

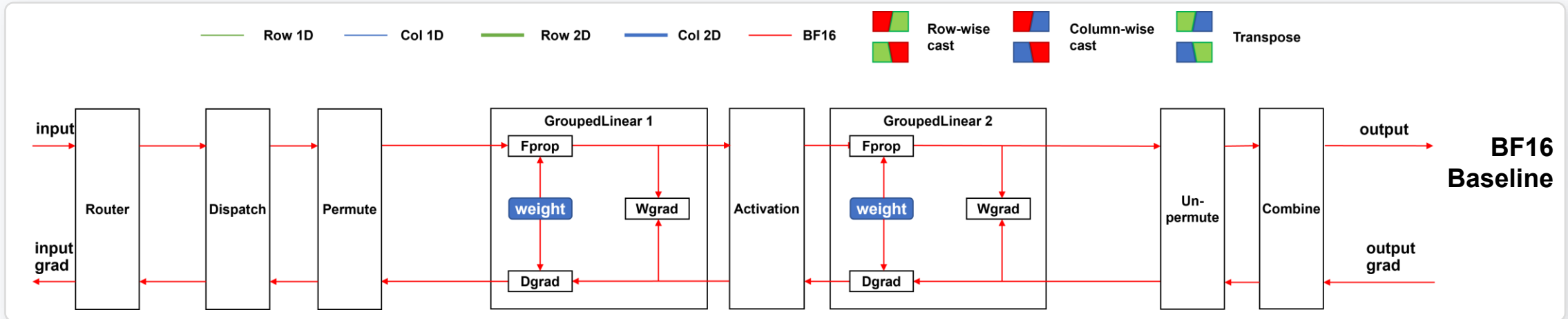
# FP8-Flow-MoE: A CASTING-FREE FP8 RECIPE WITHOUT DOUBLE QUANTIZATION ERROR

Fengjuan Wang, Zhiyi Su, Xingzhu Hu, Cheng Wang, Mou Sun

Zhejianglab

MLSys 2026

# 1. Background and Motivation: The Computational Bottlenecks of MoE



- BF16 MoE: the critical path of one MoE layer

## Background

MoE routes each token to only a small subset of experts, increasing model capacity without a proportional increase in compute

## Bottleneck

Under the BF16 baseline, communication, token reordering, and expert computation together dominate the training cost.

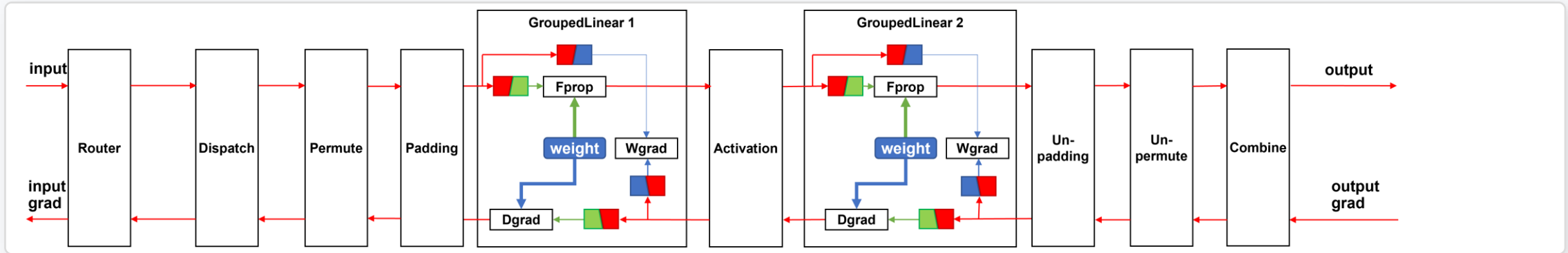
## System View

One MoE layer includes routing, dispatch, permute, expert GEMMs, activation, and combine.

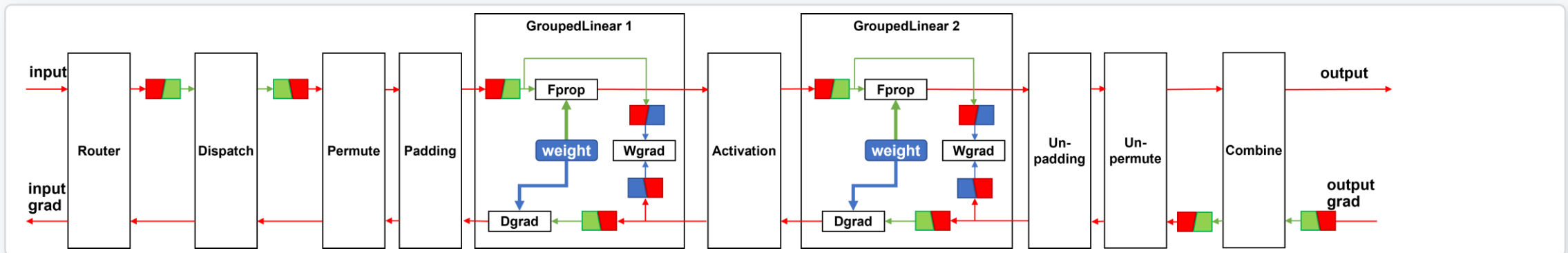
## Motivation

MoE reduces active computation, but not communication, data movement, or memory pressure; this motivates low-precision training such as FP8.

# 1. Background and Problem: Why Existing FP8-MoE Recipes Are Still Not Enough



- NVIDIA TransformerEngine Blockwise FP8 : FP8 mainly accelerates grouped GEMM, while communication and data movement remain in higher precision; Q/DQ overhead remains around grouped linear layers.

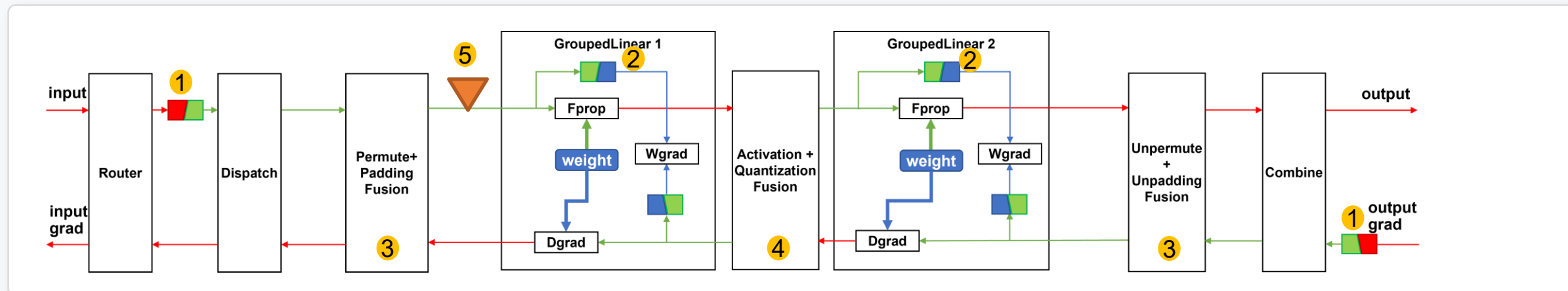


- DeepSeek-V3: FP8 is extended to communication, but one MoE layer still incurs many explicit casts. in the Wgrad path, row-wise FP8 activations must be dequantized, transposed, and requantized into column-wise FP8, which introduces double quantization error.

## 2. Method: FP8-Flow-MoE



💡 Core principle: Keep tensors in FP8 along the MoE expert path whenever possible, and retain BF16 only at numerically sensitive boundaries.



### Persistent FP8 Path

1 1 Quantize before dispatch and keep dispatch → permute → expert in FP8

### Scaling-aware Transpose

2 Convert layouts directly in FP8 and avoid double quantization error.

### Fused FP8 Operators

3 4 Fused permute\_padding, unpermuted\_unpadding, act + quant, to reduce kernel fragmentation.

### Fine-grained Recomputation

5 Recompute at the expert level with FP8 activation compression

Explicit casts

12 → 2

TGS gain

+21%

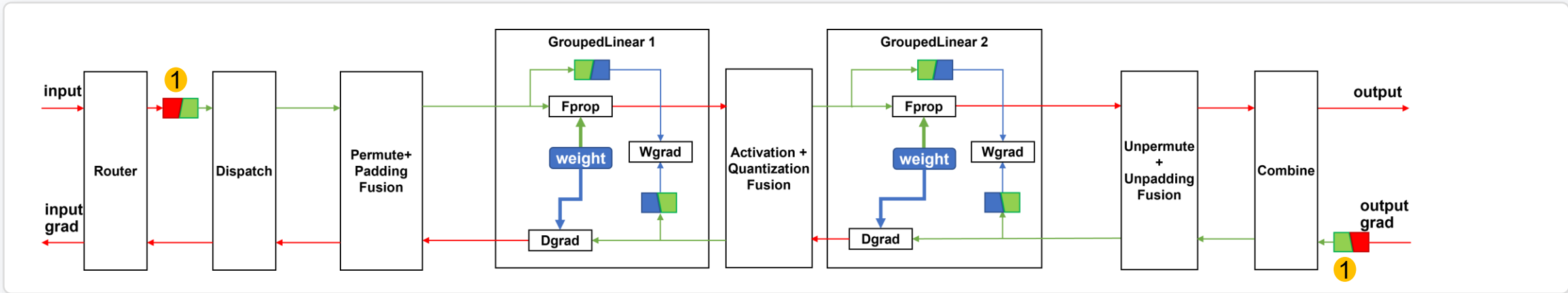
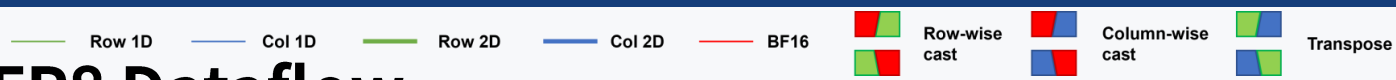
Peak memory reduction

-16.5 GB

Convergence

≈ BF16 Baseline

# 3.Key Technique I: Casting-free FP8 Dataflow



**⚡ Keep the critical MoE path in FP8 (Casting-free)**

- Quantize before dispatch, then keep dispatch → Permute → GroupedLinear(Expert) in FP8
- Benefits:** reduced cast overhead, FP8 communication, smaller FP8 buffers, and lower data-movement cost

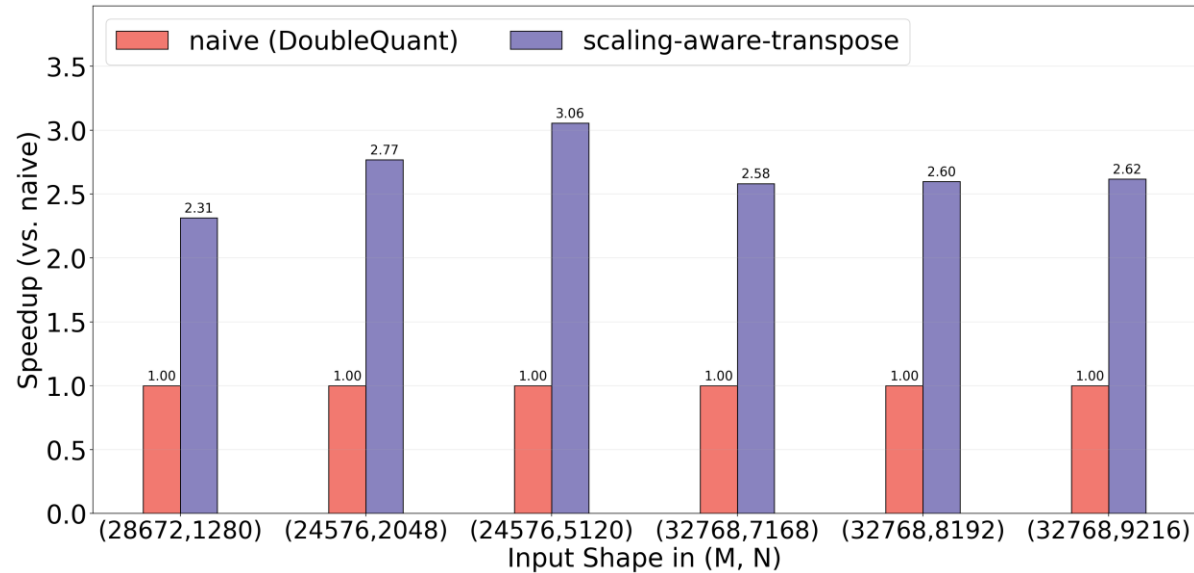
**🛡️ Retain BF16 only at two numerically sensitive boundaries**

- Grouped Linear 1 ↔ Activation: nonlinear operations are sensitive to quantization noise
- Grouped Linear 2 ↔ Combine/Reduction: reduction requires larger dynamic range and is more prone to overflow.

Keep the high-cost MoE path in FP8 whenever possible, while retaining BF16 only at necessary boundaries for stability.



## 4.Key Technique II : Scaling-aware FP8 Transpose (Micro-benchmark)



- Compared with the naive (dequantize → transpose → quantize) implementation, scaling-aware FP8 transpose operator up to **3x** speedup

## 5. Key Technique III: High-Performance FP8 Kernels (Permute + Padding)

### Fused Permute + Padding

#### ⚡ Motivation

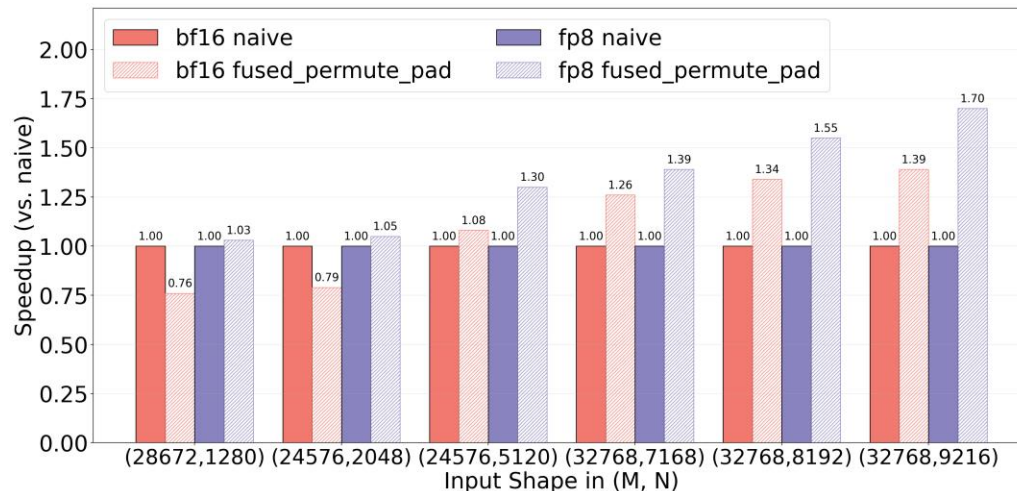
- Permute and padding are memory-bound operators and together account for a large fraction of MoE-layer latency

#### 🛡️ Benefits

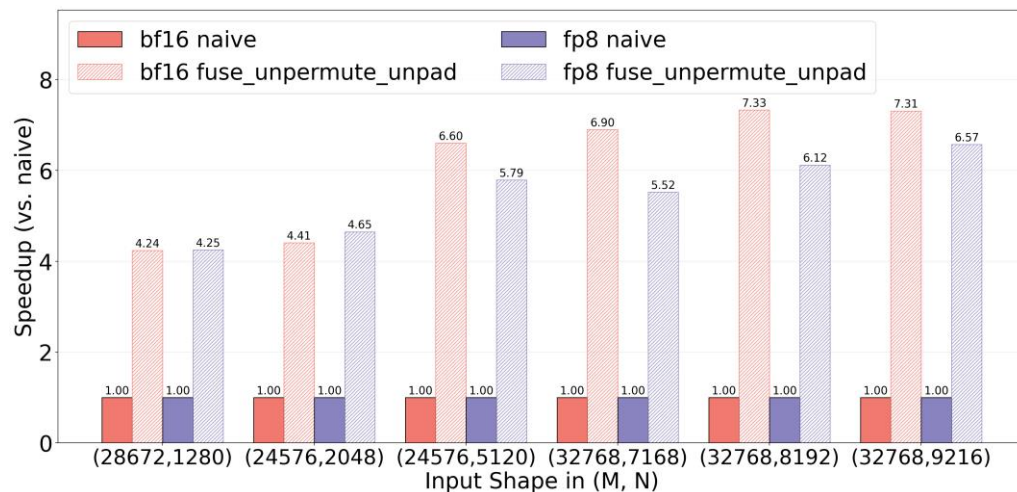
- Fusing them reduces HBM traffic, kernel launches, and intermediate buffer overhead

#### 🚀 Micro-benchmarks

- Fused permute+pad: up to 1.7× speedup
- Fused unpermute+unpad: up to 6.6× speedup
- The fusion also reduces peak memory by about ~1 GB in end-to-end training.



- Fused permute+padding case, we observe up to **1.7×** speedup



- Fused unpermute+unpadding kernel, the improvement reaches as **6.6×**

# 5. Key Technique III: High-Performance FP8 Kernels(SwiGLU + Quant)

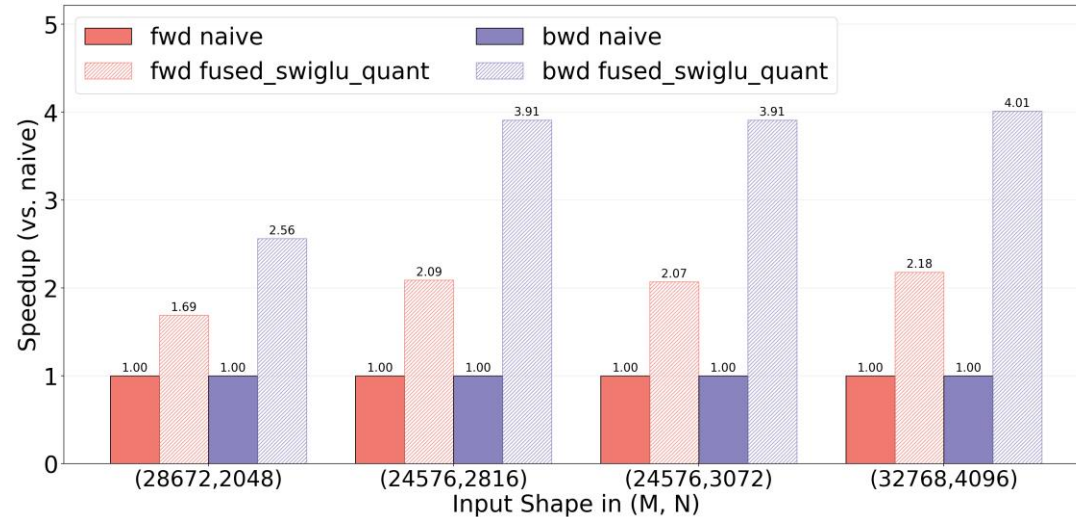
## Fused SwiGLU + Quant

### ⚡ Motivation

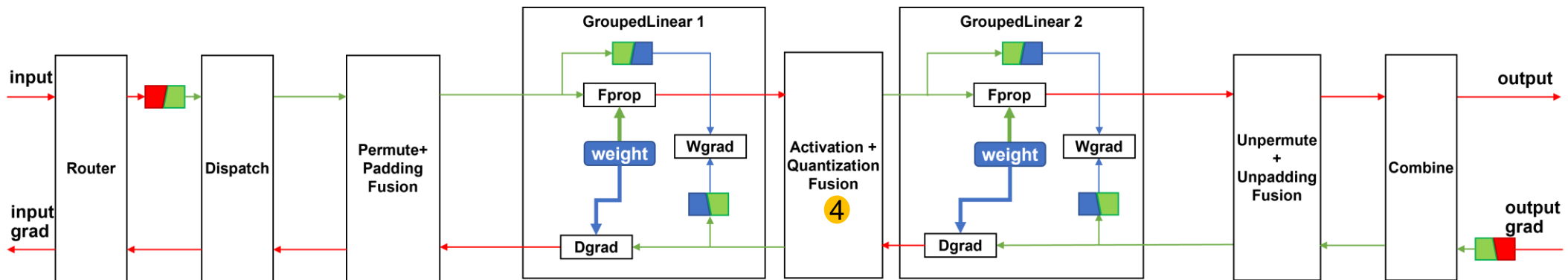
- In a naive FP8 pipeline, quantization around activation introduces extra kernels on the critical path.
- We fuse SwiGLU and quantization so that activation directly produces FP8 outputs.

### 🚀 Micro-benchmarks

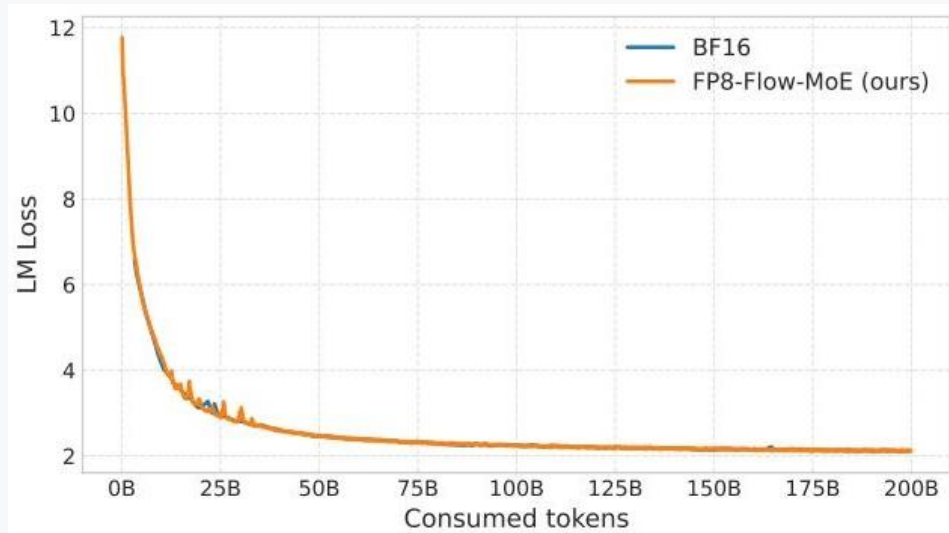
- Forward: up to 2.18× speedup
- Backward: up to 4× speedup
- This reduces kernel launch overhead and avoids extra activation round-trips.



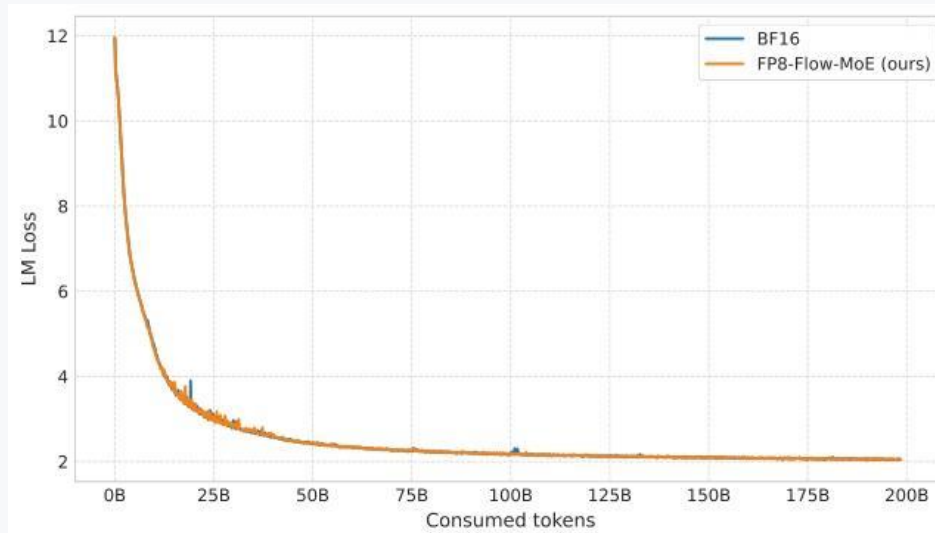
- In the forward pass, the fused kernel achieves **2.18×** speedup.
- In the backward pass, reaching **4×**



## 6. Convergence Validation



Training loss on DeepSeek-V2-Lite (16B)



Additional validation on Qwen3-30B-A3B

**Setup:** DeepSeek-V2-Lite is trained from scratch for 200B tokens; Qwen3-30B provides additional convergence validation

**Conclusion:** Train in FP8, converge like BF16; Relative loss error remains below 0.19%.

## 7. End-to-End Efficiency: higher TGS, lower memory, and better scalability at larger EP

TGS (vs BF16)  
**+16%**

TGS (vs TE Blockwise FP8)  
**+21%**

Peak Mem (vs BF16)  
**-8 GB**

Peak Mem (vs TE Blockwise FP8)  
**-16.5 GB**

Method	EP8		EP16		EP32	
	TGS	Mem	TGS	Mem	TGS	Mem
<b>BF16</b>	1,109	39	939	<b>36</b>	671	<b>43</b>
<b>Blockwise</b>	1,146	37	938	41	644	51
<b>Tensorwise</b>	1,130	38	989	42	660	53
<b>FP8-Flow-MoE</b>	<b>1,176</b>	37	<b>1,012</b>	39	<b>779</b>	49

Table 1: Throughput analysis on DeepSeek-V3 under AC=full

Method	EP8		EP16		EP32	
	TGS	Mem	TGS	Mem	TGS	Mem
<b>BF16</b>	1,178	64	1,055	71	OOM	OOM
<b>Blockwise</b>	1,178	73	1,031	77	OOM	OOM
<b>Tensorwise</b>	<b>1,231</b>	73	OOM	OOM	OOM	OOM
<b>FP8-Flow-MoE</b>	1,193	<b>56</b>	<b>1,111</b>	<b>66</b>	<b>912</b>	<b>75</b>

Table 2: Throughput analysis on DeepSeek-V3 under AC=sel (+MoE expert)

**Primary benchmark:** DeepSeek-V3 (671B), evaluated under multiple expert-parallel and checkpointing settings.

**Platform:** 32-node NVIDIA Hopper cluster, 8 GPUs/node, 80 GB/GPU

**Throughput:** up to +16% vs BF16, and +21% vs TE blockwise FP8

**Peak memory:** up to -8 GB vs BF16, and -16.5 GB vs TE FP8 baselines At EP=32, some baselines run out of memory, while FP8-Flow-MoE remains stable.

# THANKS

## Q & A

**Upstream progress: TransformerEngine merged several PRs; Megatron-LM support is in progress.**

- <https://github.com/NVIDIA/Megatron-LM/pull/2764> (Open)
- <https://github.com/NVIDIA/TransformerEngine/pull/2544> (Open)
- <https://github.com/NVIDIA/Megatron-LM/pull/2763> (Merged)
- <https://github.com/NVIDIA/Megatron-LM/pull/4038> (Merged)
- <https://github.com/NVIDIA/TransformerEngine/pull/1921> (Merged)
- <https://github.com/NVIDIA/TransformerEngine/pull/2144> (Merged)
- <https://github.com/NVIDIA/TransformerEngine/pull/2026> (Merged)
- <https://github.com/NVIDIA/TransformerEngine/pull/1866> (Merged)
- <https://github.com/NVIDIA/TransformerEngine/pull/2547> (Merged)