# DeepSeekV3-FP8 Training

XiaoTonghuan, GongPing

USTC, CHINA
**ADSLAB**
先进数据系统实验室

# Outline

- Background
- Challenges
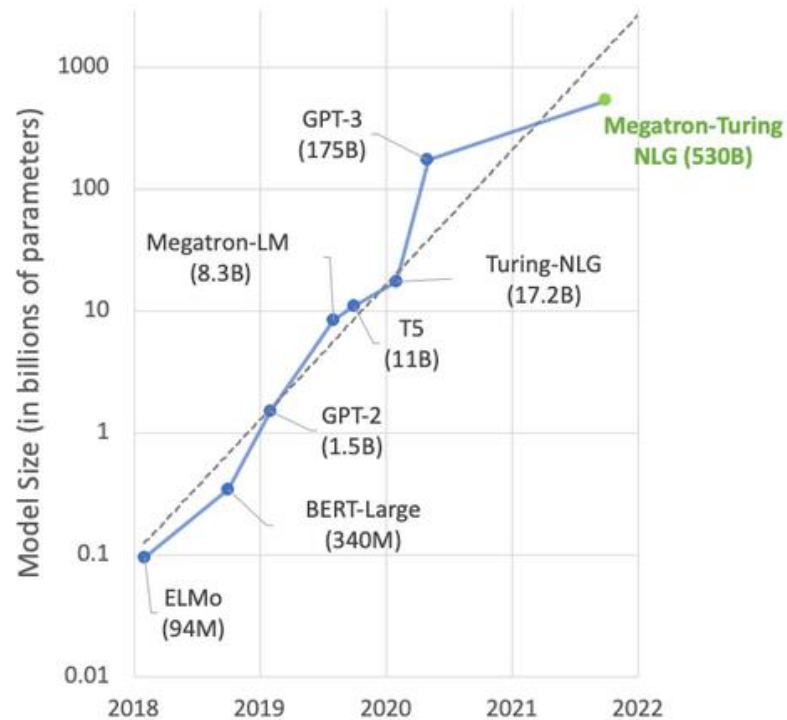- Design
- Implementation

# Outline

- **Background**
- Challenges
- Design
- Implementation

# Background

- Why FP8 training?
  - ◆ Save memory
  - ◆ computing power

# Background

- Why FP8 training?
  - ◆ Save memory



1B parameters occupy
4GB storage (32bit a parameter)

Example: GPT3(175B)
FP32: 700GB
FP8  : 175GB
Save 75% memory

# Background

- Why FP8 training?
  - ◆ computing power

NVIDIA H100 specifications (vs. NVIDIA A100)

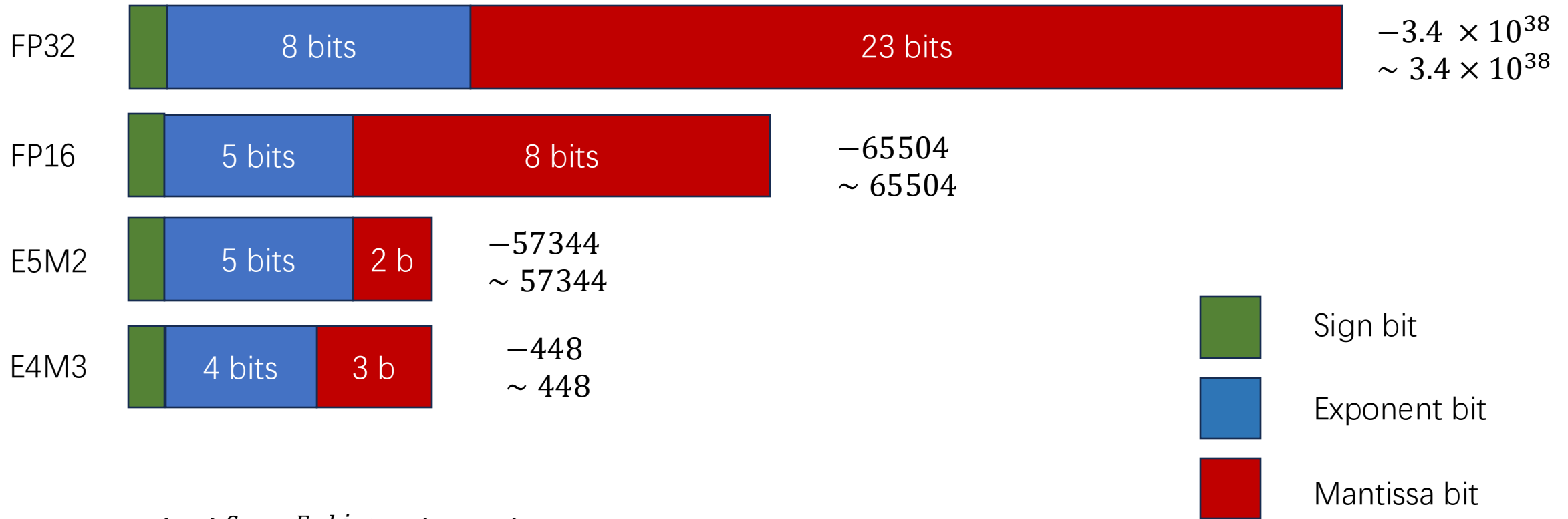| Data type | H100-SXM5 (TFLOPS) | A100-SXM4 (TFLOPS) | Difference |
|---|---|---|---|
| TF32 | 494 | 156 | 3.2x |
| BF16 | 989 | 312 | 3.2x |
| FP16 | 989 | 312 | 3.2x |
| FP8 | 1979 | - | 6.3x (vs BF16) |
| Bandwidth (GB/s) | 3350 | 2039 | 1.6x |

Table 1: FLOPS and memory bandwidth comparison between the NVIDIA H100 and NVIDIA A100. While there are 3x–6x more total FLOPS, real-world models may not realize these gains.

Computing power doubles

# Background

- Data Types
  - FP32, FP16, FP8(E5M2, E4M3)

FP32

| | 8 bits | 23 bits |

$-3.4 \times 10^{38}$
$\sim 3.4 \times 10^{38}$

FP16

| | 5 bits | 8 bits |

$-65504$
$\sim 65504$

E5M2

| | 5 bits | 2 b |

$-57344$
$\sim 57344$

E4M3

| | 4 bits | 3 b |

$-448$
$\sim 448$

Sign bit

Exponent bit

Mantissa bit

$$Value = (-1)^S \times 2^{E-bias} \times (1 + M)$$

# Background

- Data Types
  - FP32, FP16, FP8(E5M2, E4M3)



FP32 | 8 bits | 23 bits | $-3.4 \times 10^{38} \sim 3.4 \times 10^{38}$

FP16 | 5 bits | 8 bits | $-65504 \sim 65504$

E5M2 | 5 bits | 2 b | $-57344 \sim 57344$

E4M3 | 4 bits | 3 b | $-448 \sim 448$ — What we use, but Why?
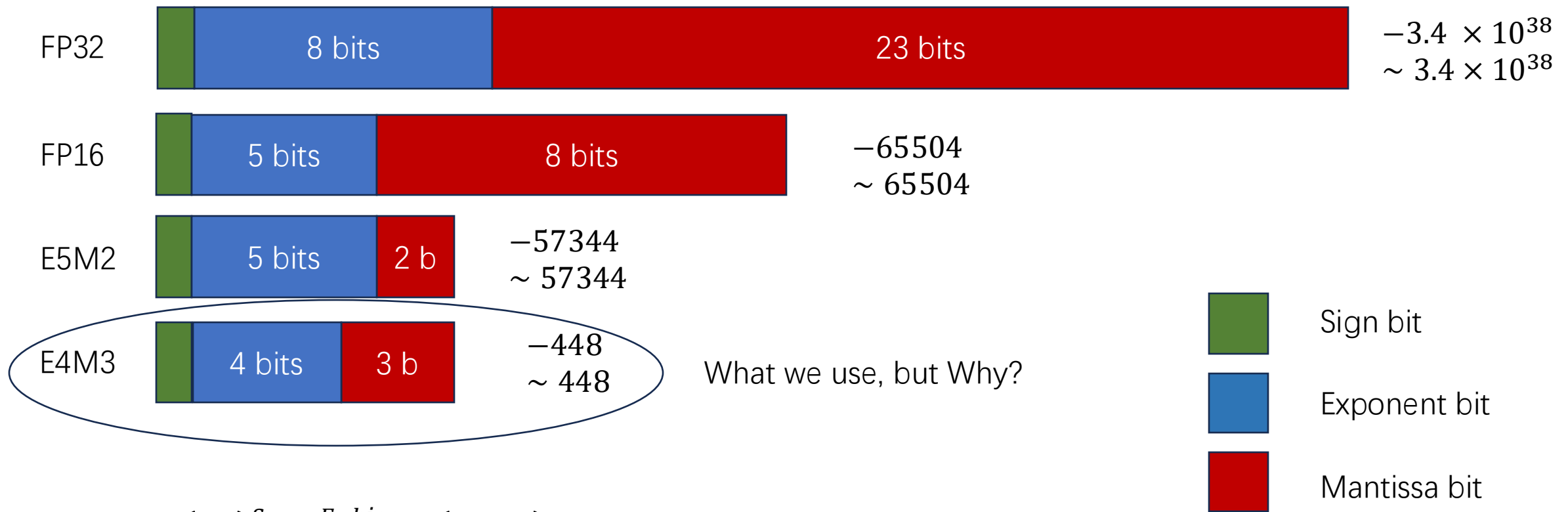
Sign bit
Exponent bit
Mantissa bit

$$Value = (-1)^S \times 2^{E-bias} \times (1 + M)$$

# Background

- Data Types
  - ◆ FP32, FP16, FP8(E5M2, E4M3)



$$Value = (-1)^S \times 2^{E-bias} \times (1+M)$$

What we use, but Why?
1. LLM normalization means narrow range of value
2. Low precision of E5M2 is too difficult to overcome

# Outline

- Background
- **Challenges**
- Design
- Implementation

# Difficulty in using FP8(E4M3)

- FP8 low precision
  - Cannot satisfy all precision requirements in training
    - Precision problem is too hard to overcome, so not all FP8
      - Where to use FP8, Where to maintain original format(BF16,FP32)
    - Narrow range cause overflow/underflow in Conversion
      - Conversion between different Floating point numbers

# Outline

- Background
- Challenges
- Design
  - How to handle mixed precision
    - Where to use FP8 ?
    - Conversion between different format
- Implementation
  - Deep GEMM
  - Mixed precision Framework

# Outline

- Background
- Challenges
- Design
  - ◆ How to handle mixed precision
    - ▪ Where to use FP8 ?
    - ▪ Conversion between different format
- Implementation
  - ◆ Deep GEMM
  - ◆ Mixed precision Framework
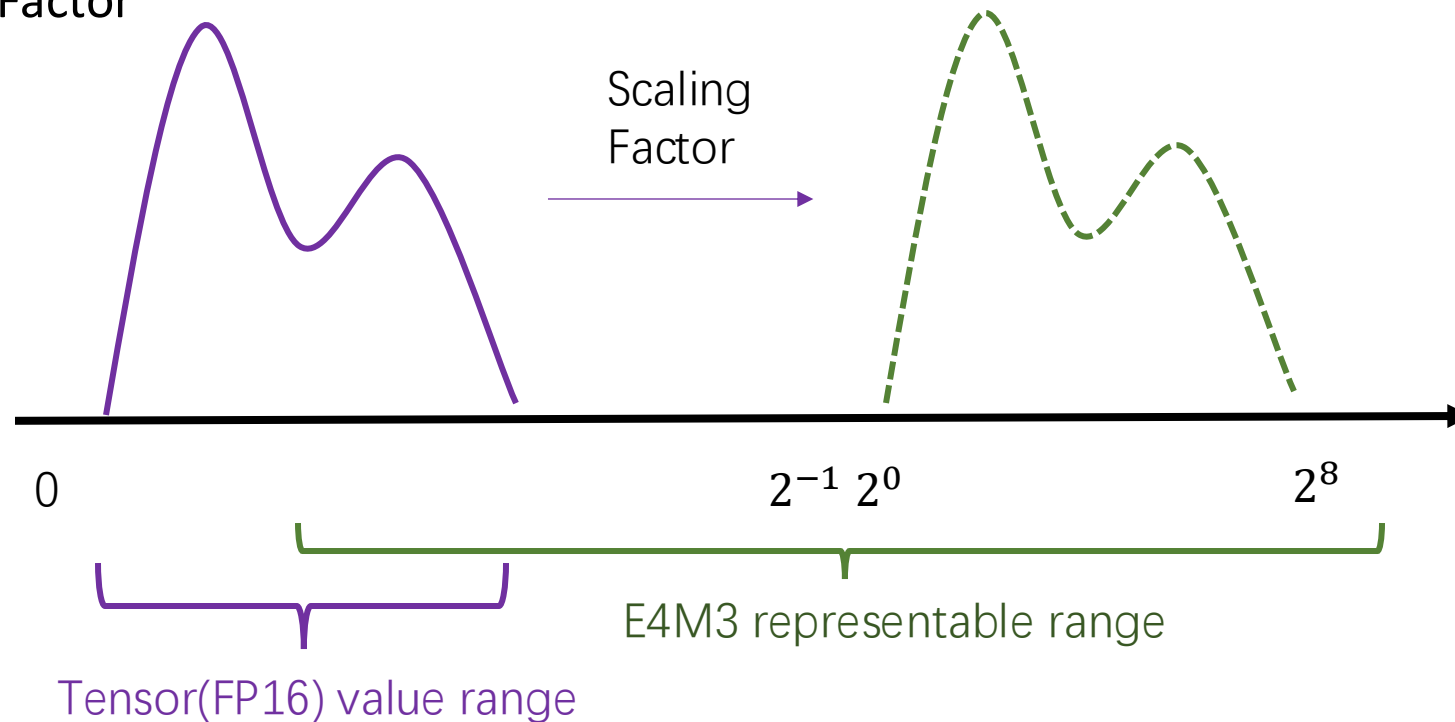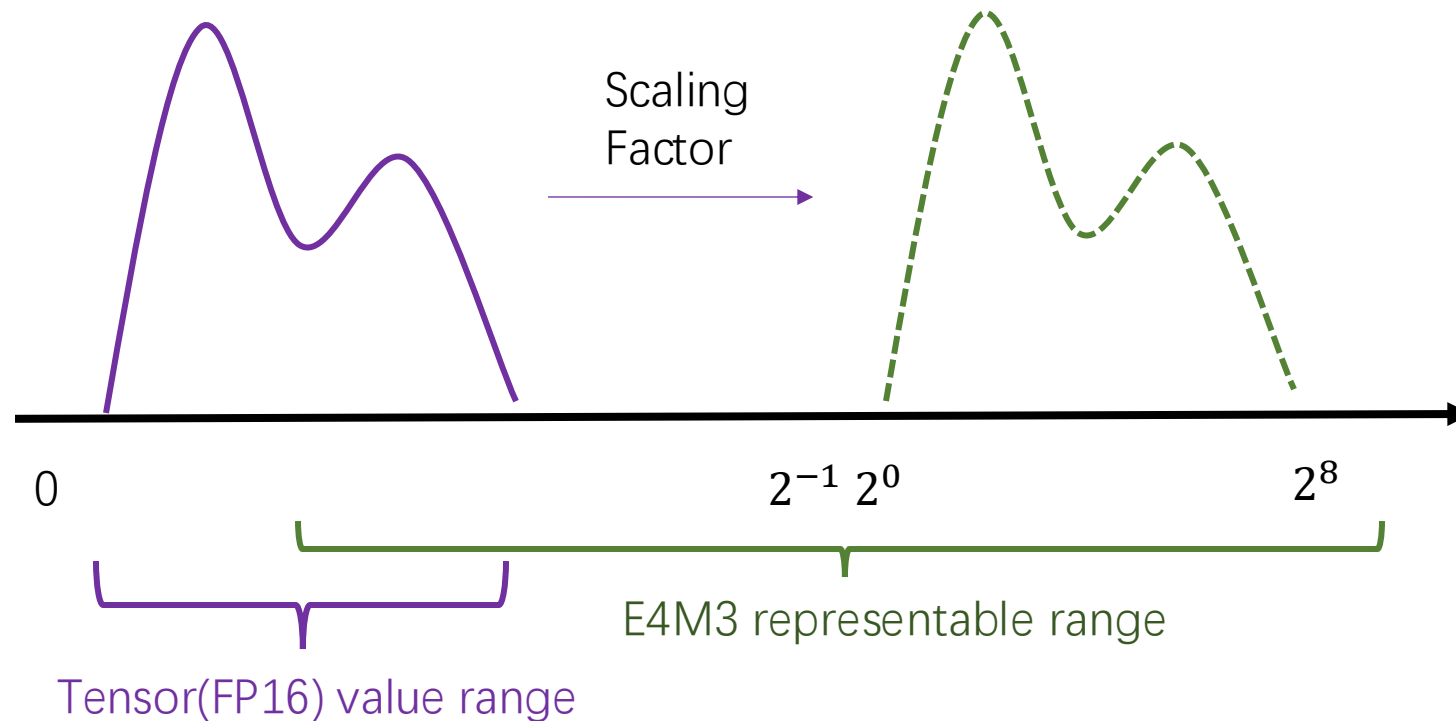
# Where to use FP8 ?

- Determine where to use FP8/BF16/FP32 by careful investigations
  - ◆ BF16/FP32
    - Embedding module(low utilization rate)
    - Output head(Low utilization rate)
    - MoE gating modules(only 1%~5% overhead of MoE)
    - normalization operators(A high precision requirement (e.g., 1e-6) )
    - attention operators(A high precision requirement )
    - Weights(Weight update use FP32, Computation use FP8)
    - weight gradients(FP32, low computing power consume)
    - optimizer states(BF16, BF16 is enough for Optimizer DeepSeek experiment proof)
  - ◆ FP8
    - MLP on MoE
    - MLP before/after Attention

# Where to use FP8 ?

- Determine where to use FP8/BF16/FP32 by careful investigations
  - ◆ BF16/FP32
    - ■ Embedding module(low utilization rate)
    - ■ Output head(Low utilization rate)
    - ■ MoE gating modules(only 1%~5% overhead of MoE)
    - ■ normalization operators(A high precision requirement (e.g., 1e-6) )
    - ■ attention operators(A high precision requirement )
    - ■ Weights(Weight update use FP32, Computation use FP8)
    - ■ weight gradients(FP32, low computing power consume)
    - ■ optimizer states(BF16, BF16 is enough for Optimizer DeepSeek experiment proof)
  - ◆ FP8
    - ■ MLP on MoE
    - ■ MLP before/after Attention

What important is the law
1. Not used in scenarios with low computational demands
2. Not used in applications requiring high precision

# Outline

- Background
- Challenges
- Design
  - How to handle mixed precision
    - Where to use FP8 ?
    - Conversion between different format
- Implementation
  - Deep GEMM
  - Mixed precision Framework

# Overflow/Underflow

- Overflow/Underflow
  - ◆ What's the problem
  - ◆ Scaling Factor



Scaling Factor

$2^{-1}$ $2^0$ $2^8$

0

Tensor(FP16) value range

E4M3 representable range

# Conversion between different format

- High precision to Low precision
  - Scaling (Make sure the value is within the FP8 range)
  - Cast(type conversion)



Scaling Factor

$0$        $2^{-1}$ $2^{0}$      $2^{8}$

E4M3 representable range

Tensor(FP16) value range

# Conversion between different format

- High precision to Low precision
  - Scaling (Make sure the value is within the FP8 range)
  - Cast (Type casting)
- Low precision to High precision
  - Type casting to a wider scope
  - multiply Scaling Factor

# Scaling Factor

- Use Dynamic Scaling Factor
  - Dynamic: Obtained during calculation
  - Scaling Factor: max(abs(x)) / MAX_E4M3

# Scaling Factor

- Use Dynamic Scaling Factor
  - ◆ Dynamic: Obtained during calculation
  - ◆ Scaling Factor: max(abs(x)) / MAX_E4M3

What is X ?
Entire tensor or part of tensor ?

# Scaling Factor

- What is X ? (Select an area to find scaling factor)



| Per-tensor | Tile-wise | Block-wise |

# Scaling Factor

- Select an area to find scaling factor

Use before Deep Seek
(e.g. NVDIA)

$N_C$

1

$N_C$

$N_C$

$N_C$

Per-tensor

Tile-wise

Block-wise

# Scaling Factor

● Select an area to find scaling factor

Deep Seek V3 Use different method for different value

$N_C$

1

$N_C$

$N_C$

Per-tensor

Tile-wise

Block-wise

# Scaling Factor

- Select an area to find scaling factor
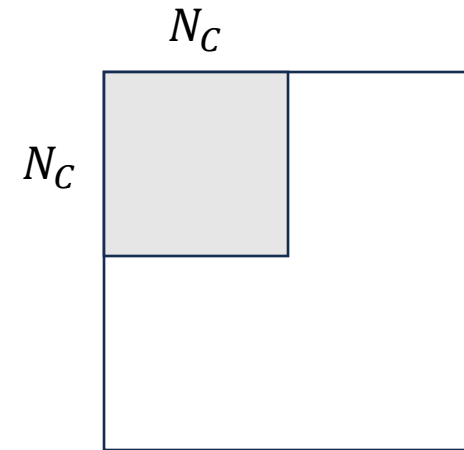  - ◆ Activation: tile-wise
  - ◆ Weight: block-wise



Per-tensor      Tile-wise      Block-wise

# Scaling Factor

- Select an area to find scaling factor
  - ◆ Activation: tile-wise
  - ◆ Weight: block-wise     Different distribution of outlier between activation and weight



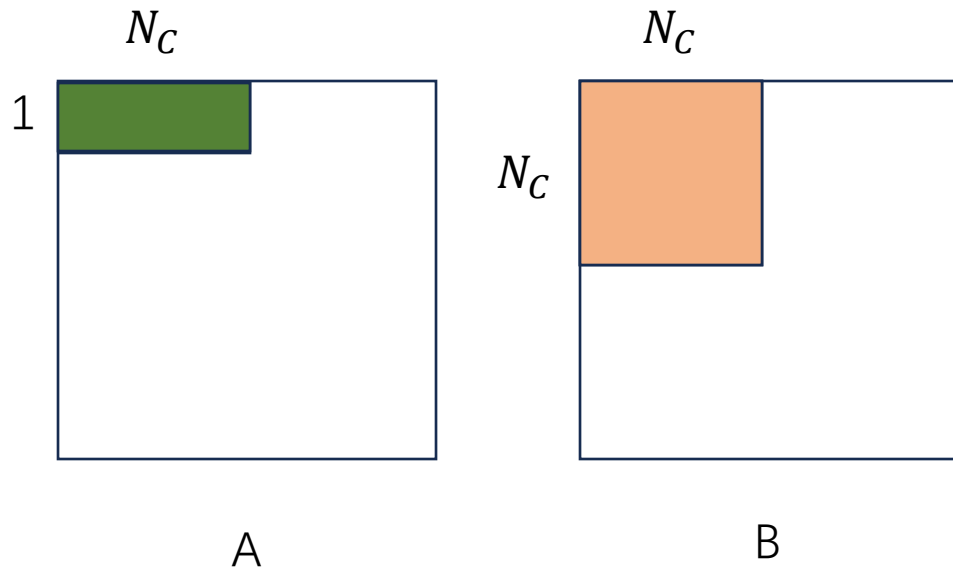Per-tensor        Tile-wise        Block-wise

# Outline

- Background

- Challenges

- Design
  - ◆ How to handle mixed precision
    - ▪ Where to use FP8 ?
    - ▪ Conversion between different format

- Implementation
  - ◆ Deep GEMM
  - ◆ Mixed precision Framework
  - ◆ Other

# GeMM

- GeMM(General Matrix-Matrix Multiplication)
- A @ B = C(FP8,FP8, FP32)
  - ◆ Note: A, B with Scaling Factor
  - ◆ Blocks of matrix process MMA
    - ■ MMA operation (wgmma instruction, FP32)
  - ◆ Merge block results
    - ■ Multiply scaling factor
    - ■ Combine partial sum
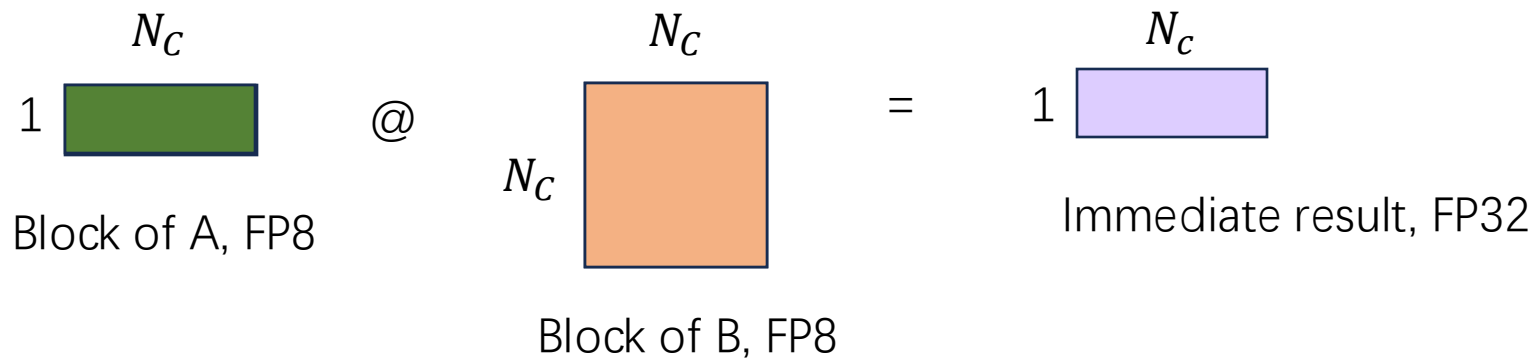  - ◆ Use two warpgroup(one execute mma operation, another merge)

# GeMM

- GeMM(General Matrix-Matrix Multiplication)
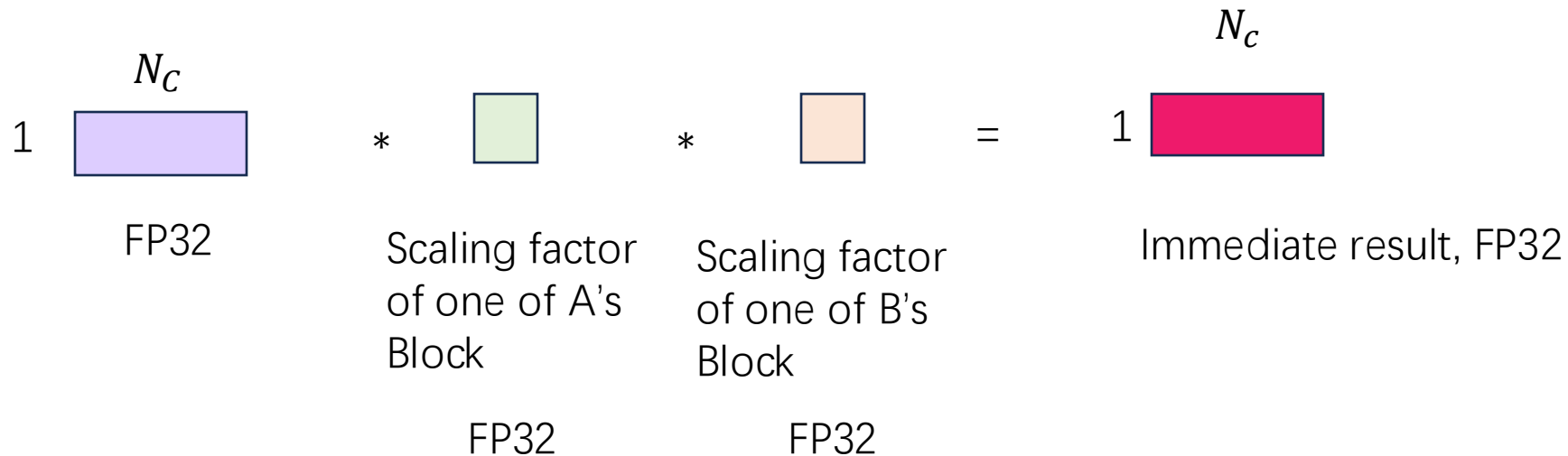- A @ B = C
  - ◆ Blocks of matrix process MMA



A

B

# GeMM

- GeMM(General Matrix-Matrix Multiplication)
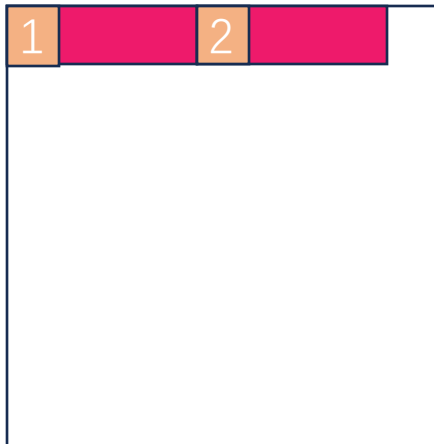- A @ B = C
  - ◆ Blocks of matrix process MMA

$N_C$

1 ▮

@

$N_C$

$N_C$

=

1 ▮

$N_C$

Block of A, FP8

Block of B, FP8

Immediate result, FP32

# GeMM

- GeMM(General Matrix-Matrix Multiplication)
- A @ B = C
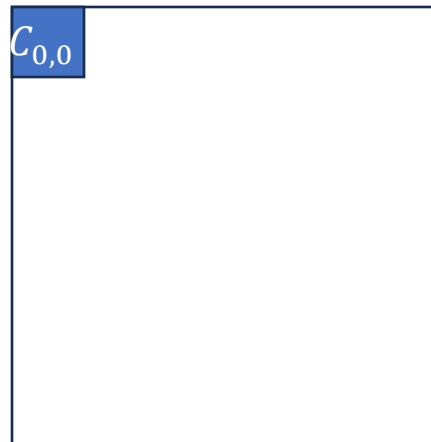  - ◆ Promotion and Merge block results
    - ■ Multiply scaling factor

$N_c$

1 ▭ FP32

*

Scaling factor of one of A's Block

FP32

*

Scaling factor of one of B's Block

FP32

=

$N_c$

1 ▭ Immediate result, FP32

# GeMM

- GeMM(General Matrix-Matrix Multiplication)
- A @ B = C
  - ◆ Promotion and Merge block results
    - Combine partial sum



Immediate result

C

$Add\ 1,2\ \dots\ to\ C_{0,0}$

# Outline

- Background & Challenges
- Design
  - Overflow/Underflow
  - Where to use FP8 ?
  - Conversion between different floating point
- **Implementation**
  - Deep GEMM
  - **Mixed precision Framework**
  - Other

# Mixed precision Framework
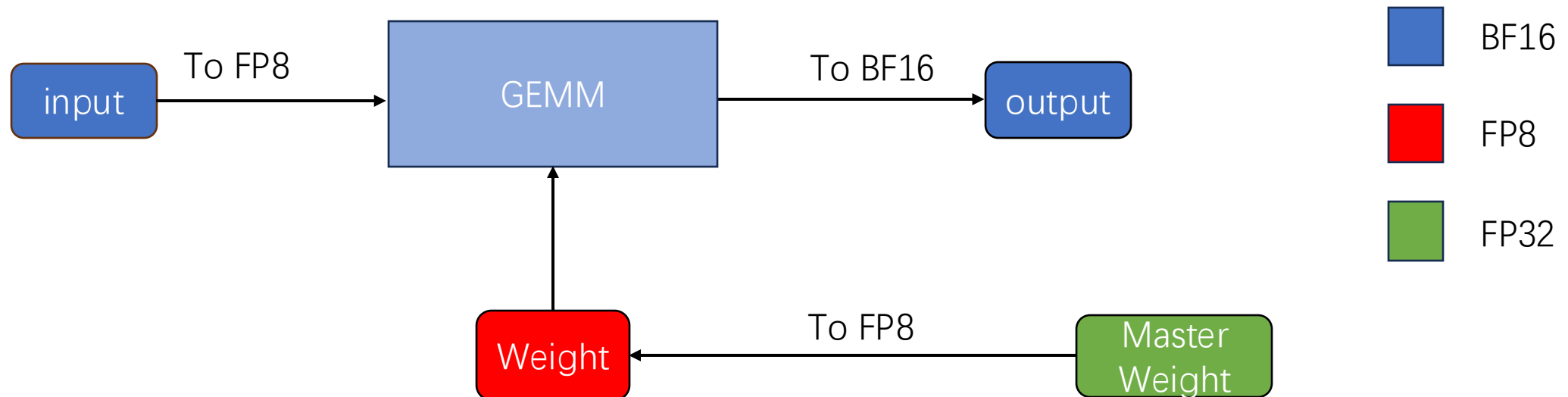
- forward pass
- activation backward pass
- weight backward pass

# Mixed precision Framework

- forward pass
- activation backward pass
- weight backward pass

This partitioning is based on the fact that backward passes rely on two key matrix multiplications.
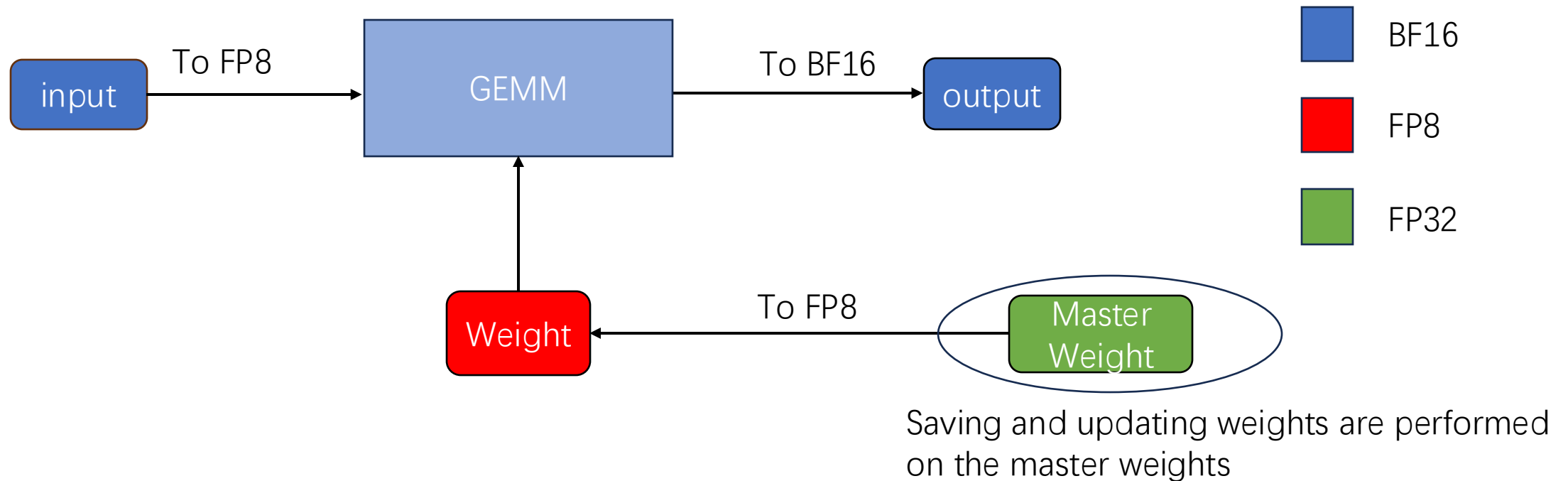
# Mixed precision Framework
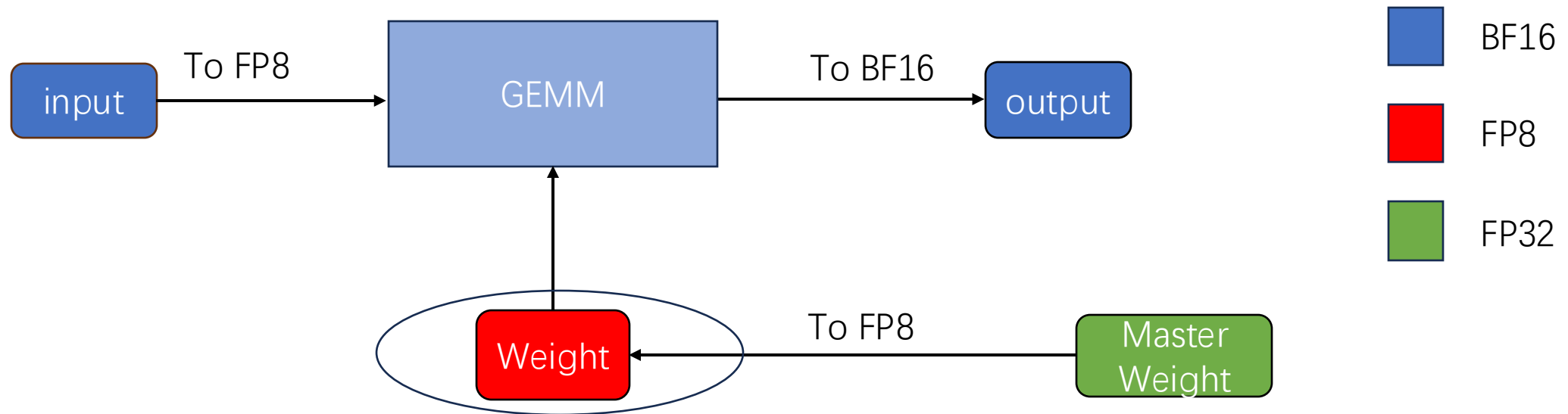
- Mixed precision Framework
  - ◆ forward pass



To FP8 → GEMM → To BF16 → output

input → GEMM

Weight → GEMM

Master Weight → To FP8 → Weight

BF16 (blue)
FP8 (red)
FP32 (green)

# Mixed precision Framework

- Mixed precision Framework
  - forward pass



Saving and updating weights are performed on the master weights

# Mixed precision Framework

- Mixed precision Framework
  - forward pass



input — To FP8 → GEMM — To BF16 → output

Weight ← To FP8 — Master Weight

BF16
FP8
FP32

FP8 weight is used for computations.

# Mixed precision Framework

- Mixed precision Framework
  - forward pass



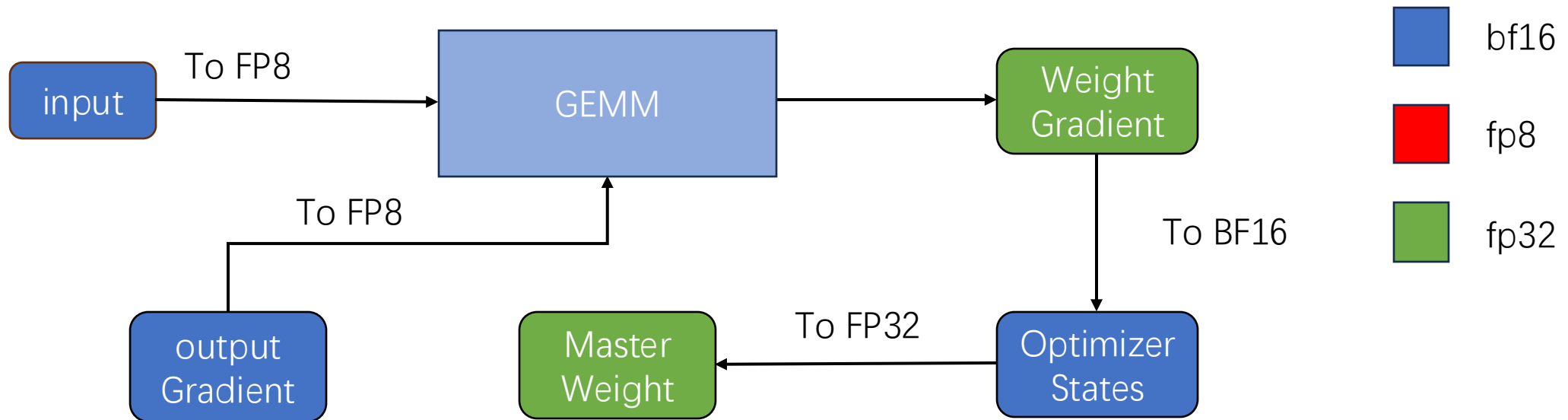Directly updating weights in FP8 can lead to precision loss and vanishing/exploding gradient issues.

# Mixed precision Framework
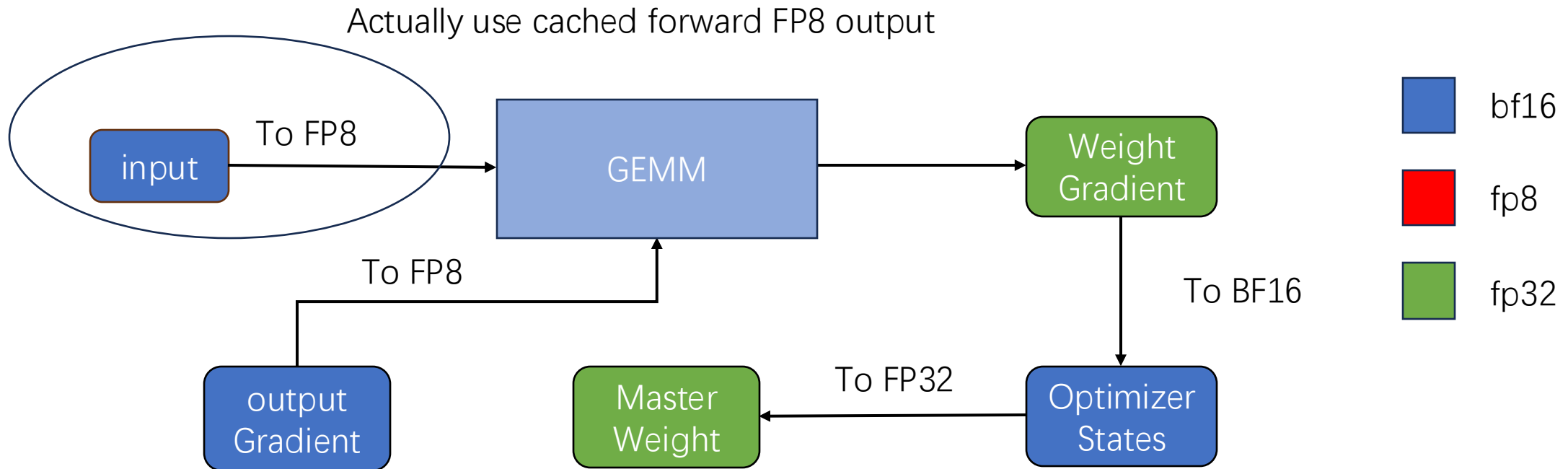
- Mixed precision Framework
  - activation backward pass

# Mixed precision Framework

- Mixed precision Framework
  - weight backward pass

# Mixed precision Framework

- Mixed precision Framework
  - weight backward pass
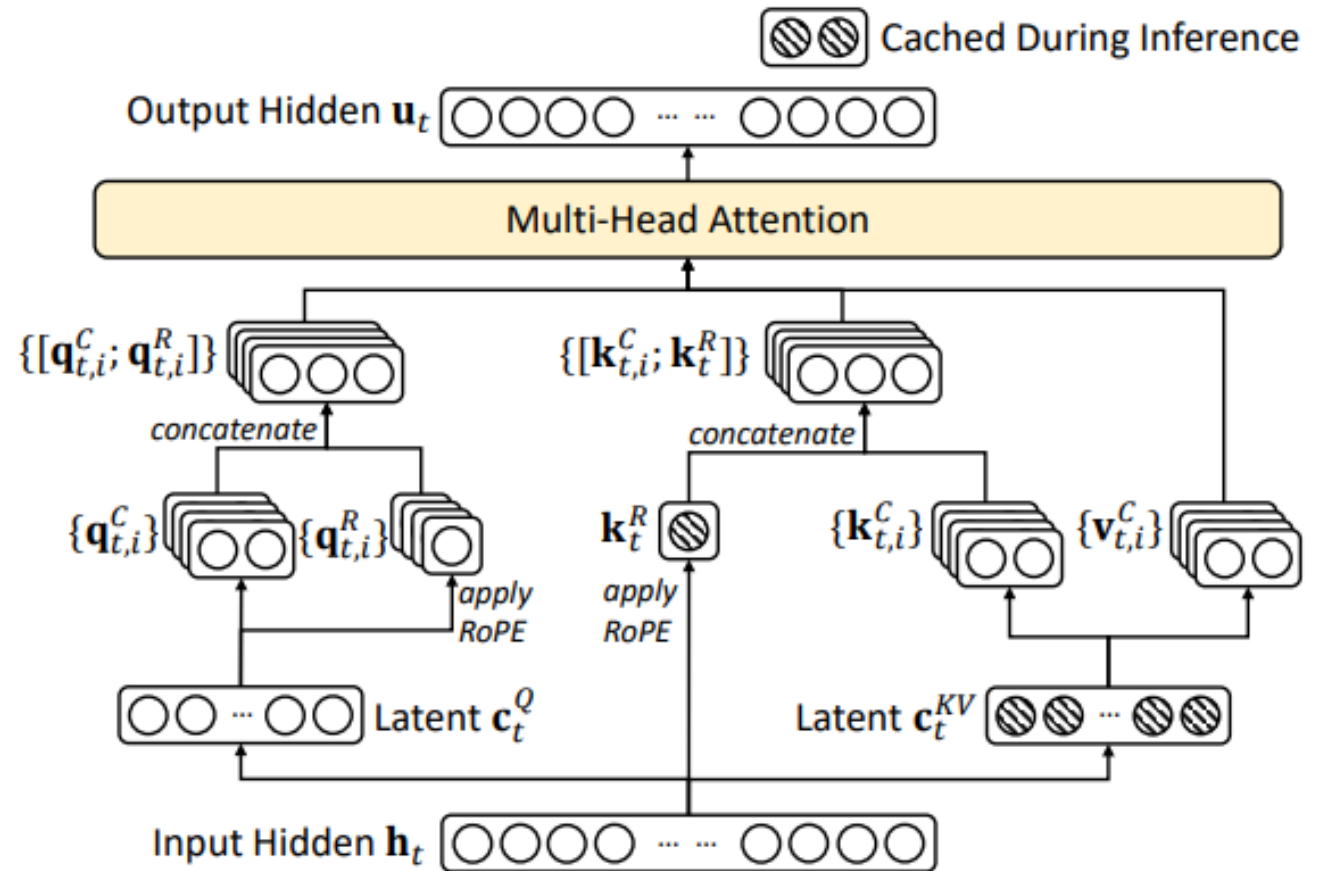
Actually use cached forward FP8 output

# Outline

- Background & Challenges
- Design
  - Overflow/Underflow
  - Where to use FP8 ?
  - Conversion between different floating point
- **Implementation**
  - Deep GEMM
  - Mixed precision Framework
  - **Other**

# Other

- Low-Precision Activation(Compared to BF16 activation)
  - When backward pass Cached FP8 Activation, special considerations are taken on several operators for low-cost high-precision training
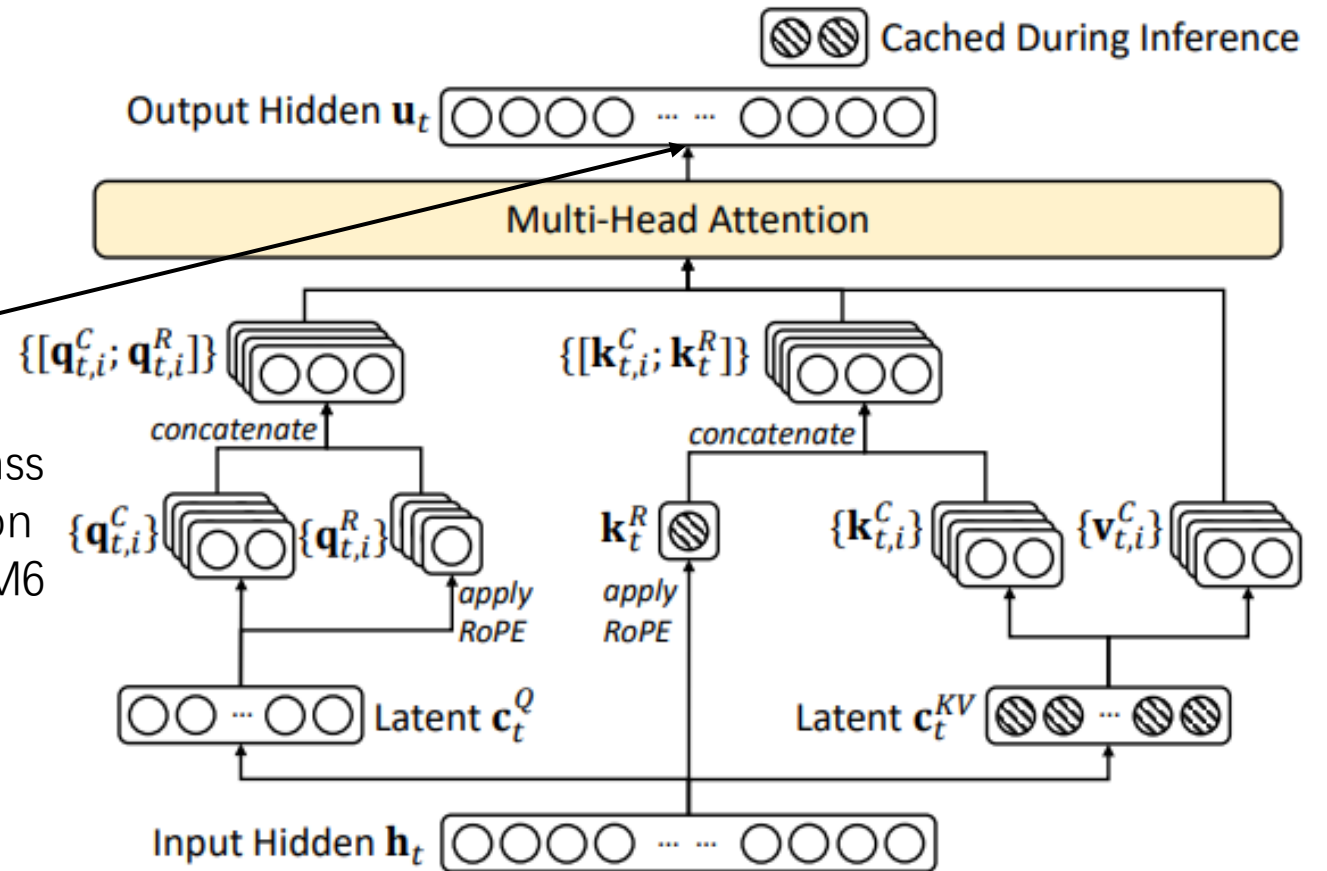  - Attention's backward pass
  - Input of SwiGLU (MoE)

# Other

- Low-Precision Activation(Compared to BF16 activation)
  - Attention's backward pass

# Other

- ● Low-Precision Activation(Compared to BF16 activation)
  - ◆ Attention's backward pass



Here use attention output in backward pass
Don't use FP8 activation, because attention
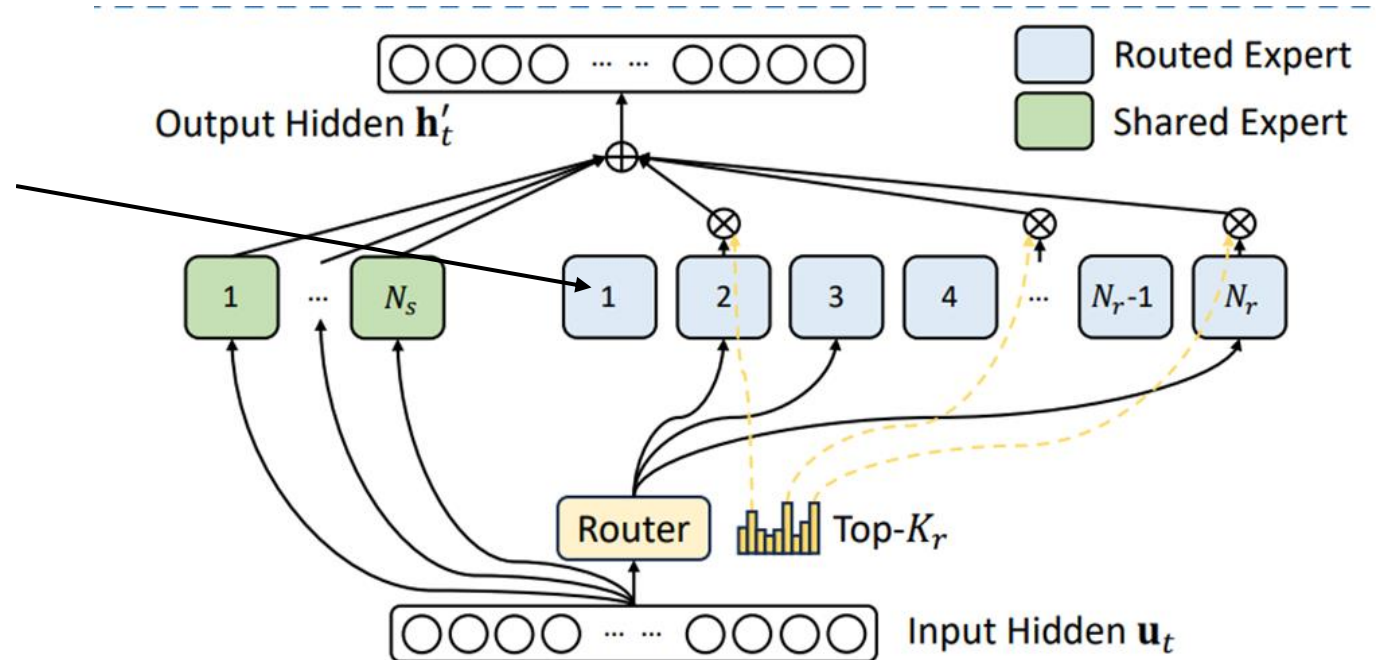need high precision, so deepseek use E9M6

# Other

- Low-Precision Activation(Compared to BF16 activation)
  - Input of SwiGLU (MoE)

Expert:

```
self.w2(F.silu(self.w1(x)) * self.w3(x))
```
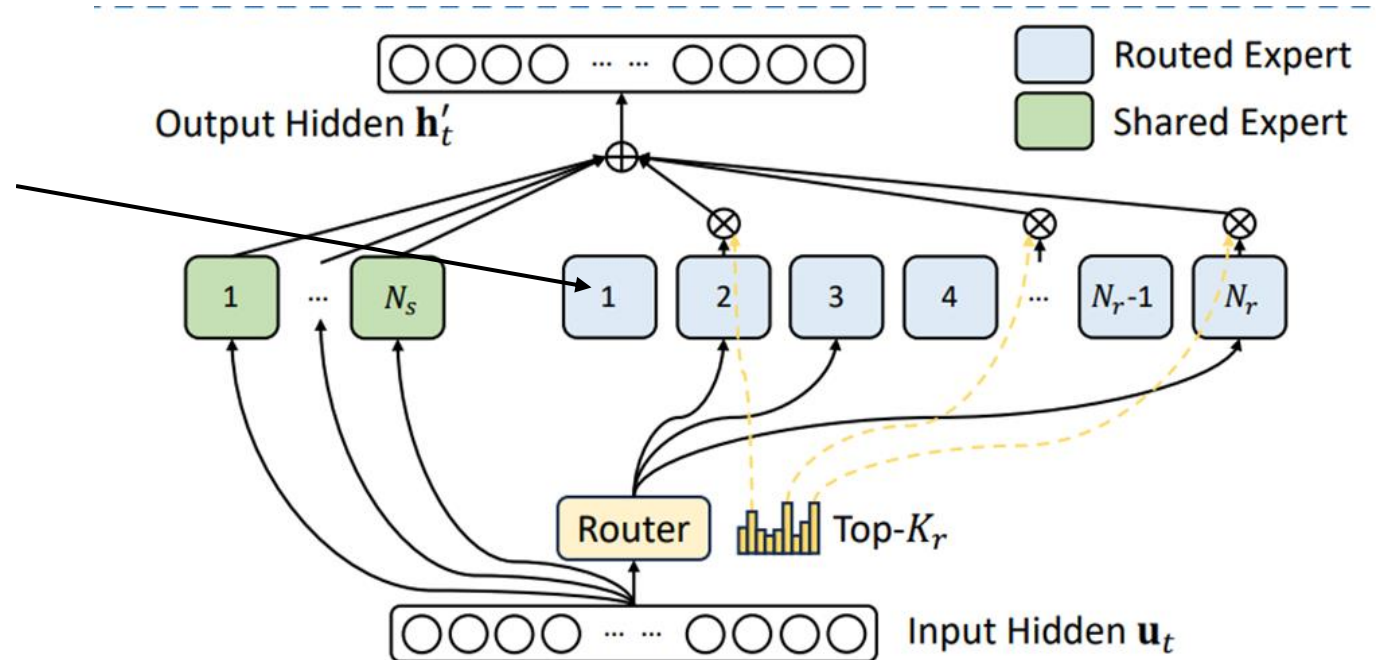
# Other

- Low-Precision Activation(Compared to BF16 activation)
  - ◆ Input of SwiGLU (MoE)

Expert:

```
self.w2(F.silu(self.w1(x)) * self.w3(x))
```

SwiGLU(x)



Output Hidden $\mathbf{h}'_t$

Input Hidden $\mathbf{u}_t$

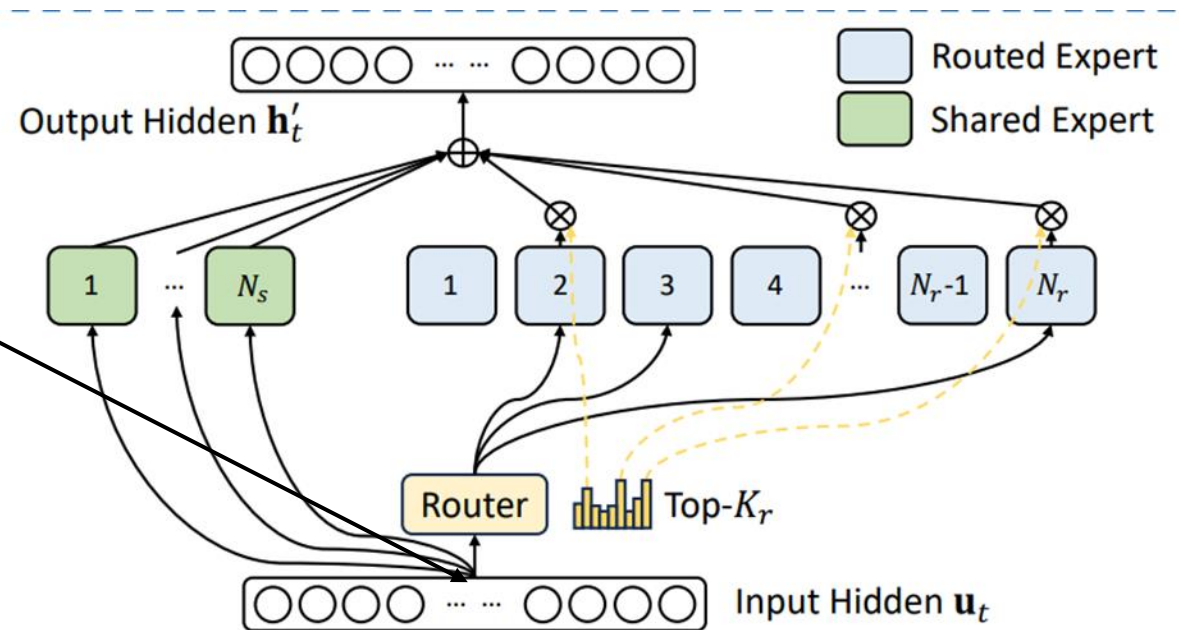Router  Top-$K_r$

Routed Expert
Shared Expert

# Other

- Low-Precision Activation(Compared to BF16 activation)
  - Input of SwiGLU (MoE)
    - cache the inputs of the SwiGLU operator and recompute its output in the backward (Saving its output incurs significant memory overhead)
    - striking a balance between memory efficiency and computational accuracy

# Other

- Low-Precision Communication
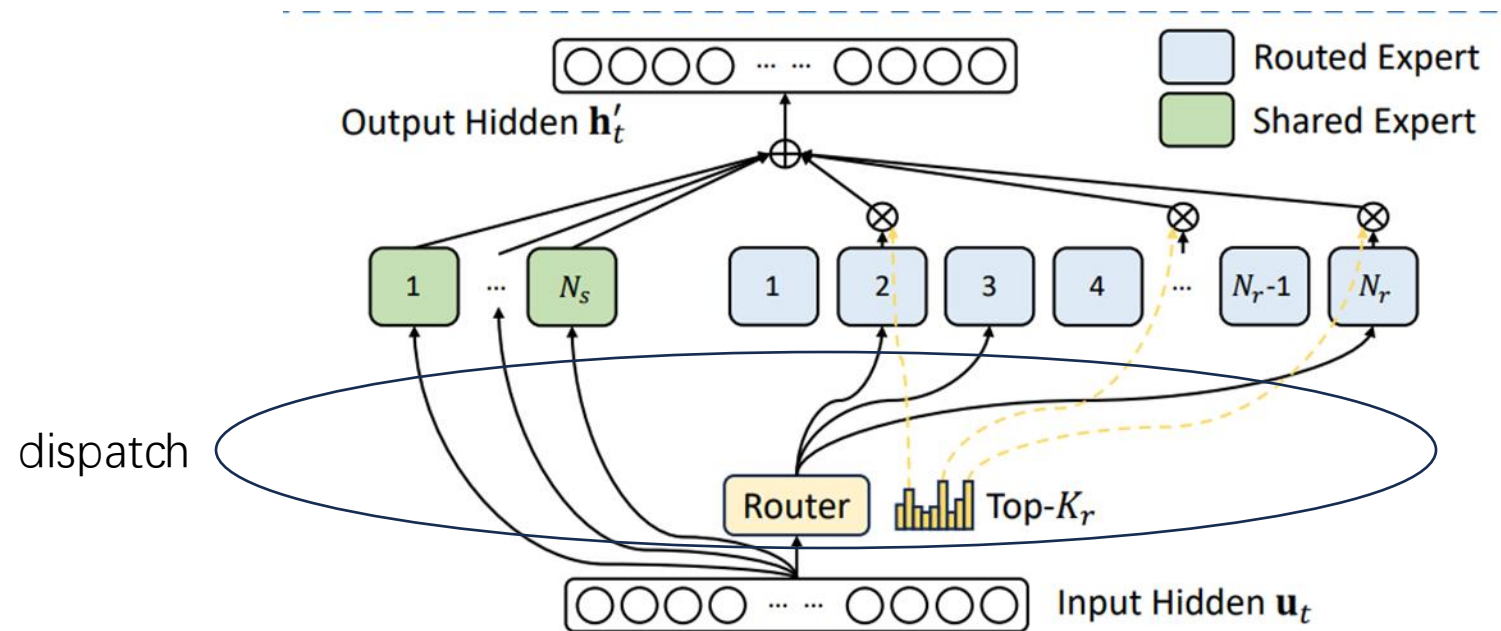  - ◆ Scale input of expert to FP8, then dispatch, which decrease communication overhead
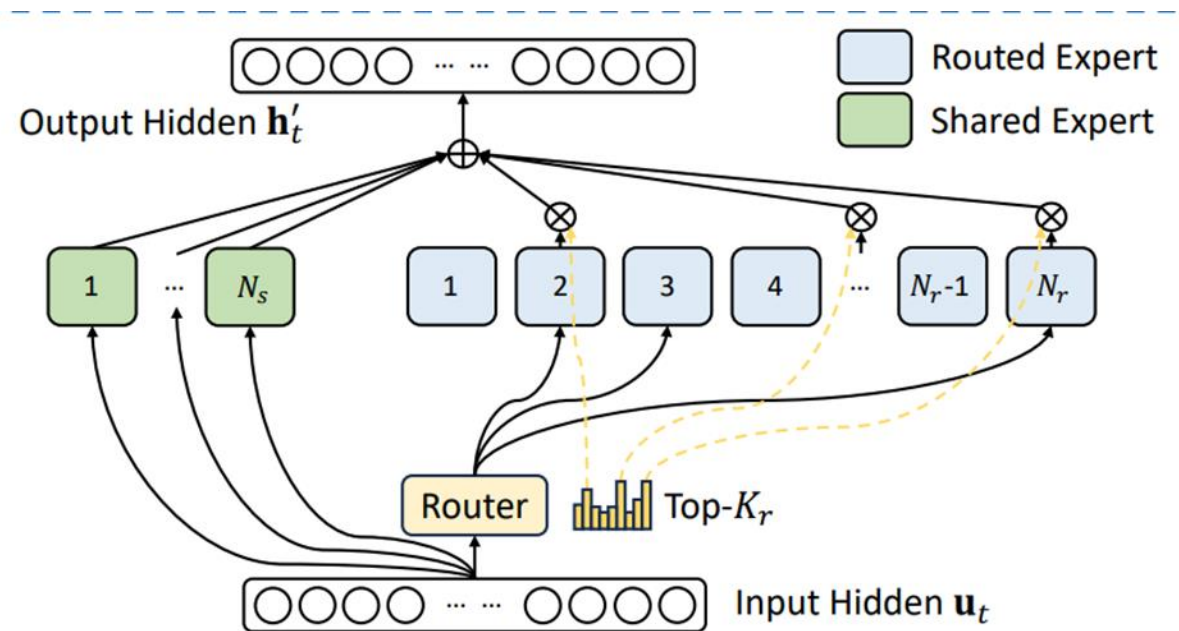
Scale BF16 to FP8

# Other

- Low-Precision Communication
  - Scale input of expert to FP8, then dispatch, which decrease communication overhead

# Other

- Low-Precision Communication
  - ◆ For combine component, retain them in BF16 to preserve training precision

# Other

- **Low-Precision Communication**
  - For combine component, retain them in BF16 to preserve training precision