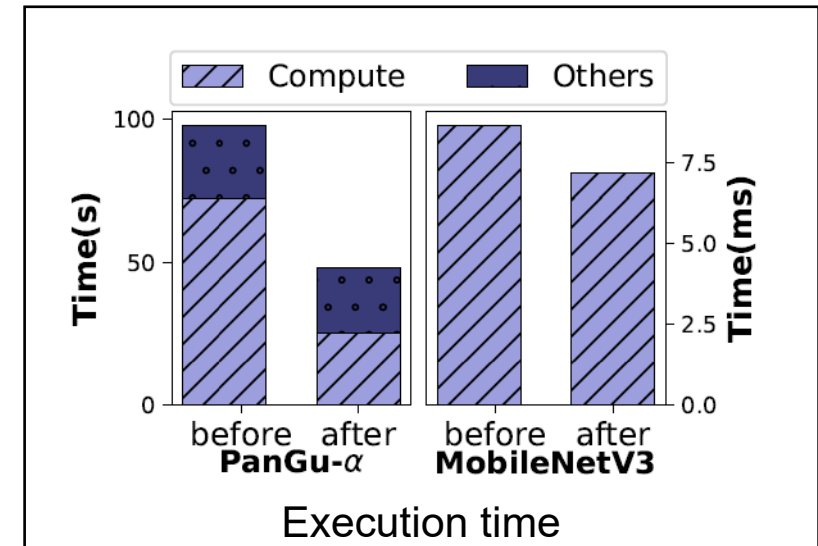
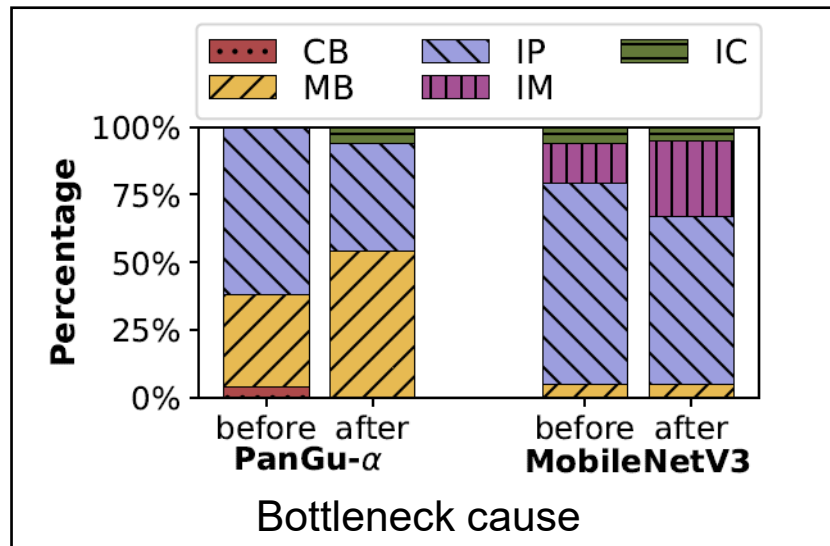
Conclusion



Evaluation on End-to-End Optimization

Device: Ascend 910 (Training); Ascend 310 (Inference)

Workloads: 100B PanGu- α (Training); MobileNetV3 (Inference)



Training

The ratio of *insufficient parallelism* reduced by **21.38%**.

The *iteration time* speedup is **2.04x**.

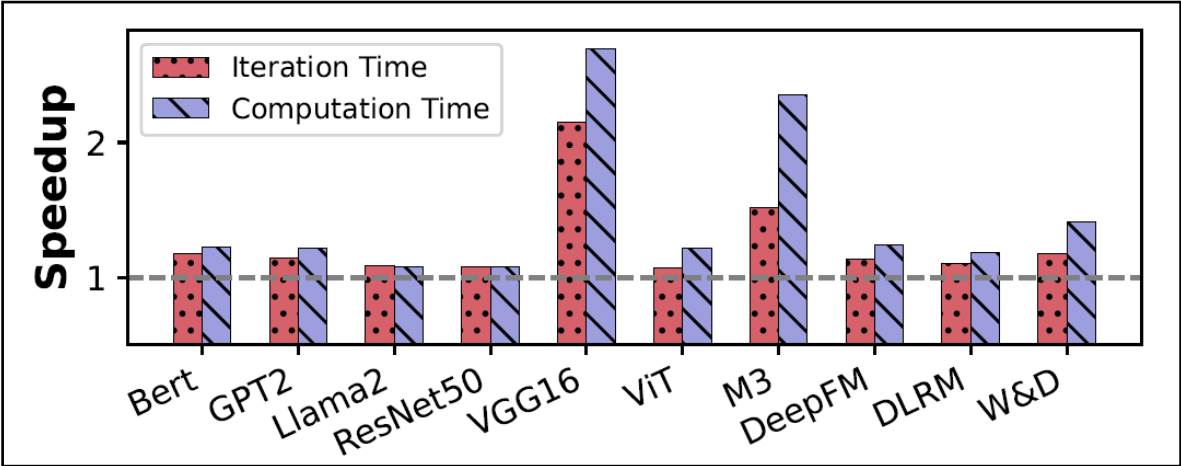
Inference

The ratio of *insufficient parallelism* reduced by **11.61%**.

The *total time* speedup is **1.21x**.

Overall Optimization Results

Type	Model	Parameter	Dataset	#NPUs
Vision	MobileNetV3(M3)	5.4M	ImageNet2012	8
	ResNet50	25.6M		
	ViT	86M		
	VGG16	138.4M		
NLP	Bert	110M	WikiText2	8
	GPT2	355M		
Recommendation	DeepFM	16.5M	Criteo	8
	Wide and Deep(W&D)	75.84M		
	DLRM	540M		
LLM	Llama 2	7B	WikiText2	8
	PanGu- α	100B	1.1TB Chinese Dataset	128



Our optimizations cover **11** different models.

Computation time speedups range from **1.08-2.7×**.

Iteration time speedups range from **1.07-2.15×**.



Outline



Introduction



System Design



Case Study



Evaluation



Conclusion



Conclusion

1. We propose a component-based roofline model and underutilization analysis to identify the operator bottlenecks on Ascend.
2. Through in-depth operator optimization case studies, we guide users on how to complete optimization.
3. Based on extensive practical optimization experiments, we share our practical insights and valuable experiences.

Future Work

1. The component-based roofline model can extend to other DSAs like TPU.
2. Depth studies of hardware architecture, especially its interaction with the software.





Thanks

Q&A

yuhangzhou@smail.nju.edu.cn

