

# FlashGEMM: Mesh-Aware Efficient GEMM for 3D-Stacked LLM Accelerators

Xin Fan<sup>1</sup>, Chen Bai<sup>1,2,†</sup>, Xin Yang<sup>2</sup>, Zhenhua Zhu<sup>1,3</sup>,  
Yanhong Wang, Zhaode Zhang, Yuan Xie<sup>1</sup>

<sup>1</sup> The Hong Kong University of Science and Technology

<sup>2</sup> Fudan University

<sup>3</sup> Tsinghua University

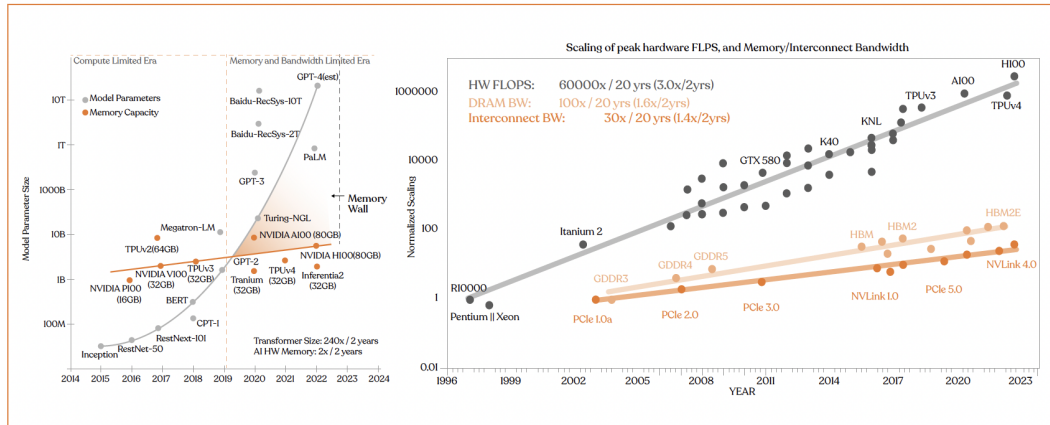
Apr. 20, 2026



- ① Introduction
- ② FlashGEMM
- ③ Evaluation
- ④ Conclusion

# Introduction

# Memory Wall in Large Language Models (LLM) Era

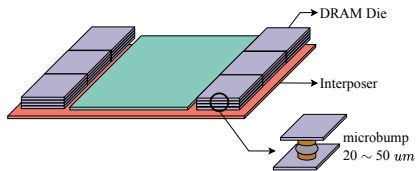


## Memory Wall <sup>1</sup>

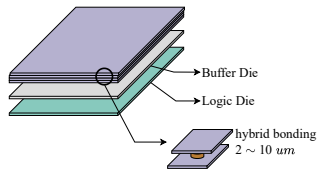
DRAM bandwidth scaling: only 1.6x per two years

This model-hardware scaling disparity creates bottlenecks in both memory capacity and bandwidth.

<sup>1</sup> MVP Ventures (n.d.). *Venture Bytes #111: AI Has a Memory Problem.*



NVIDIA H100 with HBM



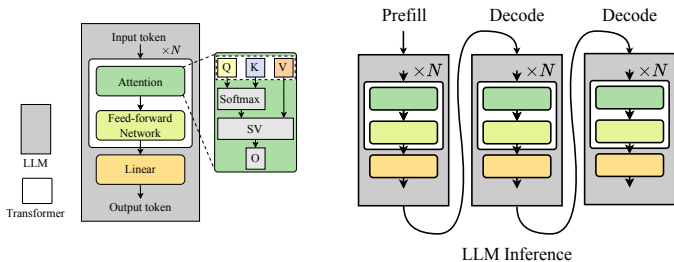
Accelerator with 3D-Stacked Memory

**HBM:** Memory dies (HBM stacks) are placed **side by side** around the logic die.

- Bandwidth density  $\propto$  logic die side length  $N$ .
- Example: NVIDIA B200 achieves up to **8 TB/s** aggregate bandwidth.

**3D-Stacked Memory:** Memory dies are stacked **vertically** on top of the logic die.

- Bandwidth density  $\propto$  logic die area  $N^2$
- Provides **4 $\times$**  higher memory bandwidth than HBM.



## Large Language Models Architecture

- LLMs primarily use transformer blocks.
- **Three key operations:** general matrix-matrix multiplication (GEMM), general matrix-vector multiplication (GEMV) and element-wise operator.
- **GEMM** is the dominant operation.

## LLM Inference Workflow

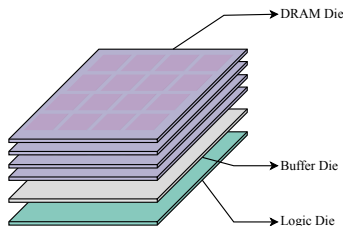
### Prefill Stage

- Processes the entire input prompt in one pass.
- Typically compute-bound due to GEMM.

### Decode Stage

- Use key-value (KV) cache to avoid recomputation.
- Process one new token in one pass.
- Memory-bound due to GEMV and its batched version, flat GEMM.

<sup>2</sup> Flat GEMM denotes  $C_{m \times n} = A_{m \times k} B_{k \times n}$  where  $m \ll k, n$ . For example,  $m = 8$ , and  $k = n = 2048$ .



## 3D-Stacked LLM Accelerators

## Problem Formulation

- Previous distributed GEMM algorithms are inefficient due to unawareness of 3D memory and mesh network-on-chip (NoC) topology.

How to customize an efficient GEMM algorithm for 3D-stacked LLM accelerators?

LLM inference is GEMM-intensive.

<sup>34567</sup>

<sup>3</sup>Robert A Van De Geijn and Jerrell Watts (1997). “SUMMA: Scalable Universal Matrix Multiplication Algorithm”. In: 9.4, pp. 255–274.

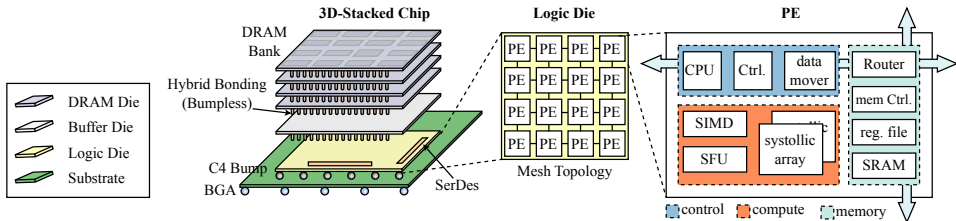
<sup>4</sup>Congjie He et al. (2025). “WaferLLM: Large Language Model Inference at Wafer Scale”. In: *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.

<sup>5</sup>Lynn Elliot Cannon (1969). *A Cellular Computer to Implement the Kalman Filter Algorithm*. Montana State University.

<sup>6</sup>Martin D. Schatz, Robert A. van de Geijn, and Jack Poulson (Jan. 2016). “Parallel Matrix Multiplication: A Systematic Journey”. In: *SIAM J. Sci. Comput.*, pp. C748–C781.

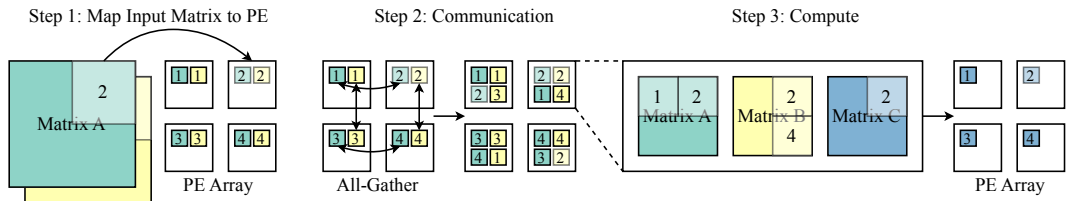
<sup>7</sup>Hyounghwook Nam, Gerasimos Gerogiannis, and Josep Torrellas (2025). “MeshSlice: Efficient 2D Tensor Parallelism for Distributed DNN Training”. In: *IEEE/ACM International Symposium on*

# FlashGEMM



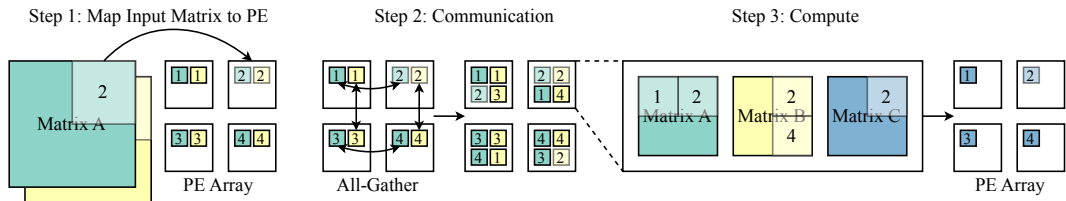
## Architecture Overview of the 3D-Stacked LLM Accelerator.

- Memory dies are stacked vertically on top of the logic die via hybrid bonding.
- Processing Elements (PEs) on logic die are connected through a 2D mesh NoC.



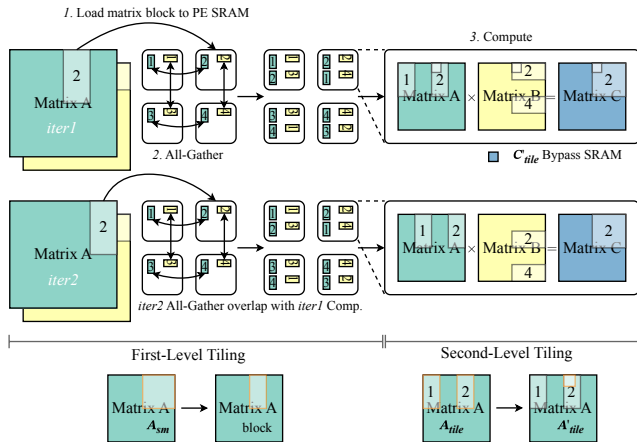
## Original Scalable Universal Matrix Multiplication Algorithm

- Poor communication-computation overlap.
- Long communication path in mesh topology.
- Static output-stationary dataflow.



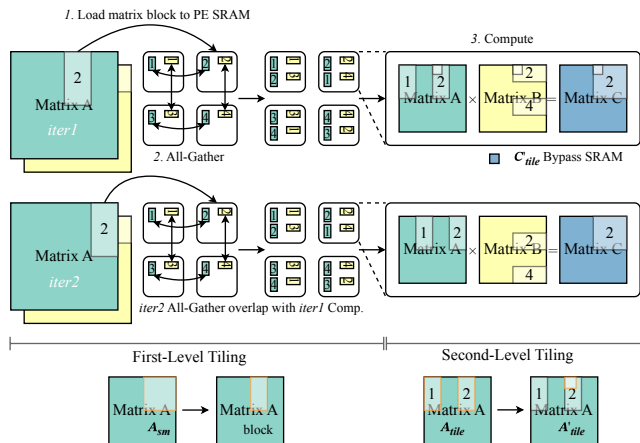
## Original Scalable Universal Matrix Multiplication Algorithm

- Poor communication-computation overlap.
- Long communication path in mesh topology.
- Static output-stationary dataflow.
- In the following, we discuss **FlashGEMM**, our efficient distributed GEMM algorithm.
  - Communication-computation overlap. → Reduce exposed communication time.
  - All-gather algorithm. → Shorten communication hops.
  - Dynamic dataflow. → Minimizes communication volumes.



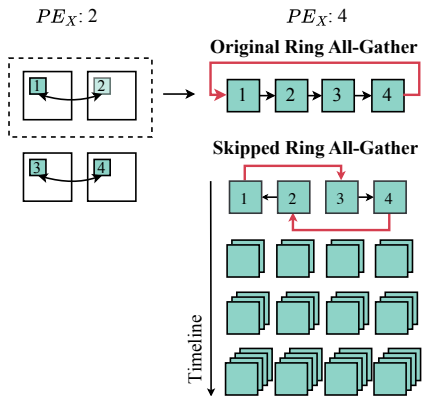
## Two-level Tiling Strategy.

- **First-level tiling:** Slices matrix along K-dimensions. Row/column all-gather of slices run in parallel with computation.
- **Second-level tiling:** Splits every SRAM-level tile into smaller, register-level tiles.



## Two-level Tiling Strategy.

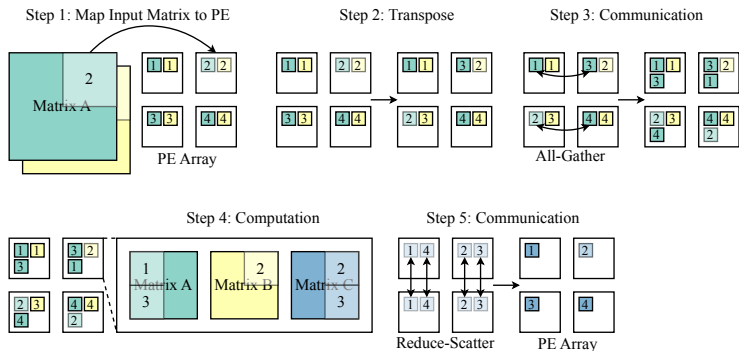
- The stationary variables, i.e.  $C$ , bypass SRAM and load directly from stacked DRAM into registers.



Skipped-Ring All Gather

## Skipped Ring All-Gather

- **All-Gather.** Every GPU collects data from all others.
- Odd-numbered (Even-numbered) PEs send data to the next odd-numbered(Even-numbered) PE.
- Eliminate the wararound path resides in orginal ring all-gather.
- The longest path is 2 hops.



## Weight-Stationary Dataflow

- In the decode stage, the input and output data volume is negligible compared to the weight matrix, owing to the KV cache.
- Adopt output-stationary and weight-stationary dynamically to minimize data movement.

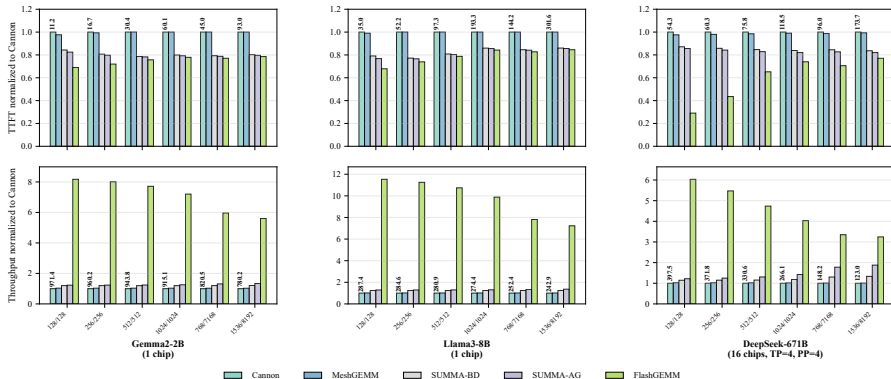
# Evaluation



## Architecture configurations

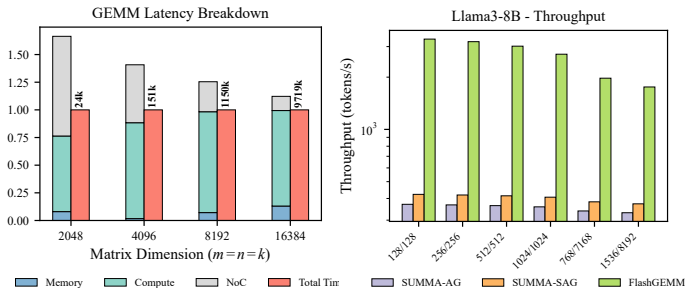
Metric	Value
Silicon area of the monolithic design	$32 \times 25 \text{ mm}^2$
# of DRAM stacks	4
Total DRAM capacity	64 GB
Aggregate DRAM bandwidth	24 TB/s
# of processing elements (PEs)	$16 \times 16$
Computing power	1024 TFLOPS@FP8
NoC bisection bandwidth	1.0 TB/s
Scale-up bandwidth	800 GB/s

# Performance Comparison



## Performance comparison of distributed GEMM algorithms on 3D-stacked LLM accelerators

- FlashGEMM delivers  $1.23\times$  TTFT reduction in prefill and  $5.44\times$  throughput improvement in decode, respectively.



**(a) Latency breakdown. (b) LLM throughput under different settings.**

- FlashGEMM achieves computation-communication overlap and delivers  $1.12\sim 1.66\times$  speedup for GEMM with varying matrix dimensions.
- FlashGEMM and SUMMA with skipped-ring all-gather deliver  $7.44\times$  and  $1.14\times$  average throughput improvements over SUMMA, respectively.

# Conclusion



- We propose FlashGEMM, an efficient matrix-multiplication algorithm for 3D-stacked LLM accelerators with mesh topology.
- Evaluations show that FlashGEMM delivers up to  $1.50\times$  improvement in time-to-first-token (TTFT) and  $7.11\times$  improvement in throughput over prior methods.

**THANK YOU!**



- Lynn Elliot Cannon (1969). *A Cellular Computer to Implement the Kalman Filter Algorithm*. Montana State University.
- Congjie He et al. (2025). “WaferLLM: Large Language Model Inference at Wafer Scale”. In: *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- MVP Ventures (n.d.). *Venture Bytes #111: AI Has a Memory Problem*.  
<https://www.mvp.vc/venture-bytes/venture-bytes-111-ai-has-a-memory-problem>. Accessed: 2025-11-13.
- Hyoungwook Nam, Gerasimos Gerogiannis, and Josep Torrellas (2025). “MeshSlice: Efficient 2D Tensor Parallelism for Distributed DNN Training”. In: *IEEE/ACM International Symposium on Computer Architecture (ISCA)*, pp. 821–834.
- Martin D. Schatz, Robert A. van de Geijn, and Jack Poulson (Jan. 2016). “Parallel Matrix Multiplication: A Systematic Journey”. In: *SIAM J. Sci. Comput.*, pp. C748–C781.
- Robert A Van De Geijn and Jerrell Watts (1997). “SUMMA: Scalable Universal Matrix Multiplication Algorithm”. In: 9.4, pp. 255–274.