



# Compute Scaling for AI

Fan Yang

Systems Research Group

Microsoft Research Aisa

In collaboration with



Imperial College  
London

# The Key to AI Success: *Scaling Laws*

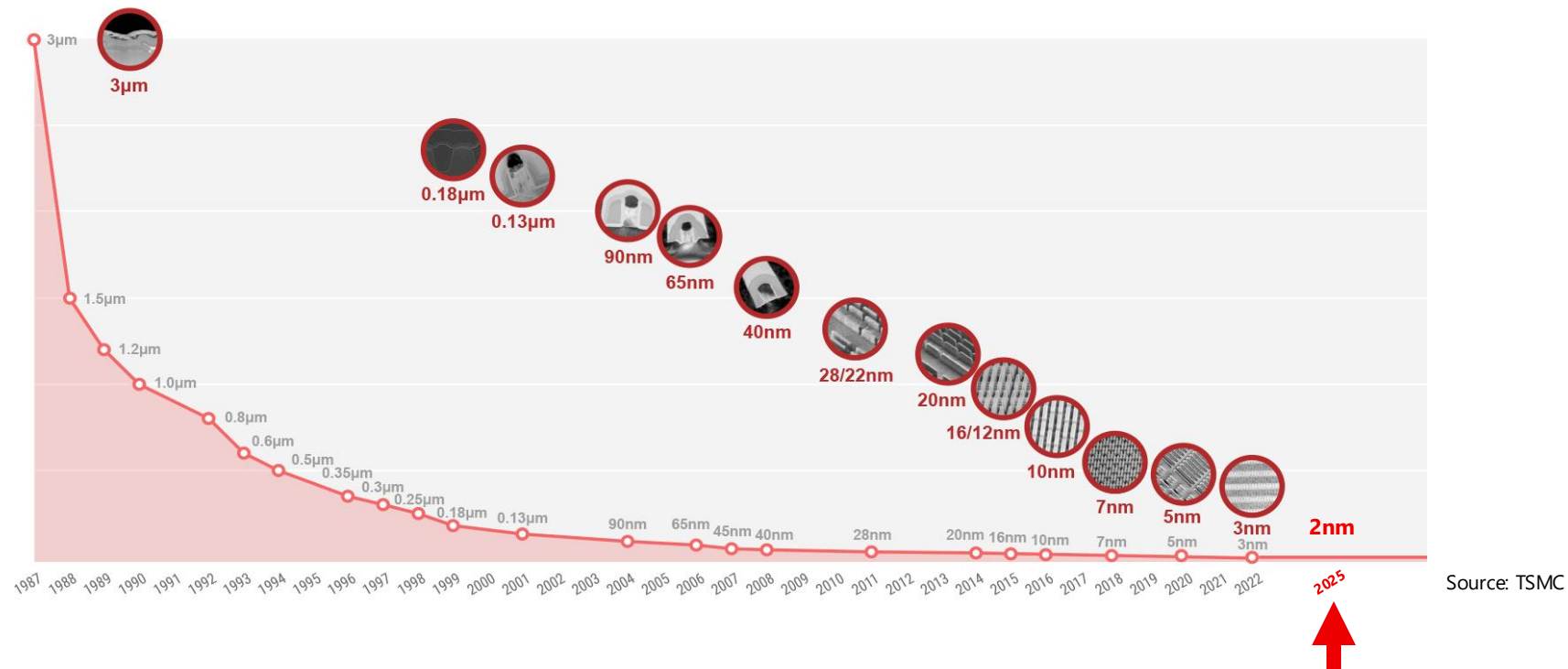
Four dimensions

- Model – size
- Data – volume
- Time – test-time scaling
- Compute – foundation



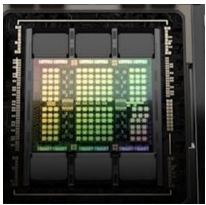
# Compute Scaling Faces Fundamental Challenges

- Semiconductor nodes approaching physical limits



# An Arms Race for Larger, More Complicated AI Chips

NVIDIA H100

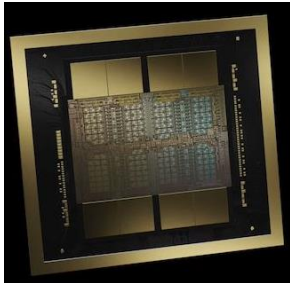


Single-Die  
814 mm<sup>2</sup>

Larger chips

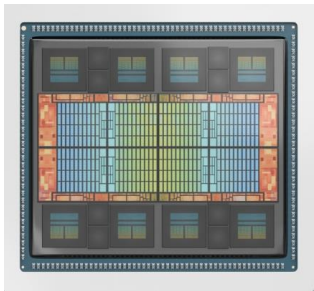


NVIDIA B200



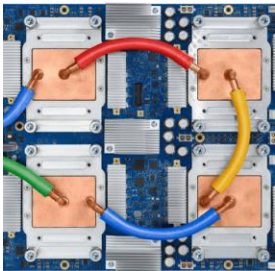
Multi-Die, 1600 mm<sup>2</sup>

AMD MI300



Multi-Die, 1017 mm<sup>2</sup>

Google TPU v3



Multi-Chip, 4 x 700 mm<sup>2</sup>



GPU architecture

CUDA core: general-purpose computation

TensorCore: matrix multiplication /w wider tile shape  
(MMA → WMMA → WGMMA)

TMA: async memory load & store

More complex hardware



# Problem – Design & Manufacturing Difficulties

- Higher chances of design flaws and/or manufacture defects

Mon 5 Aug 2024 // 13:23 UTC

**UPDATED** Nvidia is understood to be delaying shipments of its Blackwell GPUs until the first quarter of 2025, and it appears the problems may be due to the complexity of the chip-on-wafer-on-substrate (CoWoS) packaging tech that TSMC is using to manufacture the next-gen hardware.

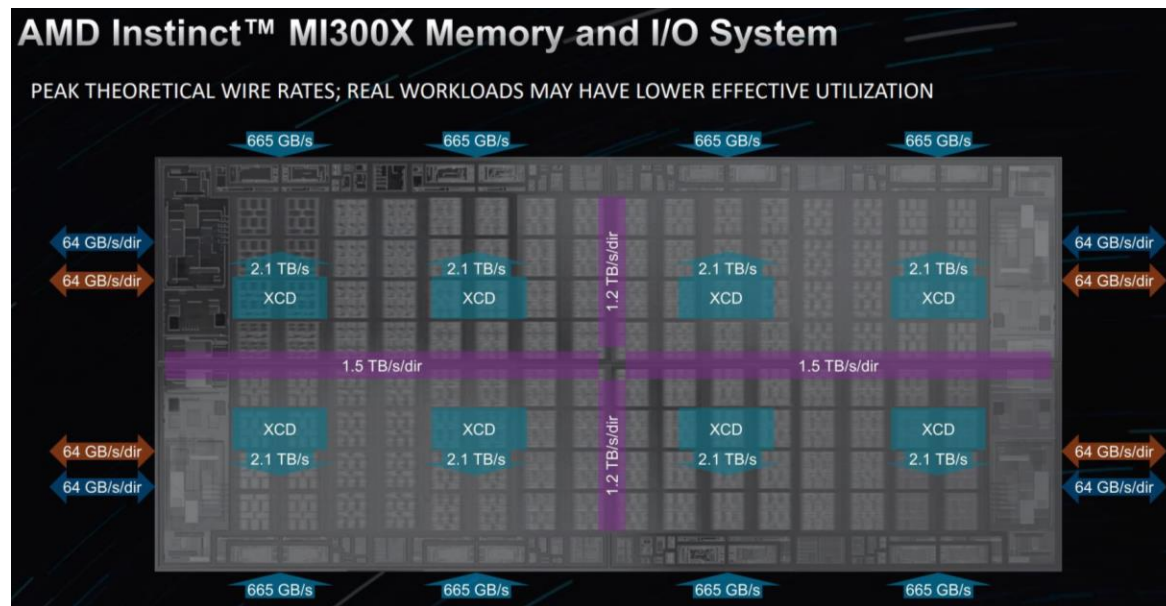
**2025 Jan 13 (Reuters)** - Nvidia's ([NVDA.O](#)) top customers are delaying orders of the AI chip leader's latest 'Blackwell' racks due to overheating issues, the Information reported on Monday.



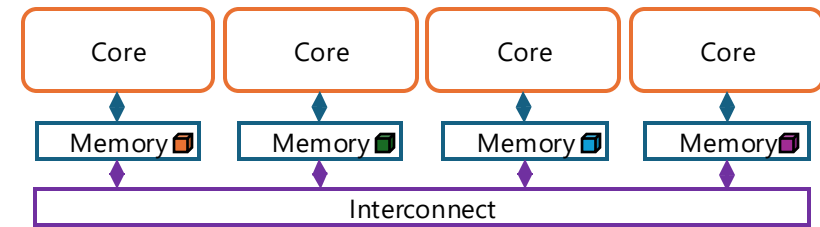


# Problem – Addressing the Uniqueness of Large Chips

- Example: harder to hide non-uniformity in larger chips



Source: AMD



AI workloads were predominantly optimized for chips with shared memory architectures



# A Path Forward Driven by Simplicity

- A codesign to *simplify* both AI models and AI chips
  - BitNet – a simpler 1-bit LLM (Bitnet.cpp)
  - LUT tensor core – a simplified tensor core for further compute scaling (ISCA'25)
- A stack *modeling* the key properties of modern AI chips
  - WaferLLM – a new system stack for wafer-scale chips (OSDI'25)



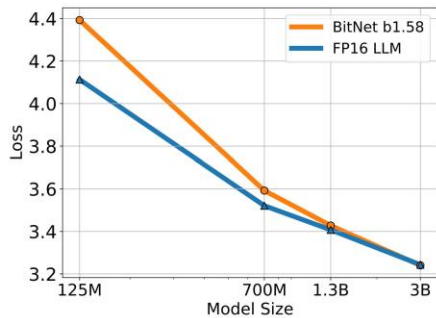
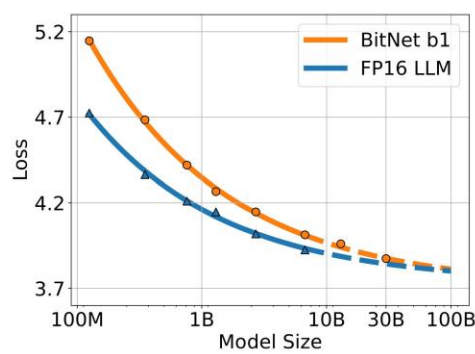
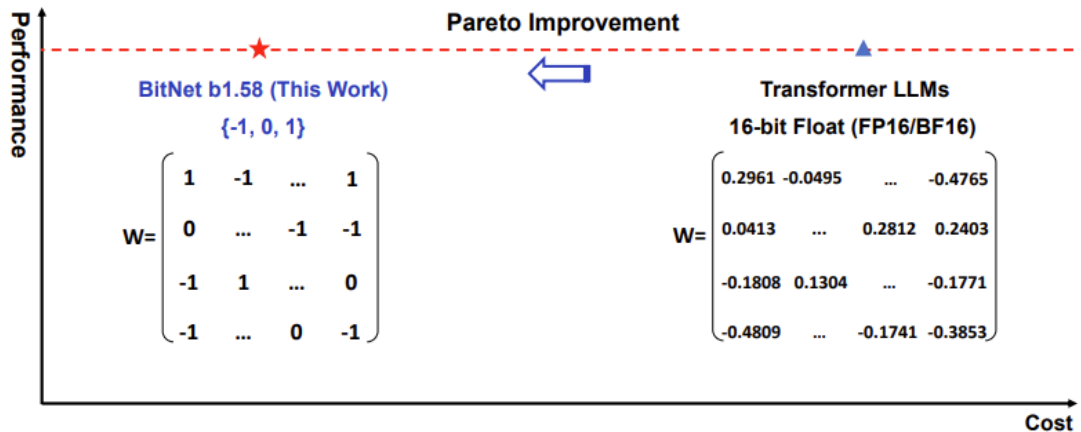
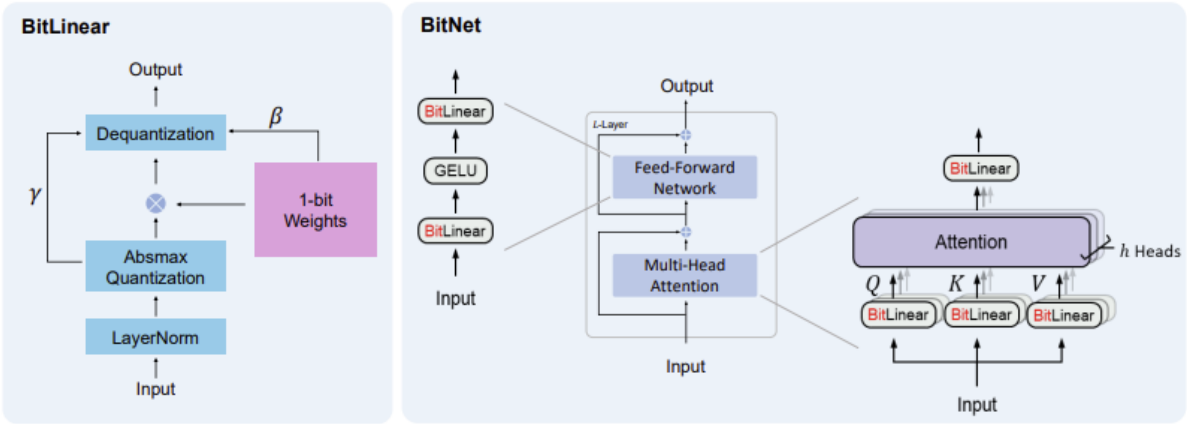
# A Path Forward Driven by Simplicity

- A codesign to *simplify* both AI models and AI chips
  - **BitNet** – a **simpler 1-bit LLM** (Bitnet.cpp)
  - LUT tensor core – a simplified tensor core for further compute scaling (ISCA'25)
- A stack *modeling* the key properties of modern AI chips
  - WaferLLM – a new system stack for wafer-scale chips (OSDI'25)





# The Era of 1-bit LLMs



Credit: BitNet team

1-bit LLMs match full-precision LLMs when the model scale is large enough



# A Path Forward Driven by Simplicity

- A codesign to *simplify* both AI models and AI chips
  - BitNet – a simpler 1-bit LLM (Bitnet.cpp)
  - **LUT tensor core – a simplified tensor core for further compute scaling** (ISCA'25)
- A stack *modeling* the key properties of modern AI chips
  - WaferLLM – a new system stack for wafer-scale chips (OSDI'25)



# LUT TensorCore: Compute Scaling by Simplicity

- 1-bit LLMs: replacing ALU w/ LUT (Lookup Table)

Appendix D

**TABLE OF NATURAL TRIGONOMETRIC FUNCTIONS**

(Note: When entering tables with an angle larger than 45 select such angle from right hand side and obtain values in column corresponding to the function at bottom of page.)

Angle	Sin	Cos	Tan	Cot	Sec	Csc
0°	.0000	1.0000	.0000	∞	1.000	∞
1	.0174	.9998	.0175	57.29	1.000	.89
2	.0349	.9994	.0349	28.64	1.001	.86
3	.0523	.9986	.0524	19.08	1.001	.82
4	.0698	.9975	.0699	14.30	1.001	.80
5	.0872	.9962	.0875	11.43	1.001	.78
6	.1045	.9945	.1051	9.514	1.001	.76
7	.1219	.9925	.1228	8.017	1.001	.74
8	.1392	.9903	.1405	6.913	1.001	.72
9	.1564	.9877	.1584	6.011	1.002	.70
10	.1736	.9848	.1763	5.311	1.002	.68
11	.1908	.9816	.1946	4.705	1.002	.66
12	.2079	.9781	.2126	4.170	1.002	.64
13	.2250	.9744	.2306	3.707	1.002	.62
14	.2419	.9703	.2483	3.306	1.002	.60
15	.2588	.9659	.2667	2.967	1.002	.58
16	.2756	.9613	.2857	2.682	1.002	.56
17	.2924	.9564	.3057	2.443	1.002	.54
18	.3090	.9513	.3267	2.238	1.002	.52
19	.3256	.9460	.3493	2.058	1.002	.50
20	.3420	.9405	.3736	1.904	1.002	.48
21	.3582	.9348	.3996	1.771	1.002	.46
22	.3743	.9289	.4272	1.656	1.002	.44
23	.3902	.9228	.4564	1.556	1.002	.42
24	.4060	.9165	.4872	1.470	1.002	.40
25	.4216	.9100	.5197	1.396	1.002	.38
26	.4371	.9033	.5540	1.332	1.002	.36
27	.4524	.8964	.5902	1.276	1.002	.34
28	.4676	.8893	.6283	1.226	1.002	.32
29	.4826	.8820	.6684	1.182	1.002	.30
30	.4975	.8746	.7115	1.143	1.002	.28
31	.5123	.8670	.7576	1.109	1.002	.26
32	.5269	.8593	.8067	1.079	1.002	.24
33	.5414	.8514	.8589	1.052	1.002	.22
34	.5557	.8434	.9142	1.027	1.002	.20
35	.5699	.8353	.9726	1.004	1.002	.18
36	.5840	.8270	1.0343	.982	1.002	.16
37	.5979	.8186	1.1004	.961	1.002	.14
38	.6117	.8100	1.1718	.940	1.002	.12
39	.6253	.8013	1.2487	.920	1.002	.10
40	.6388	.7925	1.3313	.900	1.002	.08
41	.6521	.7836	1.4198	.880	1.002	.06
42	.6653	.7745	1.5144	.860	1.002	.04
43	.6784	.7653	1.6154	.840	1.002	.02
44	.6913	.7559	1.7230	.820	1.002	.00
45	.7041	.7464	1.8375	.800	1.000	.00
	Cos	Sin	Cot	Tan	Csc	Sec

482

Lookup Table

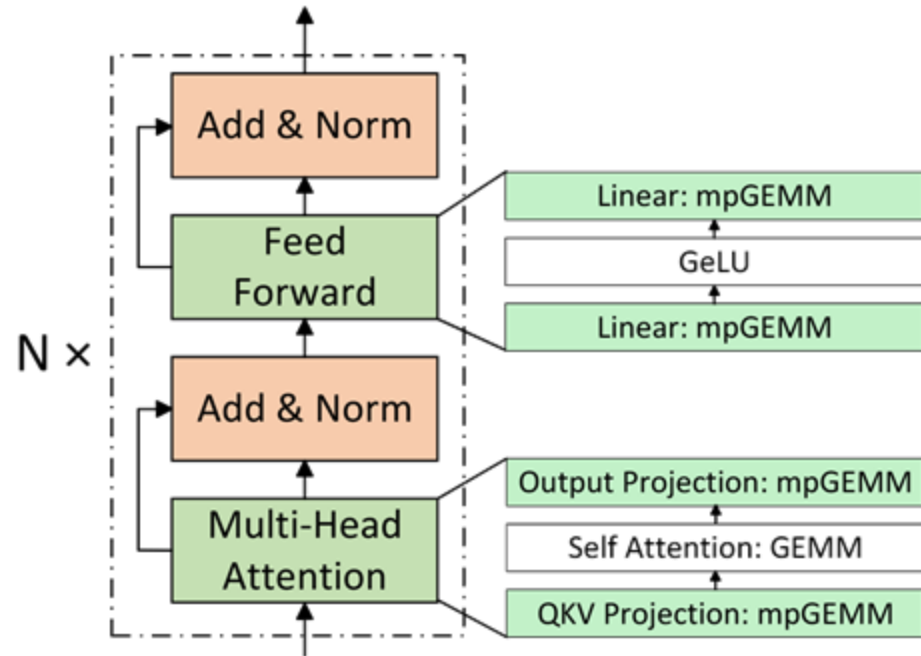
FFT Test

	time	Signal
Row 1	0.001	2.086
Row 2	0.002	2.209
Row 3	0.003	1.214
Row 4	0.004	1.632
Row 5	0.005	1.899
Row 6	0.006	-0.4744
Row 7	0.007	1.521
Row 8	0.008	0.8616
Row 9	0.009	0.5272
Row 10	0.01	2.404

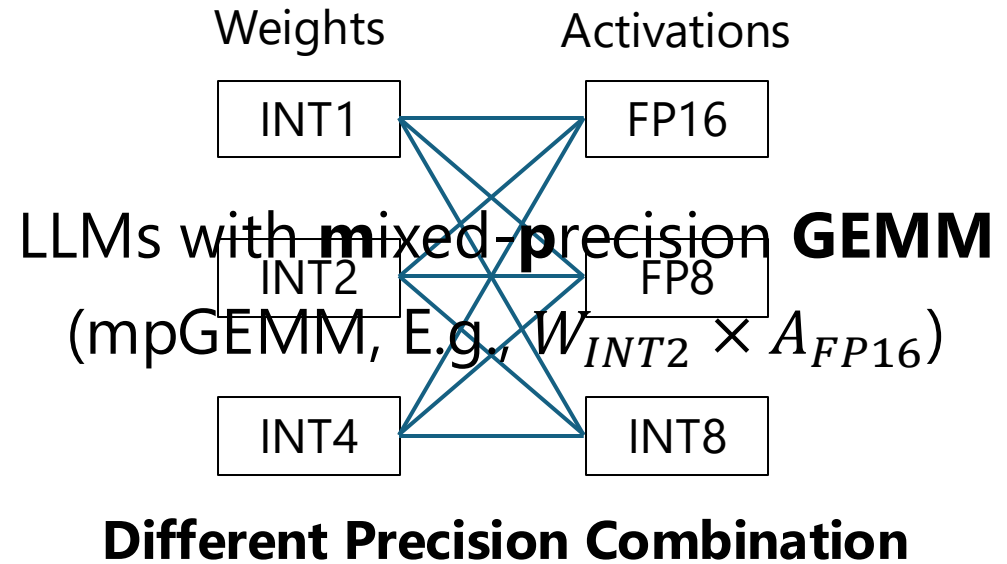
When compute couldn't scale, scientists transform them into tables/codebooks



# The Inconvenience – Not Entirely 1-Bit

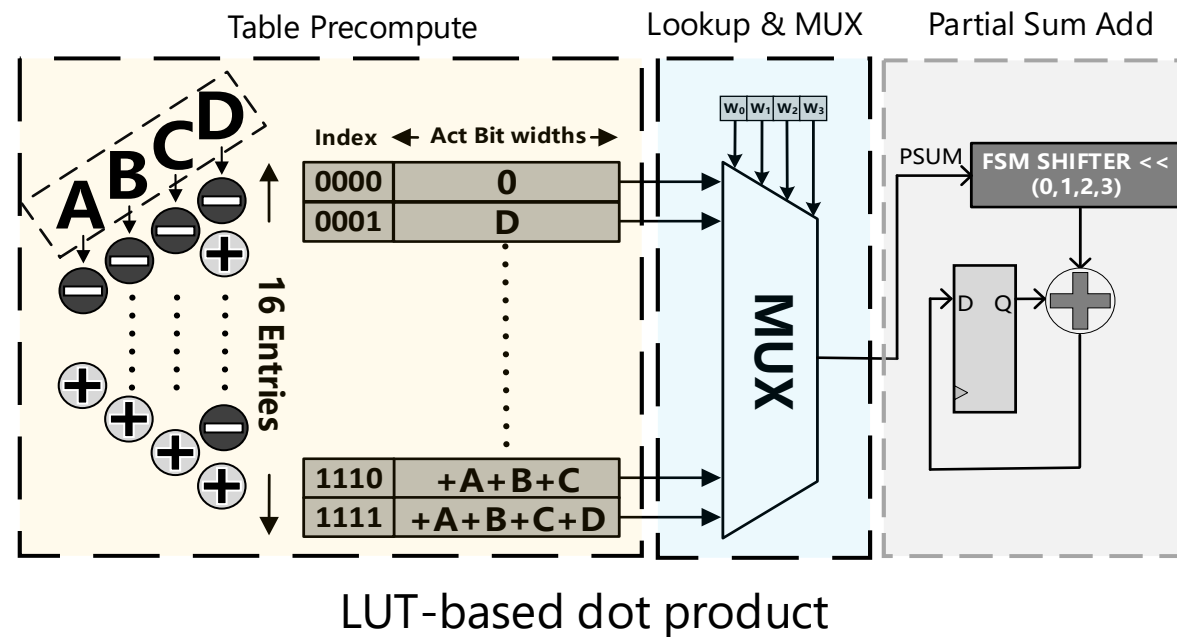


A transformer layer



# The Inconvenience – Table Overhead

- Table size still non-negligible : e.g.,  $W_{INT2} \times A_{FP16}$
- Table precompute overhead



# A Software-Hardware Codesign

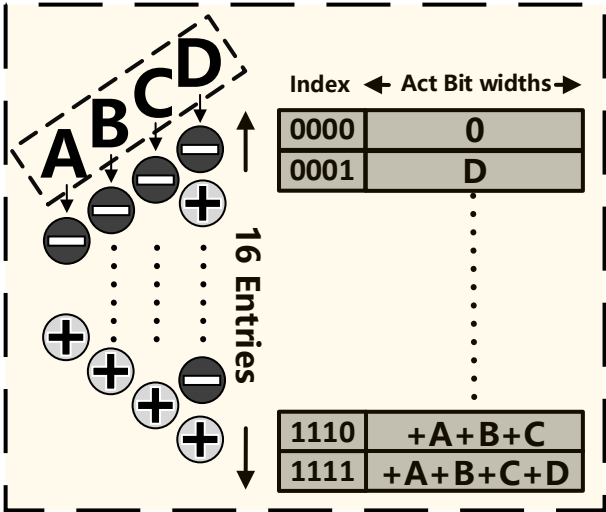
## Large Table Size

	000	001	010	011	100	101	110	111
000	000000	000000	000000	000000	000000	000000	000000	000000
001	000000	000001	000010	000011	000100	000101	000110	000111
010	000000	000010	000100	000110	001000	001010	001100	001110
011	000000	000011	000110	001001	001100	001111	010010	010101
100	000000	000100	001000	001100	010000	010100	011000	011100
101	000000	000101	001010	001111	010100	011001	011110	100011
110	000000	000110	001100	010010	011000	011110	100100	101010
111	000000	000111	001110	010101	011100	100011	101010	110001



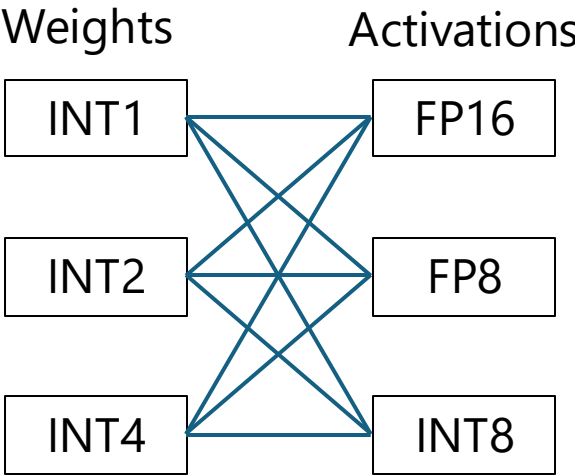
Table Symmetrization

## Table Precompute Overhead



Dedup & Fusion  
(Ladder, OSDI'24)

## Different Precision Combination

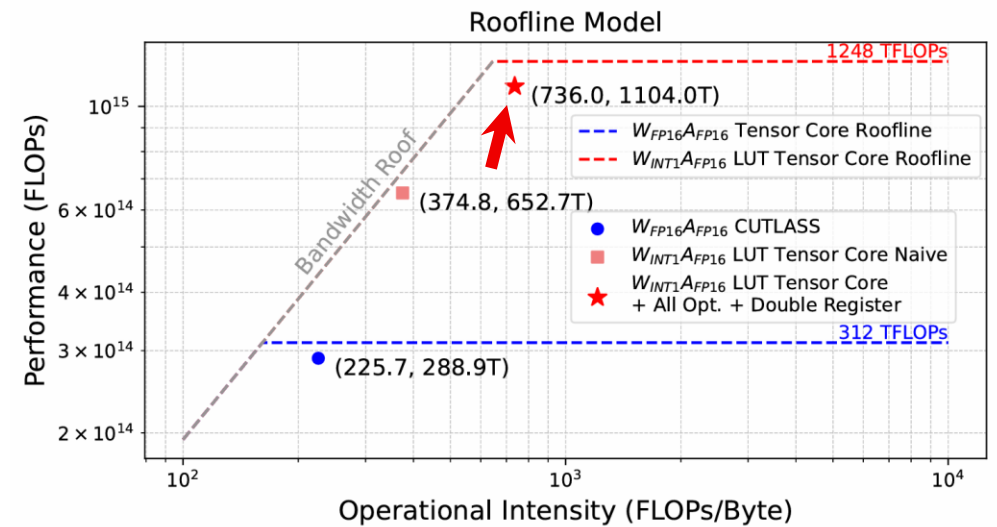
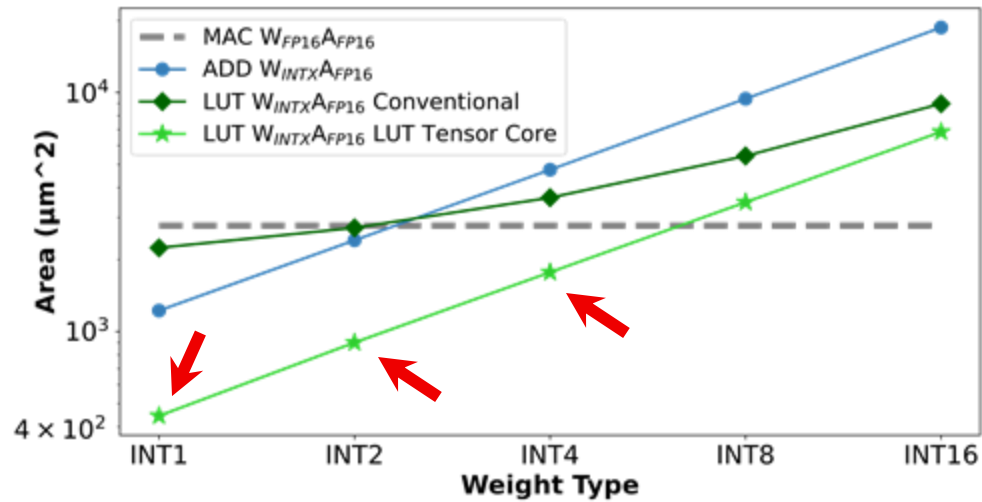


Bit-Serial Circuit





# LUT Tensor Core Scales Better at Lower Bits



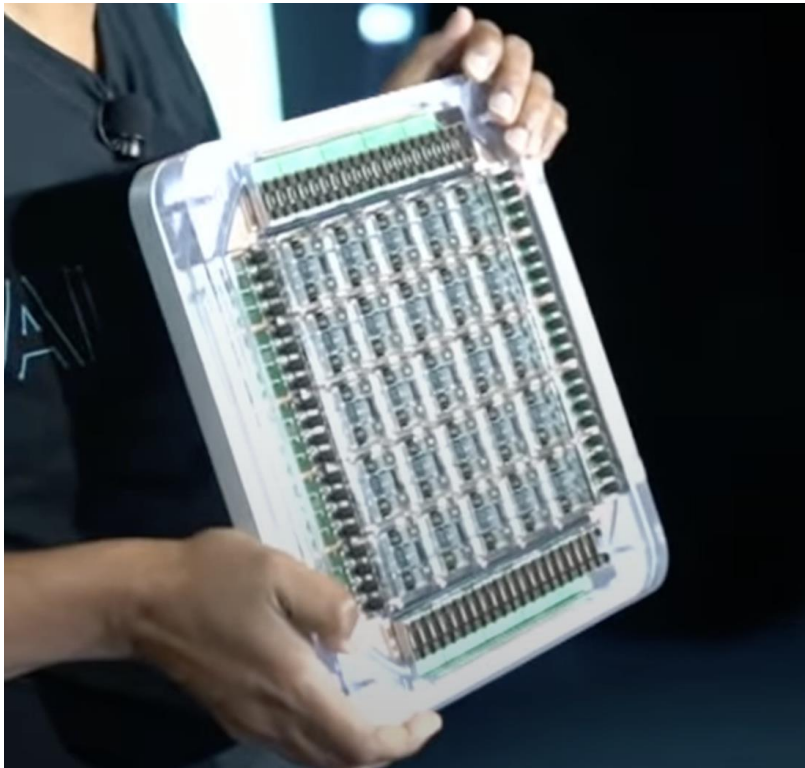
# A Path Forward Driven by Simplicity

- A codesign to *simplify* both AI models and AI chips
  - BitNet – a simpler 1-bit LLM (Bitnet.cpp)
  - LUT tensor core – a simplified tensor core for further compute scaling (ISCA'25)
- A stack *modeling* the key properties of modern AI chips
  - **WaferLLM – a new system stack for wafer-scale chips** (OSDI'25)

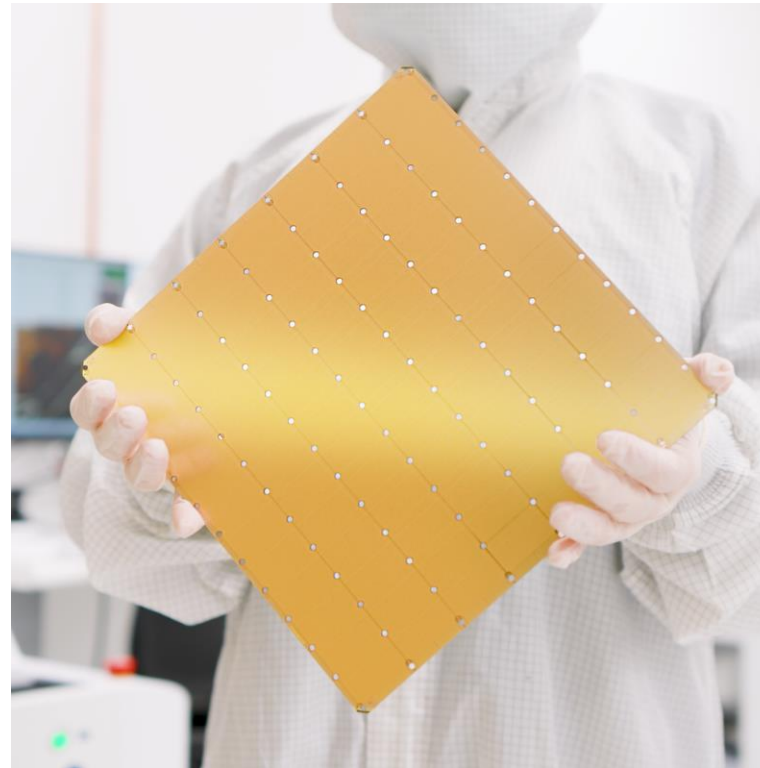


# A System Stack Modeling Modern AI Chips

## Case study – Wafer-scale chips



Tesla Dojo



Cerebras WSE-3



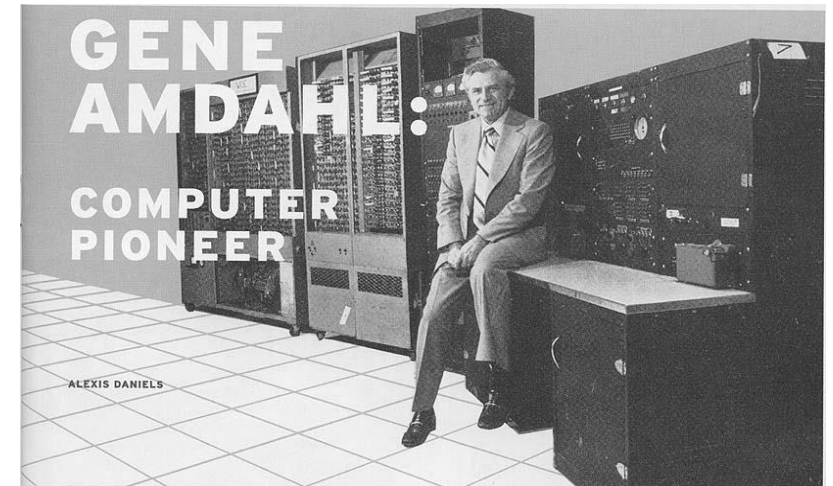
NVIDIA B200

Credit: Cerebras and Tesla



# Wafer-Scale is Not a New Dream

- Gene Amdahl shared a similar observation [1]
  - The pioneer of mainframe machines
  - The author of Amdahl's Law
- Amdahl co-founded Trilogy Systems
  - Attempted to design the first wafer-scale chips
  - The biggest investment (\$200M) in Silicon Valley in the 1980s
- Trilogy Systems failed due to
  - **Low yields at wafer-scale**
  - **Weak market demand**



[1] <https://spectrum.ieee.org/whats-better-than-40-gpubased-servers-a-server-with-40-gpus>

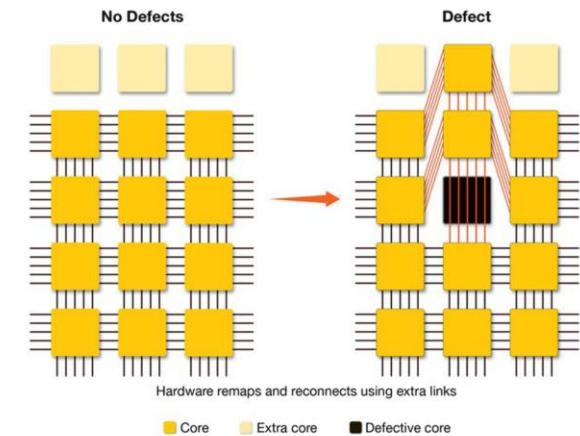


# It is about Time (40 Years Later)

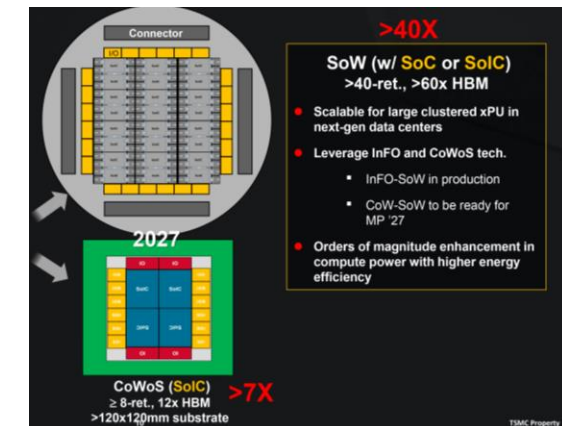
- AI chases extreme efficiency/performance
- Manufacturing improvement
  - Bypass the defective cores with redundant wires

**High yield: 93%** core active (WSE-3) vs **92%** (NV H100)

- A wave of wafer-scale computers is coming
  - **>40X** compute and bandwidth expected **by 2027** [2]
  - Advanced packaging (CoWoS), 3DIC (TSMC SoIC)



Example of hardware remapping [1]



TSMC Roadmap – System-on-Wafer[2]

[1] <https://www.cerebras.ai/blog/100x-defect-tolerance-how-cerebras-solved-the-yield-problem>

[2] <https://www.tomshardware.com/tech-industry/tsmc-to-go-3d-with-wafer-sized-processors-cow-sow-system-on-wafer-technology-allows-3d-stacking-for-the-worlds-largest-chips>



# Wafer-Scale Integration – Better Compute Scaling

	System-on-Die	System-on-Wafer
Area	Typically 858 mm <sup>2</sup>	Typically 73062 mm <sup>2</sup>
#Transistors (TSMC n3)	1 trillion	Up to 91 trillion
Interconnect	PCB/RDL/SUB/WoW	Wafer
Die-to-die efficiency	~10s pJ/bit	~0.1s pJ/bit
Die-to-die bandwidth	~ 1-10s TB/s	~ 10 - 100s TB/s
Memory Bandwidth	10s TB/s (crossbar)	10s PB/s (aggregated via mesh)
Off-chip memory	10s - 100s GB HBM	10s TB DRAM via Ethernet <b>10s TB HBM/DRAM via TSMC SoW in 2027</b>

~100x

~100x

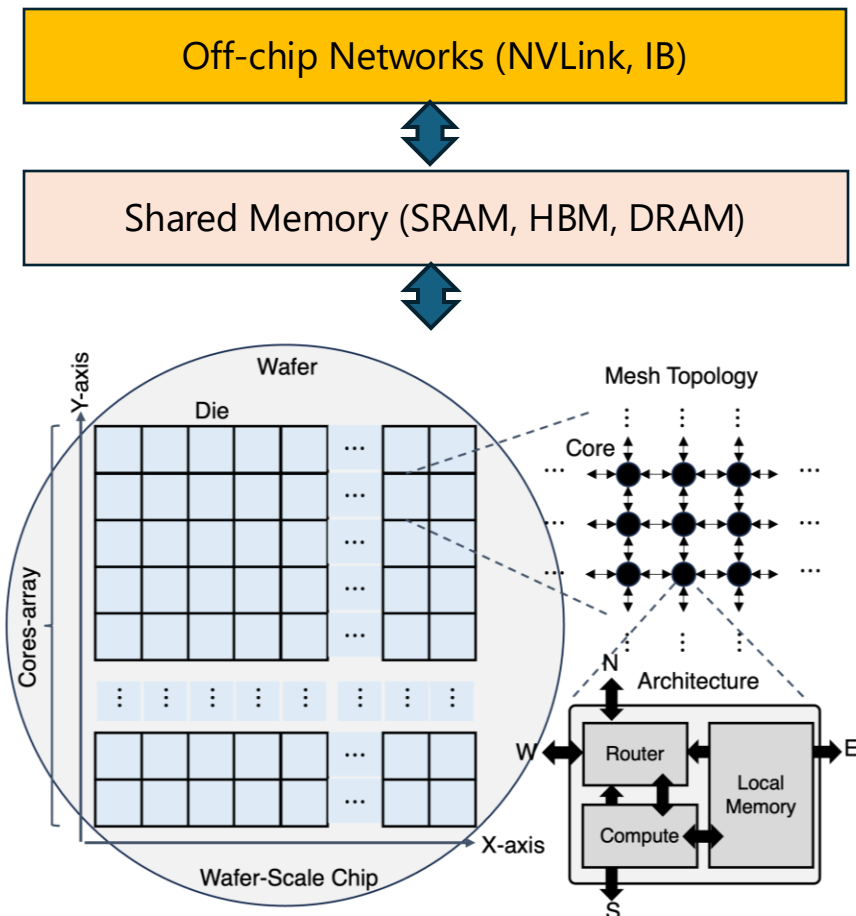
~1,000x

- **Emerging wafer-scale systems:** Cerebras, Tesla Dojo, NVIDIA and more reported by TSMC
- **Growing adoption:** Perplexity, Mixtral, Meta AI, G42, ...





# Systems Ready for Wafer-Scale Chips?



## Extensive research on scaling LLM with off-chip networks

- **Topology:** Clos, 3D-Torus
- **System:** Megatron-LM, PyTorch, JAX, TensorFlow, nnScaler
- **Multi-dimensional Parallelism:** TP, PP, DP, EP
- **Communication Operator:** Ring allreduce, All-to-All for MoE

## Extensive research on LLM with on-chip shared memory

- **Operator:** FlashAttention, MLA, PageAttention,
- **Compiler:** Ladder [OSDI'24], and T10 [SOSP'24]

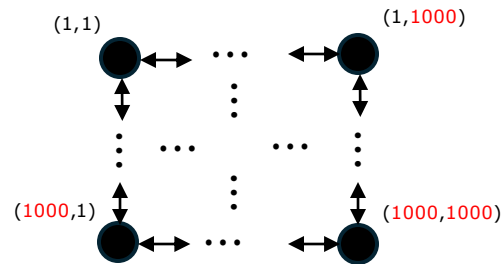
## Wafer-scale AI software remains largely unexplored

- Existing NoC research targets CPUs and small scale (up to 100s)
- Suffer severe communication bottlenecks
- ...



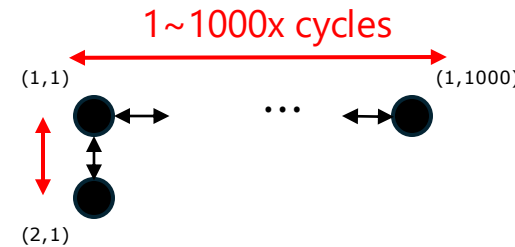
# PLMR – A Simplicity-Driven Model for Wafer-Scale Chips

## 1. Million-scale Parallelism (PLMR)



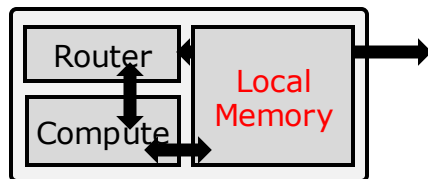
From hundreds of parallelism in a crossbar to millions of parallelism in a mesh

## 2. Highly non-uniform access Latency (PLMR)



From shared memory and small NUMA to large-scale non-uniform memory

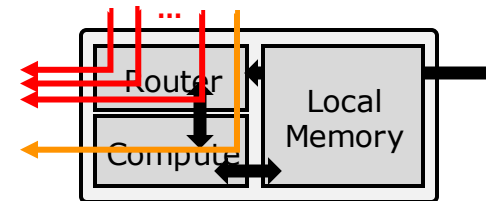
## 3. Constrained local Memory (PLMR)



From coarse-grained tile pipeline to fine-grained tile pipeline

100s KB – 1s MB

## 4. Constrained Routing resources (PLMR)



From centralized routing to decentralised NoC routing

Only support 10s routing entries

→ Routing on NoC

→ Routing on Compute Engine

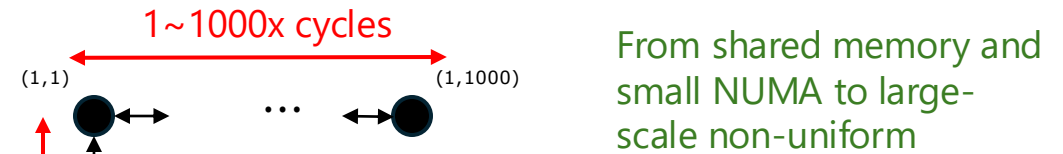


# PLMR – A Simplicity-Driven Model for Wafer-Scale Chips

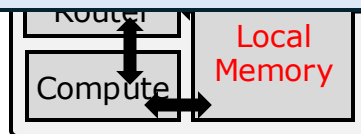
## 1. Million-scale Parallelism (PLMR)



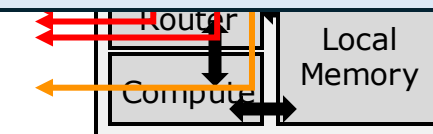
## 2. Highly non-uniform access Latency (PLMR)



**PLMR model** - The **key shift** from *shared-memory architectures* to *on-chip large-scale, distributed memory systems*



From coarse-grained tile pipeline to fine-grained tile pipeline

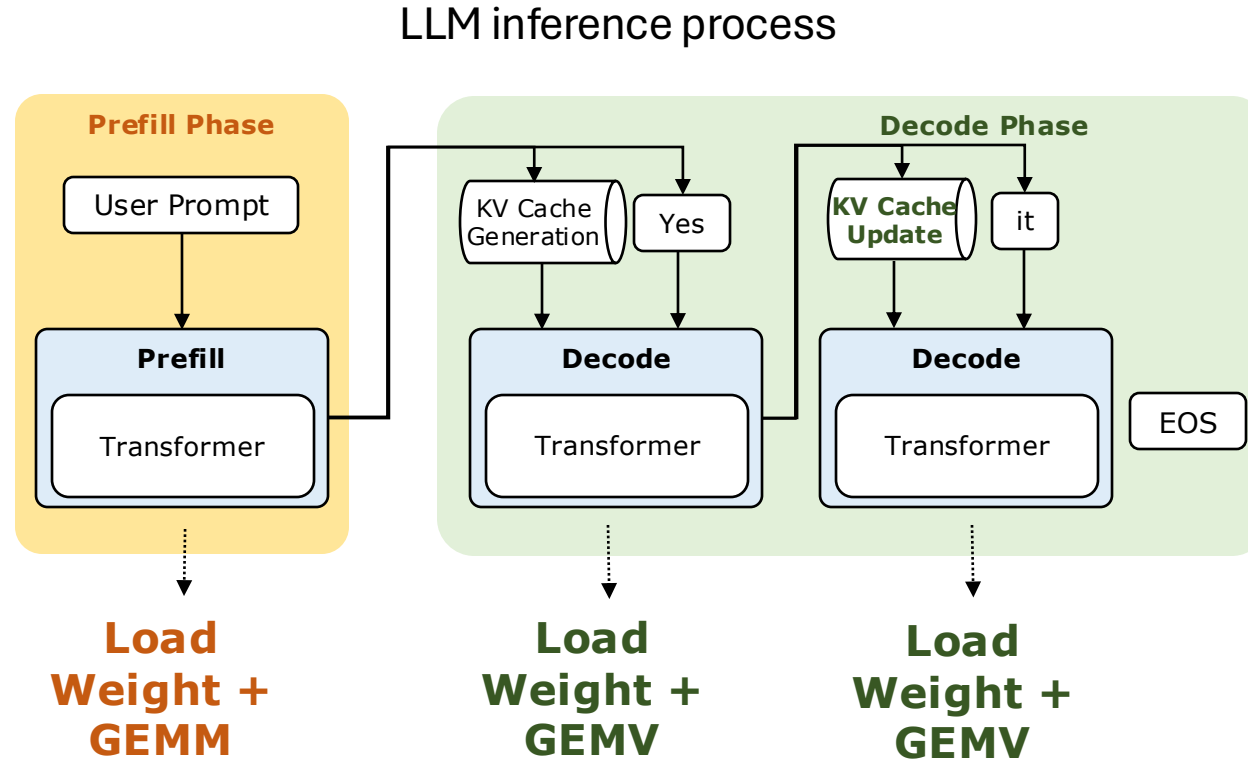


From centralized routing to decentralised NoC routing

→ Routing on NoC  
→ Routing on Compute Engine



# LLM Inference on Wafer-Scale Chips



**PLMR Compliant**



# WaferLLM: World-First Wafer-Scale LLM Inference System

## Goals

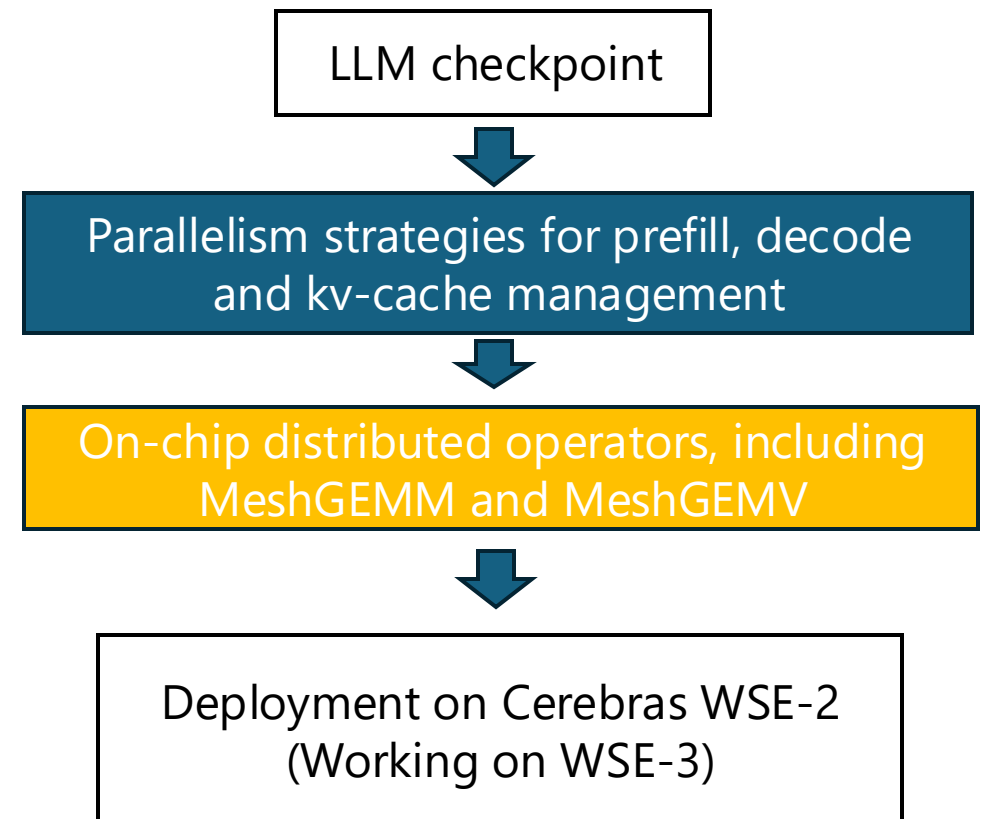
- **Design the entire stack guided by PLMR**
- Generalise across hardware backends

## Contributions

- New prefill parallelism strategies
- New decode parallelism strategies
- New KV-cache algorithm – Shift-based update
- New GEMM algorithm - MeshGEMM
- New GEMV algorithm - MeshGEMV

**First LLM inference system to reach **2700 token/s per request****

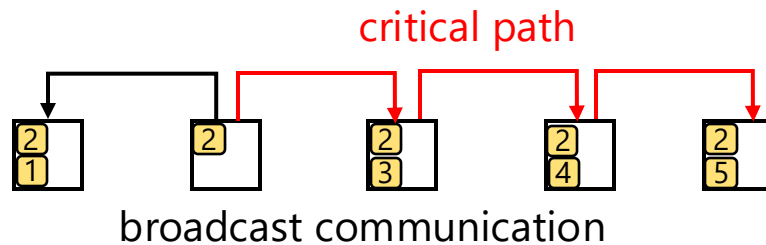
## Current WaferLLM Architecture



# PLMR Compliant MeshGEMM

Prefill is bottlenecked by GEMM, which requires each submatrix to traverse all row (or column) cores, constrained by properties L and R.

**SUMMA**  
(SOTA in All-to-All / Clos)



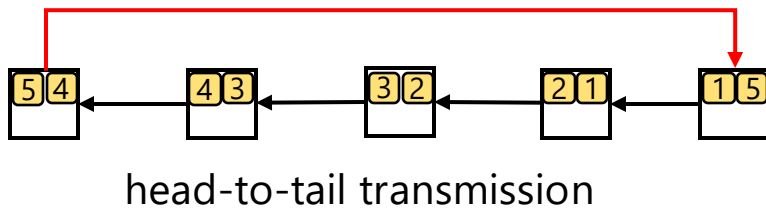
Latency (L)

$O(N)$

Routing (R)

$O(N)$

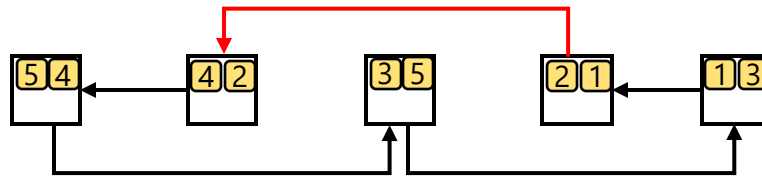
**Cannon**  
(SOTA in 2D Torus)



$O(N)$

$O(1)$

**MeshGEMM**  
(Ours)



$O(1)$

$O(1)$

Evaluation results:

1. **1.3–2×** faster GEMM with matrix sizes from 2K to 8K
2. Reduces communication overhead by **2–5×**.

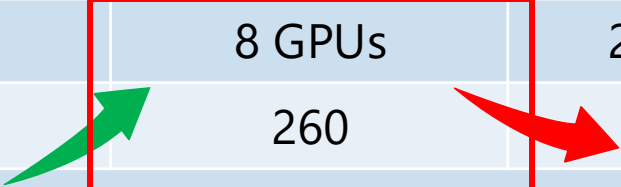




# Comparison with SOTA LLM Inference Systems

We compare **WaferLLM on real Cerebras WSE-2 chip (TSMC 7nm)** with **SGLang/vLLM on NVLink/IB-connected A100 GPU (TSMC 7nm)** in **performance** and **energy efficiency**

Decode (4K in, 4K out, BSZ=1)	LlaMA3-8B		
	1 GPU	8 GPUs	2x8 GPUs
SGLang (A100) Token/s per request	78	260	164
WaferLLM (WSE-2) Tokens/s per request		<b>2700</b>	
A100/WSE-2 Energy Ratio	0.92	2.22	7.02



- **6-20x** faster than SoTA off-chip solutions on LLM model size range from 8B to 70B
- **2-2.5x** more energy efficiency than GPU interconnect – currently the only one on the market beyond NVLink
- **Outperforms best-case off-chip scaling in both speed and efficiency**



# Summary

- To further scale compute for AI
  - Going after simplicity that enables true scalability
  - A new AI stack following the PLMR model to address the system challenges

