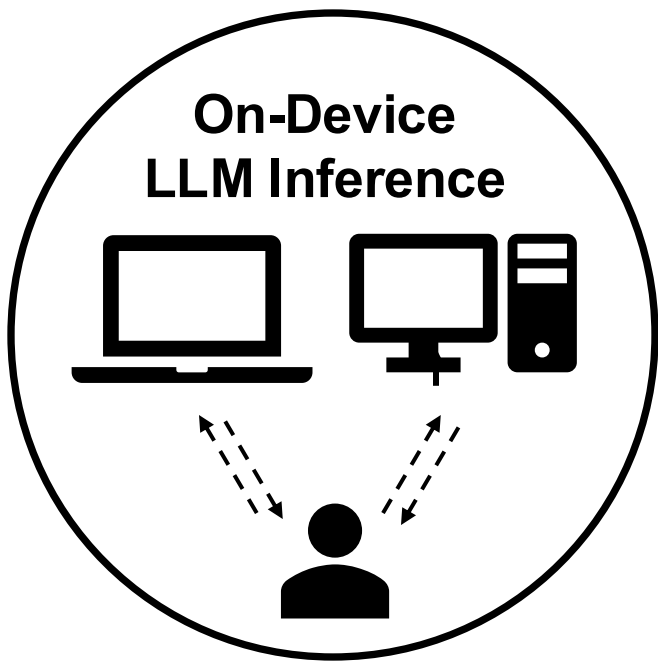# DecDEC: A Systems Approach to Advancing Low-Bit LLM Quantization

**Yeonhong Park\*, Jake Hyun\*, Hojoon Kim, Jae W. Lee**
**Seoul National University**

\* These authors equally contributed to this work

# Preliminary: On-Device LLM Inference
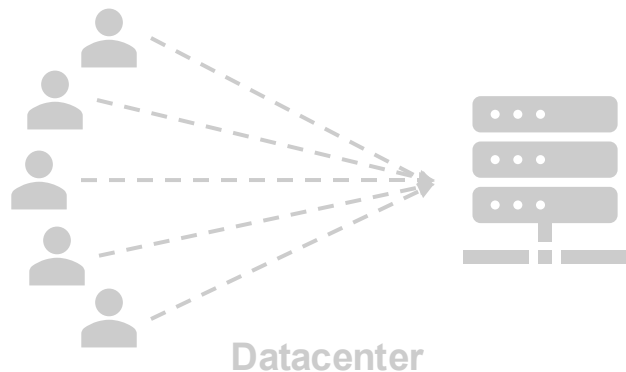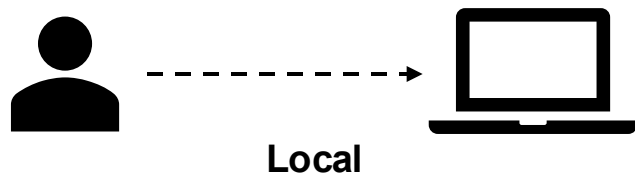


**On-Device LLM Inference**
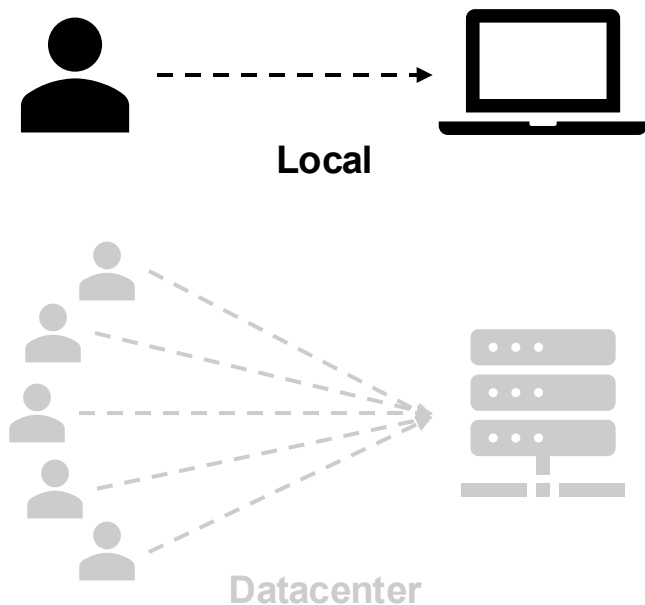
Privacy

Availability

Low Latency

# Preliminary: On-Device LLM Inference

Single-Query Inference
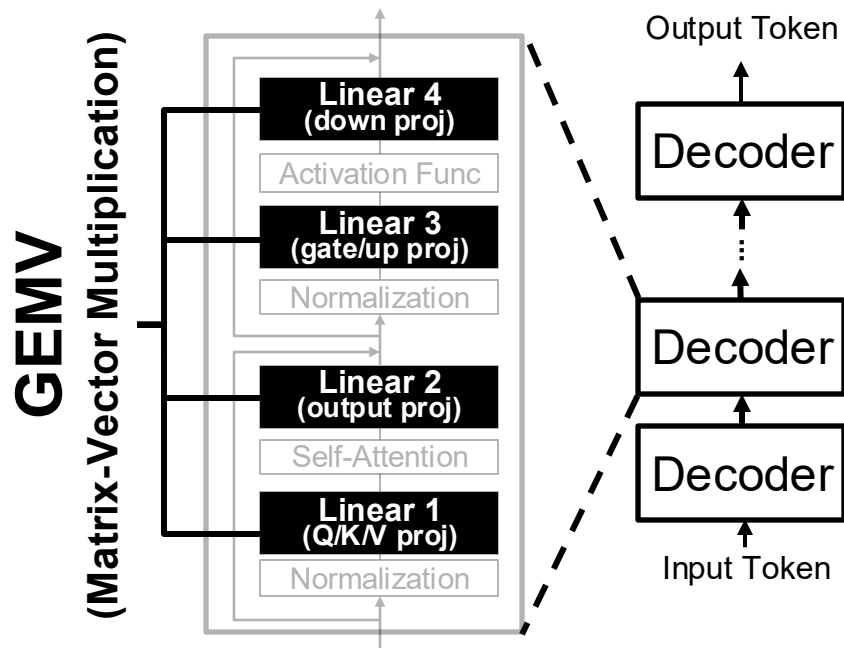


Local

Datacenter

# Preliminary: On-Device LLM Inference

## Single-Query Inference

**Local**

**Datacenter**

## Memory-Bound

**GEMV (Matrix-Vector Multiplication)**

- Linear 4 (down proj)
- Activation Func
- Linear 3 (gate/up proj)
- Normalization
- Linear 2 (output proj)
- Self-Attention
- Linear 1 (Q/K/V proj)
- Normalization

Output Token

Decoder

Decoder

Decoder

Input Token

# Preliminary: Weight-Only Quantization

**Weight**
**(INT4)**   **Activation**
**(FP16)**

Dequant

FP16 ALU

Output

**Weight-Only
Quantization**

* Save only memory
  bandwidth

* No loss of activation
  precision

* Good for single-batch
  inference (on-device)

Weight
(INT4)   Activation
(INT4)

INT4 ALU

Output

**Weight-Activation
Quantization**

* Save both compute
  and memory bandwidth

* Good for large-batch
  inference (data-center)

# Upsides and Downsides of Quantization



**GPU Memory Usage**

**Latency**

**Accuracy**

# Research Goal: (Almost) Free Accuracy Recovery



**GPU Memory Usage**

FP16   INT3/4   **Goal**

NO INCREASE

**Latency**

FP16   INT3/4   **Goal**

NO INCREASE

**Accuracy**

FP16   INT3/4   **Goal**

Accuray Recovery

# Research Goal: (Almost) Free Accuracy Recovery

# Q: What's the Cost?

GPU Memory Usage

Latency

Accuracy

# Research Goal: (Almost) Free Accuracy Recovery

Accuray Recovery

# Q: What's the Cost?
## A: CPU Memory

FP16   INT3/4   Goal
**GPU Memory Usage**

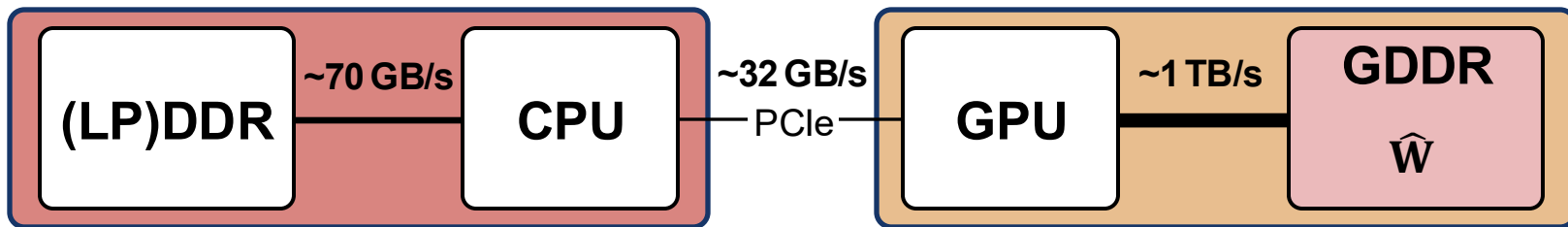FP16   INT3/4   Goal
**Latency**

FP16   INT3/4   Goal
**Accuracy**

# Suggestion: CPU-Augmented Quantized LLM Inference

# Suggestion: CPU-Augmented Quantized LLM Inference

# Suggestion: CPU-Augmented Quantized LLM Inference



| | | | |
|------|------|------|------|
| -0.2 | 0.2 | 0 | 0.1 |
| -0.2 | 0.7 | 0 | 0.3 |
| -0.3 | -0.1 | 0.2 | -0.5 |
| 0 | -0.3 | 0.4 | 0.1 |

$R$

| | | | |
|------|------|------|------|
| -3.2 | 0.2 | -1.0 | 3.1 |
| 2.8 | 3.7 | -2.0 | 1.3 |
| 1.7 | -1.1 | 1.2 | 0.5 |
| -4.0 | -2.7 | 0.4 | 2.1 |

$W$

| | | | |
|------|------|------|------|
| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

$\widehat{W}$

$$R = W - \widehat{W}$$

# Suggestion: CPU-Augmented Quantized LLM Inference



| (LP)DDR **R** | ~70 GB/s | CPU | ~32 GB/s PCIe | GPU | ~1 TB/s | GDDR $\widehat{\mathbf{W}}$ |

Fetch **R** for every GEMV

$$(\widehat{\mathbf{W}} + \mathbf{R}) * \mathbf{x} = \mathbf{W} * \mathbf{x}$$

# Suggestion: CPU-Augmented Quantized LLM Inference

(LP)DDR ~70 GB/s | CPU | ~32 GB/s PCIe | GPU | ~1 TB/s | GDDR

$R$ | | | | | $\widehat{W}$

**Can fully recover full-precision weight, but may incur prohibitive slowdown**

Fetch $R$ for every GEMV

$$(\widehat{W} + R) * x = W * x$$

# Suggestion: CPU-Augmented Quantized LLM Inference

**(LP)DDR** — ~70 GB/s — **CPU** — ~32 GB/s / PCIe — **GPU** — ~1 TB/s — **GDDR**

$R$                            $\widehat{W}$

Selectively fetch $\mathbf{R \odot M}$ for every GEMV

$$(\widehat{\mathbf{W}} + \mathbf{R \odot M}) * \mathbf{x} = \mathbf{W} * \mathbf{x}$$

$\mathbf{M}$: binary mask

# Key Research Question

How to determine a subset of residuals to fetch ($M$)?

# Key Research Question

**How to determine a subset of residuals to fetch ($M$)?**

A good $M$ should:

\* Select the most impactful portions

For RTX-4090, PCIe BW : GPU BW = 32 GB/s : 1 TB/s

$\cong$ **1: 30**

# Key Research Question

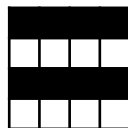**How to determine a subset of residuals to fetch ($\mathbf{M}$)?**

A good $\mathbf{M}$ should:

\* Select the most impactful portions

    For RTX-4090, PCIe BW : GPU BW = 32 GB/s : 1 TB/s

                          $\cong$ **1: 30**

\* Be structured for efficient processing

# Opportunity: Not All Channels are Equally Important



$$X \quad * \quad \widehat{W} \quad \leftarrow \quad W$$

| 1.1 | -0.2 | -3.0 | 0.7 |

X

| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

$\widehat{W}$

| -3.2 | 0.2 | -1.0 | 3.1 |
| 2.8 | 3.7 | -2.0 | 1.3 |
| 1.7 | -1.1 | 1.2 | 0.5 |
| -4.0 | -2.7 | 0.4 | 2.1 |

W

# Opportunity: Not All Channels are Equally Important



$$X \quad * \quad \widehat{W} \quad \Longleftarrow \quad W$$

| 1.1 | -0.2 | -3.0 | 0.7 |

| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

| -3.2 | 0.2 | -1.0 | 3.1 |
| 2.8 | 3.7 | -2.0 | 1.3 |
| 1.7 | -1.1 | 1.2 | 0.5 |
| -4.0 | -2.7 | 0.4 | 2.1 |

# Opportunity: Not All Channels are Equally Important

| 1.1 | -0.2 | -3.0 | 0.7 |

X

*

| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

$\widehat{W}$

←

| -3.2 | 0.2 | -1.0 | 3.1 |
| 2.8 | 3.7 | -2.0 | 1.3 |
| 1.7 | -1.1 | 1.2 | 0.5 |
| -4.0 | -2.7 | 0.4 | 2.1 |

W

# Opportunity: Not All Channels are Equally Important

| 1.1 | -0.2 | -3.0 | 0.7 |

*

| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

←

| -3.2 | 0.2 | -1.0 | 3.1 |
| 2.8 | 3.7 | -2.0 | 1.3 |
| 1.7 | -1.1 | 1.2 | 0.5 |
| -4.0 | -2.7 | 0.4 | 2.1 |

X          $\widehat{W}$          W

# Opportunity: Not All Channels are Equally Important

| | | | |
|---|---|---|---|
| 1.1 | -0.2 | -3.0 | 0.7 |

\*

| | | | |
|---|---|---|---|
| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

←

| | | | |
|---|---|---|---|
| -3.2 | 0.2 | -1.0 | 3.1 |
| 2.8 | 3.7 | -2.0 | 1.3 |
| 1.7 | -1.1 | 1.2 | 0.5 |
| -4.0 | -2.7 | 0.4 | 2.1 |

$X$　　　　　$\widehat{W}$　　　　　$W$

# Opportunity: Not All Channels are Equally Important

Activation Outlier

| 1.1 | -0.2 | -3.0 | 0.7 |

$X$

*

| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

$\widehat{W}$

Salient Channel

| -3.2 | 0.2 | -1.0 | 3.1 |
| 2.8 | 3.7 | -2.0 | 1.3 |
| 1.7 | -1.1 | 1.2 | 0.5 |
| -4.0 | -2.7 | 0.4 | 2.1 |

$W$

# Opportunity: Not All Channels are Equally Important

Compensate for errors in **salient channels**

| 1.1 | -0.2 | -3.0 | 0.7 |

\*

| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

\+

| -0.3 | -0.1 | 0.2 | 0.5 |

X                 $\widehat{W}$                 $\mathbf{R} \odot \mathbf{M}$

# Challenge: Dynamic Nature of Activation Outliers

Decoding Step $i$



X          *          $\widehat{W}$

# Challenge: Dynamic Nature of Activation Outliers

Decoding Step $i + 1$

| 4.1 | -1.2 | 0.8 | 0.3 |
|---|---|---|---|

\*

| -3 | 0 | -1 | 3 |
|---|---|---|---|
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

X                    $\widehat{W}$

# Challenge: Dynamic Nature of Activation Outliers

Decoding Step $i + 2$

| | | | |
|---|---|---|---|
| -3 | 0 | -1 | 3 |
| 3 | 3 | -2 | 1 |
| 2 | -1 | 1 | 1 |
| -4 | -3 | 0 | 2 |

| -0.1 | 0.8 | 1.1 | -3.7 |
|---|---|---|---|

\* 

X        $\widehat{W}$

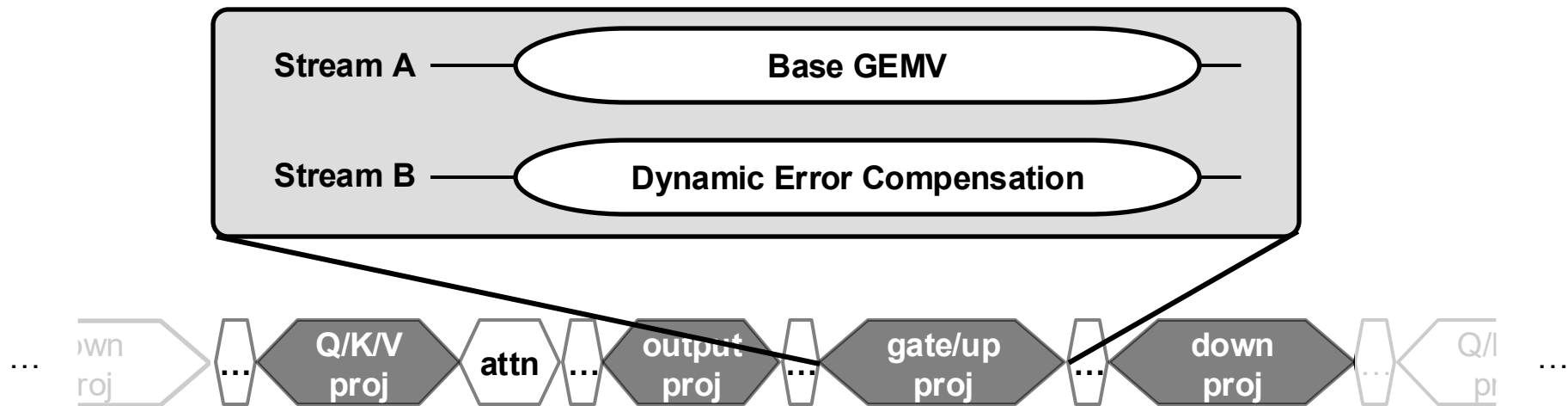# Challenge: Dynamic Nature of Activation Outliers



*Distribution of activation outliers (<u>top 5%</u>) across 100 decoding steps*
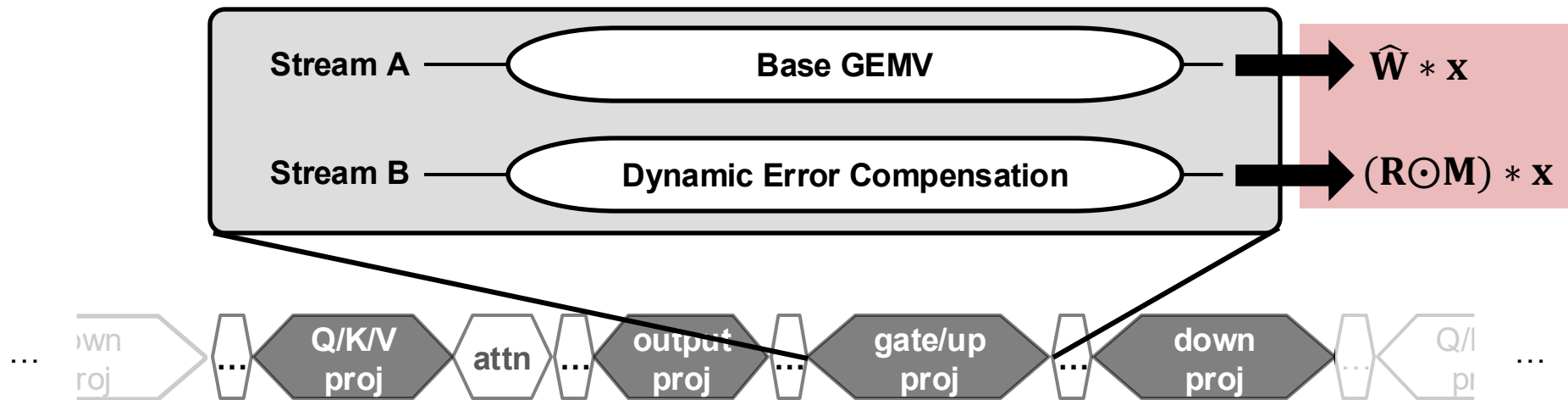
# DecDEC: Decoding with Dynamic Error Compensation

Inference system for quantized LLMs that performs **Dec**oding with **D**ynamic **E**rror **C**ompenstation



Augment each linear layer with **dynamic error compensataion (DEC)**
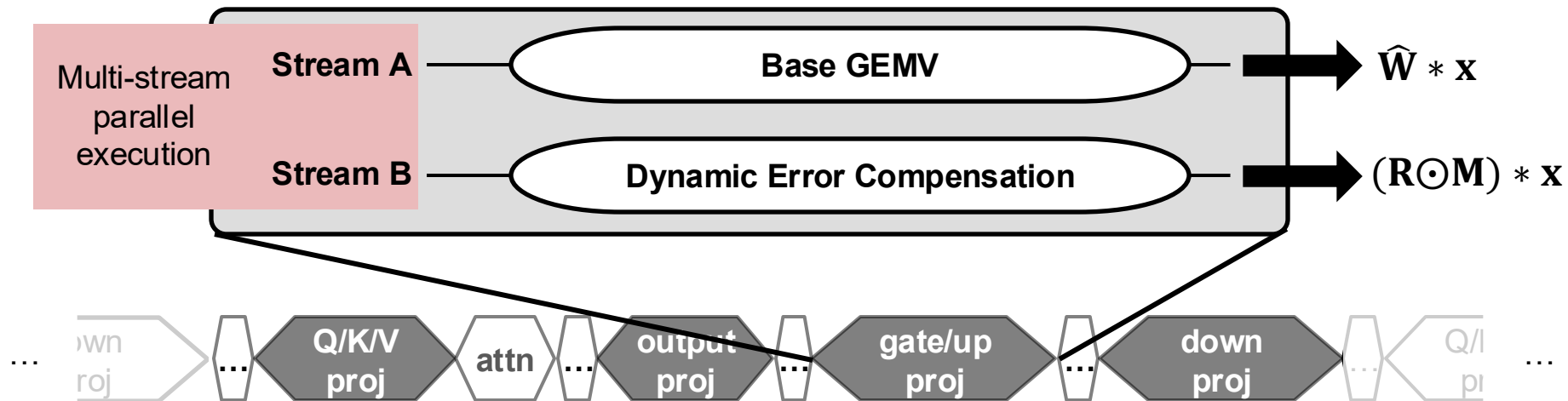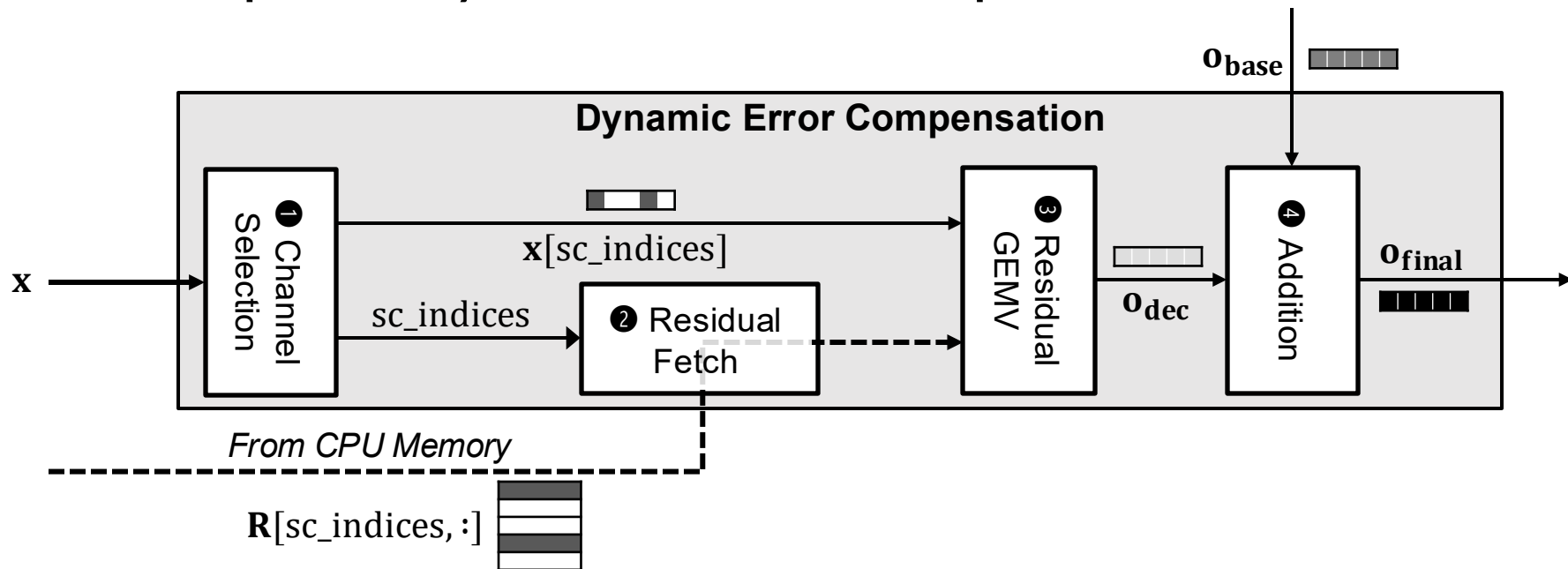
# DecDEC: Decoding with Dynamic Error Compensation

Inference system for quantized LLMs that performs **Dec**oding with **D**ynamic **E**rror **C**ompenstation

Stream A — Base GEMV — $\widehat{\mathbf{W}} * \mathbf{x}$

Stream B — Dynamic Error Compensation — $(\mathbf{R} \odot \mathbf{M}) * \mathbf{x}$

... down proj ... Q/K/V proj attn ... output proj ... gate/up proj ... down proj ... Q/K proj ...

Augment each linear layer with **dynamic error compensataion (DEC)**

# DecDEC: Decoding with Dynamic Error Compensation

Inference system for quantized LLMs that performs **Dec**oding with **D**ynamic **E**rror **C**ompenstation
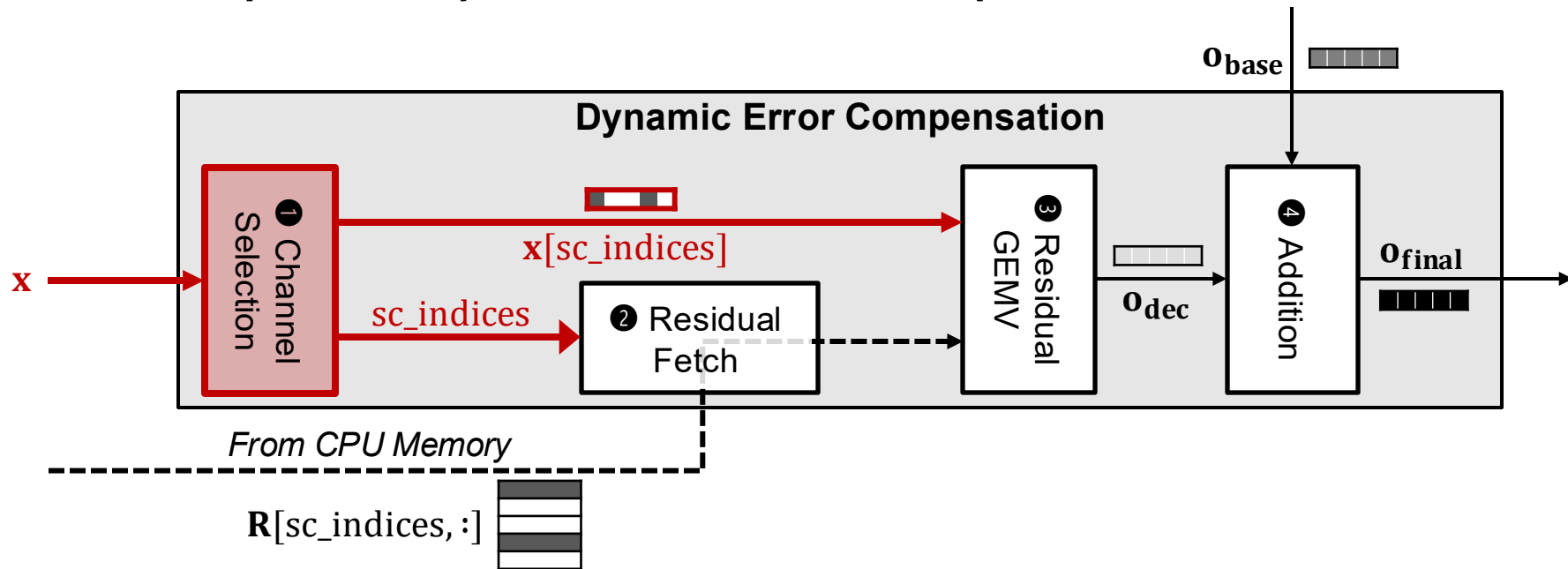


Augment each linear layer with **dynamic error compensataion (DEC)**
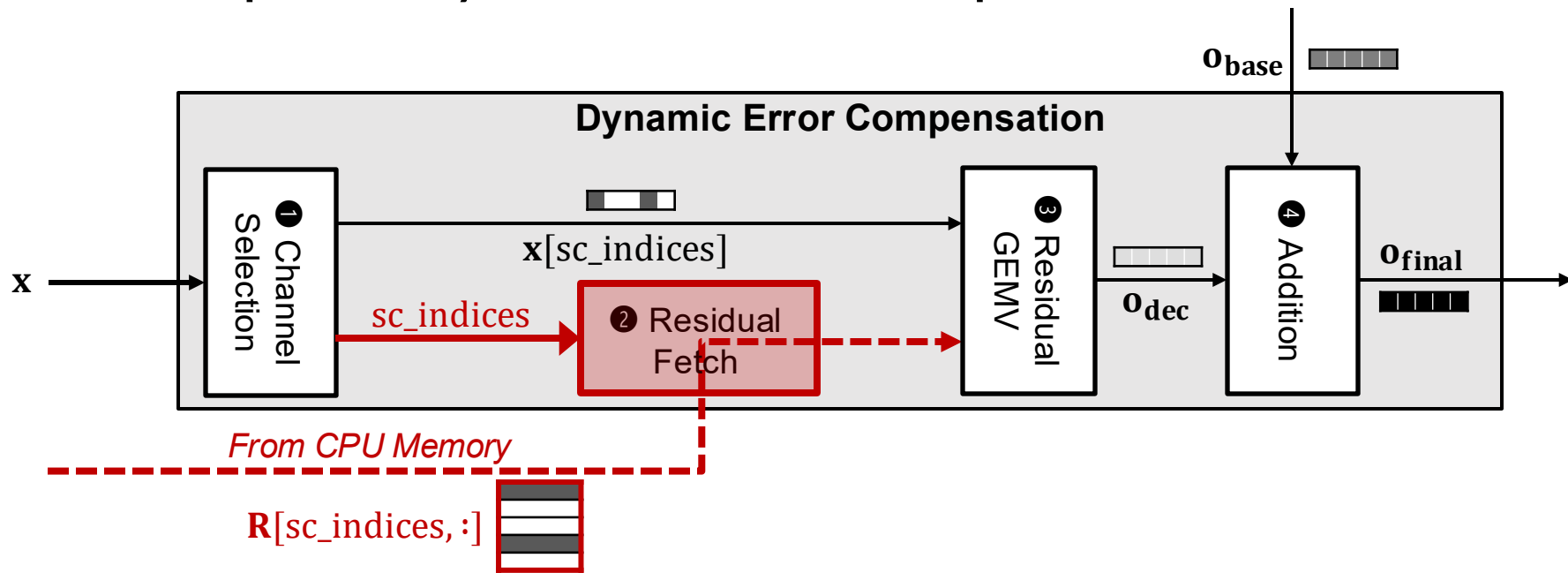
# Four Steps of Dynamic Error Compensation
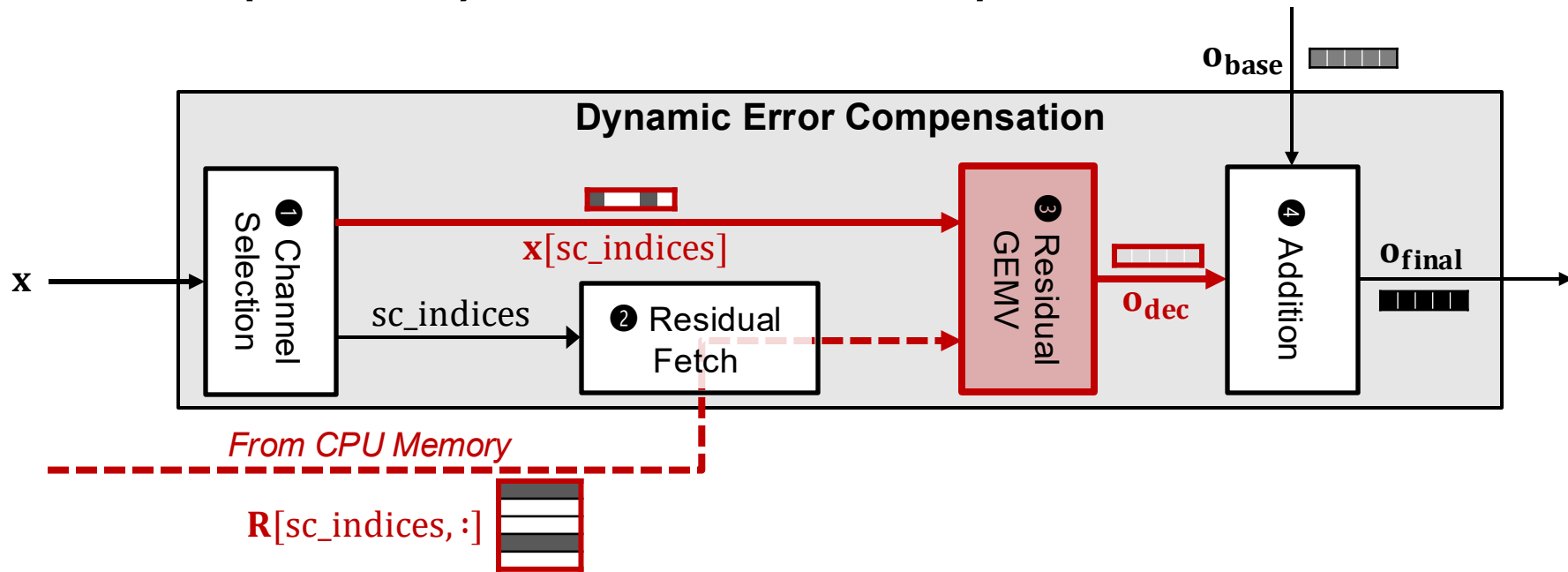
# Four Steps of Dynamic Error Compensation



① Perform Top-K on the input activation vector ($\mathbf{x}$)
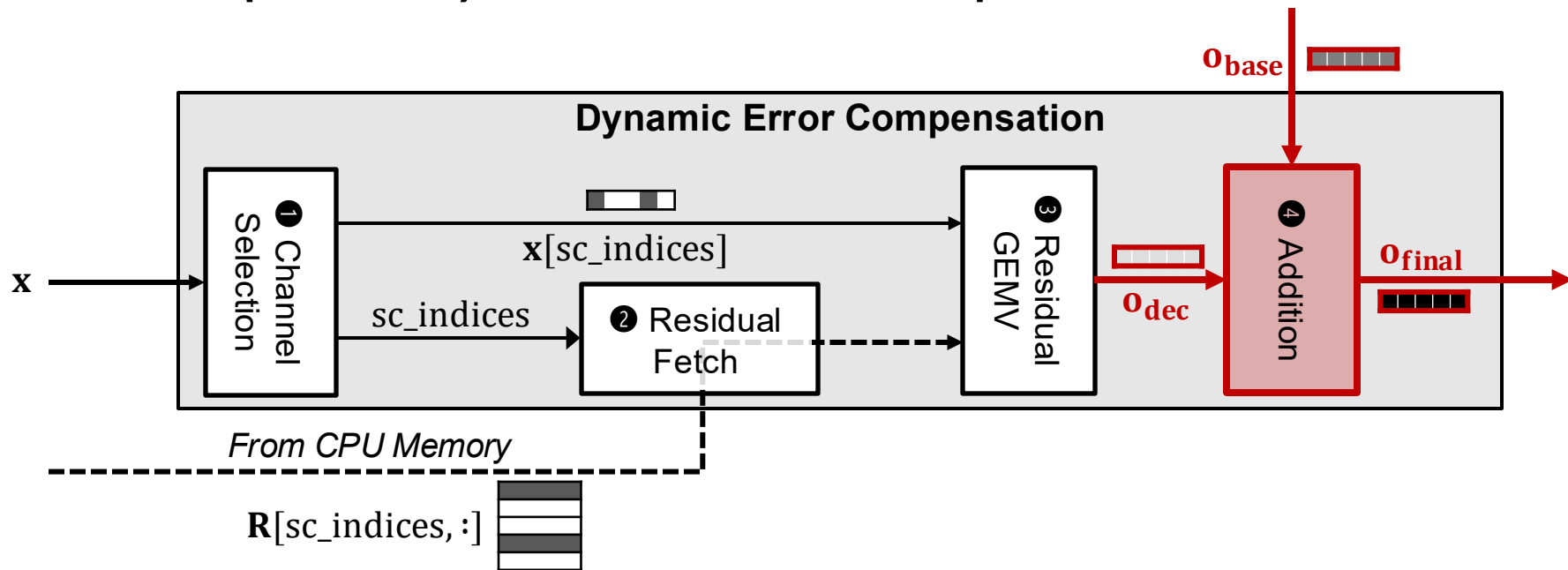
# Four Steps of Dynamic Error Compensation



② Fetch residuals for the selected channels from CPU ($\mathbf{R}[\text{sc\_indices}, :]$)

# Four Steps of Dynamic Error Compensation



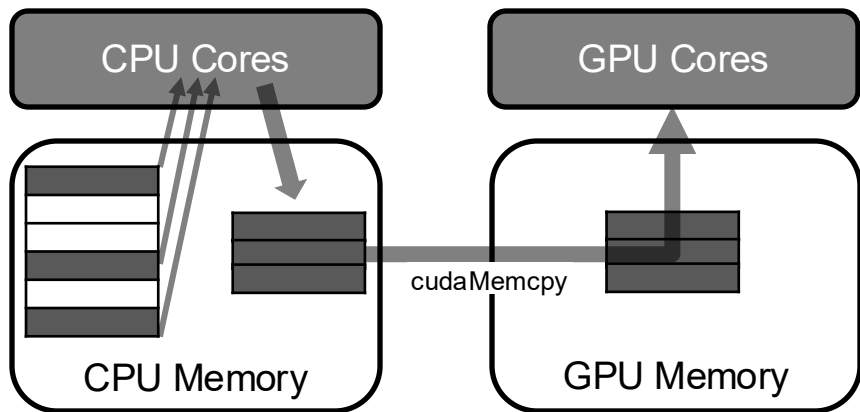③ Multiply the input activation by the selected residuals

# Four Steps of Dynamic Error Compensation



④ Add base GEMV result ($\mathbf{o_{base}}$) & residual GEMV result ($\mathbf{o_{dec}}$)

# Key Implementation Point: Zero-Copy Residual Fetch

**DMA-based Approach**

CPU Cores

GPU Cores

cudaMemcpy

CPU Memory

GPU Memory

👍 High BW util for bulk transfer

👎 Long latency

**Zero-Copy-based Approach**

CPU Cores

GPU Cores

CPU Memory

*No GPU memory BW consumption*

GPU Memory

👍 Low latency
👍 Fine-grained, cacheline-sized data access

👎 GPU core consumption

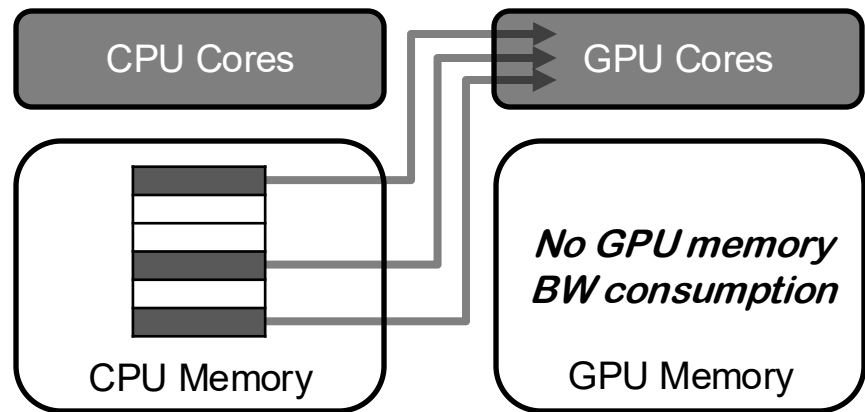# Key Implementation Point: Zero-Copy Residual Fetch

**DMA-based Approach**

## Zero-Copy is the answer!

* Need low-latency, fine-grained data access
* GPU core consumption is OK,
  while BW contention is undesirable
  (Base GEMV is memory-bound)

**Zero-Copy-based Approach**

| CPU Cores | GPU Cores |
| --- | --- |

*No GPU memory BW consumption*

CPU Memory          GPU Memory

👍 Low latency
👍 Fine-grained, cacheline-sized data access

👎 GPU core consumption

Y. Park et al., "DecDEC: A Systems Approach to Advancing Low-Bit LLM Quantization", OSDI'25

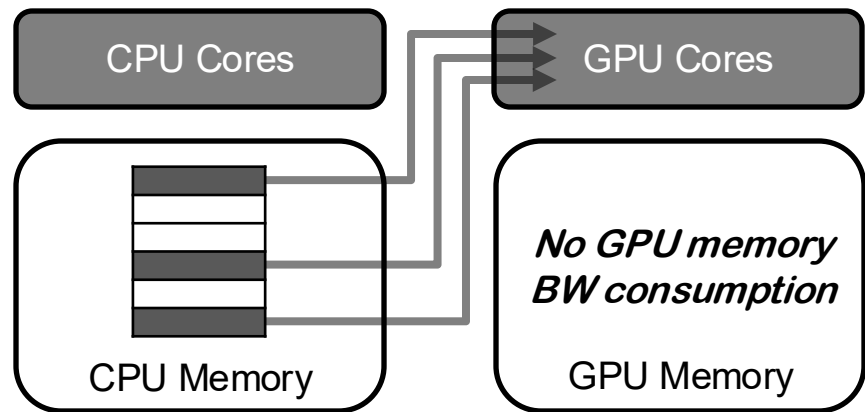# Key Implementation Point: Zero-Copy Residual Fetch

**DMA-based Approach**

**Zero-Copy is the answer!**

\* Need low-latency, fine-grained data access
\* GPU core consumption is OK,
while BW contention is undesirable
(Base GEMV is memory-bound)

**Zero-Copy-based Approach**

CPU Cores → GPU Cores

*No GPU memory BW consumption*

CPU Memory

GPU Memory

👍 Low latency
👍 Fine-grained, cacheline-sized data access

👎 GPU core consumption

Y. Park et al., "DecDEC: A Systems Approach to Advancing Low-Bit LLM Quantization", OSDI'25  40

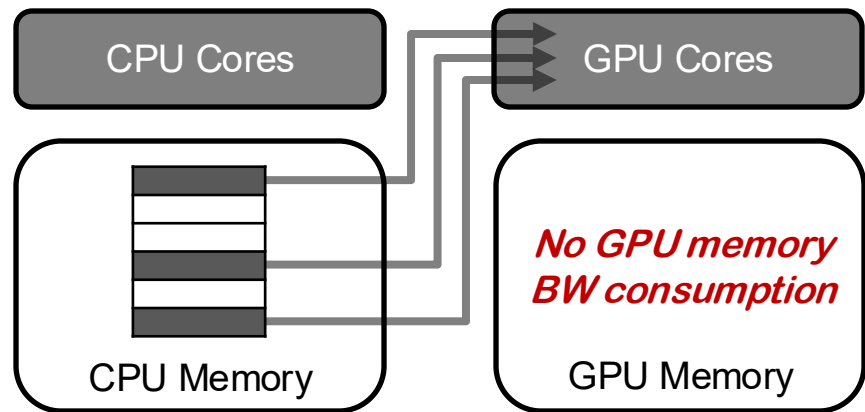# Key Implementation Point: Zero-Copy Residual Fetch

**DMA-based Approach**

## Zero-Copy is the answer!

\* Need low-latency, fine-grained data access

\* GPU core consumption is OK,
while BW contention is undesirable
(Base GEMV is memory-bound)

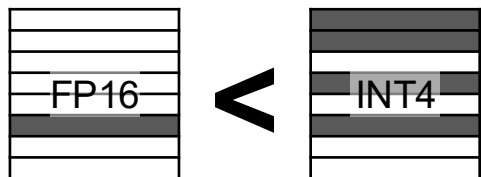**Zero-Copy-based Approach**

CPU Cores → GPU Cores

CPU Memory

*No GPU memory BW consumption*

GPU Memory

👍 Low latency
👍 Fine-grained, cacheline-sized data access

👎 GPU core consumption

# Other Optimizations

- 4-bit residual quantization

FP16   <   INT4

Introduces approximation errors,
but enables <u>fetching more channels</u>

- GPU kernel optimizations

❶ Channel Selection → ❷ Residual Fetch → ❸ Residual GEMV → ❹ Addition ➡ **DEC**
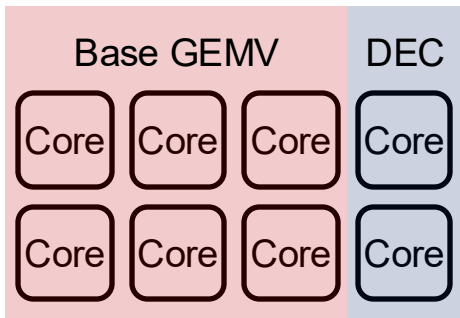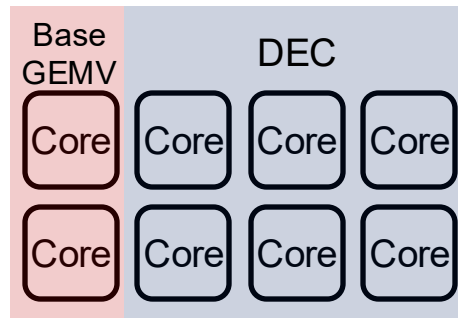
GPU-friendly approximiate
Top-K selection

Full kernel fusion

# DecDEC Parameter Tuner

- Two system parameters should be carefully tuned
  - ① $n_{tb}$: # of thread blocks (≈ # of GPU cores) to allocate to DEC
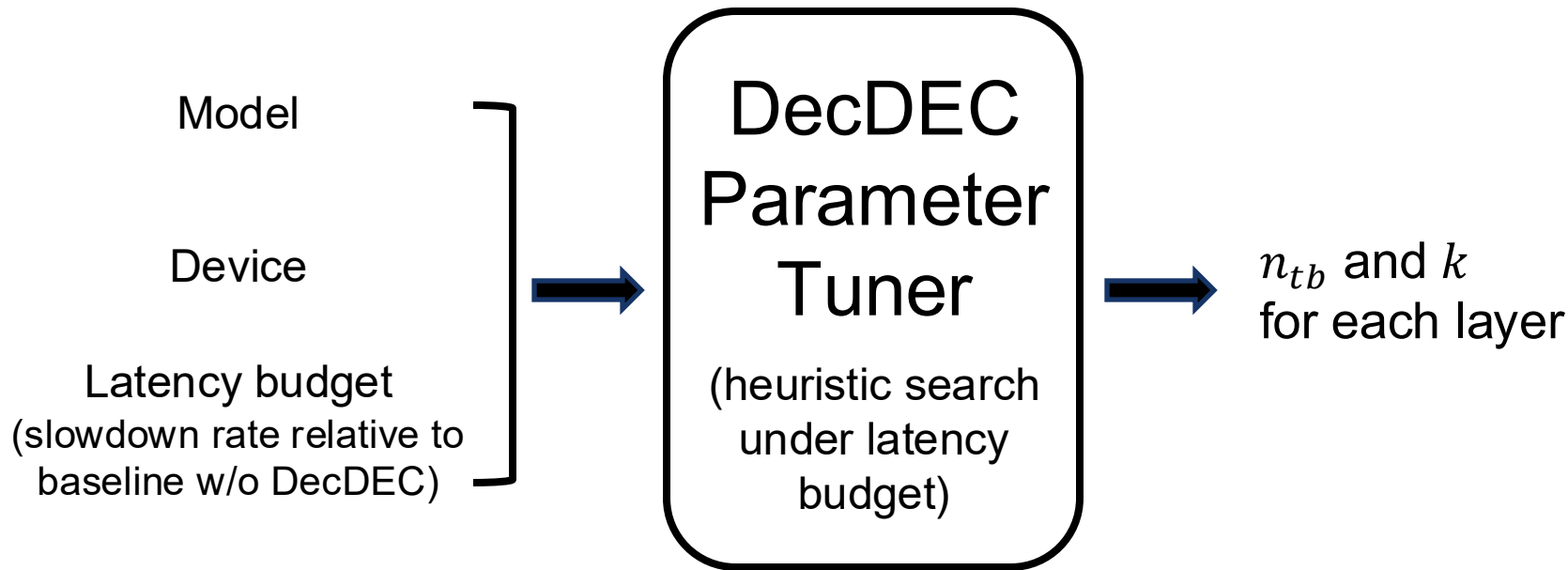


Too small $n_{tb}$ (i.e., 2) underutilizes PCIe BW
Too large $n_{tb}$ (i.e., 6) slows down the base GEMV

  - ② $k$: # channels to fetch
    → Larger is better, up to the point it incurs too high latency overhead

# DecDEC Parameter Tuner

Model

Device

Latency budget
(slowdown rate relative to
baseline w/o DecDEC)

→ **DecDEC Parameter Tuner**

(heuristic search under latency budget)

→ $n_{tb}$ and $k$ for each layer

# Evaluation

Y. Park et al., "DecDEC: A Systems Approach to Advancing Low-Bit LLM Quantization", OSDI'25

# Evaluation

- **Key Result #1:**
  DEC is effectively overlapped with base GEMV using appropriate $n_{tb}$ and $k$

- **Key Result #2:**
  DecDEC significantly improves quality with limited latency overhead

# Evaluation

- **Key Result #1:**
  DEC is effectively overlapped with base GEMV using appropriate $n_{tb}$ and $k$

- **Key Result #2:**
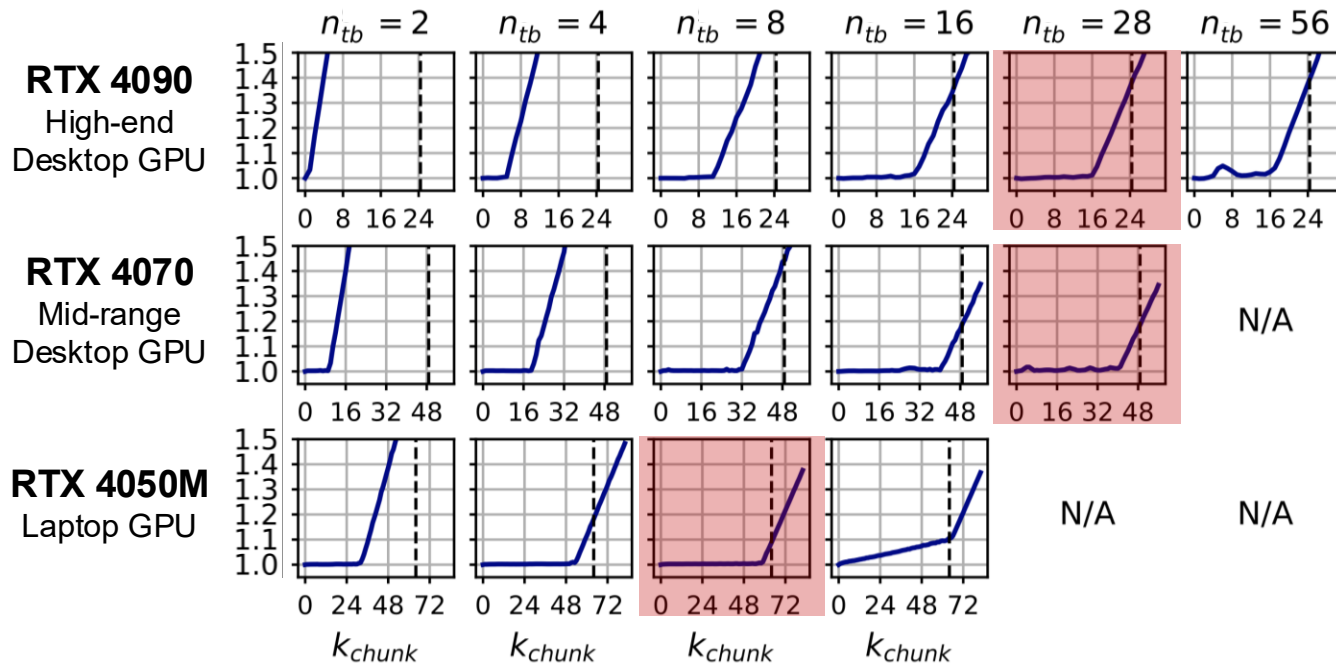  DecDEC significantly improves quality with limited latency overhead

# DecDEC Kernel Evaluation

3-bit uniform quantization / 4096 x 28672



With appropriate $n_{tb}$, DecDEC incurs almost no overhead
up to a certain value of $k$ ($k_{chunk}$: $k$ per 1024 channels)

# Evaluation

- **Key Result #1:**
  DEC is effectively overlapped with base GEMV using appropriate $n_{tb}$ and $k$

- **Key Result #2:**
  DecDEC significantly improves quality with limited latency overhead

# End-to-End Evaluation

Perplexity vs. Latency under 2.5%, 5%, 10%, and 20% target slowdown rate



| GPU Mem. | Latency | Quality |
|----------|---------|---------|
| 3-bit | ~3-bit (1.7% ↑) | Better than 3.5-bit |