

Tender: Accelerating Large Language Models via **Tensor Decomposition** and **Runtime Requantization**

Jungi Lee*, Wonbeom Lee*, Jaewoong Sim
Seoul National University

* Equal Contribution



Outline

- **Motivation**

- Challenges in Efficient LLM Inference
- Limitations of Prior Works

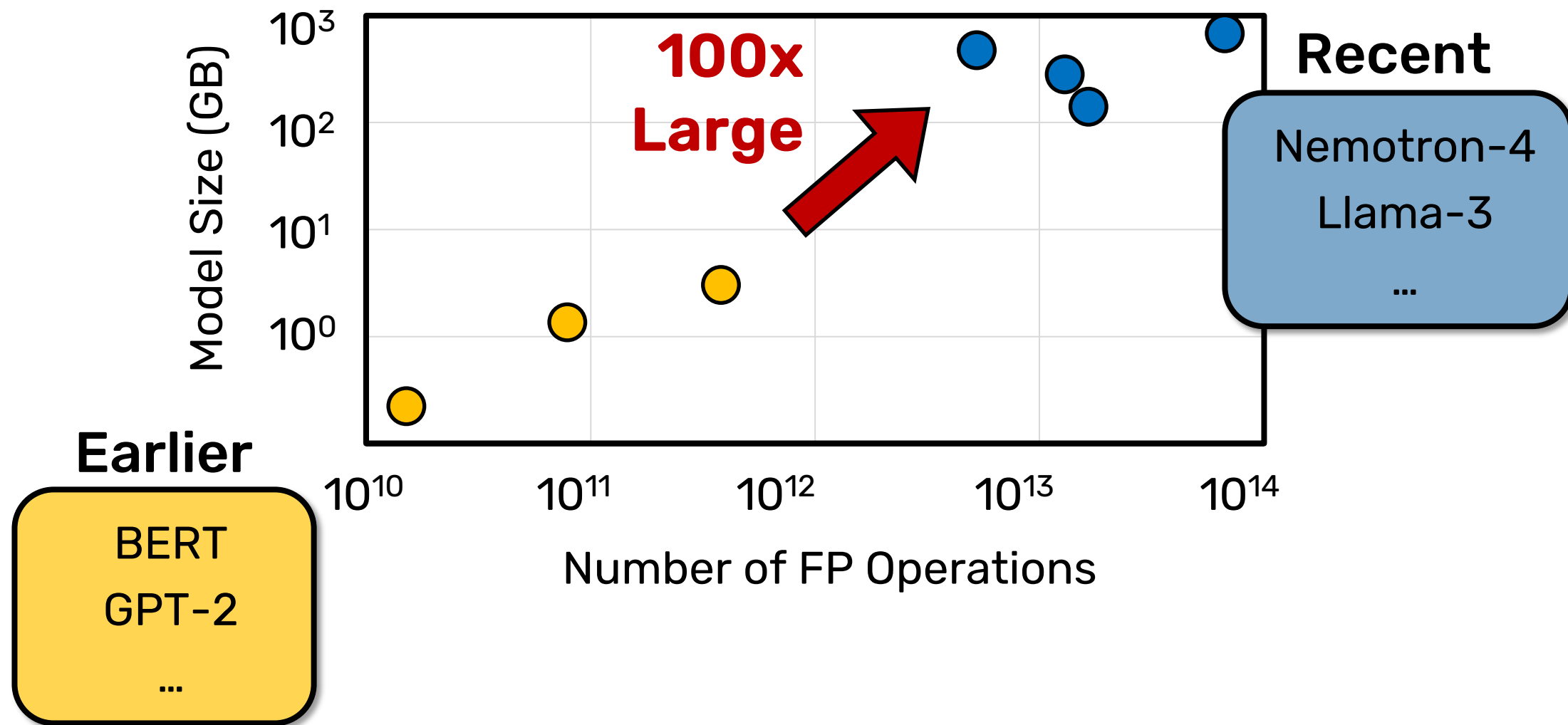
- **Tender: Algorithm-Hardware Co-design for Efficient LLM Inference**

- Tensor Decomposition
- Rescaling Operation

- **Evaluation**

- **Conclusion**

Challenges in LLM Inference



Challenges in LLM Inference



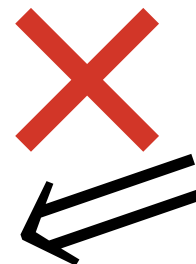
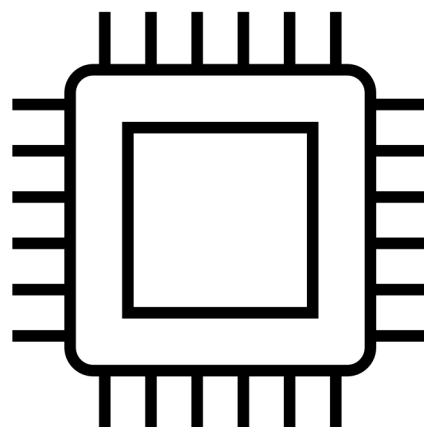
How do we serve LLMs efficiently?

Challenges in LLM Inference

Quantize **Both**
Weights & Activations



Integer pipeline



Recent

Nemotron-4
Llama-3

...

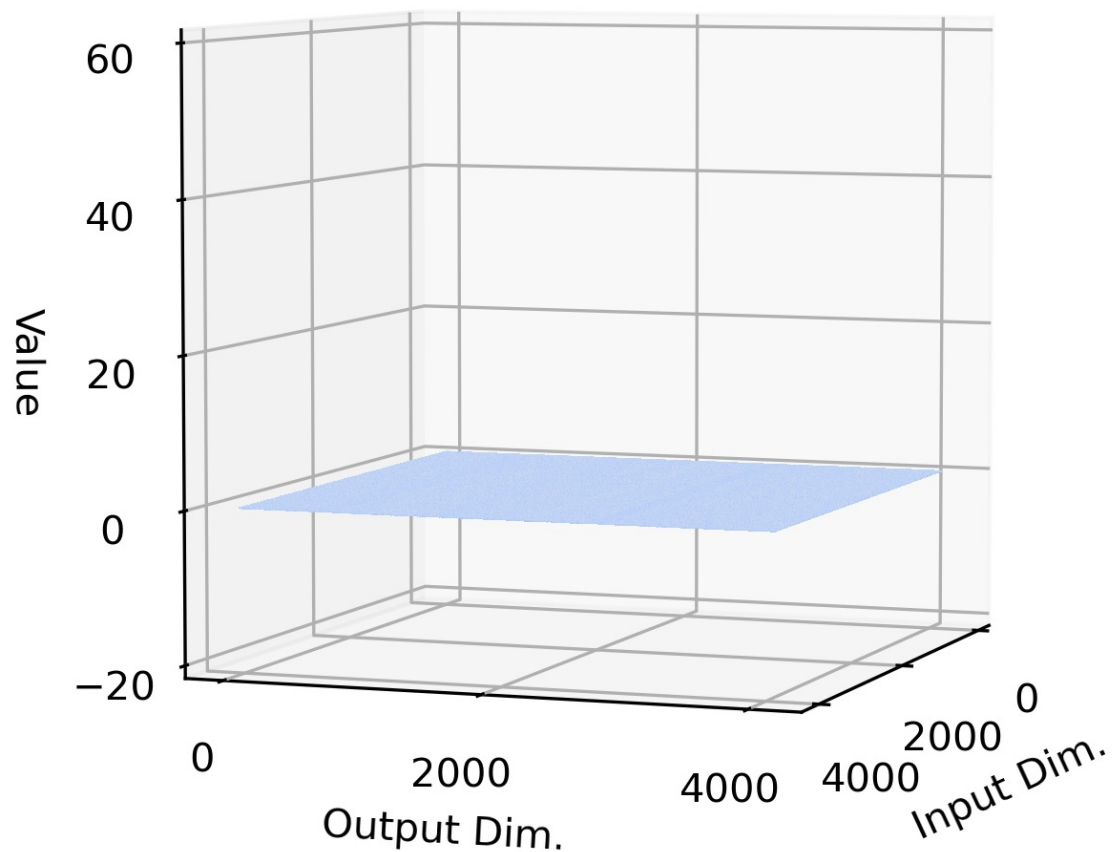
Earlier

BERT
GPT-2

...

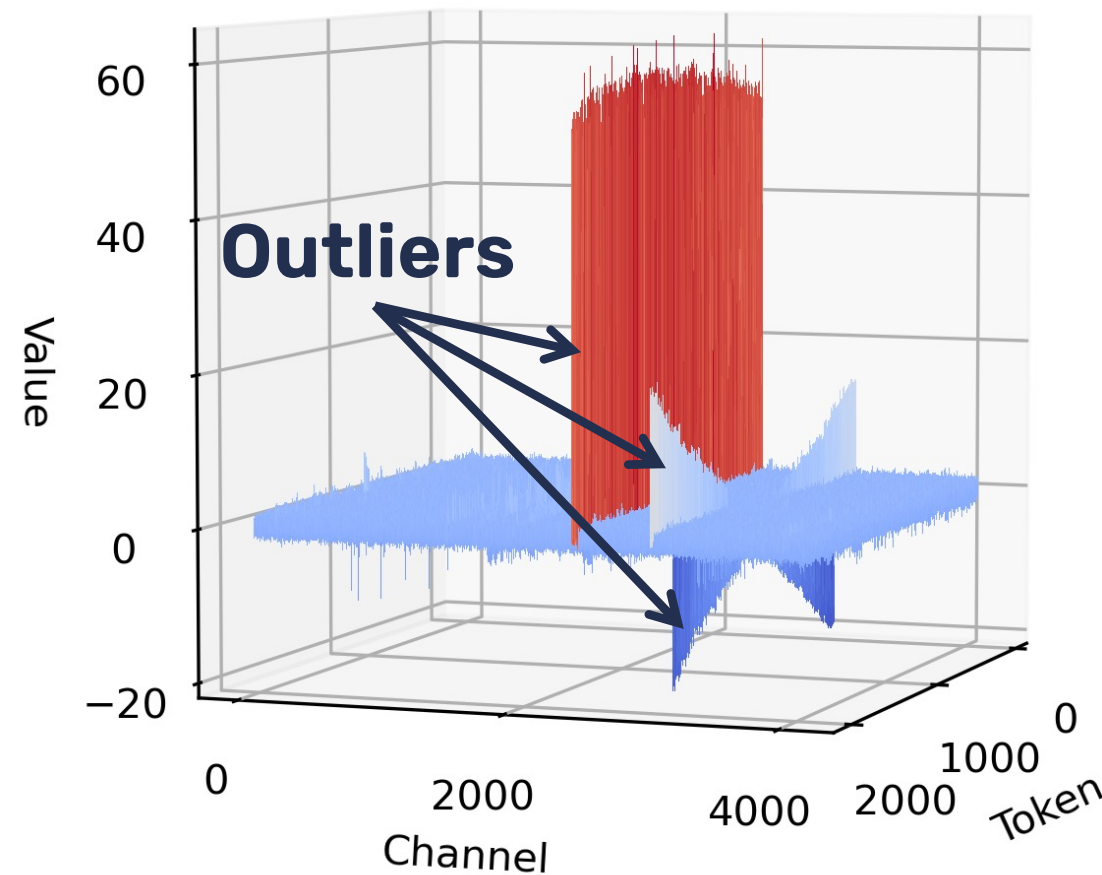
Activation Outliers in LLMs

Weight



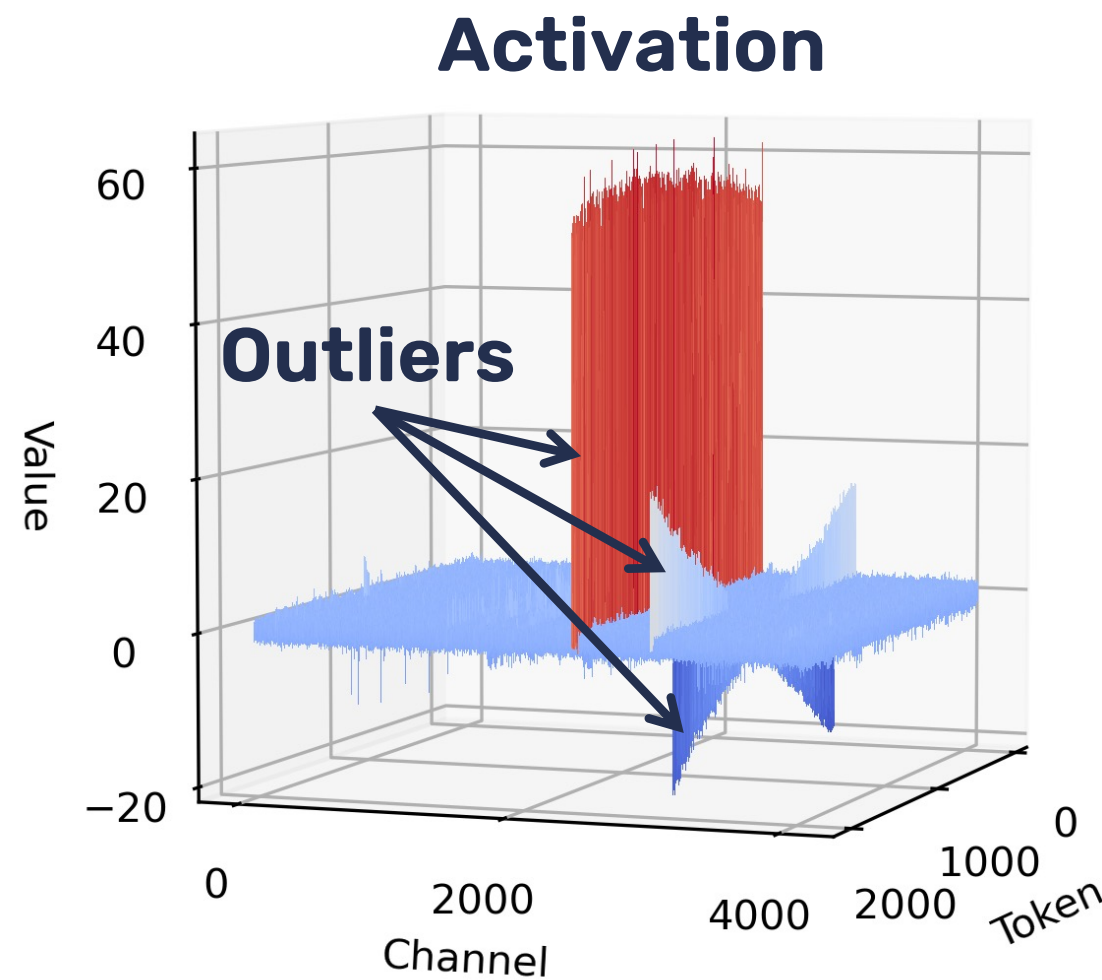
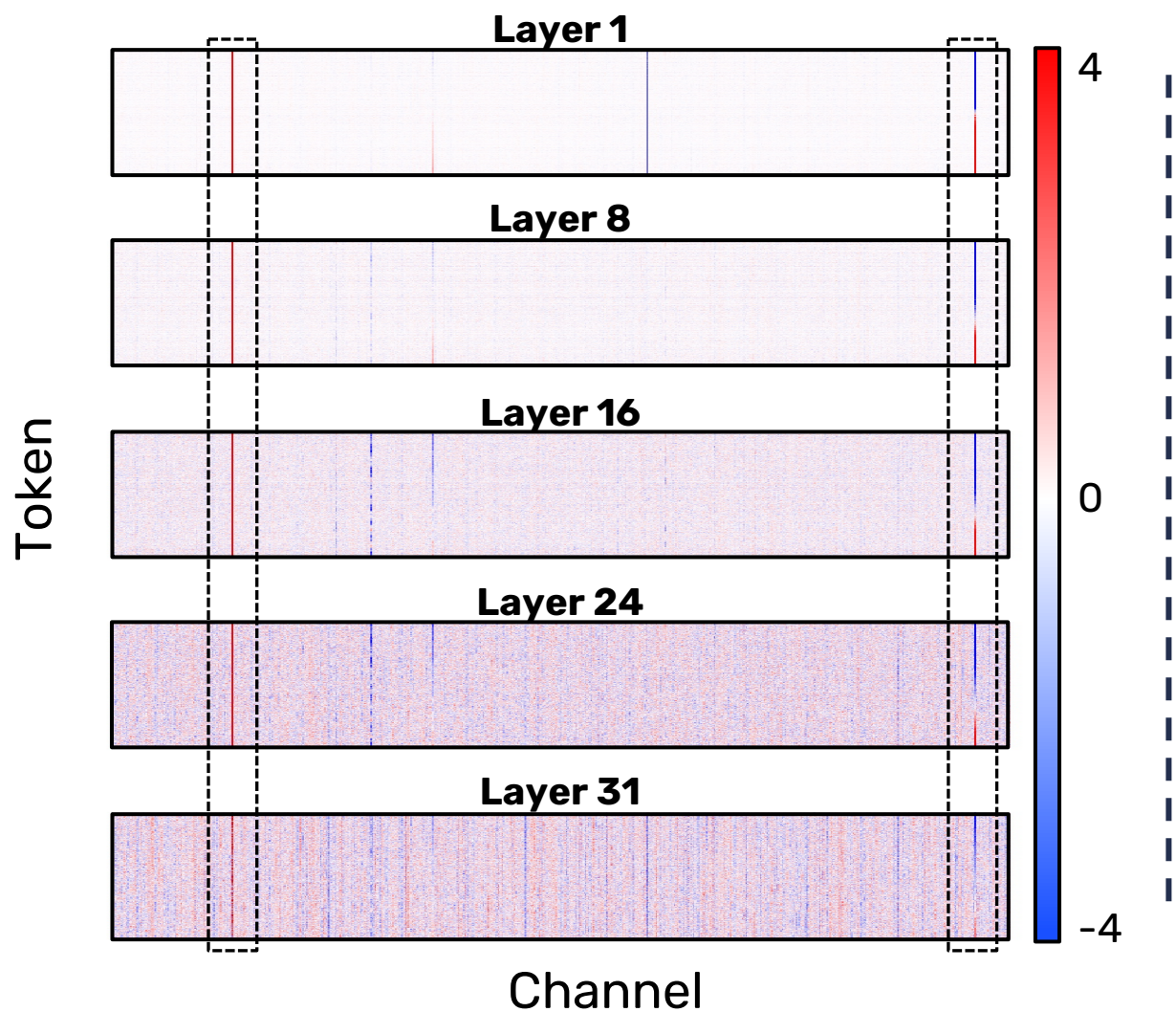
EASY to quantize!

Activation



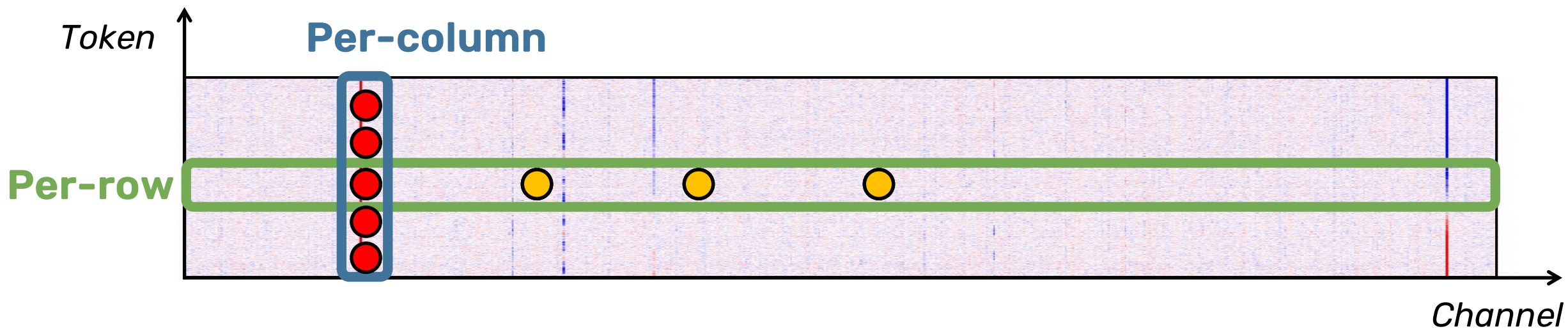
HARD to quantize!

Activation Outliers in LLMs

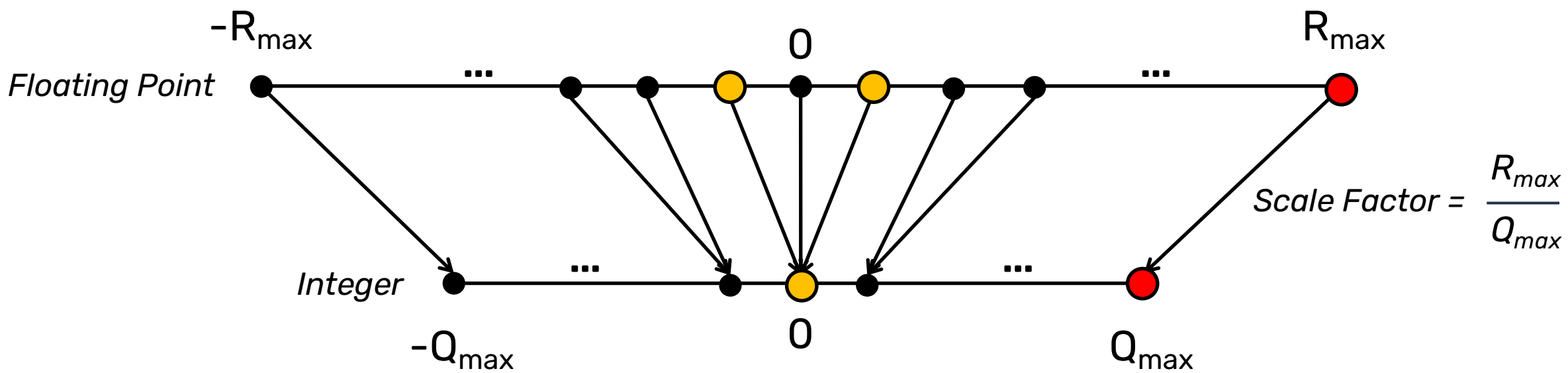
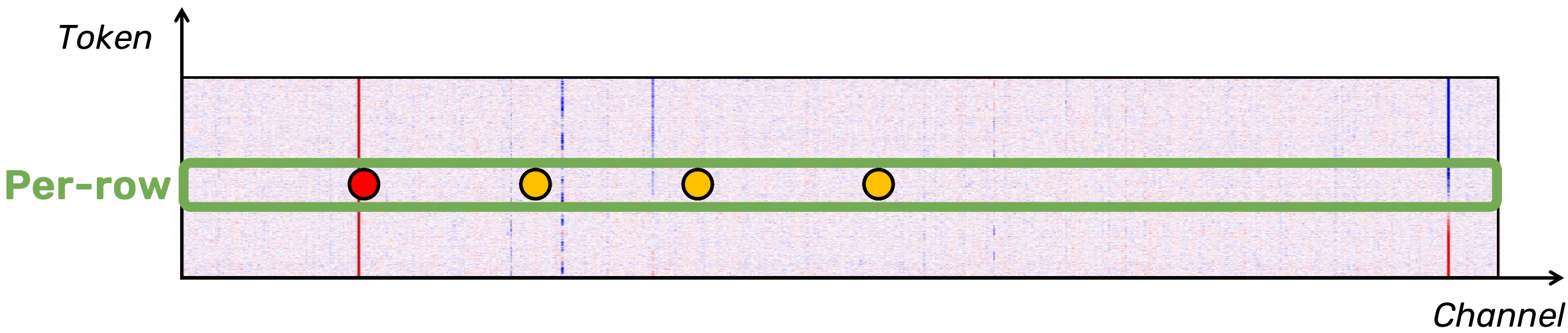


HARD to quantize!

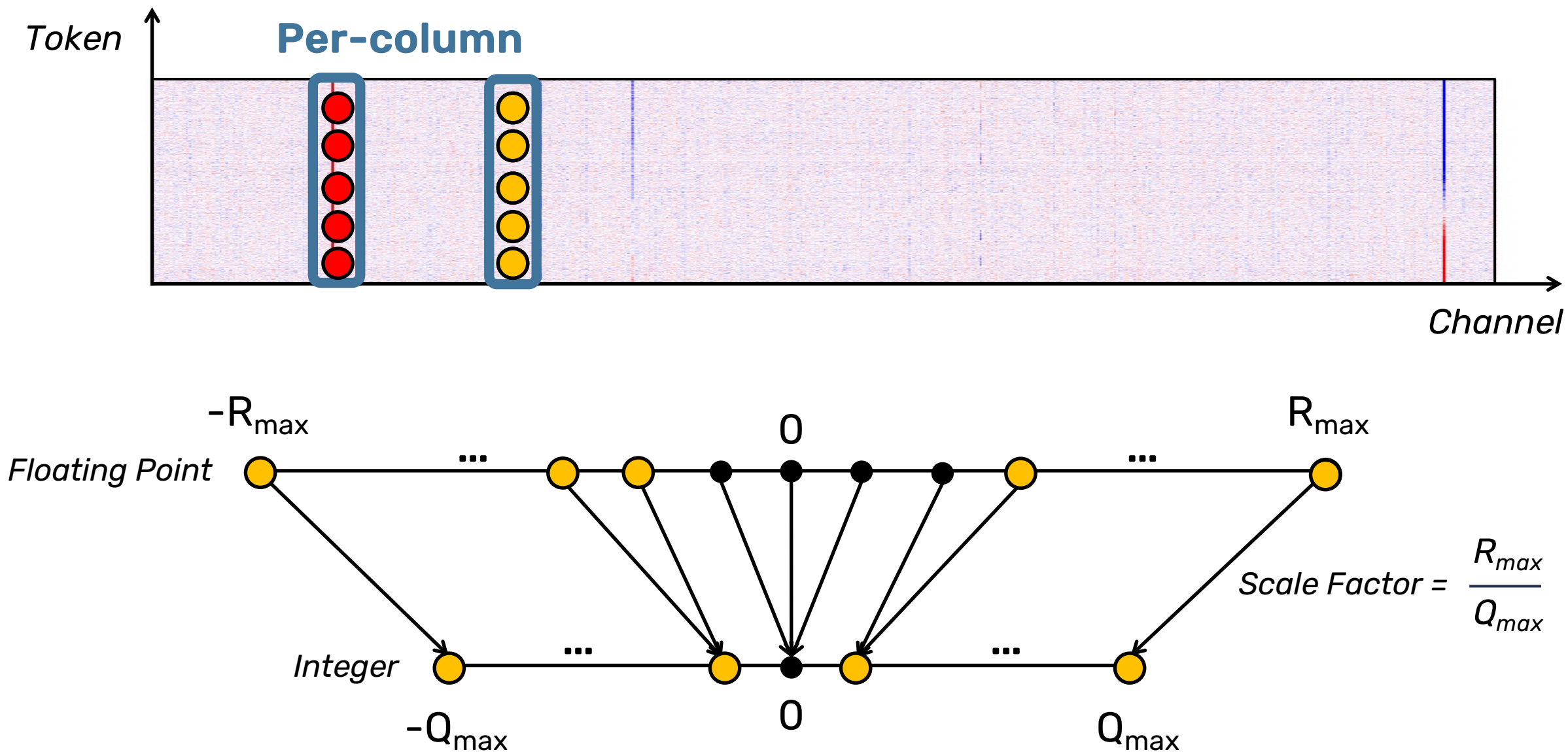
Activation Outliers in LLMs



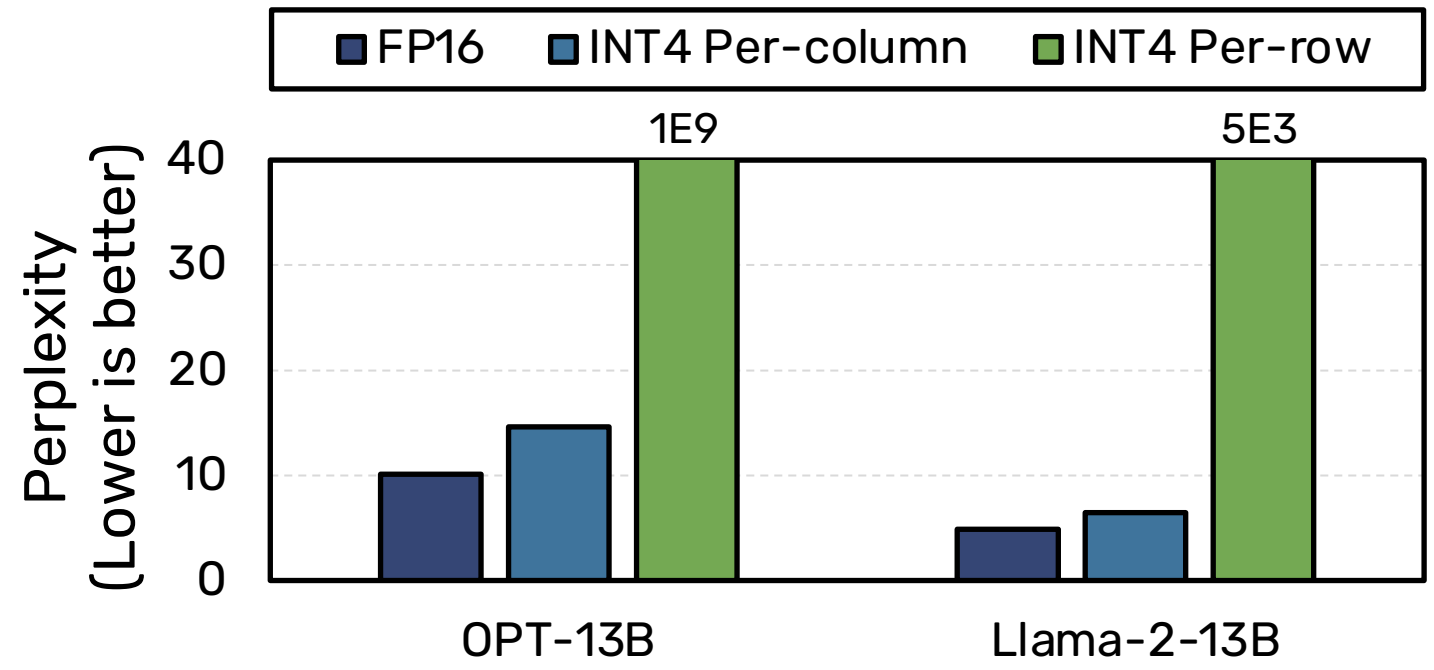
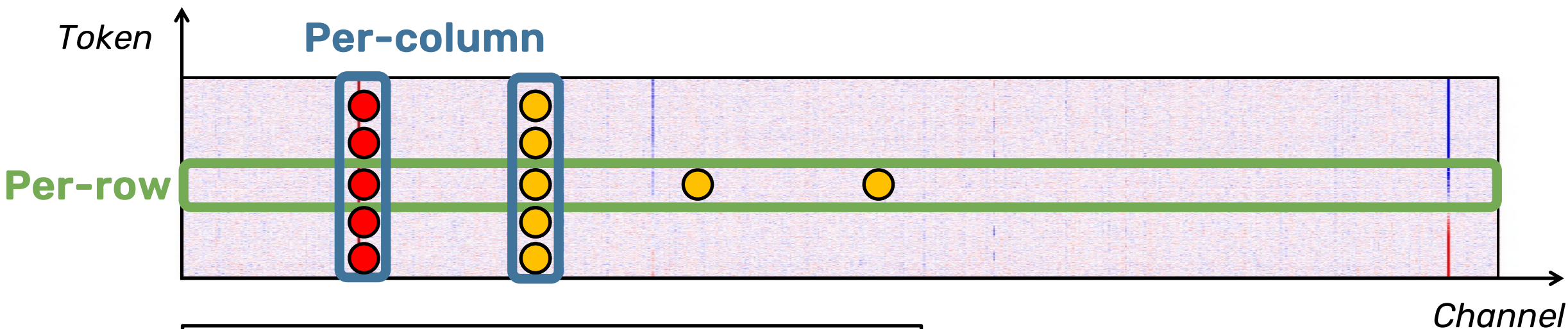
Activation Outliers in LLMs



Activation Outliers in LLMs

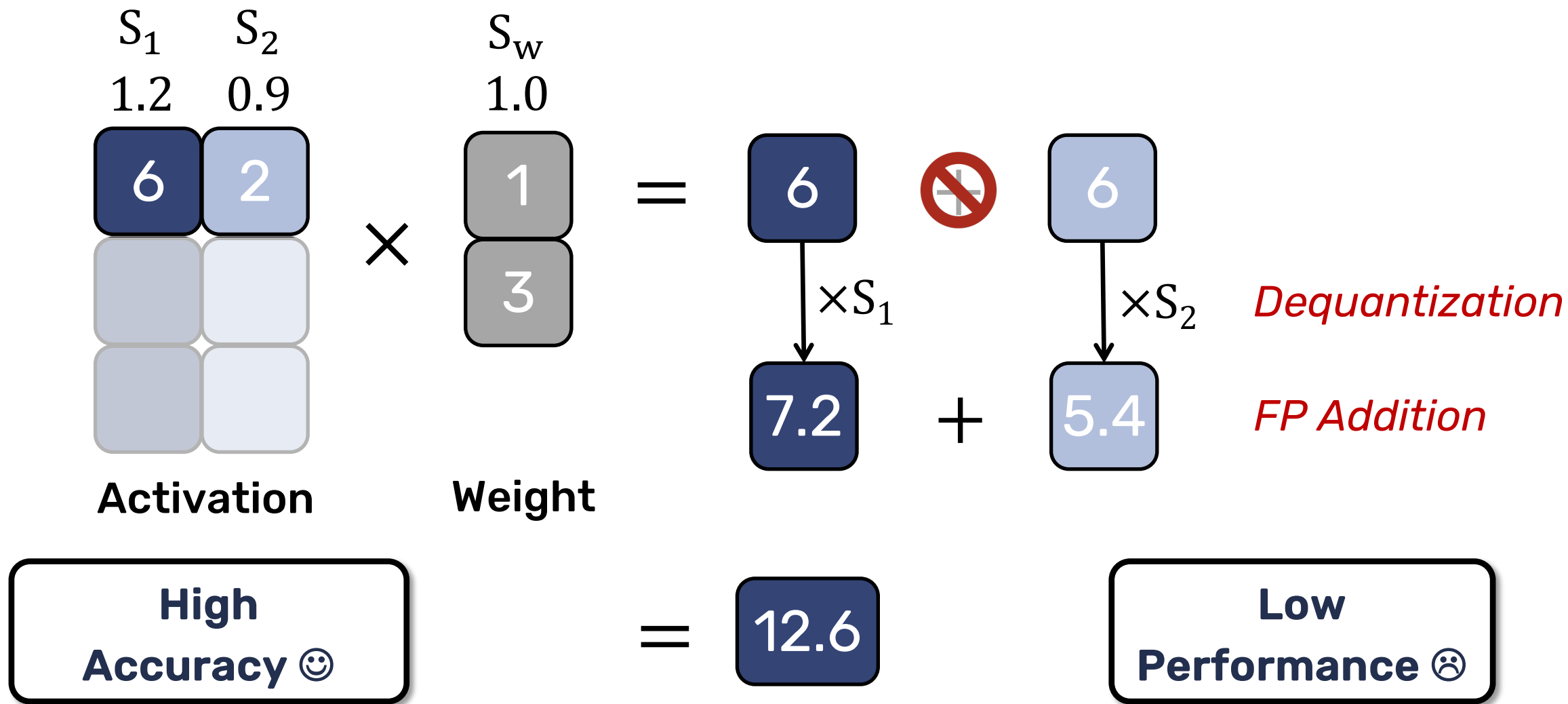


Activation Outliers in LLMs



**Per-column isolates outliers
→ Smaller Error**

Performance of Per-column Quantization



Performance of Per-column Quantization



How to split channels in activations without floating point operations?

Activation

Weight

High Accuracy 😊

=

12.6

Low Performance 😞

Limitations of Prior Works

Mixed Precision

LLM.int8() [NeurIPS'22]

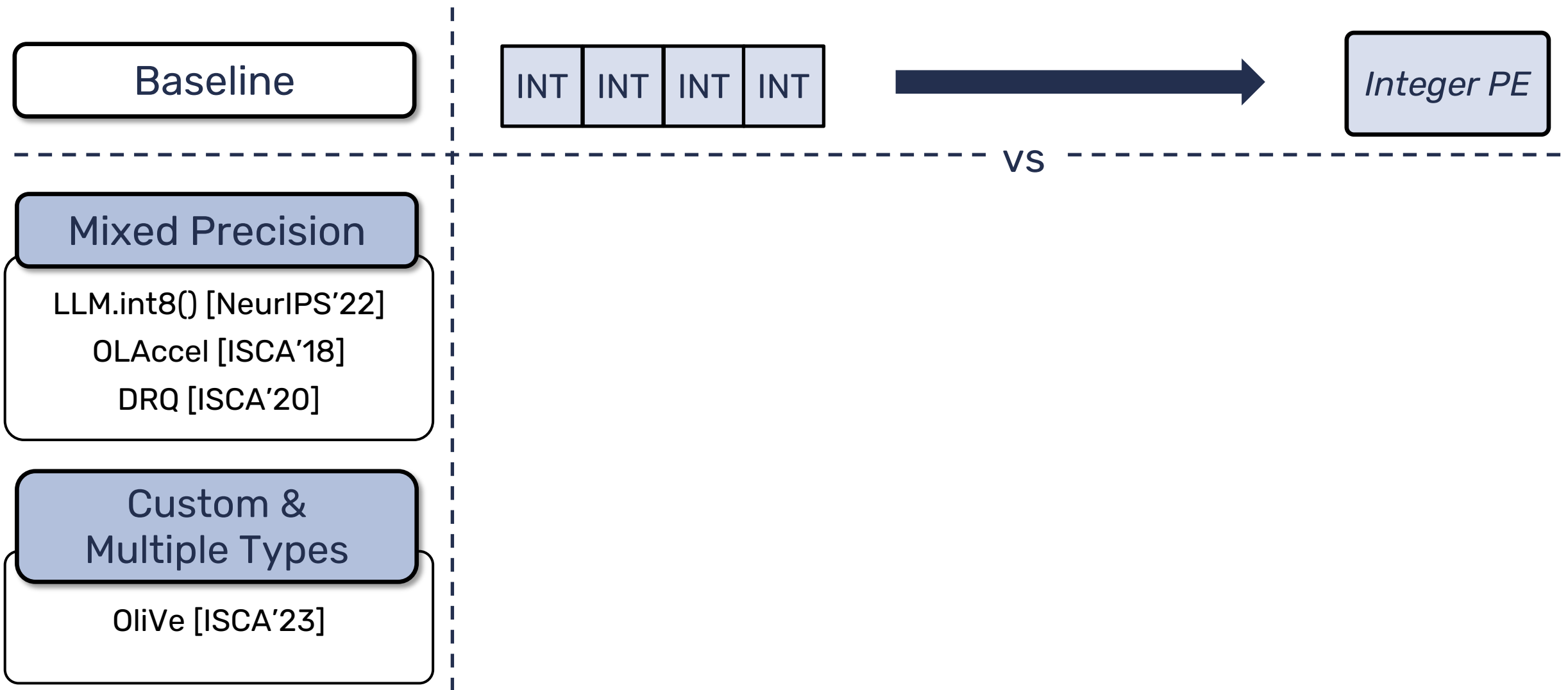
OLAccel [ISCA'18]

DRQ [ISCA'20]

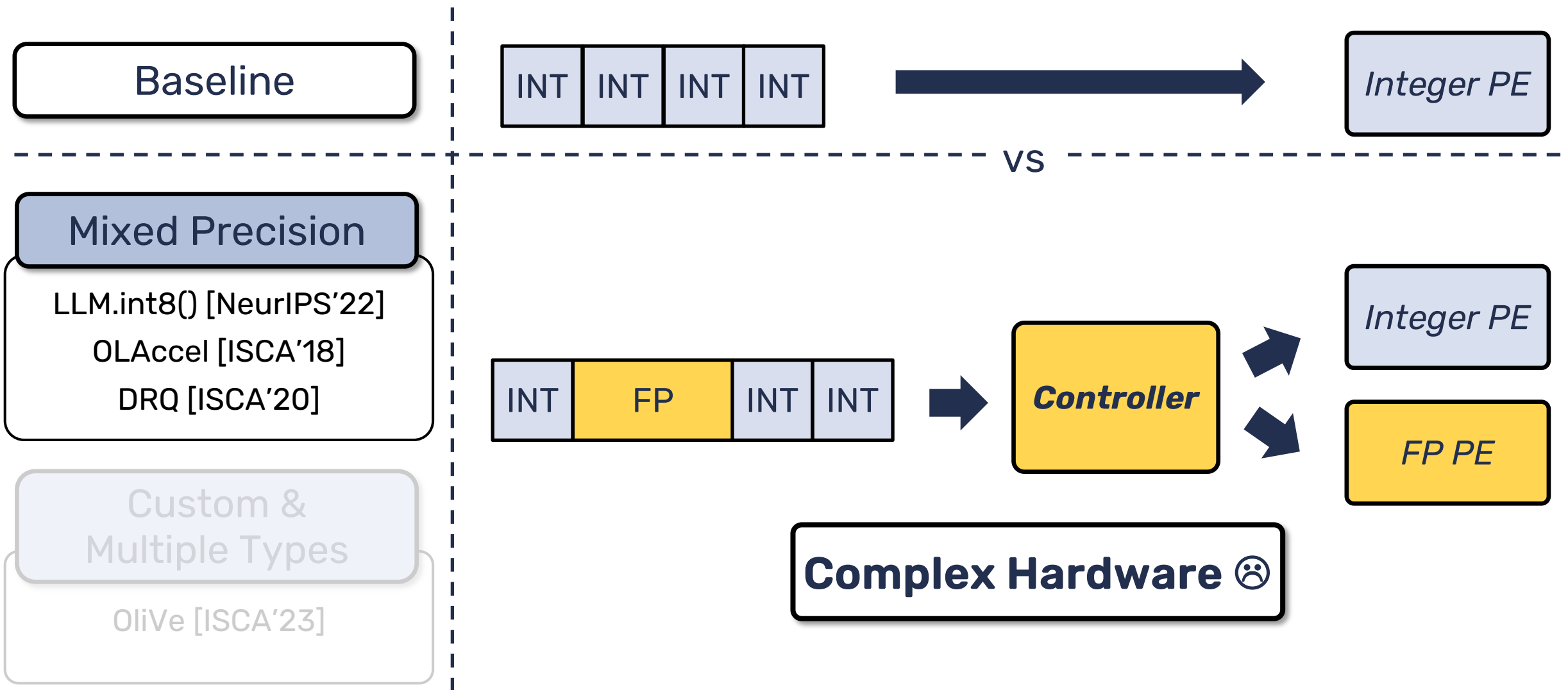
Custom & Multiple Types

OliVe [ISCA'23]

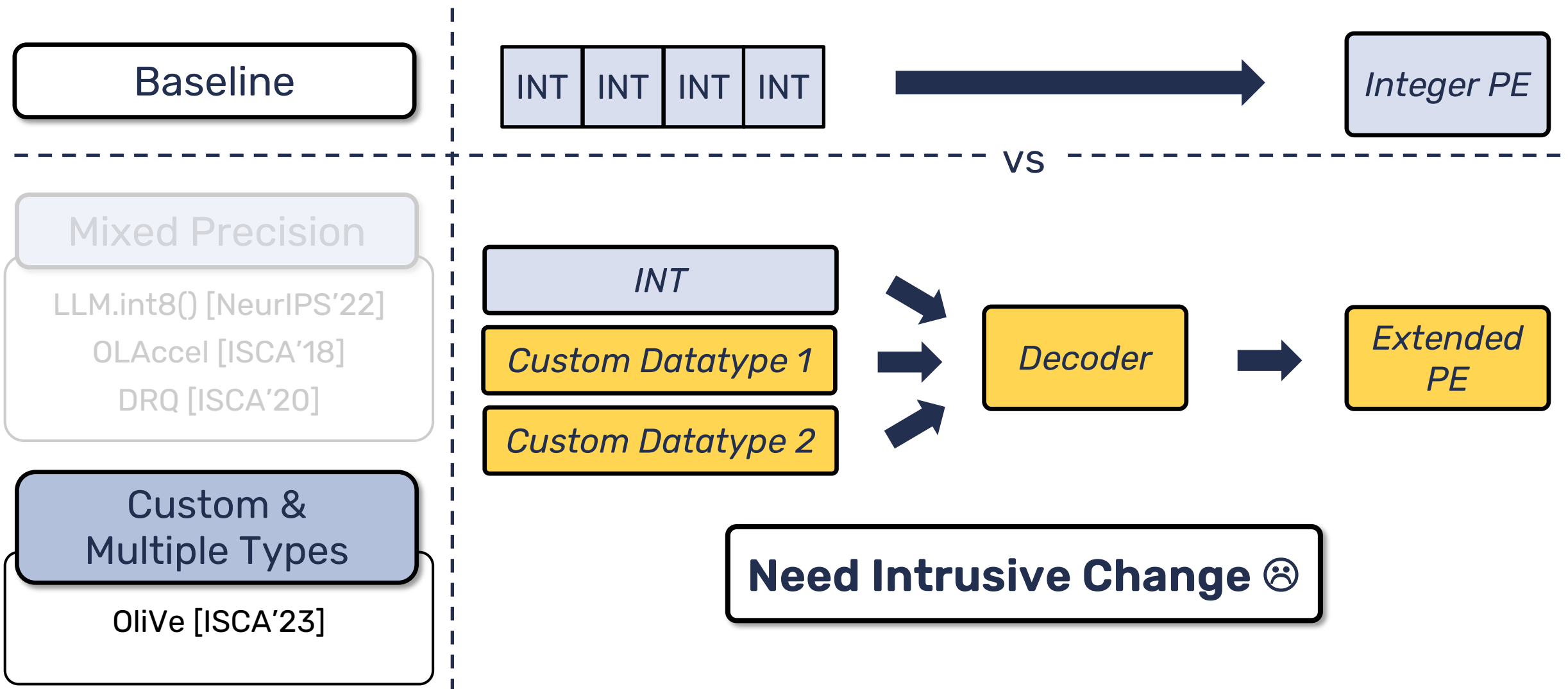
Limitations of Prior Works



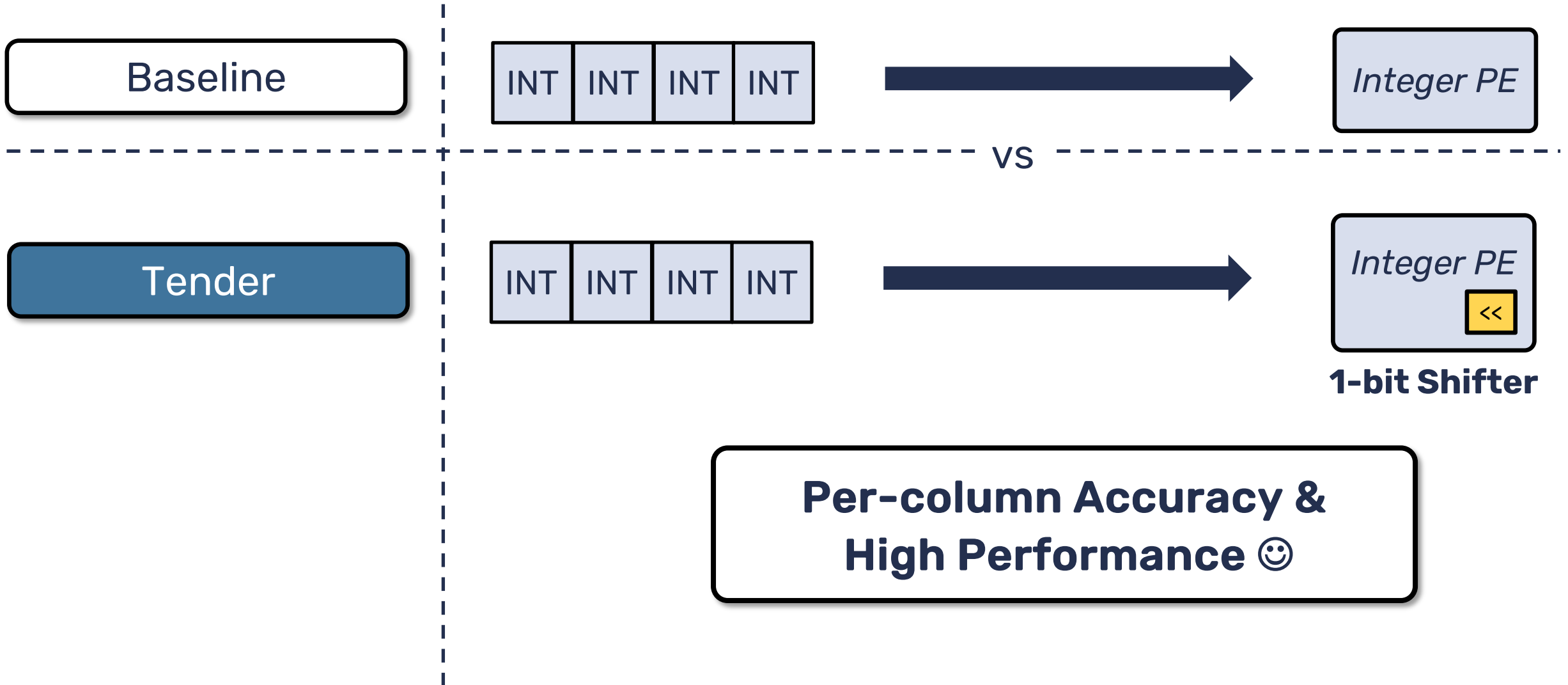
Limitations of Prior Works



Limitations of Prior Works



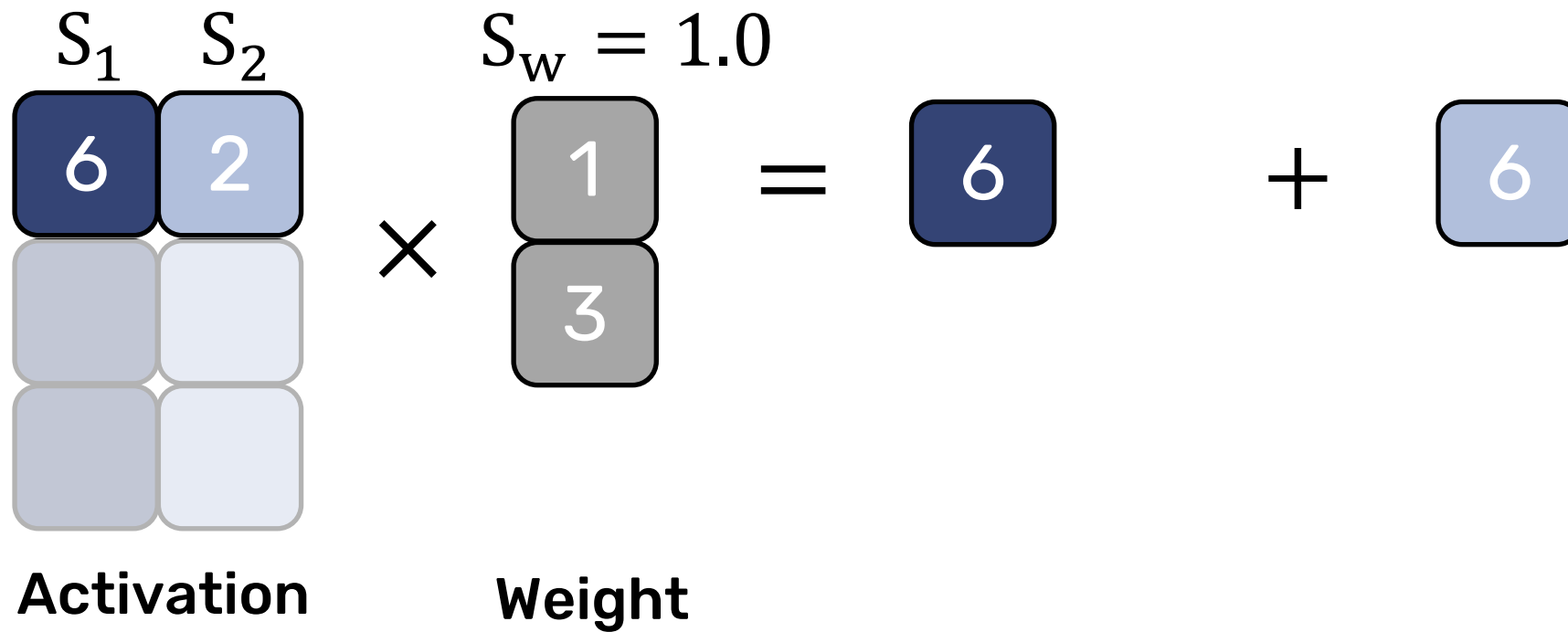
Tender Overview



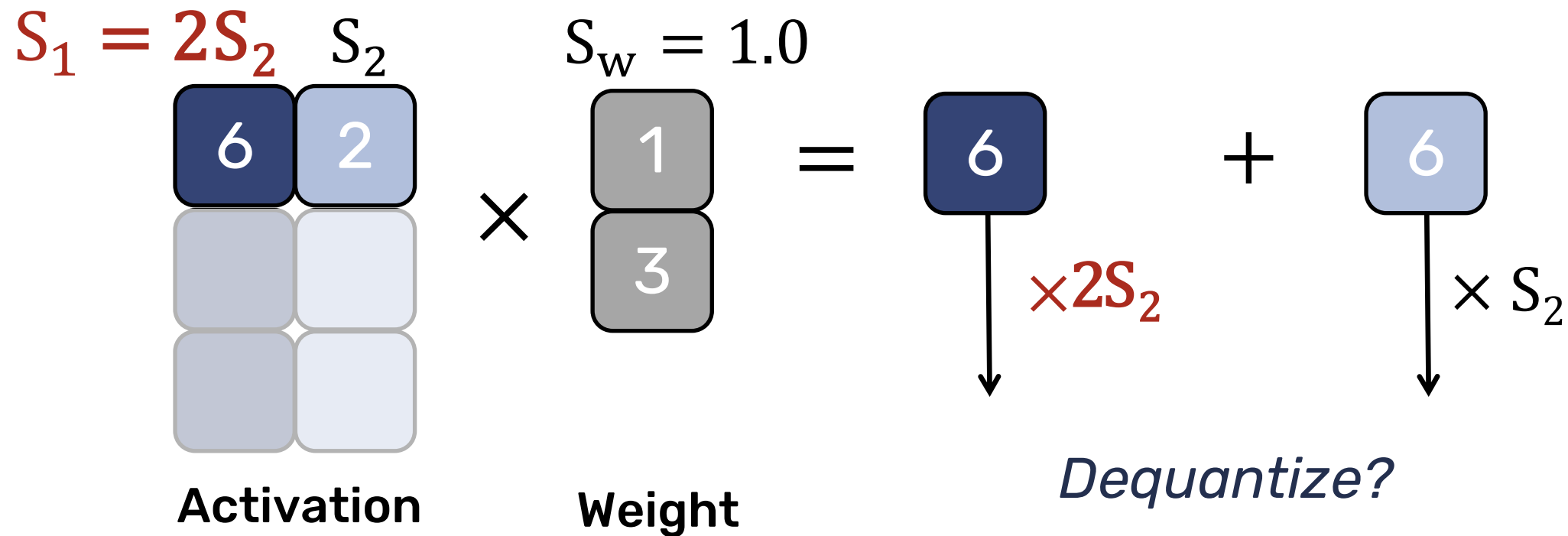
Outline

- **Motivation**
 - Challenges in Efficient LLM Inference
 - Limitations of Prior Works
- **Tender: Algorithm-Hardware Co-design for Efficient LLM Inference**
 - Tensor Decomposition
 - Rescaling Operation
- **Evaluation**
- **Conclusion**

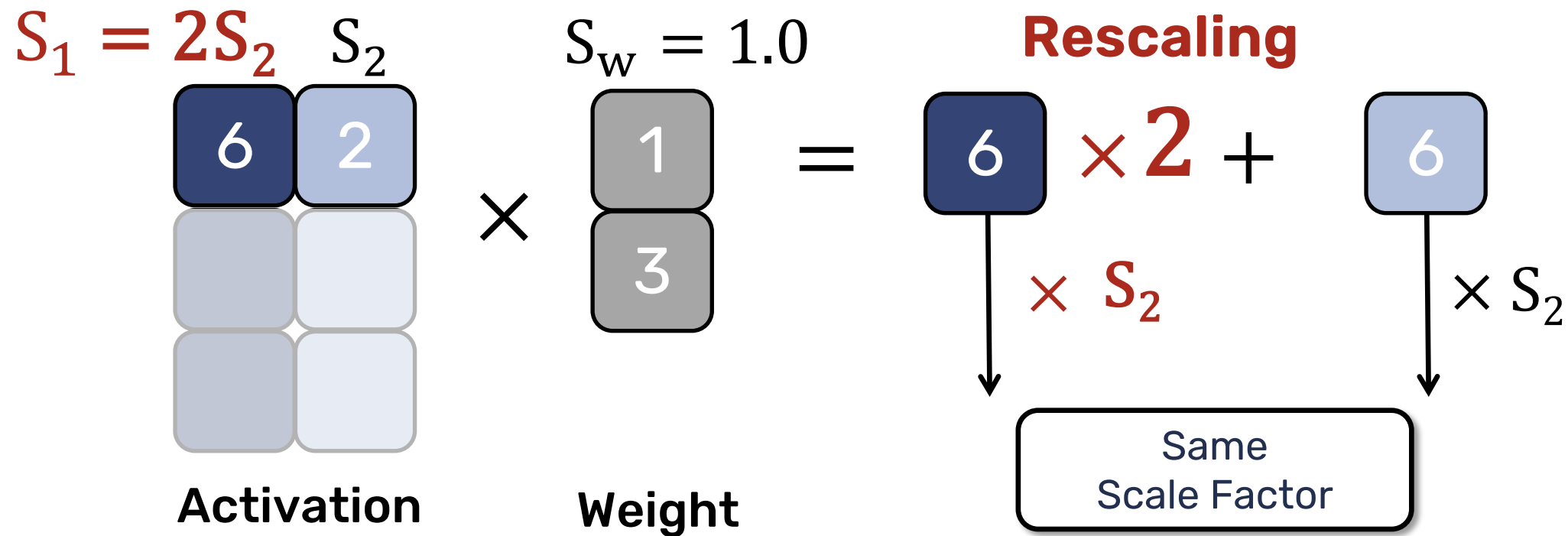
Key Insight



Key Insight



Key Insight



Key Insight

$$\begin{array}{c} S_1 = 2S_2 \\ \begin{array}{|c|c|} \hline 6 & 2 \\ \hline \hline \hline \end{array} \\ \text{Activation} \end{array} \times \begin{array}{c} S_w = 1.0 \\ \begin{array}{|c|} \hline 1 \\ \hline 3 \\ \hline \end{array} \\ \text{Weight} \end{array} = \begin{array}{c} \text{Rescaling} \\ \begin{array}{|c|} \hline 6 \times 2 \\ \hline \end{array} + \begin{array}{|c|} \hline 6 \\ \hline \end{array} \\ = \begin{array}{|c|} \hline 12 \\ \hline \end{array} + \begin{array}{|c|} \hline 6 \\ \hline \end{array} \\ = \begin{array}{|c|} \hline 18 \\ \hline \end{array} \end{array}$$

Key Insight

$$S_1 = 2S_2 \quad S_2 \quad S_w = 1.0 \quad \text{Rescaling}$$

6 2 1 = 6 × 2 + 6

The ratio between scale factors is an integer
→ Enables computing in integer pipeline

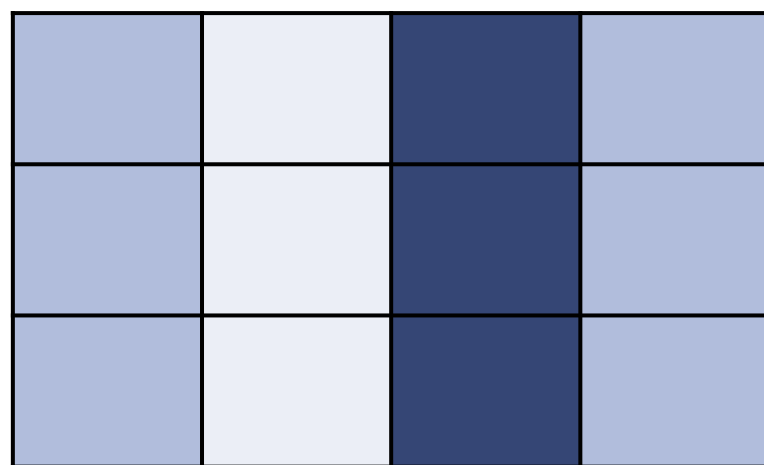
Activation

Weight

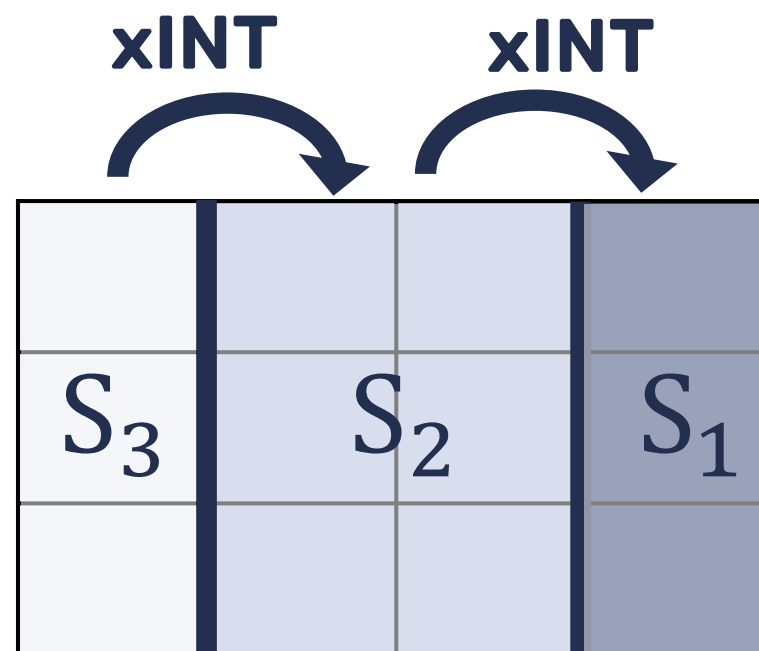
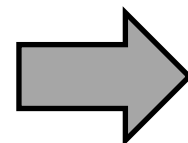
Scale Factor

= 18

Tender: Tensor Decomposition

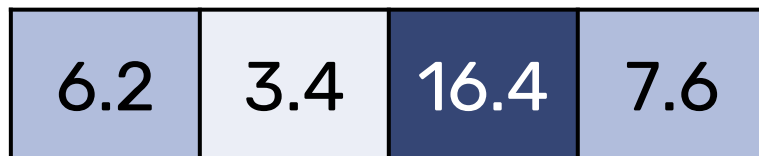


Activation

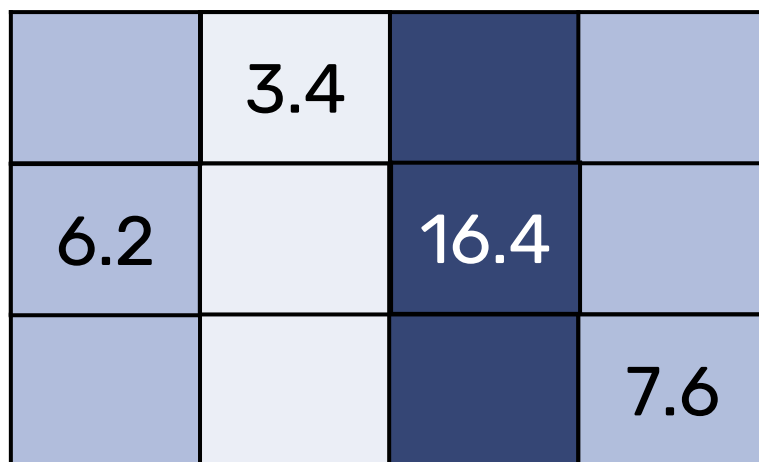


Tender: Tensor Decomposition

**Channel
Max.**

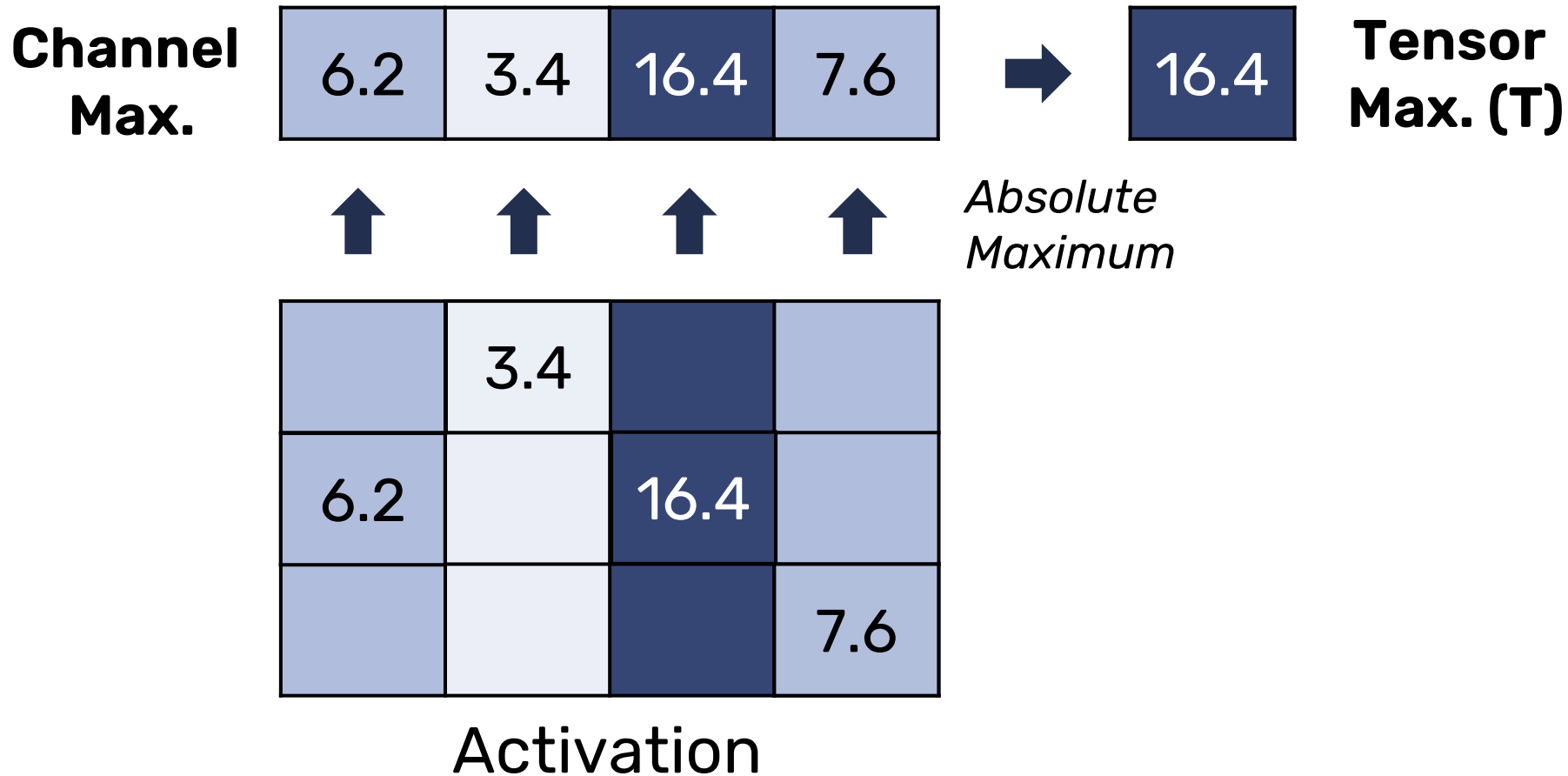


*Absolute
Maximum*

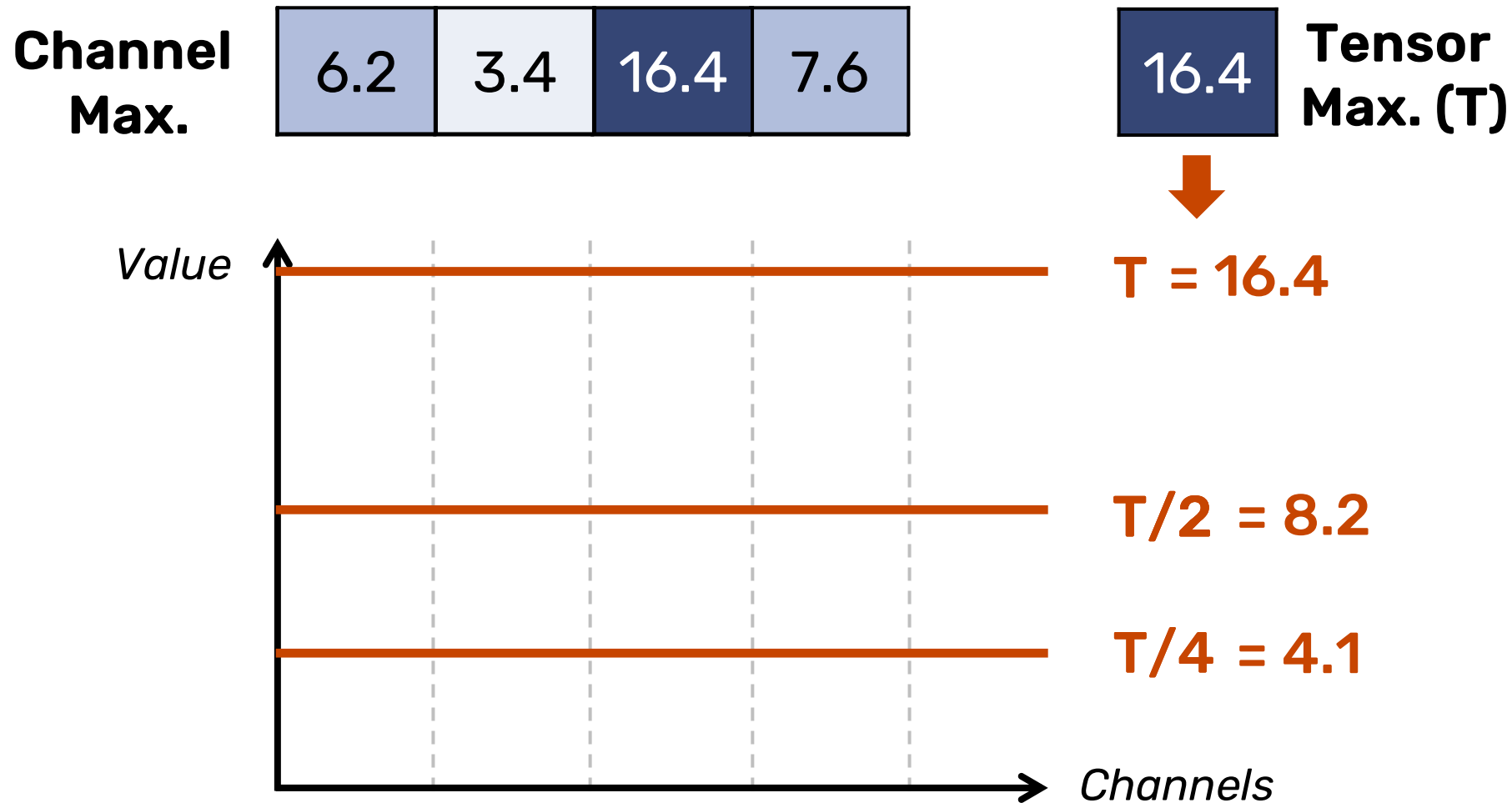


Activation

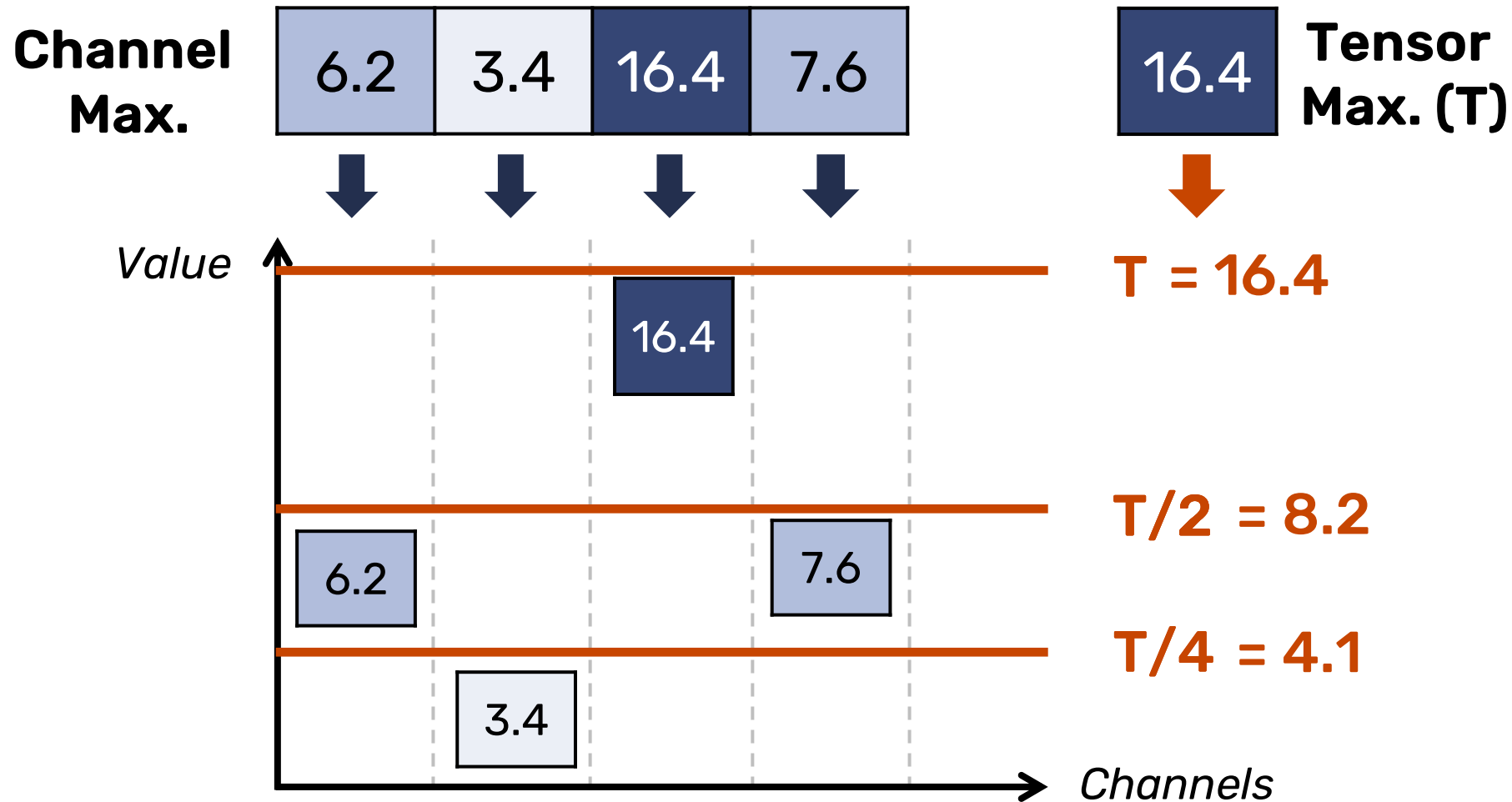
Tender: Tensor Decomposition



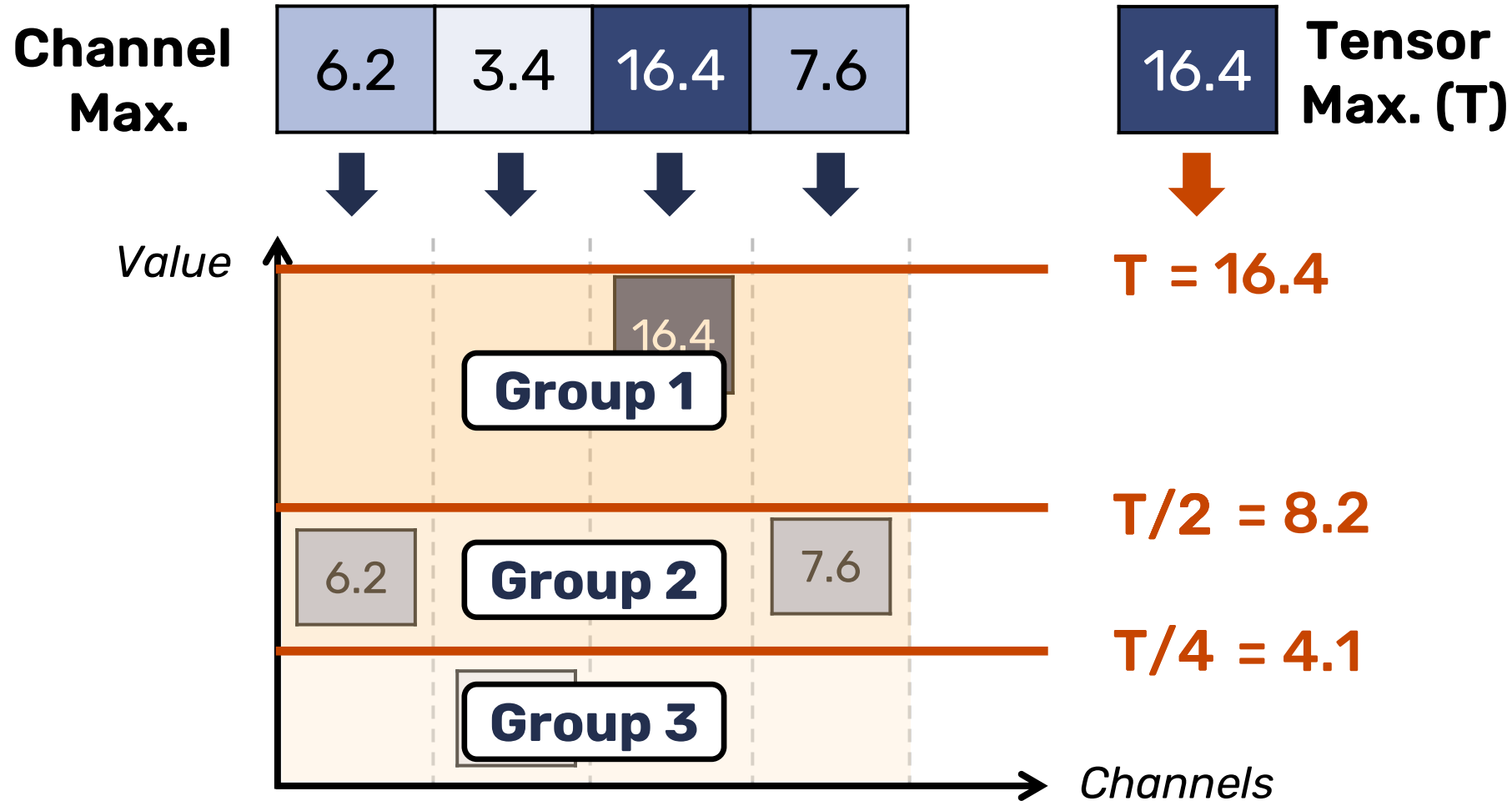
Tender: Tensor Decomposition



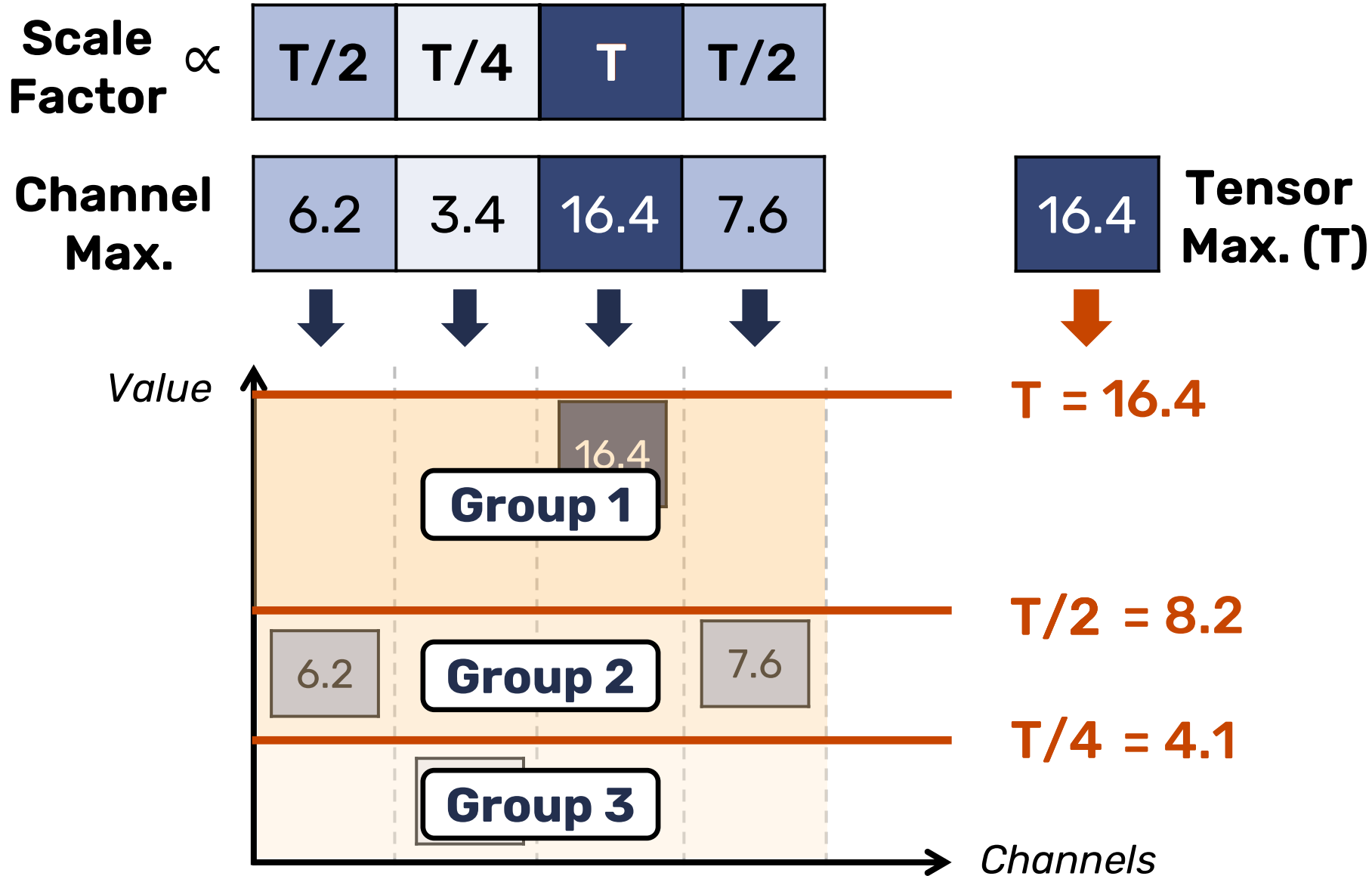
Tender: Tensor Decomposition



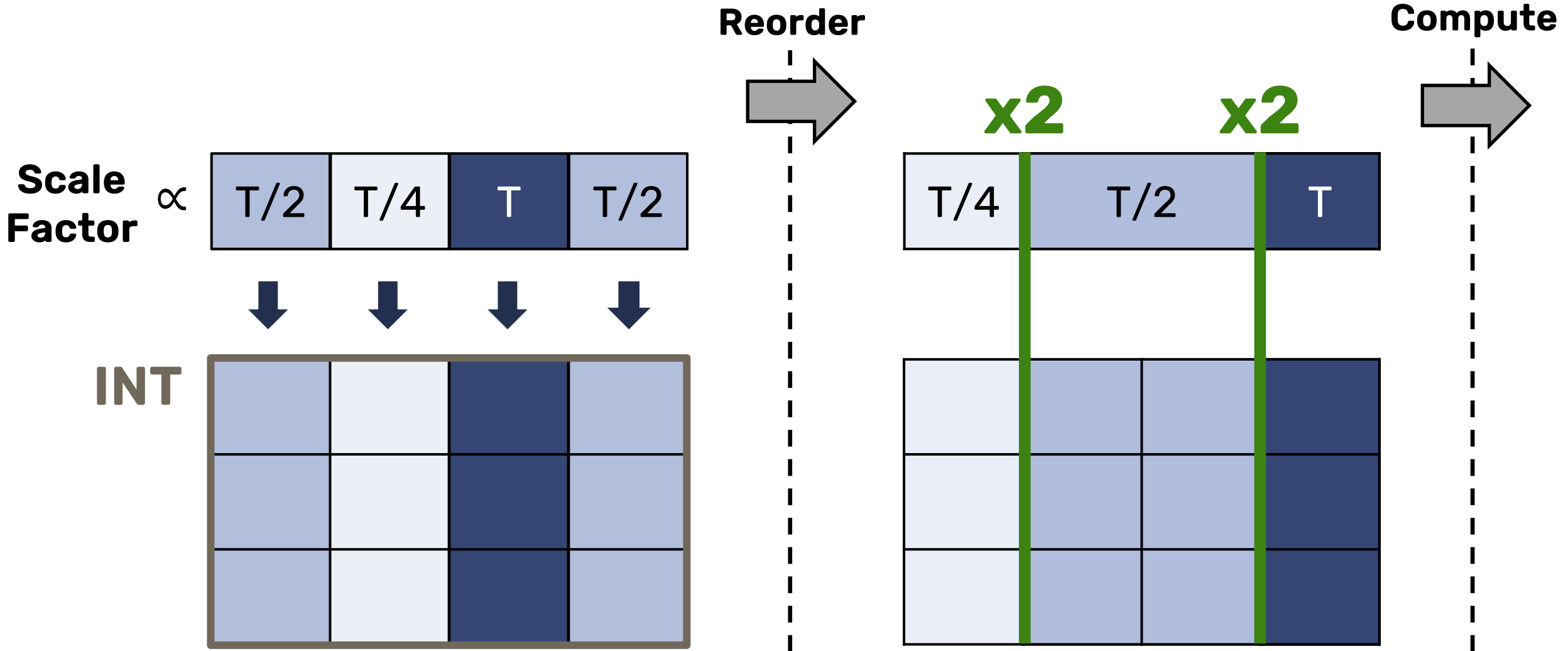
Tender: Tensor Decomposition



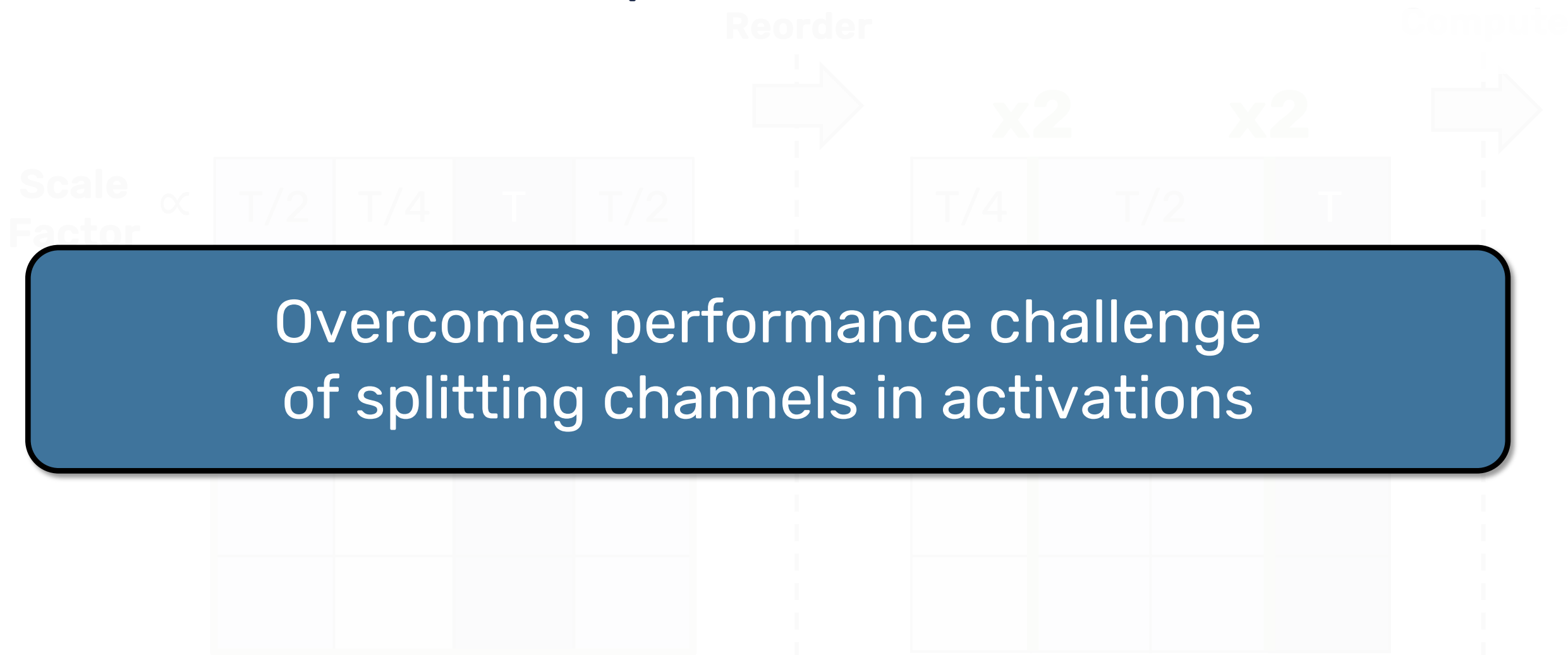
Tender: Tensor Decomposition



Tender: Tensor Decomposition



Tender: Tensor Decomposition



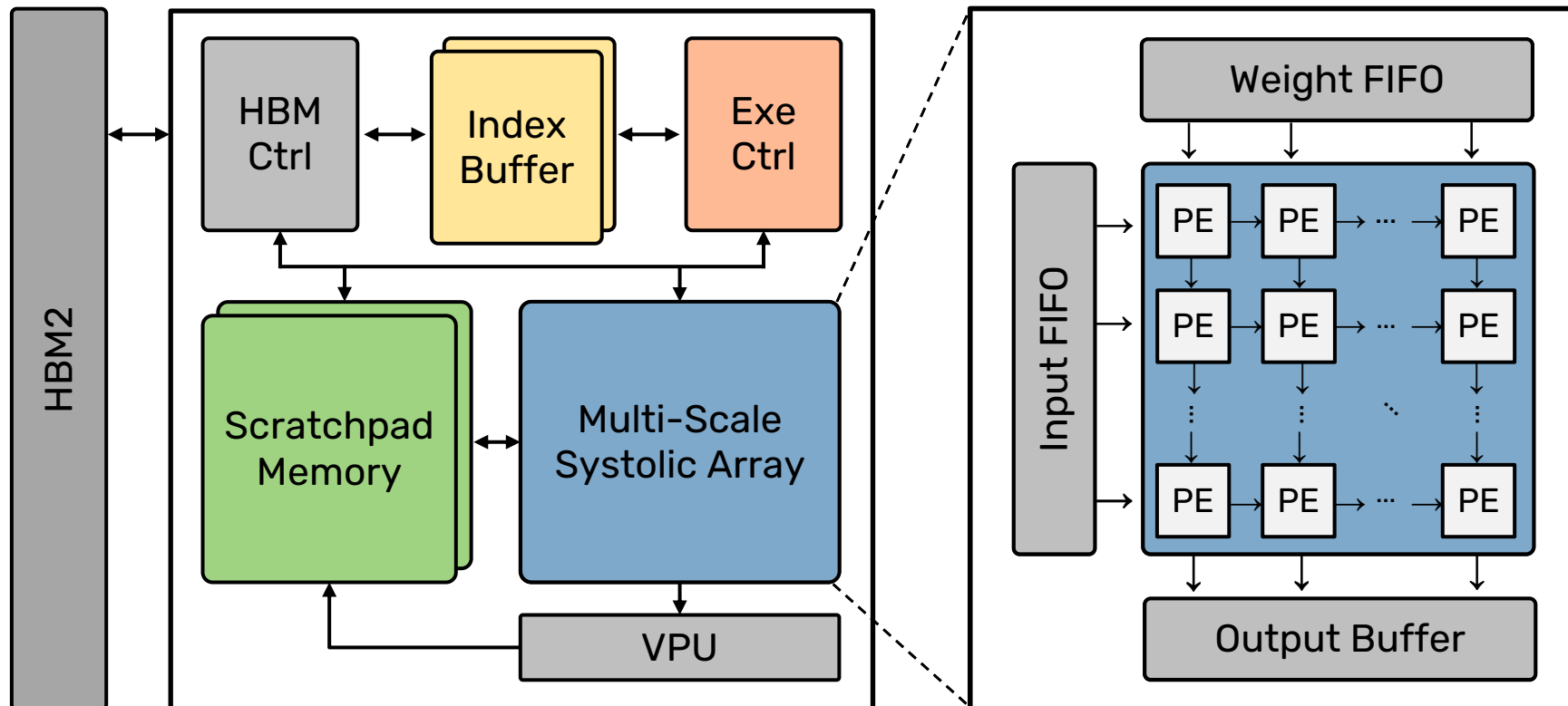
Tender: Architecture Overview

Execution Controller

- Column reordering

Multi-Scale Systolic Array (MSA)

- Computation with Rescaling



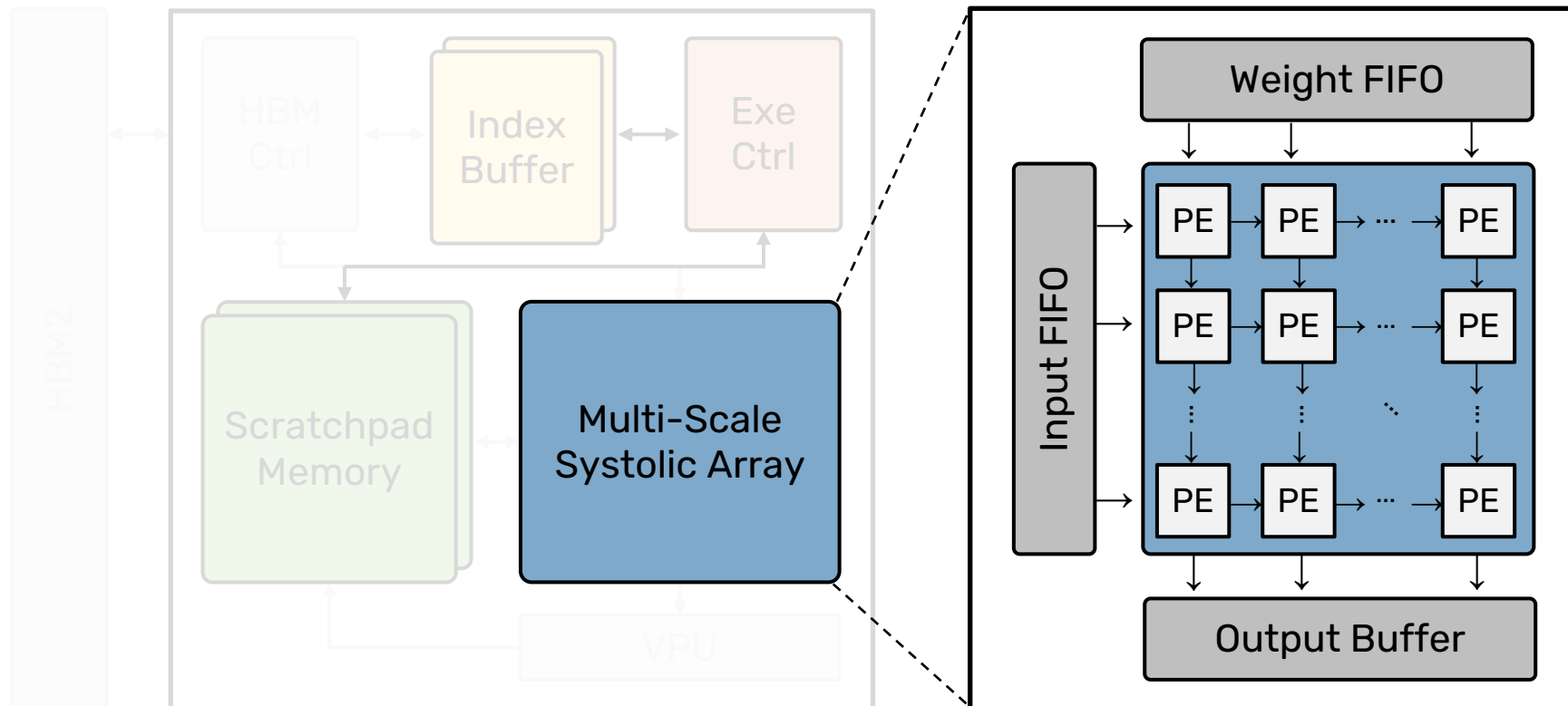
Tender: Architecture Overview

Execution Controller

- Column reordering

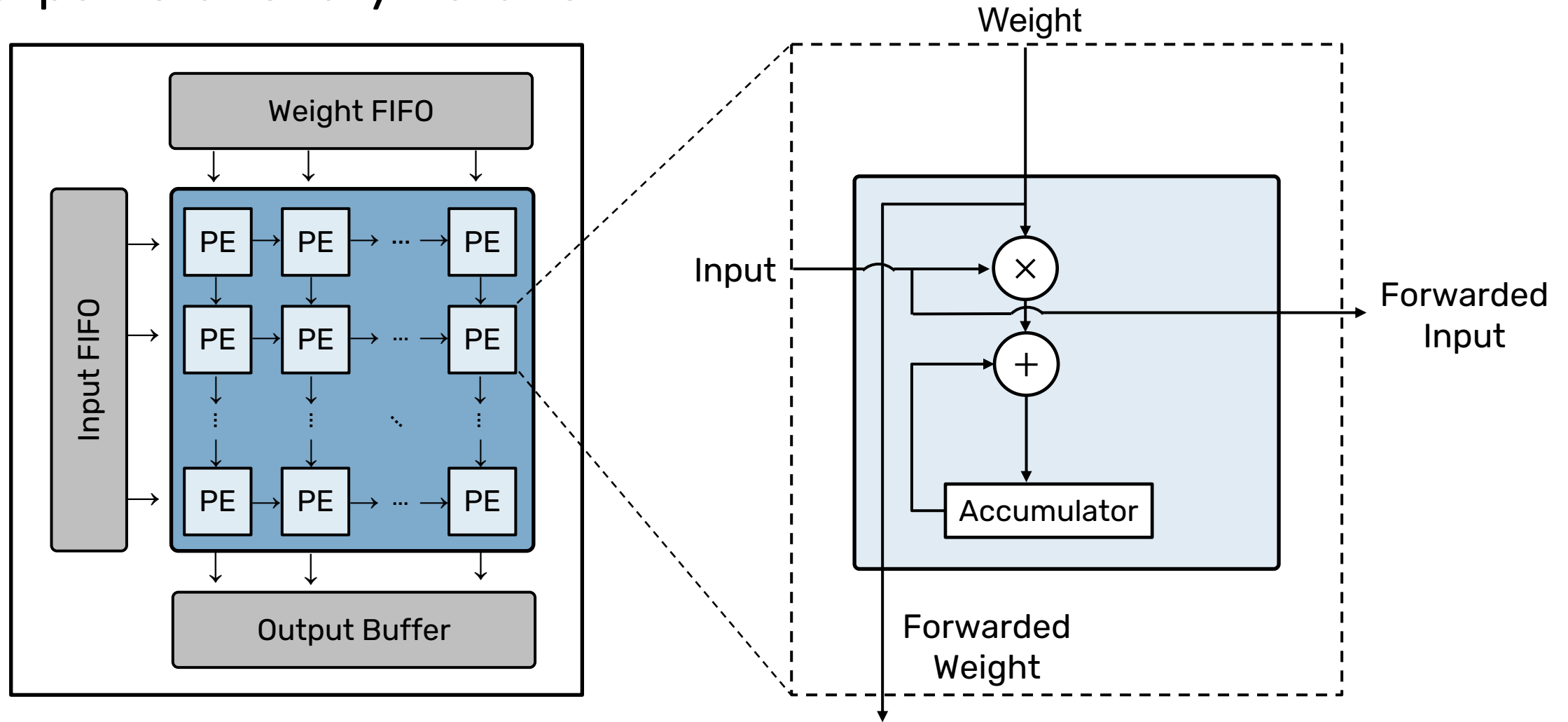
Multi-Scale Systolic Array (MSA)

- Computation with Rescaling



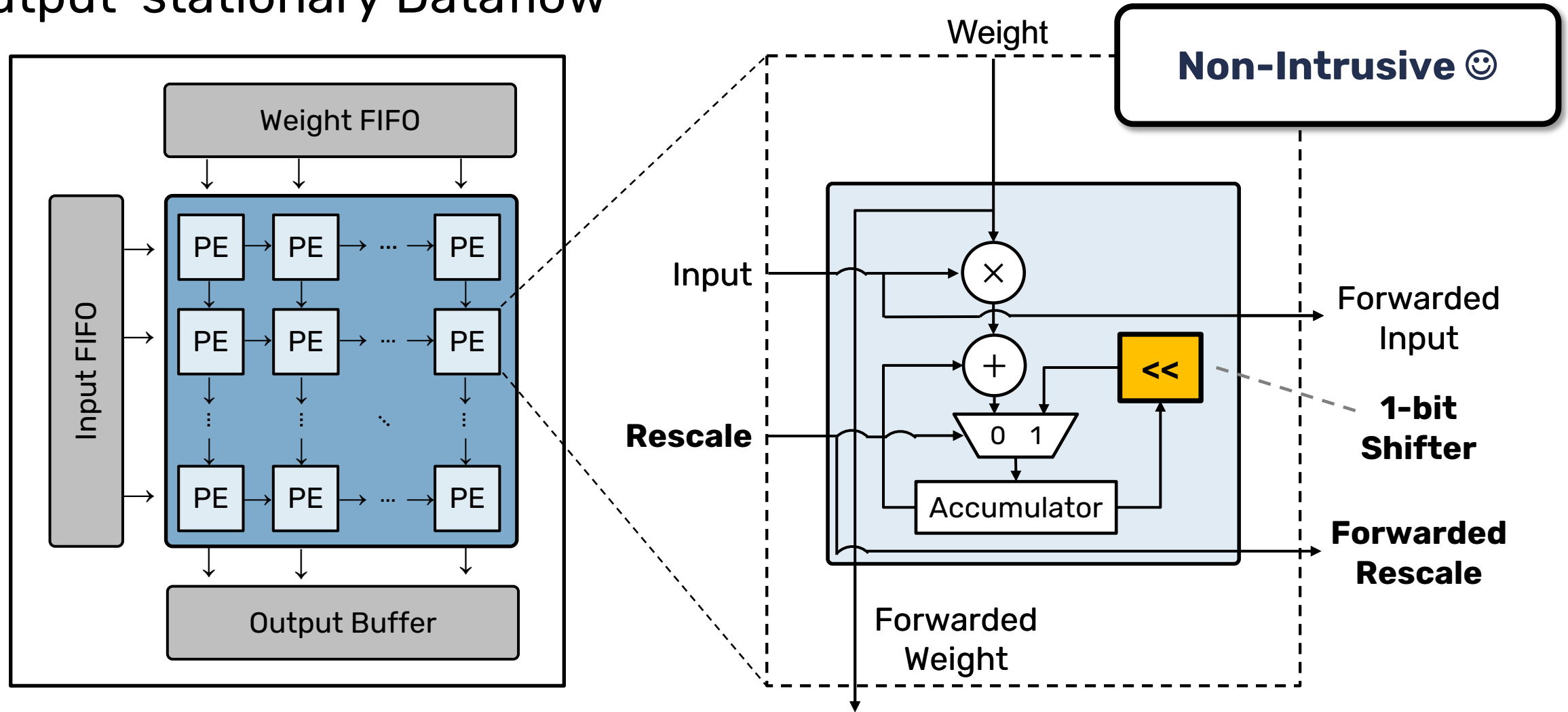
Multi-Scale Systolic Array (MSA)

Output-stationary Dataflow

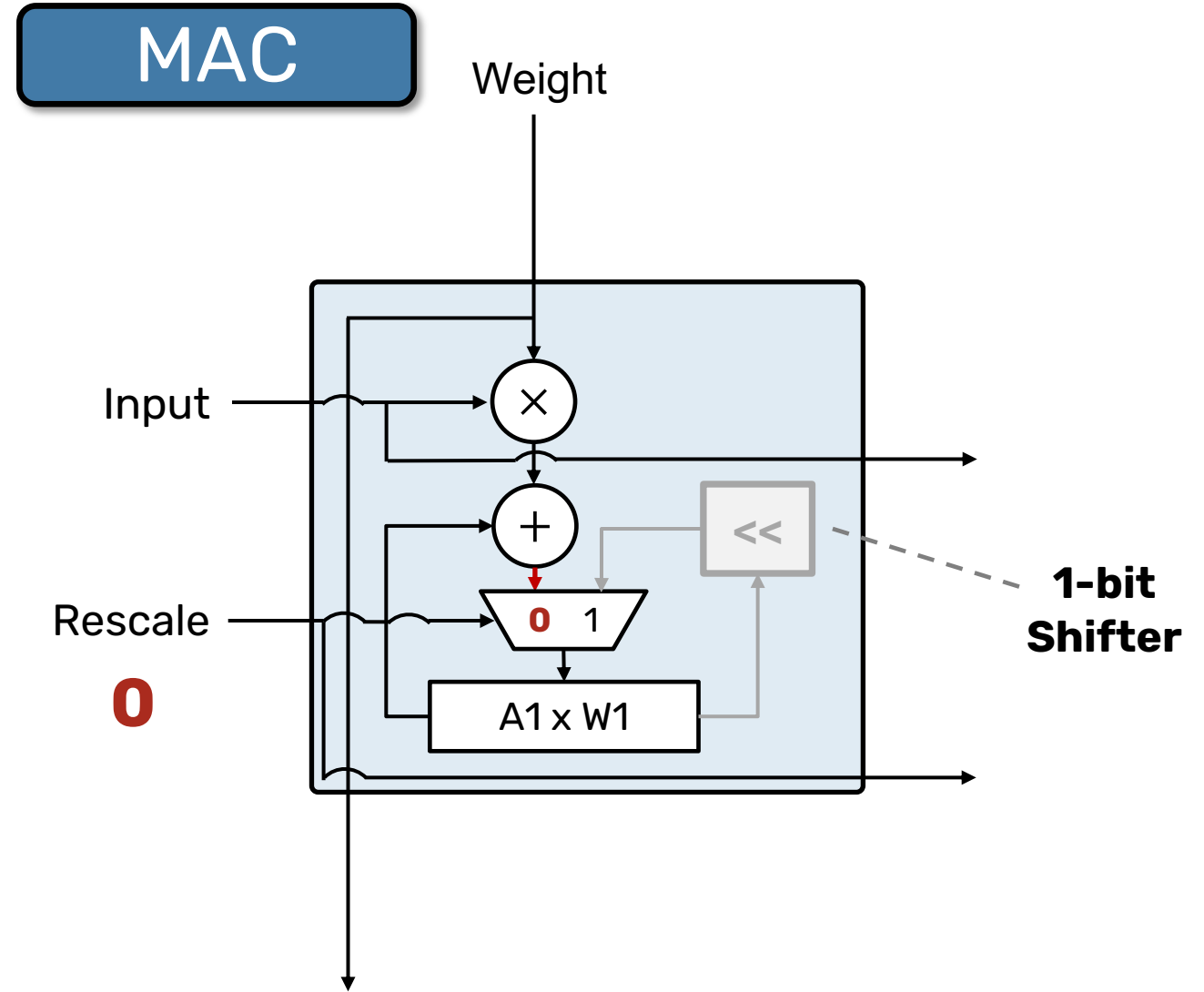
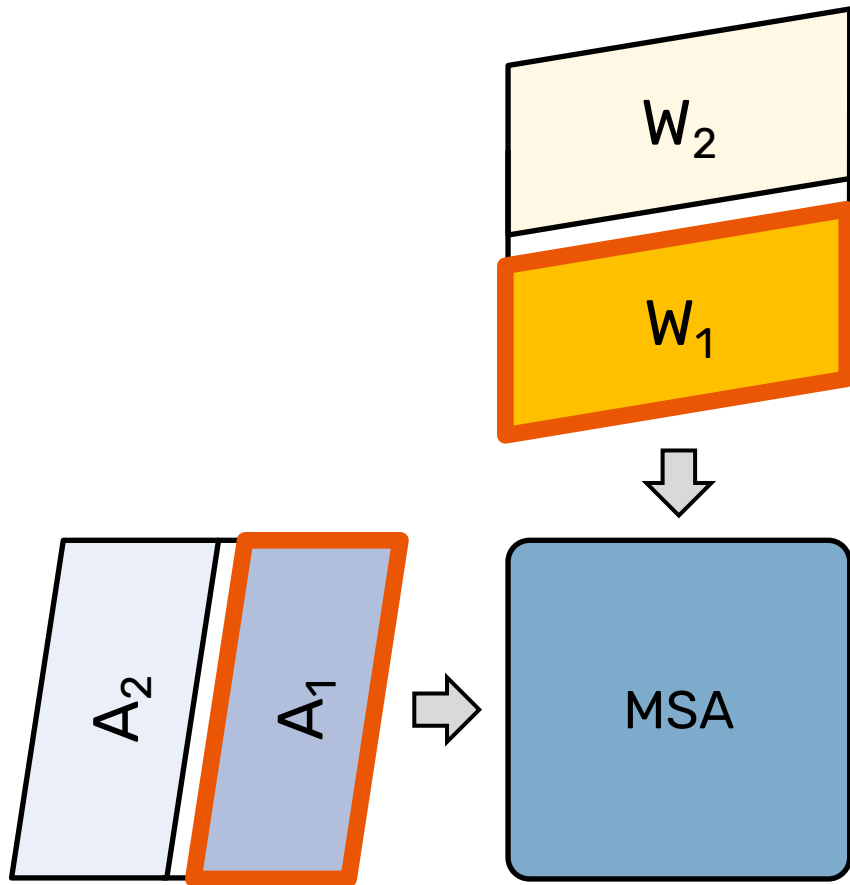


Multi-Scale Systolic Array (MSA)

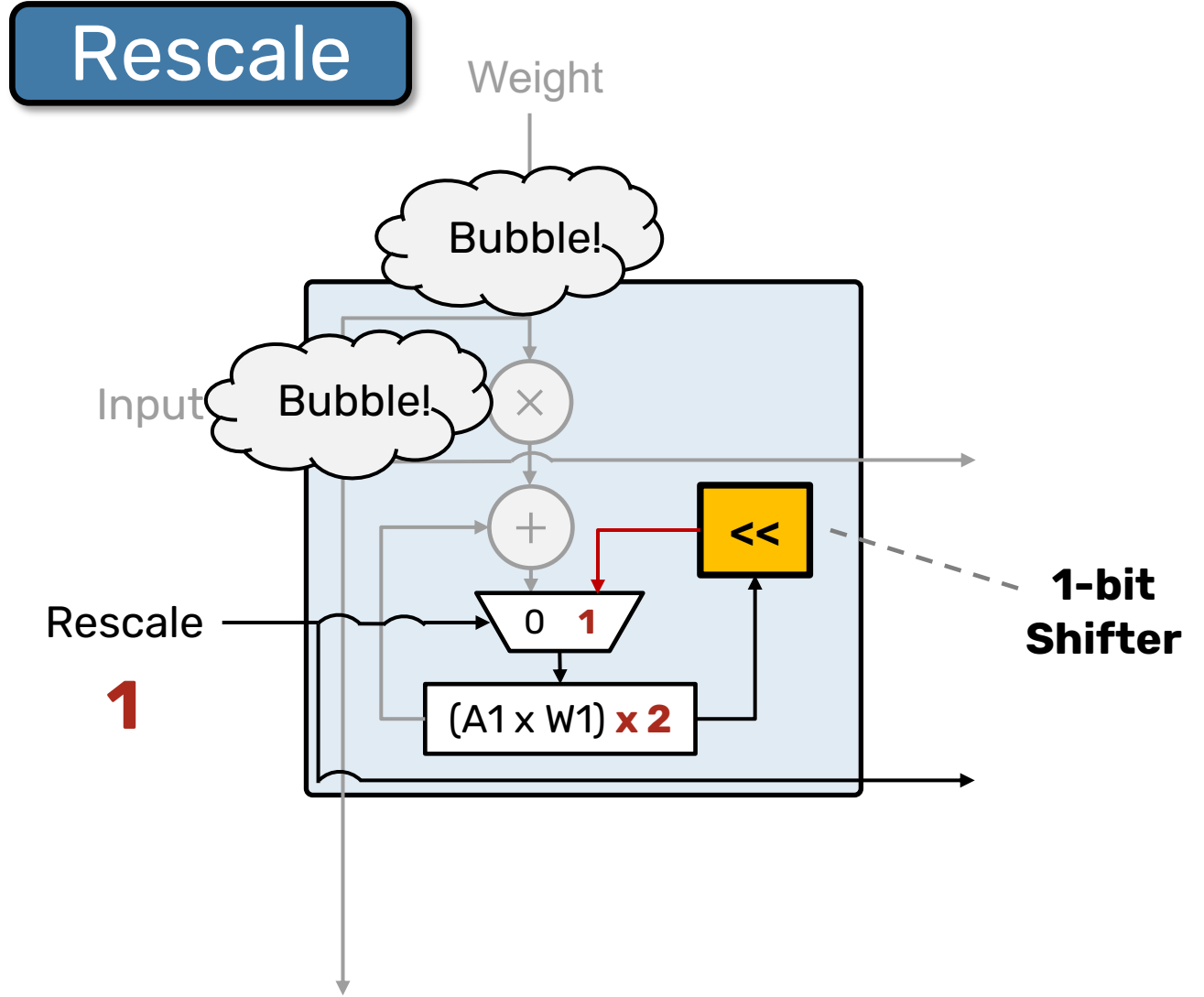
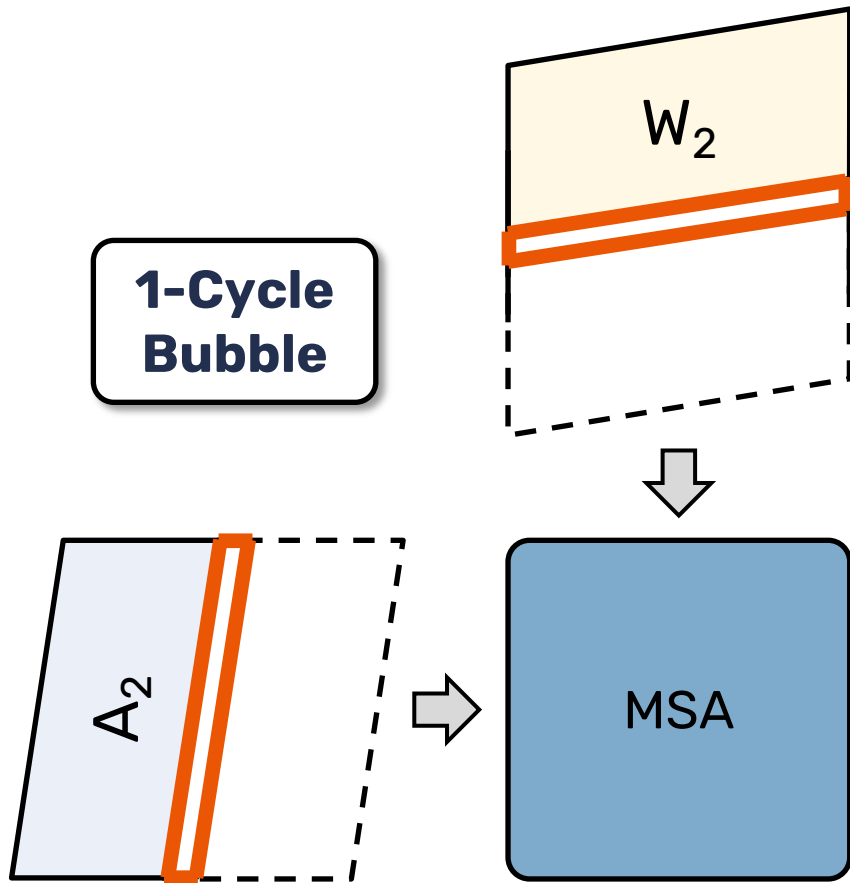
Output-stationary Dataflow



Multi-Scale Systolic Array (MSA)



Multi-Scale Systolic Array (MSA)



Outline

- Motivation
 - Challenges in Efficient LLM Inference
 - Limitations of Prior Works
- **Tender**: Algorithm-Hardware Co-design for Efficient LLM Inference
 - Tensor Decomposition
 - Rescaling Operation
- Evaluation
- Conclusion

Methodology

Models

- OPT, LLaMA, and Llama-2

Datasets

- WikiText-2 and Penn Treebank

Accuracy

- Hugging Face Library

Performance

- RTL: 28nm technology
- Cycle-level simulator

Baselines

Accuracy	
SmoothQuant	Column-wise scaling
ANT	Adaptive & Custom Types
OliVe	Adaptive & Custom Types

Performance	
OLAccel	Input - Mixed Precision
ANT	Input - Exponent & Integer
OliVe	Input - Exponent & Integer

Quantization Results

Perplexity results using *WikiText-2* dataset

* Lower is better

Precision	Scheme	OPT-66B	Llama-2-70B
FP16	Base	9.34	3.32
	SmoothQuant	9.87	17.30
INT8	OliVe	9.43	50.94
	Tender	9.43	3.48

Isolation of outliers

Quantization Results

Perplexity results using *WikiText-2* dataset

* Lower is better

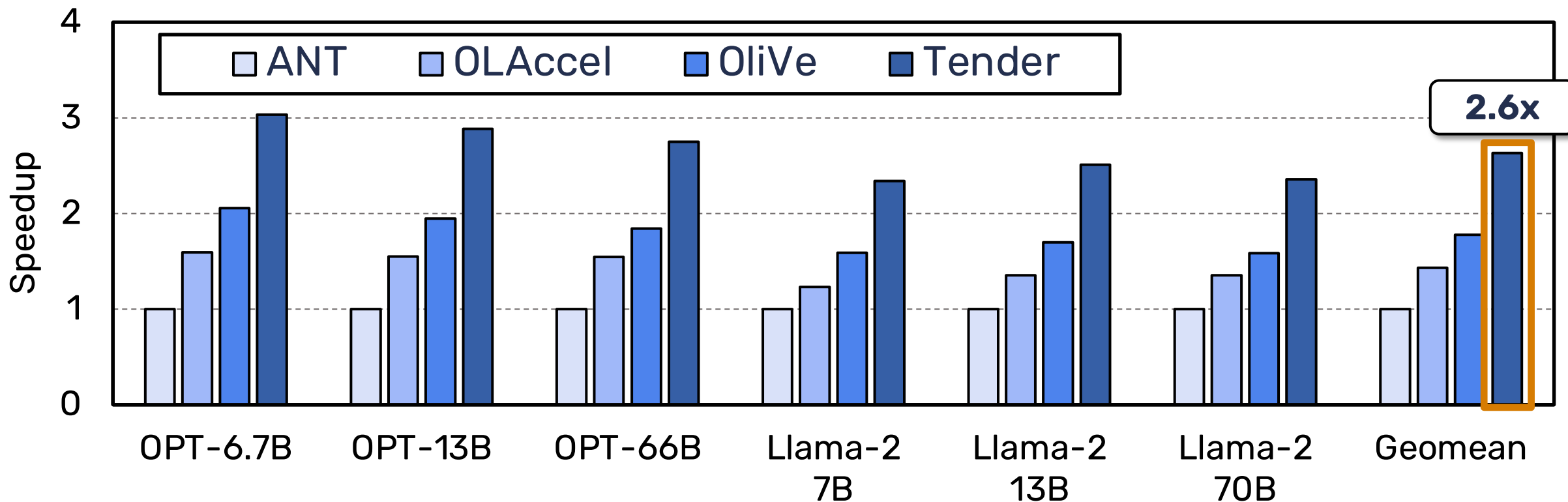
Precision	Scheme	OPT-66B	Llama-2-70B
FP16	Base	9.34	3.32
	SmoothQuant	9.87	17.30
INT8	OliVe	9.43	50.94
	Tender	9.43	3.48
INT4	SmoothQuant	6E+4	7E+4
	OliVe	6E+3	99.91
	Tender	12.38	13.43

Isolation of outliers

More **important**
in **low-precision**

Performance

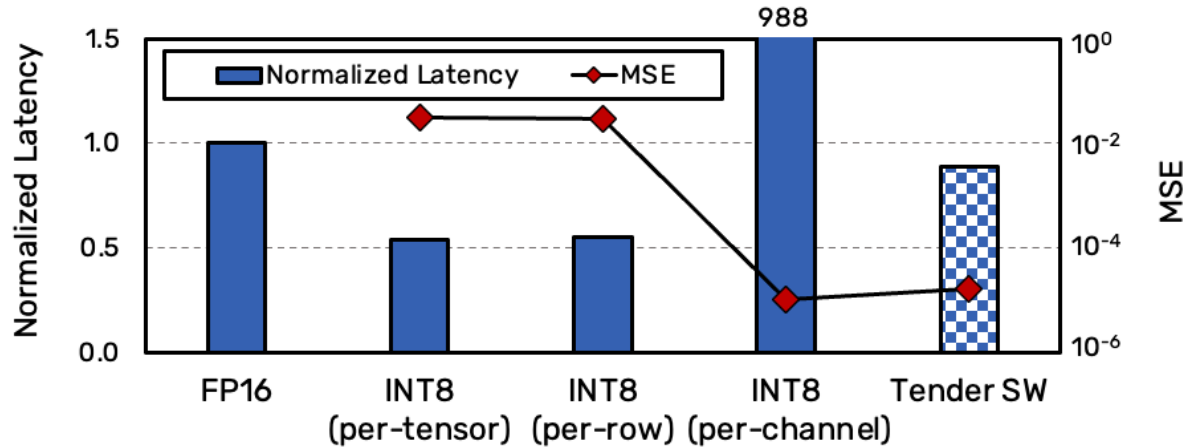
LLM inference speedup



→ With **higher accuracy**, Tender achieves **higher performance**

More Details in Our Paper

- GPU Implementation of Tender
- Tender on weight-stationary dataflow
- Hardware support for reordering
- Comparison with BFP variants
 - MSFP and MX formats
- Area & Energy Efficiency
- Others...



Conclusion

Problem

- Outliers make an efficient serving of LLM challenging
- Complex and intrusive design of prior works

Solution: Tender, efficient low-bit integer-based LLM inference accelerator

- Tensor decomposition while considering accuracy and performance
- Rescaling only requires a 1-bit shifter and 1-cycle latency

Result

- Tender achieves up to an average of **2.6x speedup** over the baseline with **substantially higher accuracy** 😊

Thank You!

Tender

Accelerating Large Language
Models via **Tensor Decomposition**
And **Runtime Requantization**

Jungi Lee (jungli.lee@snu.ac.kr)