# EXION:
# Exploiting Inter- and Intra-Iteration Output Sparsity for Diffusion Models

HPCA' 2025

Jaehoon Heo, Adiwena Putra, Jieon Yoon, Sungwoong Yune, Hangyeol Lee, Ji-Hoo
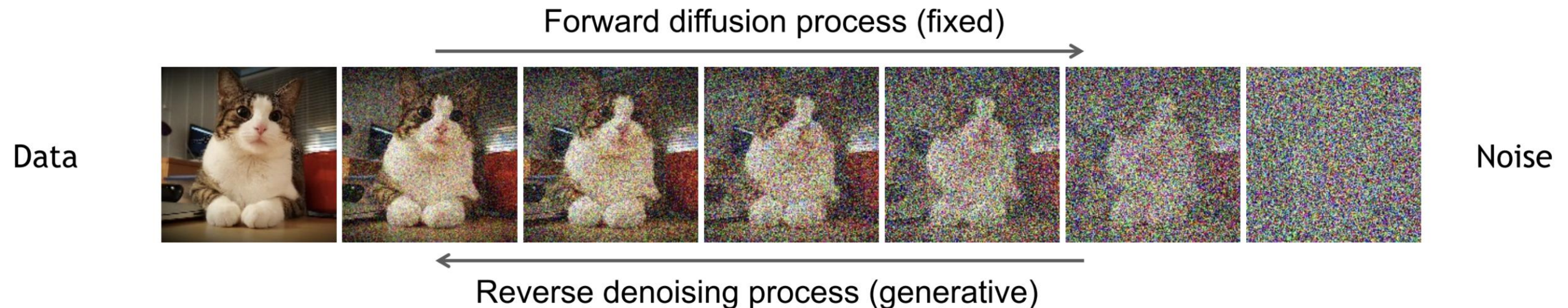
*KAIST Daejeon, Republic of Korea*

Yuge Cheng, 23/05/2025

# Preliminaries

## Denoising Diffusion Models: Learning to generate by denoising

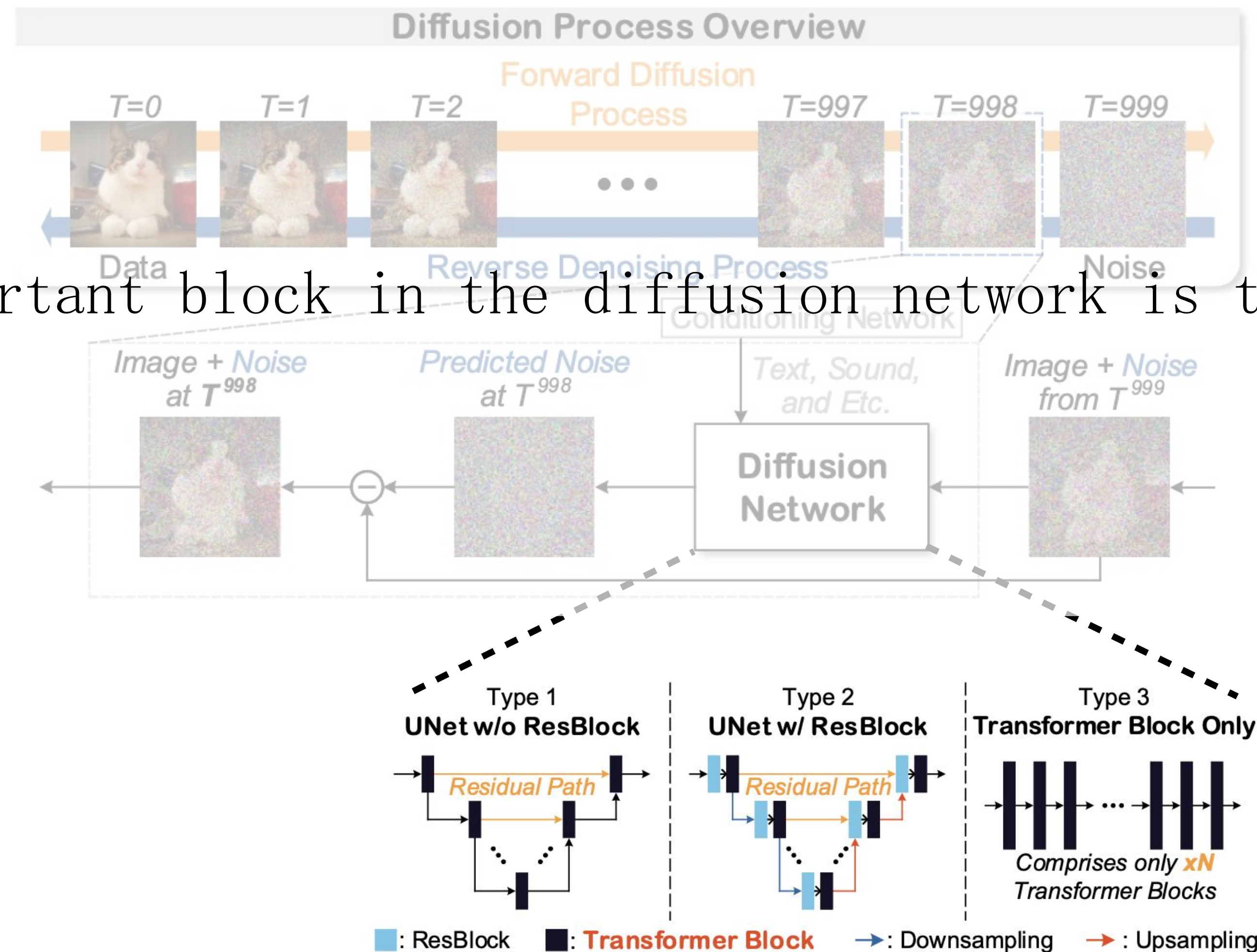Denoising diffusion models consist of two processes:

- Forward diffusion process that gradually adds noise to input

- Reverse denoising process that learns to generate data by denoising
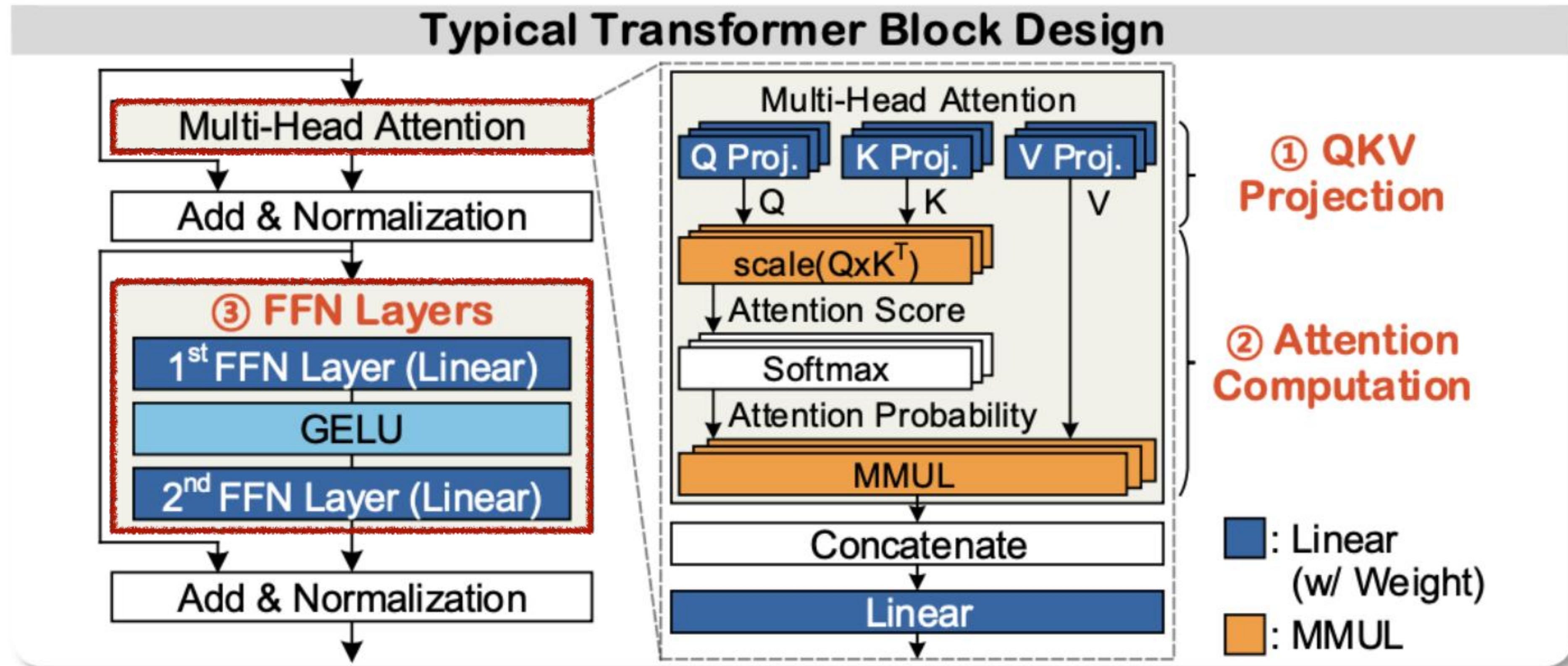
# Preliminaries

## Denoising Diffusion Models: Learning to generate by denoising



The most important block in the diffusion network is the transformer block !

# Preliminaries

Transformer Block Design



## Typical Transformer Block Design

Multi-Head Attention

Add & Normalization

③ FFN Layers
- 1st FFN Layer (Linear)
- GELU
- 2nd FFN Layer (Linear)

Add & Normalization

Multi-Head Attention

Q Proj.  K Proj.  V Proj.

Q  K  V

scale(QxK^T)

Attention Score

Softmax

Attention Probability

MMUL

Concatenate

Linear

① QKV Projection

② Attention Computation

: Linear (w/ Weight)

: MMUL

# Problems

Better outputs come with *higher* energy consumption and *longer* latency :(

- Each generation requires numerous denoising iterations

- Each iteration evaluates numerous Transformer Blocks

- Each transformer block entails a large number of operations

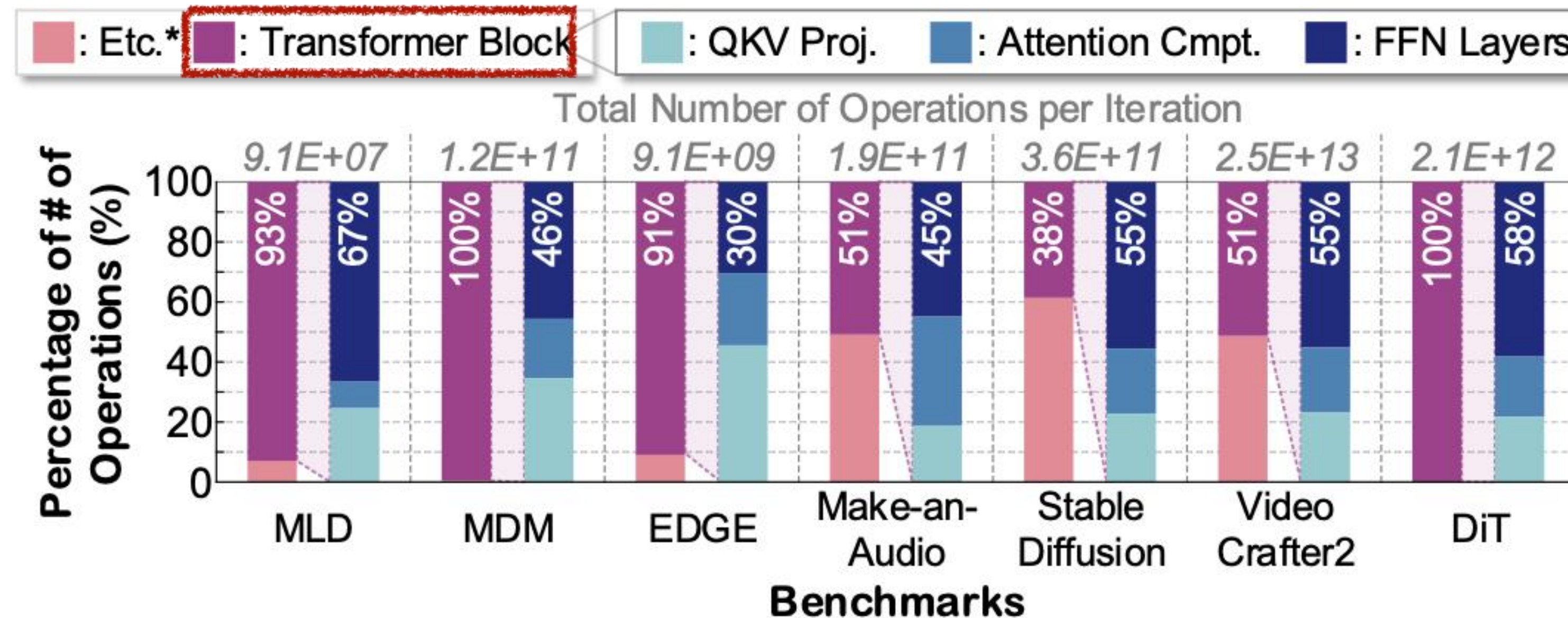| Metric | StyleGAN-XL | Stable Diffusion |
|---|---|---|
| Energy (Joules) | ~65.5 J | 1546.7 J |
| Latency (sec) | ~1.2 s | 11.8 s |

X23.6 more energy

X 9.8 slower

*Experiment results on NVIDIA's RTX 6000 Ada*

# Problems

Better outputs come with *higher* energy consumption and *longer* latency :(

Fig: Number of Operations Breakdown



- the <u>transformer block</u> accounts for the highest ratio

- the <u>FFN layers</u> are generally the most compute-intensive

# Overview

EXION: Exploiting Inter- and Intra-Iteration Output Sparsity for Diffusion Models

Key ideas: exploiting the unique <u>inter-</u> and <u>intra-iteration</u> <u>output sparsity</u>

# Overview

EXION: Exploiting Inter- and Intra-Iteration Output Sparsity for Diffusion Models

Key ideas

1. Software optimizations

   - FFN-Reuse across different iterations

   - Sparse attention computation via eager prediction

2. ConMerge: data compaction mechanism
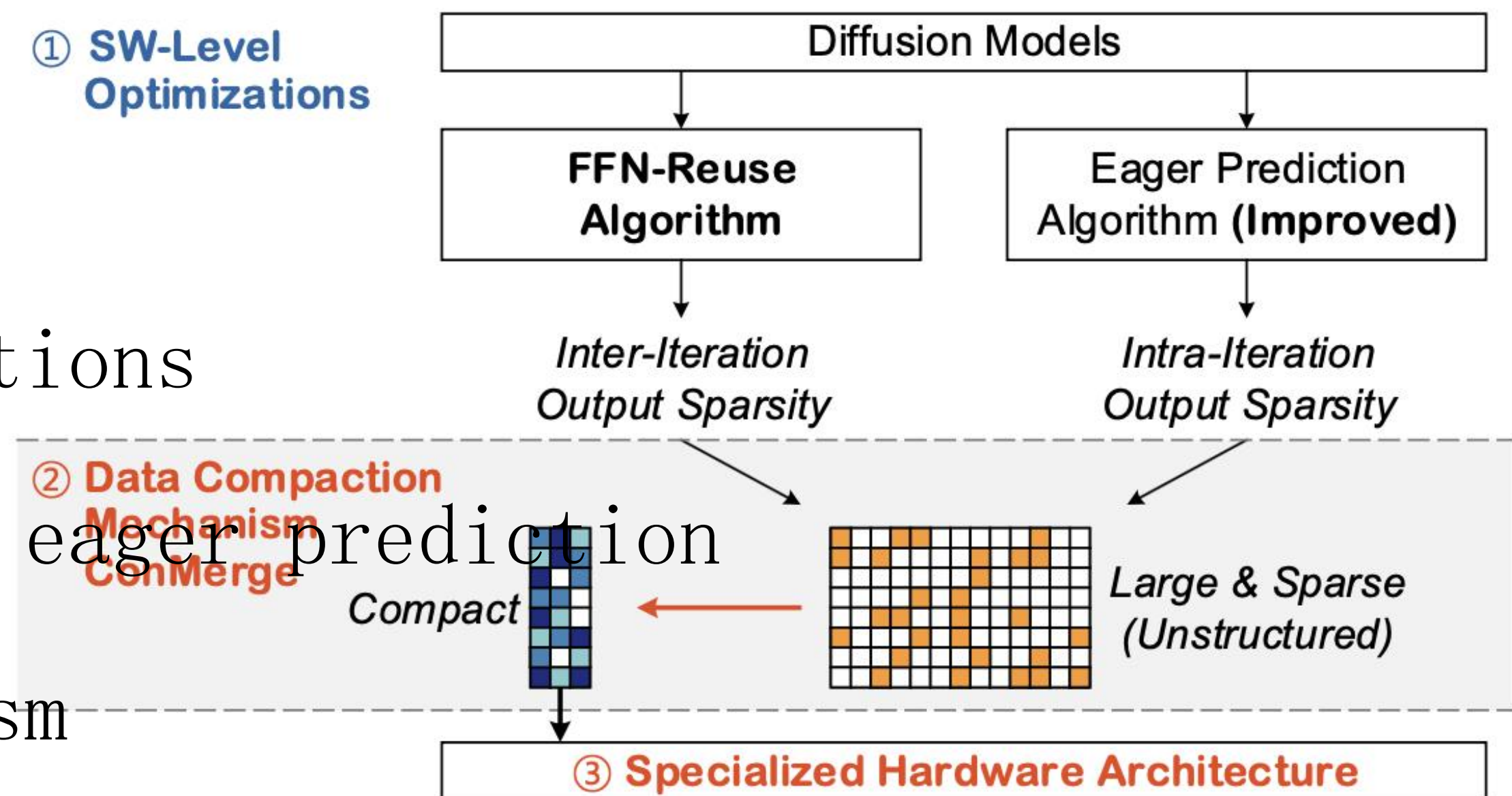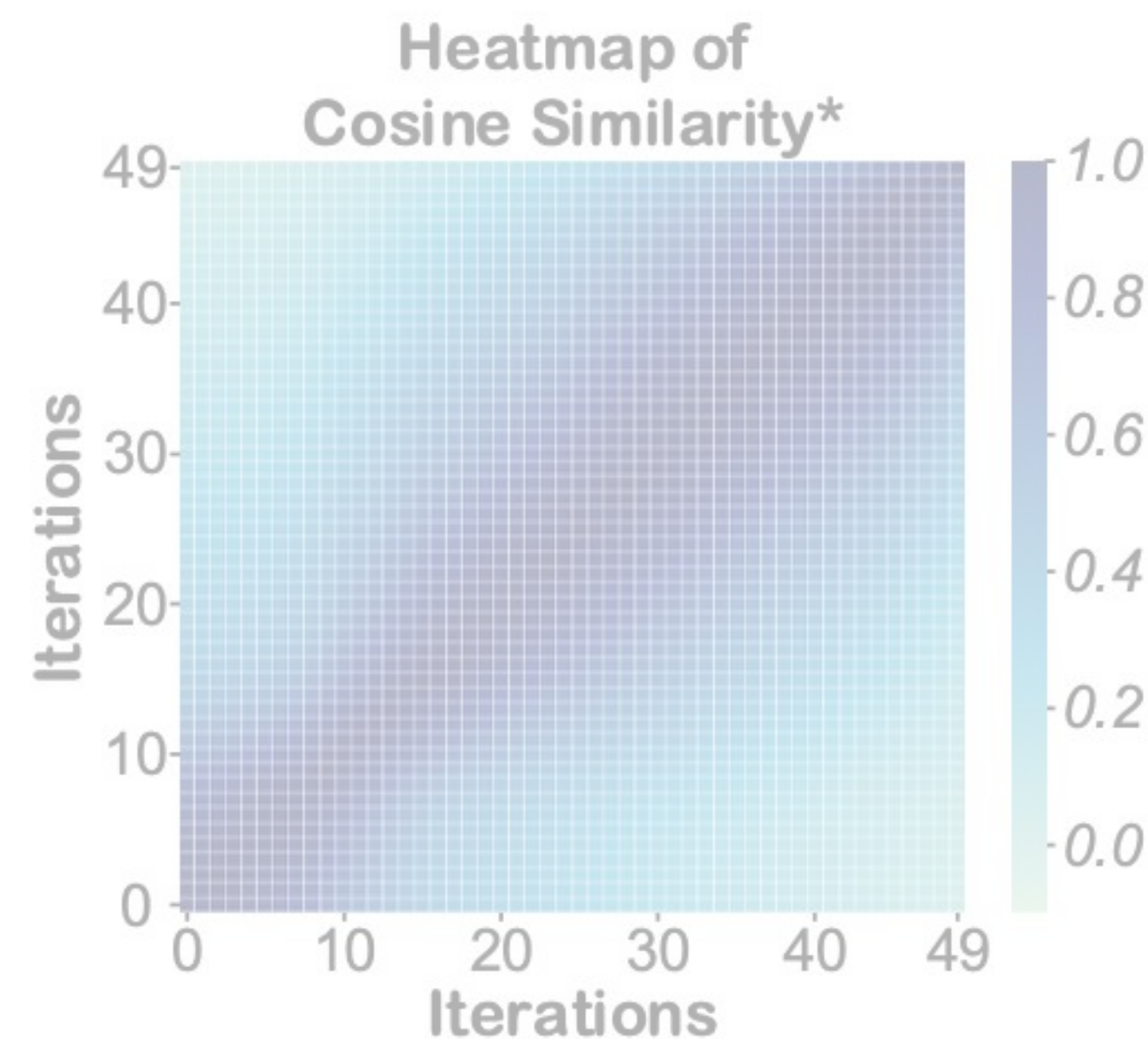
3. Specialized hardware architecture



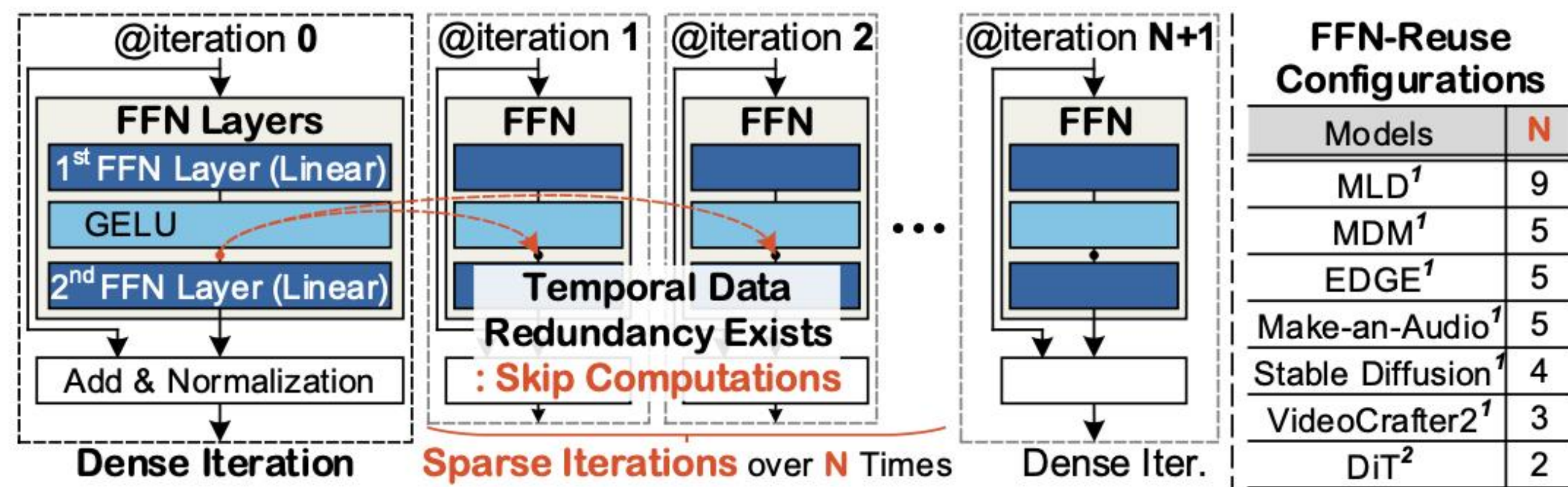Fig. 1. Overview of EXION Accelerator

# Method

## Software optimization 1. FFN-Reuse (inter iteration)

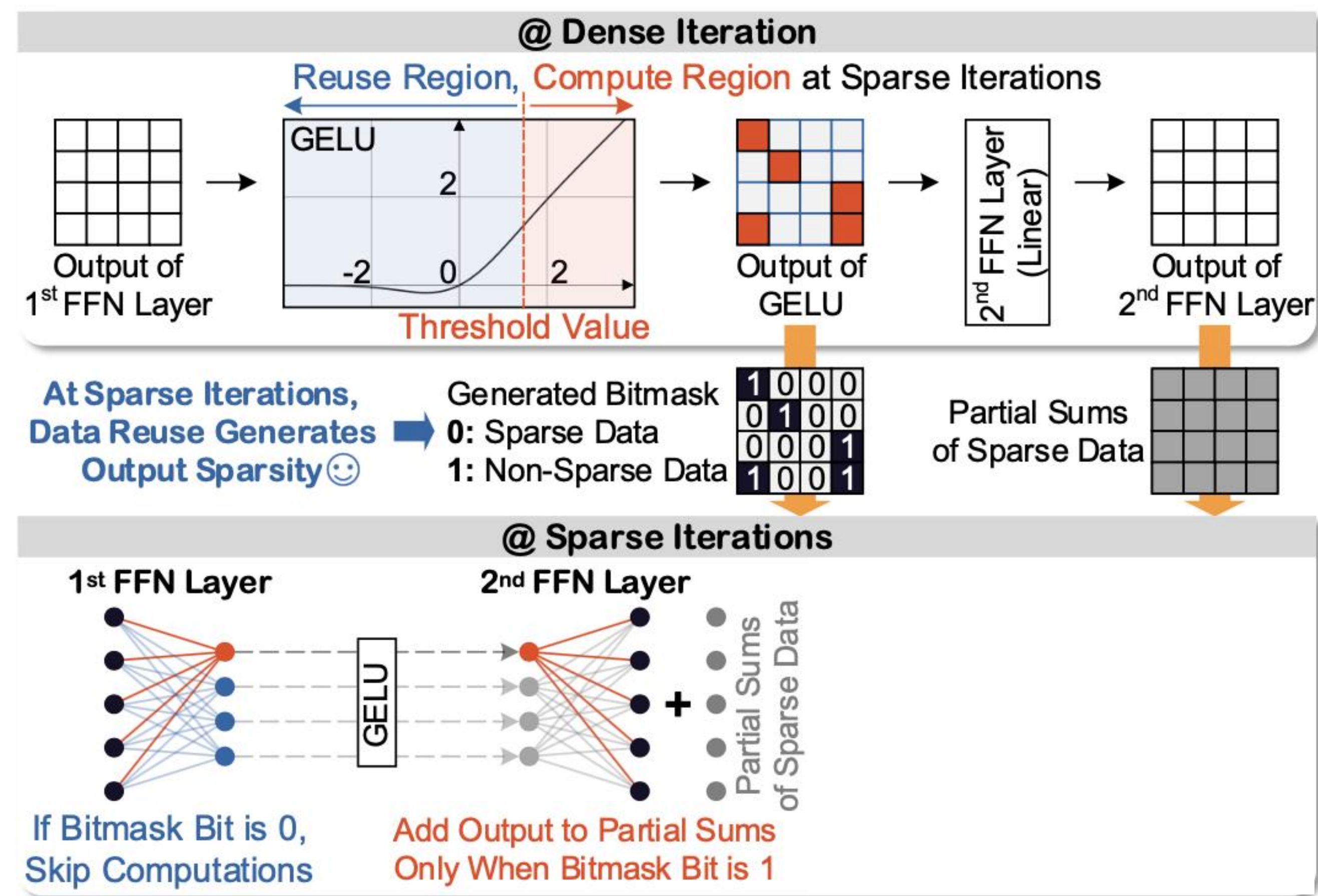Observation: temporal data redundancy exists across different iterations in



Heatmap of
Cosine Similarity*

*: Cosine similarity of 2nd block's GELU
output in FFN layers across different
iterations in DiT model

| Models | N |
|---|---|
| MLD[1] | 9 |
| MDM[1] | 5 |
| EDGE[1] | 5 |
| Make-an-Audio[1] | 5 |
| Stable Diffusion[1] | 4 |
| VideoCrafter2[1] | 3 |
| DiT[2] | 2 |

# Method

## Software optimization 1. FFN-Reuse (inter iteration)

Fig: FFN-Reuse Algorithm for Inter-iteration Output Sparsity



**@ Dense Iteration**

Reuse Region, Compute Region at Sparse Iterations

GELU

Threshold Value

Output of 1st FFN Layer → Output of GELU → 2nd FFN Layer (Linear) → Output of 2nd FFN Layer

At Sparse Iterations, Data Reuse Generates Output Sparsity☺ ⟹ Generated Bitmask
0: Sparse Data
1: Non-Sparse Data

```
1 0 0 0
0 1 0 0
0 0 0 1
1 0 0 1
```

Partial Sums of Sparse Data

**@ Sparse Iterations**

1st FFN Layer        2nd FFN Layer

GELU

Partial Sums of Sparse Data

If Bitmask Bit is 0, Skip Computations

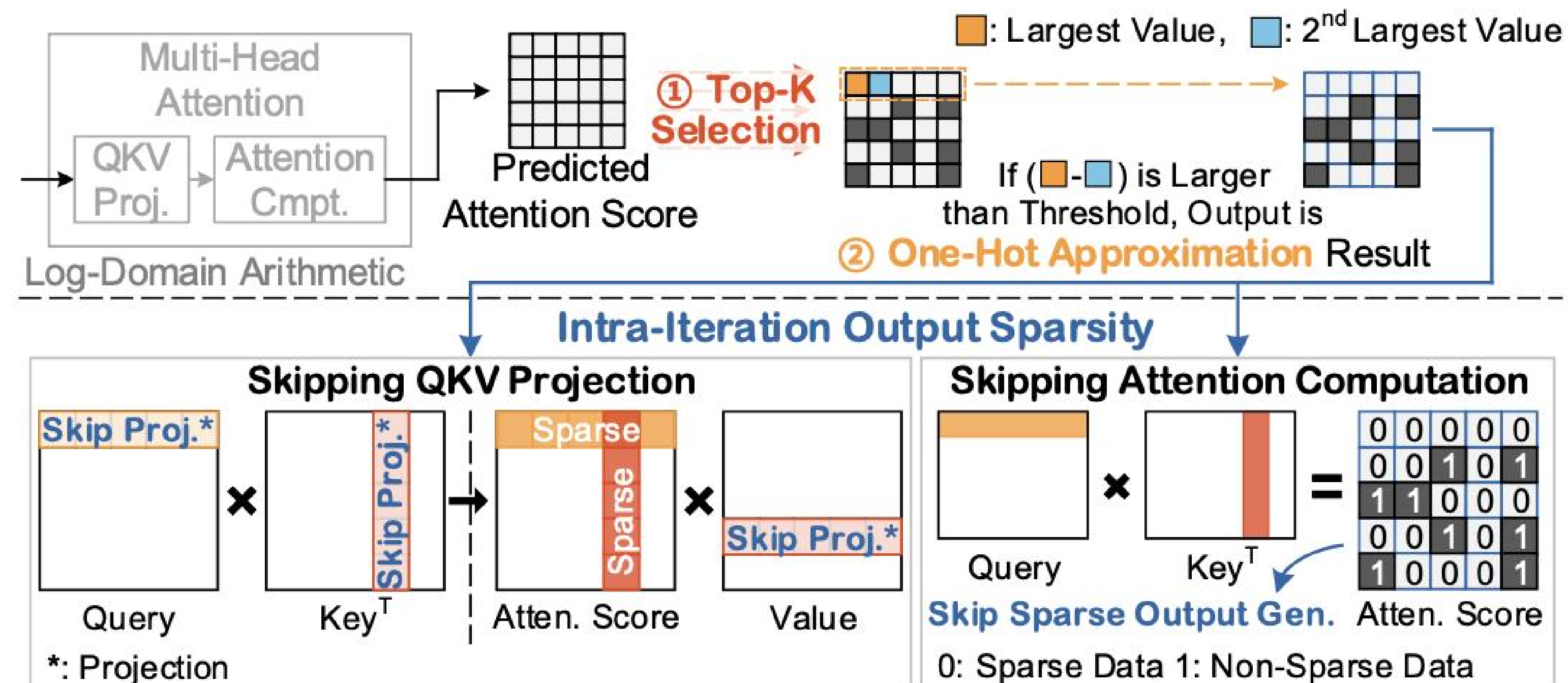Add Output to Partial Sums Only When Bitmask Bit is 1

*1*: Total 50 iterations, *2*: Total 100 iterations, *3*: FFN layers' average reduction percentage

# Method

## Software optimization 2. Early Prediction (intra iteration)

Goal: predict attention score to skip unnecessary computations

- Skip QKV projection

- Skip attention computation $(Q * K^T)$

# Method

## Software optimization 2. Early Prediction (intra iteration)

Question: how to estimate the attention early *before* the QKV generation?

**FACT: FFN-Attention Co-optimized Transformer Architecture with Eager Correlation Prediction**

Yubin Qin*
qyb20@mails.tsinghua.edu.cn
Tsinghua University
Beijing, China

Yang Wang*
wangyang_imec@mail.tsinghua.edu.cn
Tsinghua University
Beijing, China

Dazheng Deng
ddz20@mails.tsinghua.edu.cn
Tsinghua University
Beijing, China

Zhiren Zhao
zhaozr21@mails.tsinghua.edu.cn
Tsinghua University
Beijing, China

Xiaolong Yang
yangxl21@mails.tsinghua.edu.cn
Tsinghua University
Beijing, China

Leibo Liu
liulb@mail.tsinghua.edu.cn
Tsinghua University
Beijing, China

Shaojun Wei
wsj@mail.tsinghua.edu.cn
Tsinghua University
Beijing, China

Yang Hu
hu_yang@mail.tsinghua.edu.cn
Tsinghua University
Beijing, China

Shouyi Yin[†]
yinsy@mail.tsinghua.edu.cn
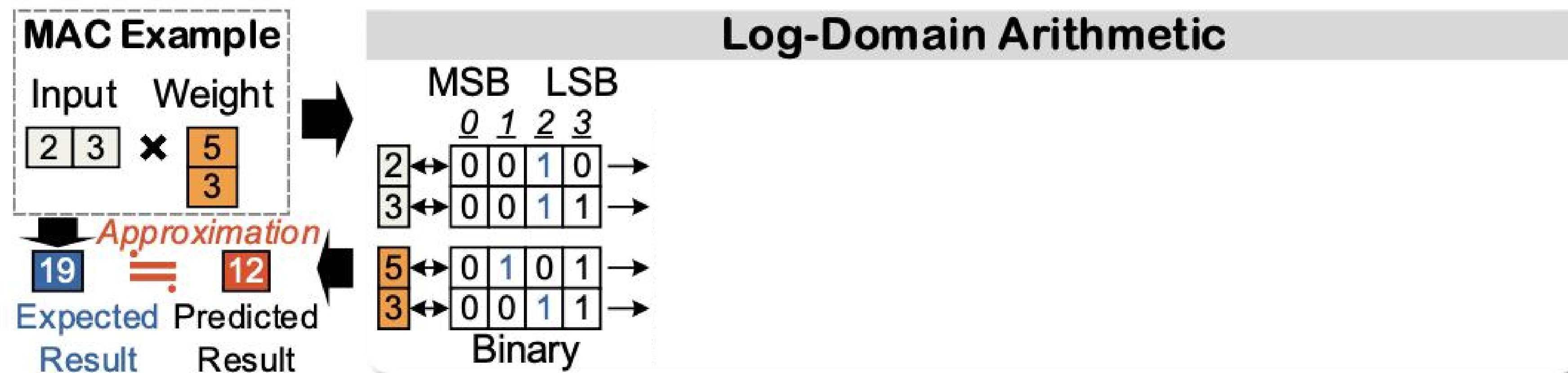Tsinghua University
Beijing, China

ISCA' 2023

# Method

## Software optimization 2. Early Prediction (intra iteration)

Key ideas: use a log-based multiplication-free prediction

- First transforms the data into log-domain by using a *leading-one detector*

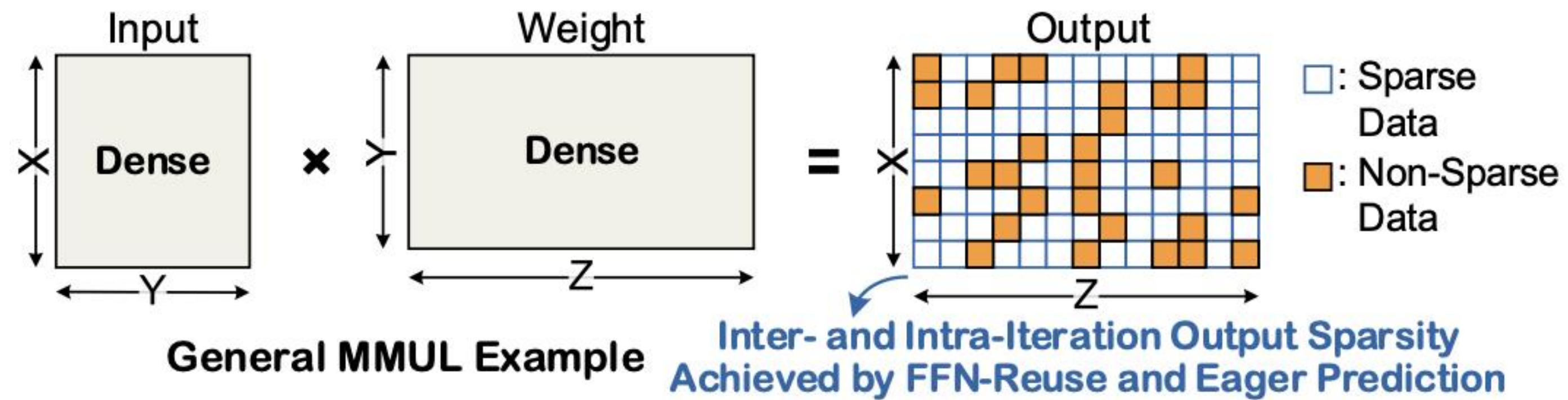- Then substitutes the costly multiplication with low-power shift-and-add op



Ps: an INT-type number $\alpha$ can be decomposed as: $\alpha = Sign \times 2^{(W-LO-1)} \times M$

the multiplication of two INT-type numbers: $\alpha \times \beta = XOR(Sign_\alpha, Sign_\beta)$

$$\times 2^{(W_\alpha + W_\beta - (LO_\alpha + LO_\beta) - 2)}$$

$$\times (M_\alpha \times M_\beta)$$

# Method

## Inter- and Intra-Iteration Output Sparsity



**General MMUL Example**

**Inter- and Intra-Iteration Output Sparsity Achieved by FFN-Reuse and Eager Prediction**

□ : Sparse Data
■ : Non-Sparse Data

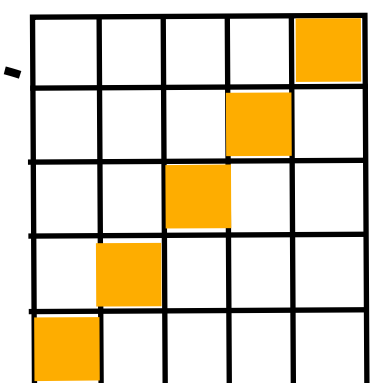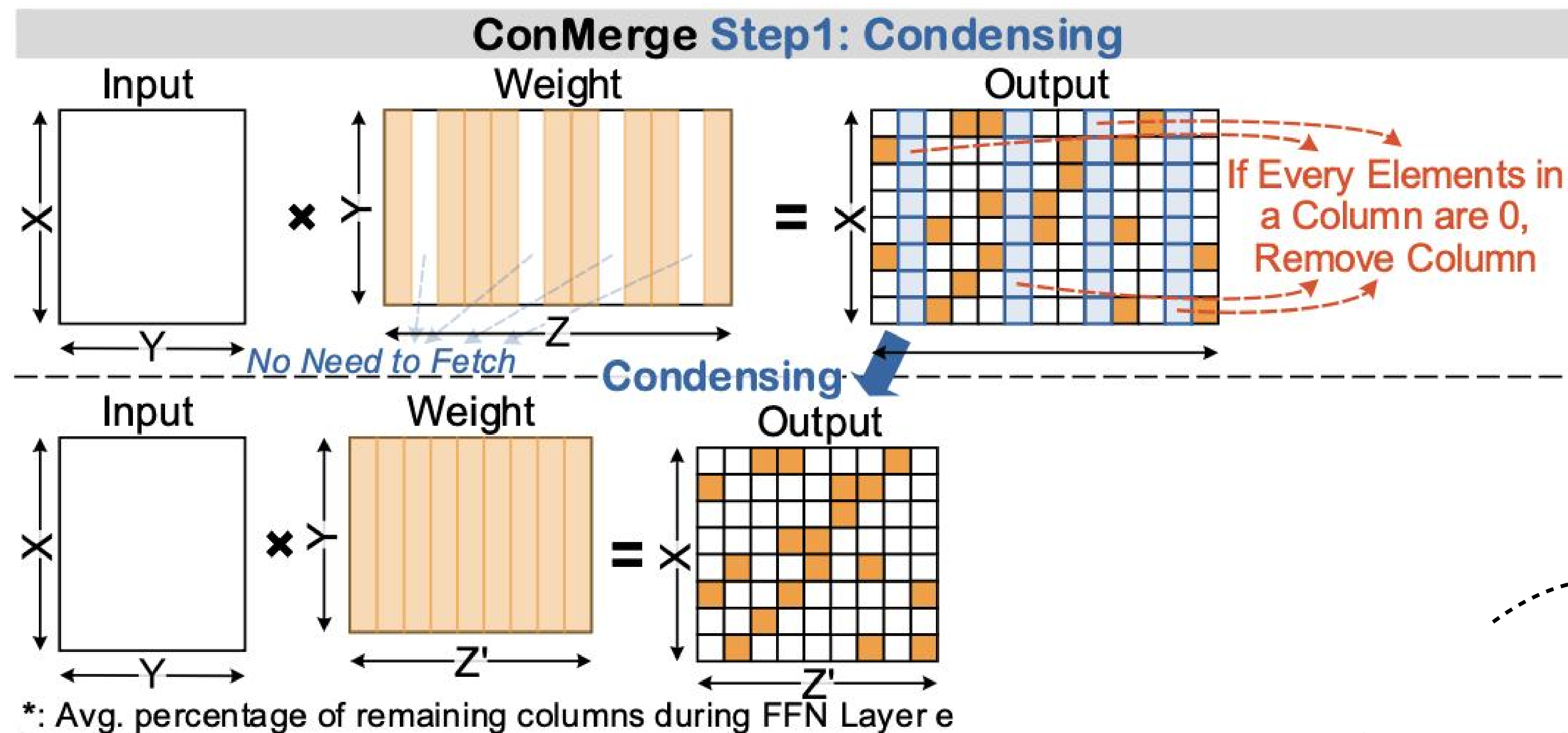Problem: Conventional HW, such as GPUs, cannot utilize it to reduce energy

# Method

## ConMerge: data compaction mechanism

Goal: to **condense** and **merge** large&sparse matrices into small&compact forms
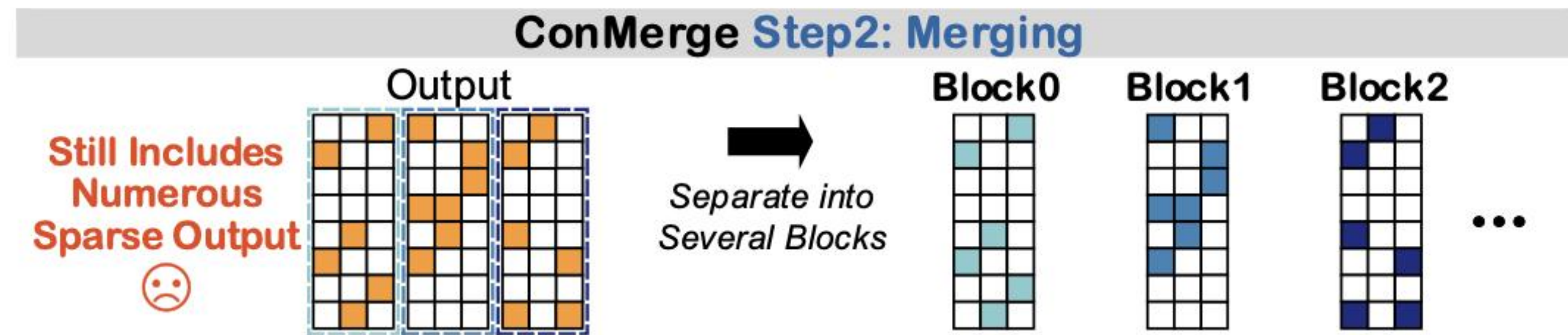
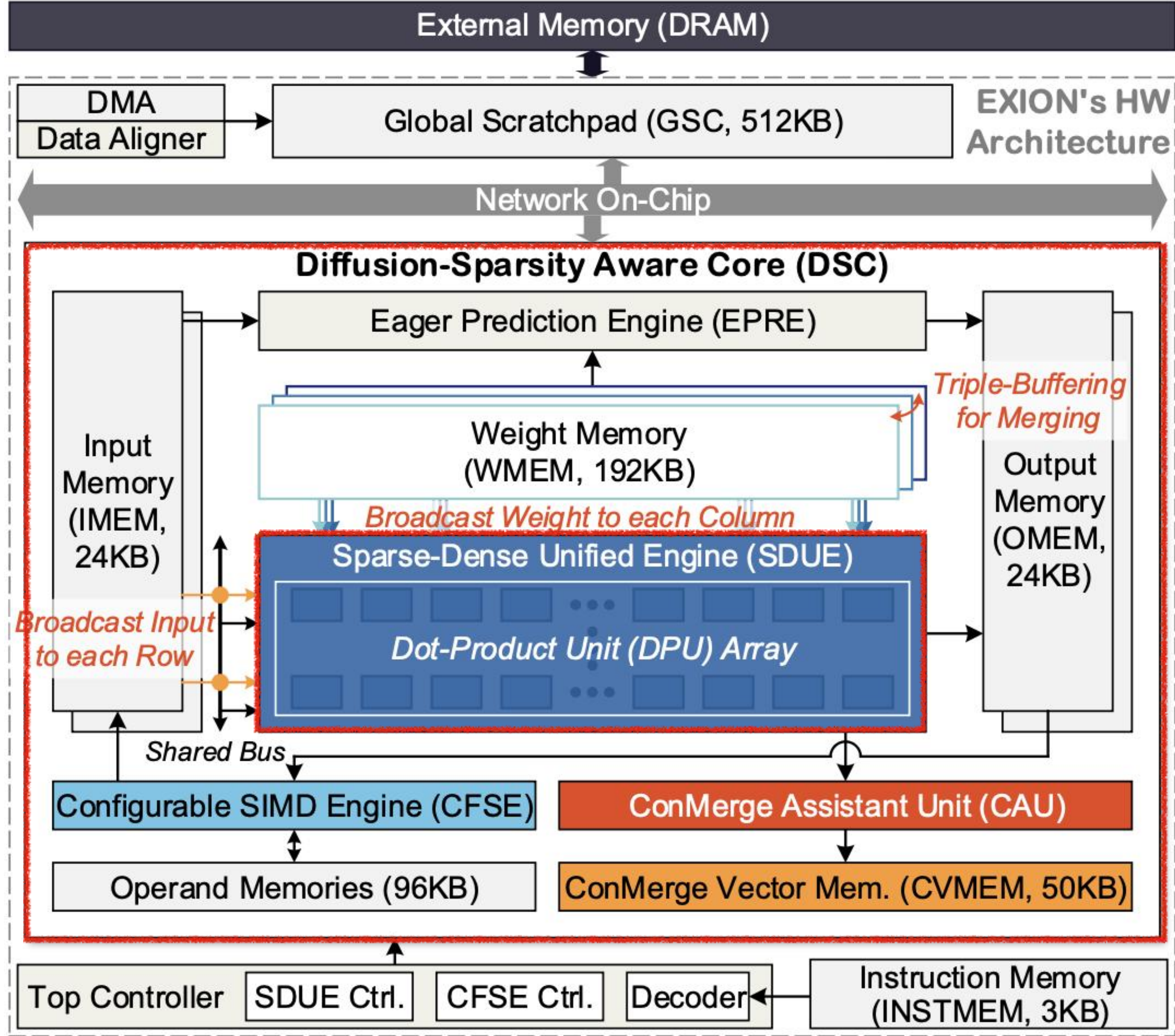# Method

ConMerge: data compaction mechanism

- condense



ConMerge Step1: Condensing

Input × Weight = Output

If Every Elements in a Column are 0, Remove Column

No Need to Fetch

Condensing

*: Avg. percentage of remaining columns during FFN Layer e

# Method

ConMerge: data compaction mechanism

- merge

# Method
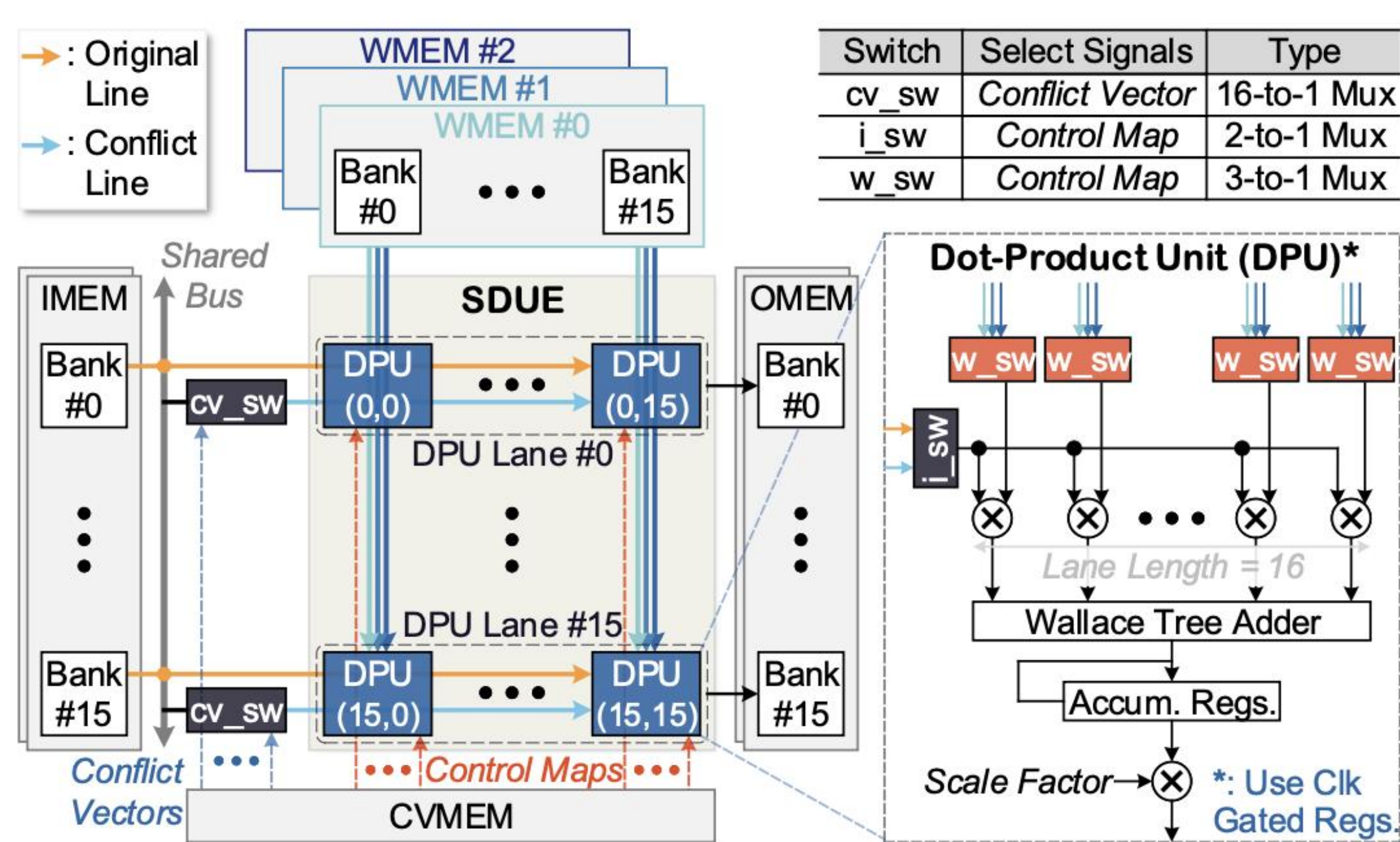
## Hardware co-design: Overview



DSC: run the diffusion model's dense and sparse iterations

Fig: Hardware architecture overview of EXION

# Method

## Hardware co-design: Sparse-Dense Unified Engine (SDUE)



cv_switch: conflict vector switch

i_sw: input switch

w_sw: weight switch

SDUE can compute:

- the normal dense output matrix

- and also the ConMerged block

# Method

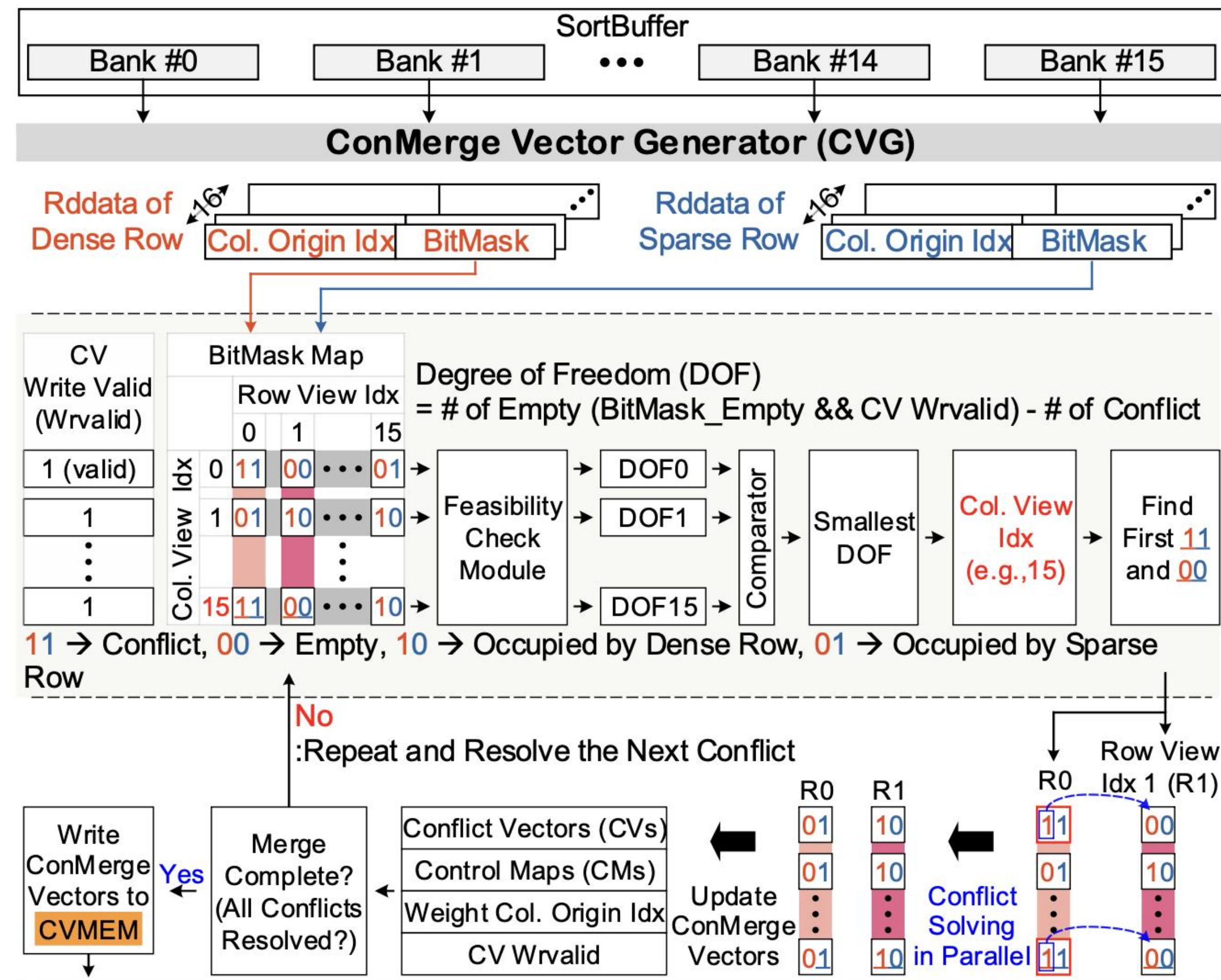## Hardware co-design: ConMerge Assistant Unit (CAU)

- Sorting strategies for fast merging

# Method

## Hardware co-design: ConMerge Assistant Unit (CAU)

- Detailed Merging Process in ConMerge Vector Generator

# Evaluation

## Experimental setup

- Workloads selected: seven different diffusion models

  text-to-motion (MLD and MDM)

  music-to-motion (EDGE)

  text-to-image (Stable Diffusion)

  class-to-image (DiT)

  text-to-audio (Make-an-Audio)

  text-to-video (VideoCrafter2)

- Hardware specifications of GPUs

| | Edge GPU | Server GPU |
|---|---|---|
| | NVIDIA Jetson Orin Nano | NVIDIA RTX 6000 Ada |
| Throughput | 40.0 TOPS | 91.1 TFLOPS[1] |
| Memory Bandwidth | 68 GB/s | 960 GB/s |
| Power Consumption | -15W | -300W |

# Evaluation

## Accuracy Evaluation

| Models | MLD | MDM | EDGE | Make-an-Audio | Stable Diffusion | DiT | VideoCrafter2 |
|---|---|---|---|---|---|---|---|
| Task | Text-to-Motion | Text-to-Motion | Music-to-Motion | Text-to-Audio | Text-to-Image | Image Generation | Text-to-Video |
| Dataset | HumanML3D | HumanML3D | AIST++ | AudioCaps | COCO 2014 | ImageNet 2012 | ECTV |
| Total Iterations | 50 | 50 | 50 | 50 | 50 | 100 | 50 |

<table>
<tr><td colspan="8" align="center"><b>Output Sparsity Achieved by EXION's Software-Level Optimizations</b></td></tr>
<tr><td rowspan="2">FFN-Reuse</td><td>Inter-Iter. Sparsity</td><td>95%</td><td>95%</td><td>95%</td><td>97%</td><td>97%</td><td>80%</td><td>70%</td></tr>
<tr><td>$N^1$</td><td>9</td><td>5</td><td>5</td><td>5</td><td>4</td><td>2</td><td>3</td></tr>
<tr><td>$EP^2$</td><td>Intra-Iter. Sparsity<br>($q\_th^3$, $k^4$)</td><td>30%<br>(q_th=0.3, k=0.7)</td><td>95%<br>(q_th=0.3, k=0.05)</td><td>50%<br>(q_th=0.9, k=0.5)</td><td>80%<br>(q_th=0.7, k=0.2)</td><td>20%<br>(q_th=0.8, k=0.8)</td><td>95%<br>(q_th=0.15, k=0.05)</td><td>50%<br>(q_th=2, k=0.5)</td></tr>
</table>

### Accuracy Evaluation Metric

| Applied Methods | FID (↓) w/ $GT^5$ | R-Precision (↑) | FID (↓) w/ $GT^5$ | PSNR w/ Vanil.$^6$(↑) | PFC (↓) | Beat Align Score (↑) | FAD (↓) w/ $GT^5$ | PSNR w/ Vanil.$^6$(↑) | FID (↓) w/ $GT^5$ | IS (↑) | PSNR w/ Vanil.$^6$(↑) | FID (↓) w/ $GT^5$ | IS (↑) | PSNR w/ Vanil.$^6$(↑) | VQA _A(↑) | IS (↑) | PSNR w/ Vanil.$^6$(↑) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Vanilla Model | 0.393 | 0.754 | 0.406 | - | 1.352 | 0.218 | 4.618 | - | 26.63 | 33.11 | - | 10.63 | 265.73 | - | 59.68 | 17.06 | - |
| FFN-Reuse | 0.401 | 0.752 | 0.467 | 19.22 | 1.389 | 0.214 | 4.932 | 27.14 | 26.05 | 33.04 | 18.20 | 14.42 | 265.73 | 15.99 | 58.92 | 16.83 | 28.25 |
| FFN-Reuse+$EP^2$ | 0.410 | 0.745 | 0.968 | 17.84 | 2.241 | 0.195 | 4.975 | 26.09 | - | - | 14.06 | - | - | 14.60 | - | - | 27.86 |
| FFN-Reuse+$EP^2$+Quant.$^7$ | 0.410 | 0.744 | 1.080 | 17.67 | 2.411 | 0.193 | 5.131 | 25.72 | - | - | 13.94 | - | - | 14.57 | - | - | 27.56 |

*1*: # of sparse iterations between two dense iterations, *2*: EP w/ TS LOD, *3*: Threshold for the difference between the largest value and the 2nd largest value, *4*: Top-k selection ratio (i.e., k=0.5 selects 50% of the data), *5*: Ground truth, *6*: Compute PSNR compared to the vanilla model, *7:* After post-training quantization (INT mixed precision, 12b for SDUE/EPRE and 16/32b for CFSE)
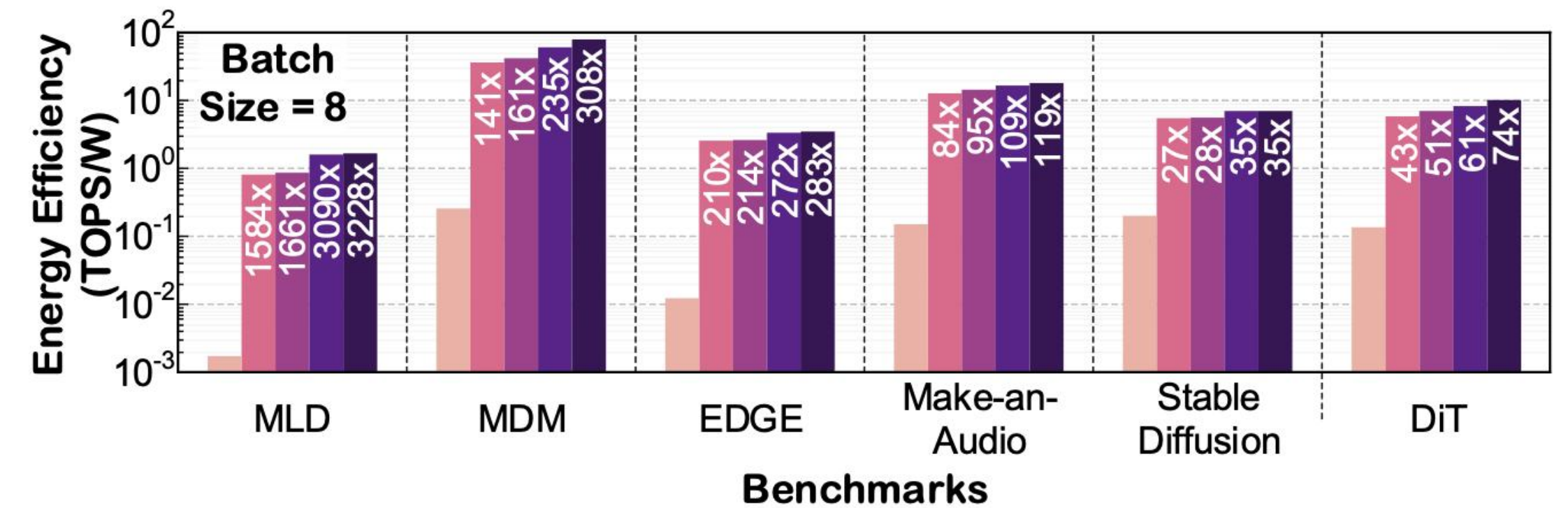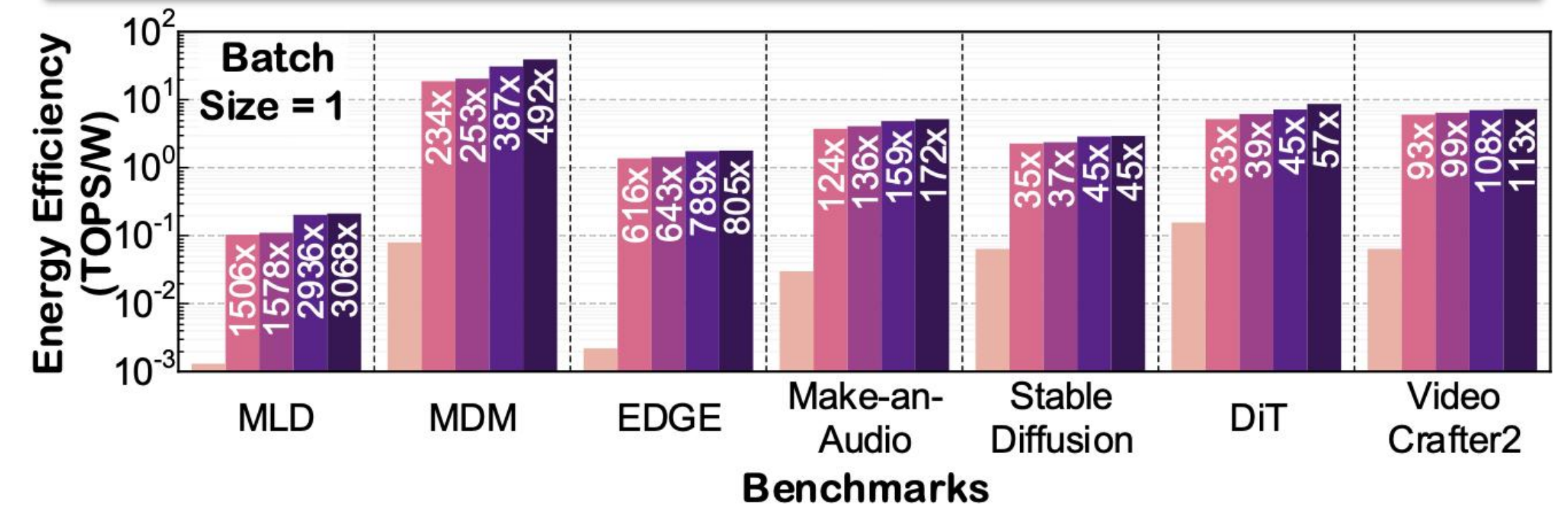
# Evaluation

## Performance Evaluation: Energy efficiency

# Evaluation

## Performance Evaluation: Latency

# Thanks !

Q&A

# Previous Solutions

## Software-based approaches

Reduce the large number of iterations (inference steps), e.g. distillation;

Cache and reuse block results;

- Problems: harms accuracy; some of them even require retraining

## Hardware accelerators

In each iteration: optimizing QKV projection and attention computation

- Problems: no significantly energy&latency reduce of the overall diffusion