# Wafer-Scale AI: GPU Impossible Performance

Sean Lie, Co-founder and CTO, Cerebras Systems

Hot Chips 2024

**cerebras**

# Cerebras Systems

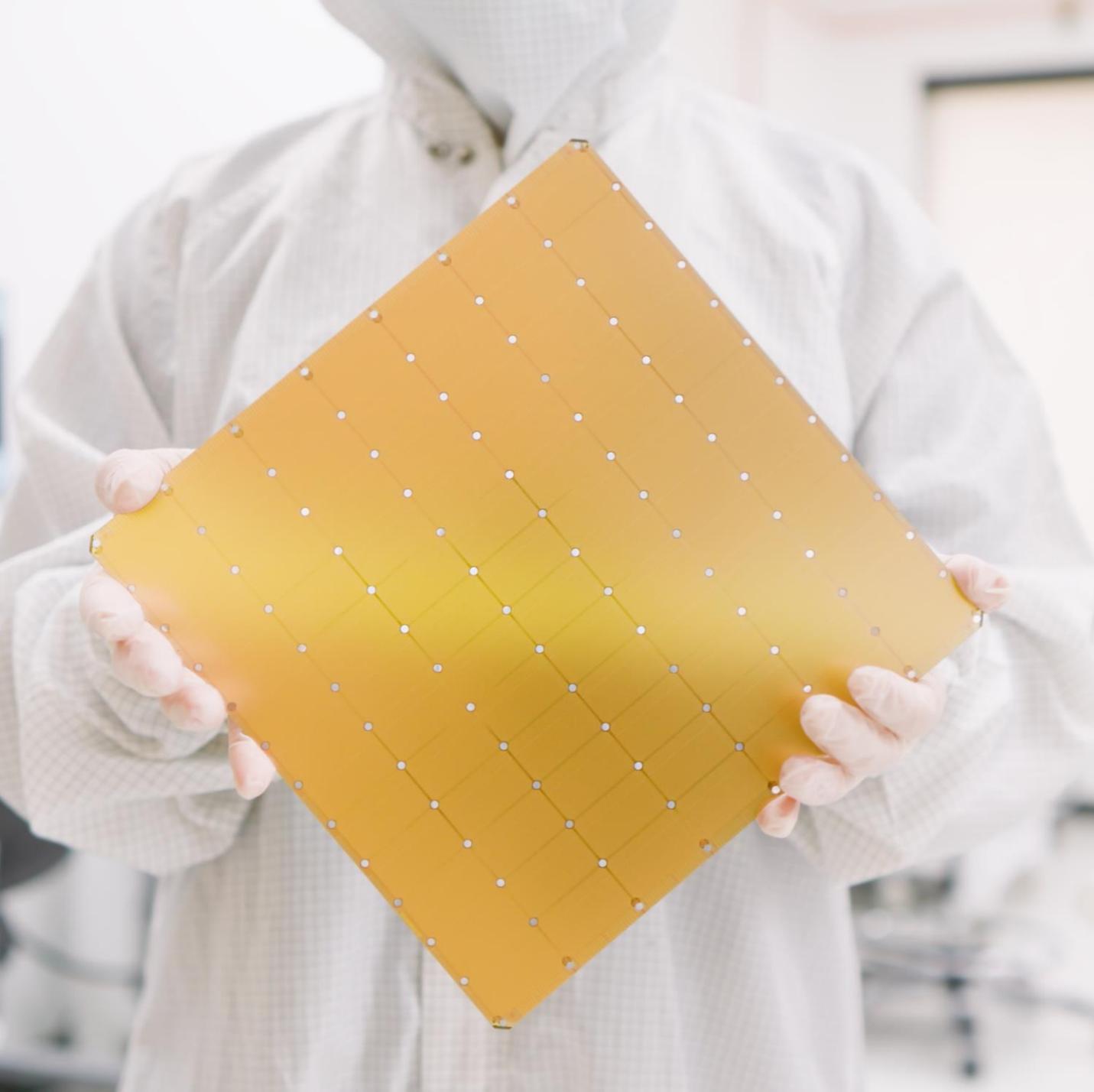**Founded in 2016**

**400 Employees**

**Offices**
Silicon Valley | San Diego | Toronto | Bangalore

**Customers**
North America | Asia | Europe

# Cerebras Wafer-Scale Engine

## The largest chip ever produced

**46,225 mm²** silicon

**4 trillion** transistors

**900,000** AI cores

**125 Petaflops** of AI compute
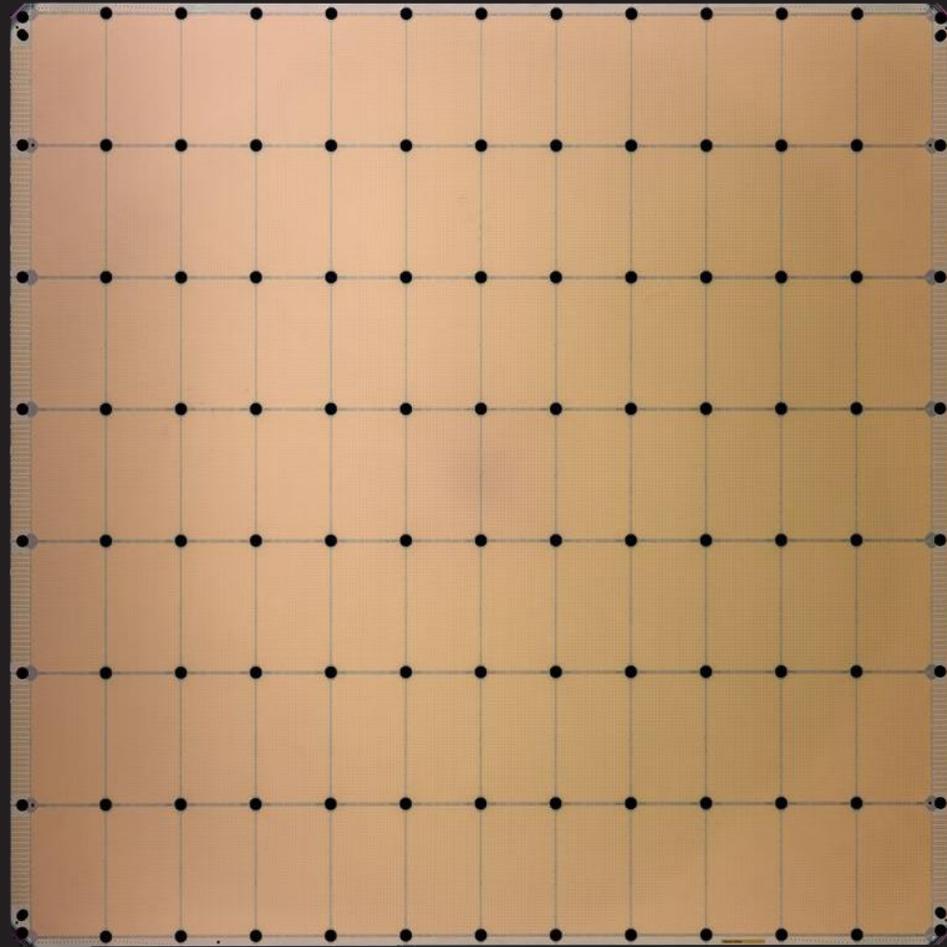
**44 Gigabytes** of on-chip memory

**21 PByte/s** memory bandwidth

**214 Pbit/s** fabric bandwidth

**5nm** TSMC process

# Cerebras Wafer-Scale Engine Versus the H100



**Cerebras WSE-3**
4 Trillion Transistors
46,225 mm$^2$ Silicon

**Largest GPU**
80 Billion Transistors
814 mm$^2$ Silicon

# Cerebras CS-3

**Condor Galaxy 1 - 4 Exaflops**
Santa Clara, California

**Condor Galaxy 2 - 4 Exaflops**
Stockton, California

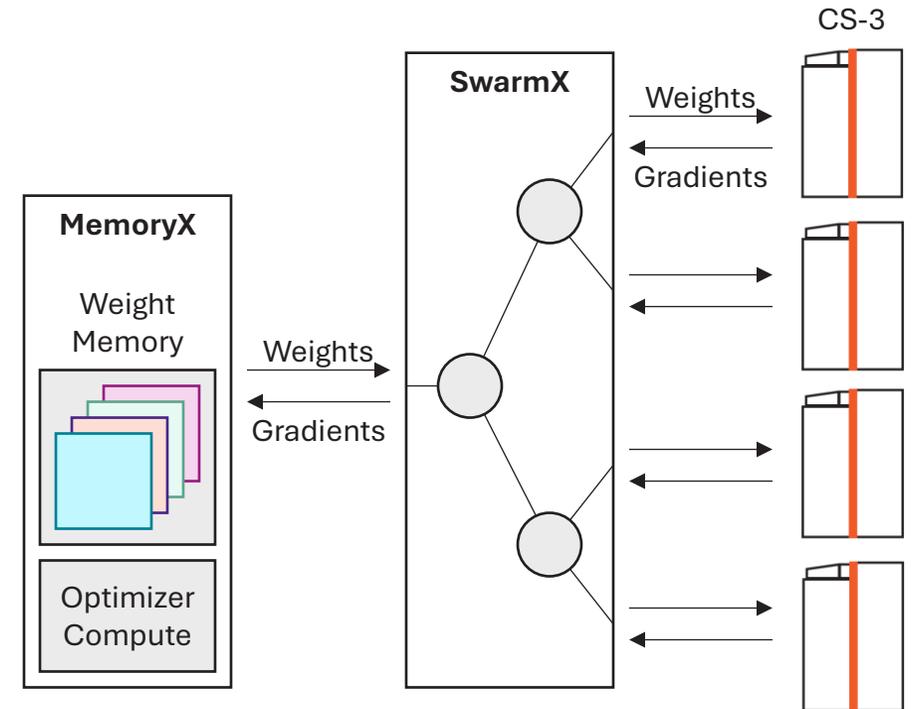**Condor Galaxy 3-5 - 20 Exaflops**
Dallas, Texas

**Condor Galaxy 6-9  -  32 Exaflops**
Minneapolis, MN

# Designed end-to-end for large-scale **Training**

**Co-designed cluster architecture to scale using data-parallelism only**

- WSE large enough to run even the largest models on a single chip

- Avoid hybrid model parallelism complexity

- MemoryX store streams weights to CS-3s

- SwarmX fabric performs broadcast/reduce

- Multi-system scaling with the same execution model as single system

**The only architecture with**

**Exaflop-scale training performance**

**But programs like a single device**

# Training SOTA large models everyday

Sample of open-source models trained on Cerebras

From multi-lingual LLMs to healthcare chatbots to code models



**BTLM-3B-8K**
3B PARAMETERS • 8K CONTEXT
7B Performance in a 3B Model
Open Source. Trained on Cerebras

**CrystalCoder**
7B PARAMETERS • 1.3T TOKENS
Coding + English. The most open source & reproducible model in the world.
Open Source. Trained on Cerebras

**cerebras**

**Jais**
13B & 30B
State of the art Arabic + English models
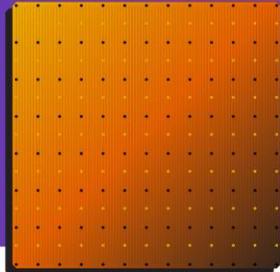Open Weights. Trained on Cerebras

**Med42**
FINED-TUNED LLAMA2-70B
Medical Q&A LLM Scores 72% on USMLE
Trained on Cerebras

**gigaGPT**
GPT-3 in 565 LINES OF CODE
Cerebras implementation of nanoGPT
Open Source. Trained on Cerebras

**SlimPajama**
627BTOKEN DATASET
Extensively deduplicated dataset with twice the perf/token
Open Source

**Cerebras-GPT**
111M–13B PARAMETERS
First family of GPT models released under Apache 2.0
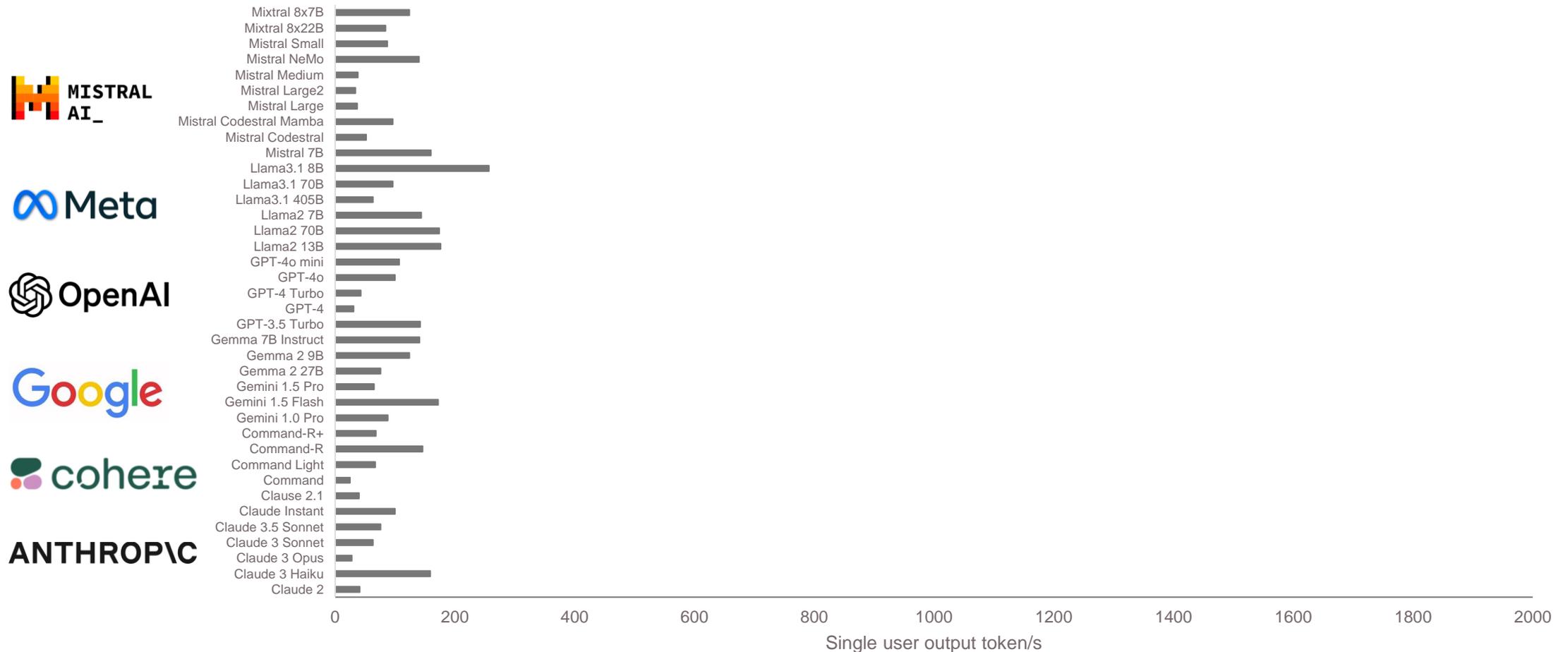Open Source. Trained on Cerebras

But also designed for **inference...**

# The Generative Inference Problem

# Generative inference today is *really slow*



Generative Inference Speed Landscape

Single user output token/s
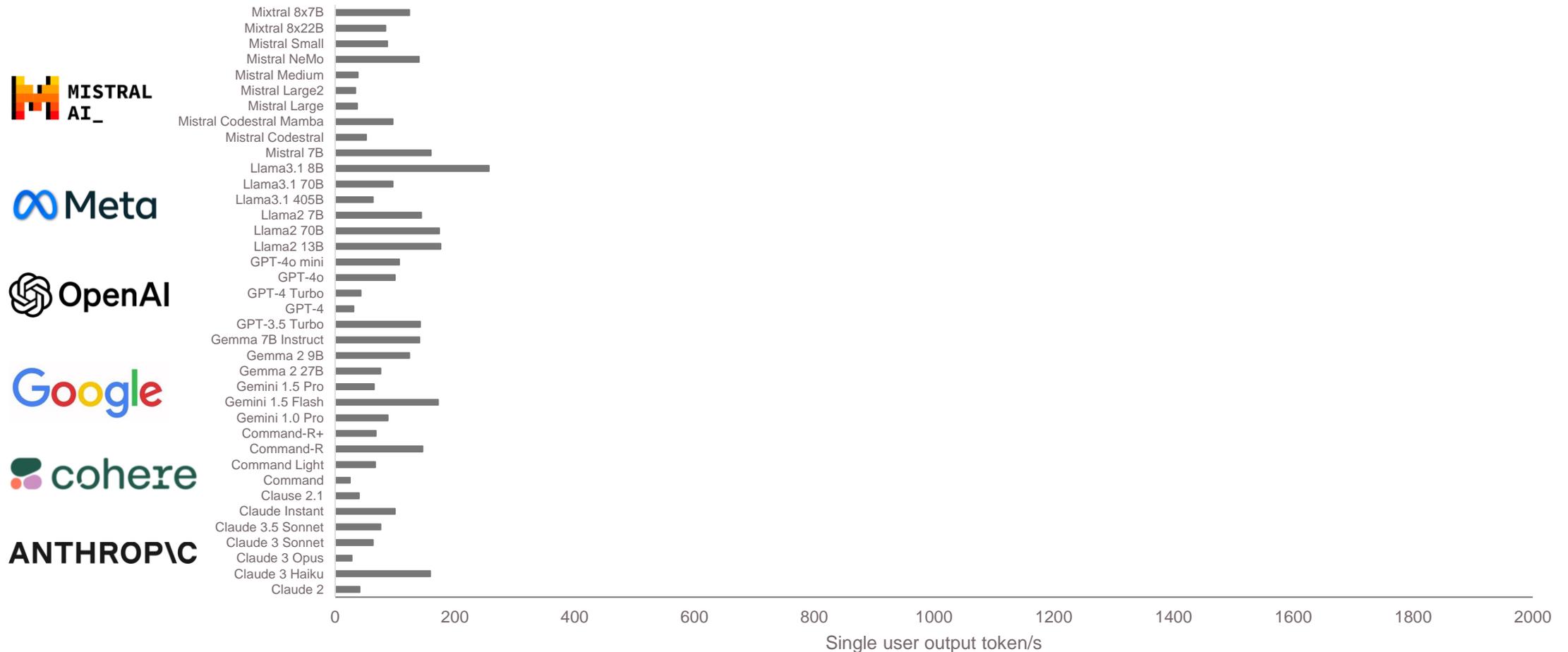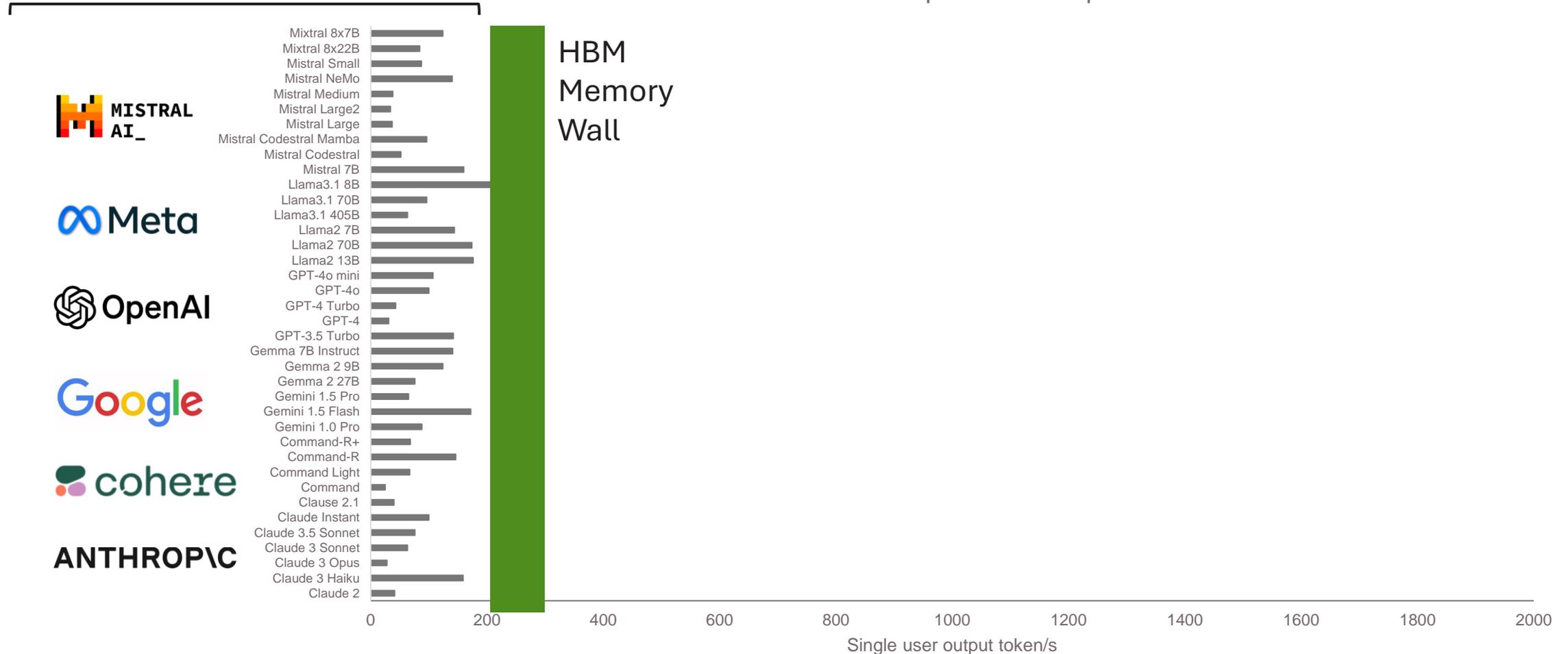
# Generative inference today is *really slow*

Different ML model architectures
Different hardware: H100, MI300, TPU, …
**Similar performance. Why?**
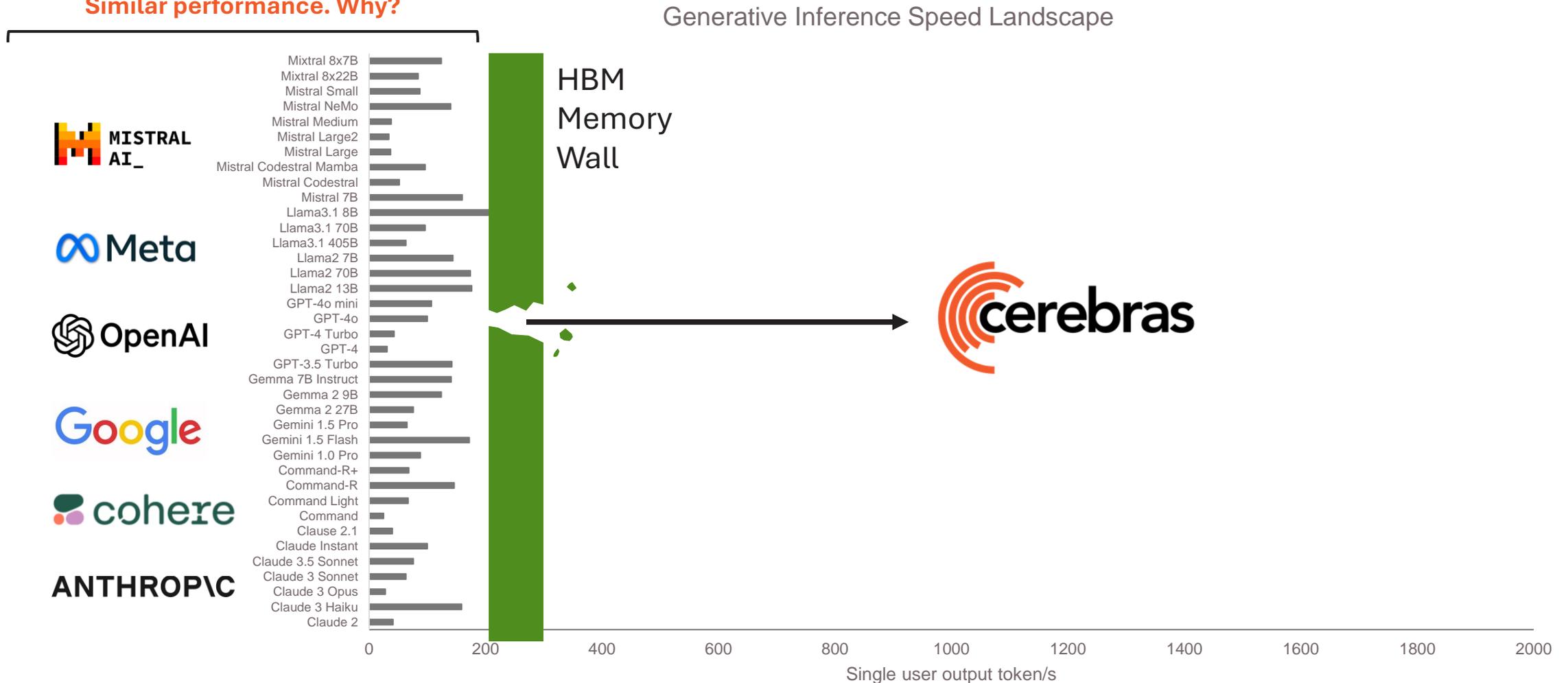
Generative Inference Speed Landscape



Single user output token/s

# Generative inference today is *really slow*

Different ML model architectures
Different hardware: H100, MI300, TPU, …
**Similar performance. Why?**

Generative Inference Speed Landscape



HBM
Memory
Wall

Single user output token/s

# Generative inference today is *really slow*

Different ML model architectures
Different hardware: H100, MI300, TPU, …
**Similar performance. Why?**

Generative Inference Speed Landscape



HBM
Memory
Wall

Single user output token/s

# Cerebras CS-3

# DGX-H100

## Llama3.1 8B

# The fastest inference on the planet on Llama3.1-8B
## 20x faster than hyperscale clouds



1,800

750

Single user
output token/s

72    79    93    164    165    225    257

Replicate    Azure    aws    OctoAI    perplexity    together.ai    Fireworks AI    groq    cerebras

What **20x speed** enables...

# Today GenAI applications are **promising but...**



**They are slow**

**Limited user engagement**

**Still primitive**

# 20x speed eliminates the wait
# Unlocks AI agent and copilot capabilities

**20x Speed enables**

- **20x more** interactive response for higher user engagement

- **20x more** model calls for chain-of-thought reasoning and more accurate responses

**Industry trying to move to *agentic* frameworks that need many LLM calls**

- **GraphRAG (MS Research) :** 5 calls per page, 100s per user request

- **GIST (DeepMind) :** 5 calls per page, several calls per user request

- **ReAct :** 2 calls per turn, 8 calls per user request

This will lead to the most powerful, most sophisticated, most engaging applications

The need for speed is even more evident with larger models

# The fastest inference on the planet on Llama3.1-70B

## 20x faster than hyperscale clouds
## 5x faster than the fastest DGX-H100 solution

Single user output token/s

| Azure | aws | perplexity | Fireworks AI | databricks | OctoAI | replicate | together.ai | groq | cerebras |
|-------|-----|------------|--------------|------------|--------|-----------|-------------|------|----------|
| 20 | 50 | 52 | 54 | 57 | 58 | 66 | 86 | 250 | 450 |

Source Artificial Analysis, Llama3 & 3.1 Results.

# The GPU Impossible

**5x faster** than the <u>fastest</u> GPU solution.

No number of GPUs can do this.

# Memory Bandwidth

# GenAI inference is a **memory bandwidth** problem

- Generating 1000 tokens takes 1000 serial passes through the model
- Each pass requires reading all model parameters from memory
- Low memory bandwidth is the bottleneck for generation performance



Time

# Why Wafer Scale Matters

| | Cerebras WSE-3 | Nvidia H100 | Cerebras Advantage |
|---|---|---|---|
| Chip size | 46,225 mm$^2$ | 814 mm$^2$ | **57x** |
| Cores | 900,000 | 16,896 FP32 + 528 Tensor | **52x** |
| On-chip memory | 44 Gigabytes | 0.05 Gigabytes | **880x** |
| Memory bandwidth | 21 Petabytes/s | 0.003 Petabytes/s | **7,000x** |

Completely removes the memory bandwidth bottleneck
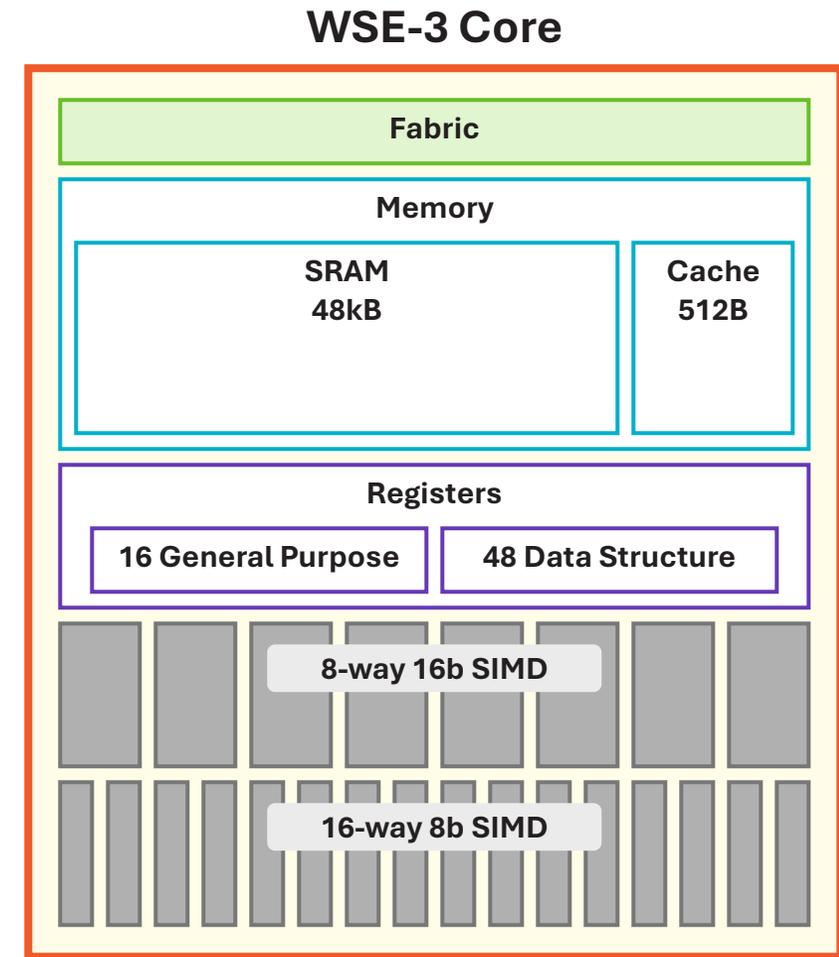
# WSE-3 Core

**Tightly coupled compute and memory**
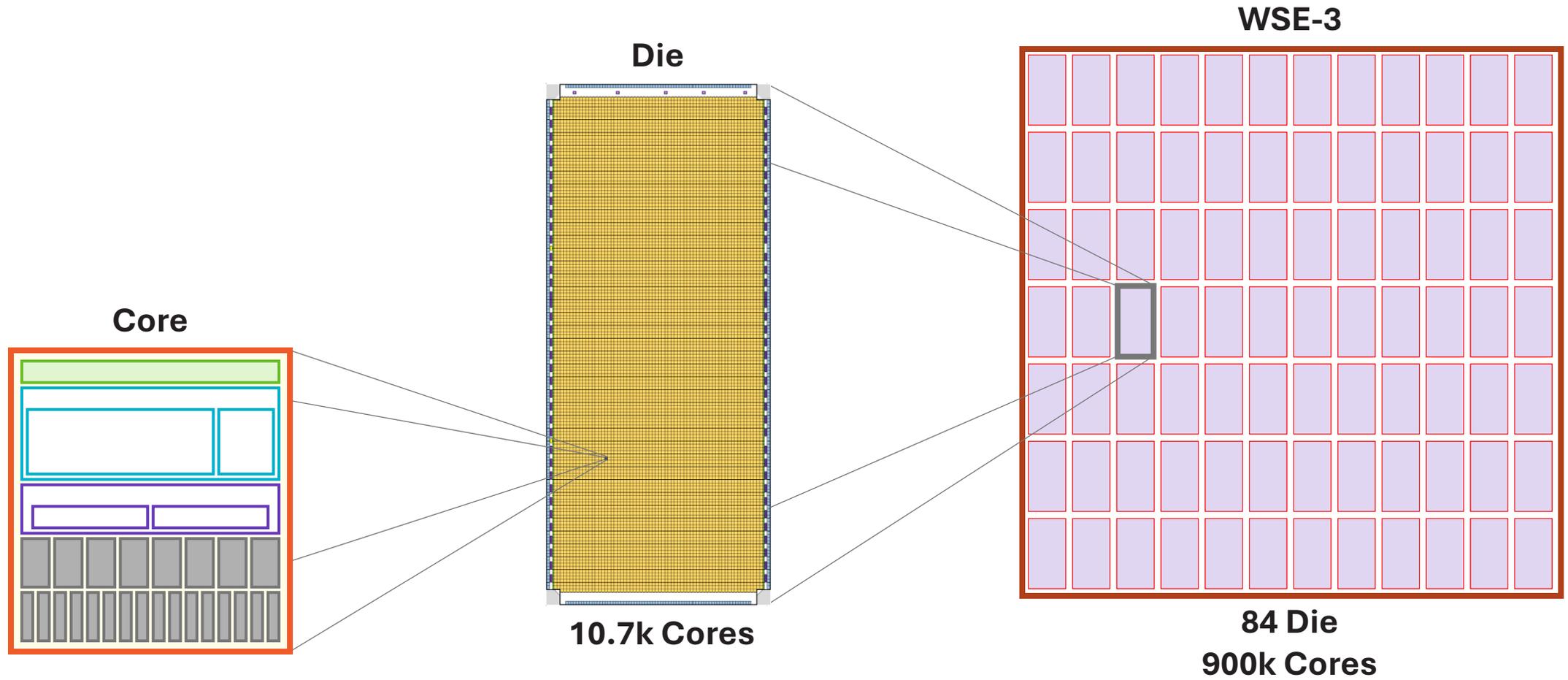
**High performance for AI compute**

- High performance tensor operations
  - 8-way SIMD for 16b data (FP/BF16)
  - 16-way SIMD for 8b data (Fixed/INT8)
- Fast non-linear functions instructions
- Fine-grained dataflow scheduling
- Native unstructured sparsity acceleration

**High bandwidth memory and cache**

- 48kB SRAM per core
- 512B local cache per core
- Full bandwidth for full SIMD performance

**WSE-3 Core**

| Fabric |
| --- |

| Memory | |
| --- | --- |
| SRAM 48kB | Cache 512B |

| Registers | |
| --- | --- |
| 16 General Purpose | 48 Data Structure |

8-way 16b SIMD

16-way 8b SIMD

# From Small Core to Massive Wafer

**Core**

**Die**

**WSE-3**

**10.7k Cores**
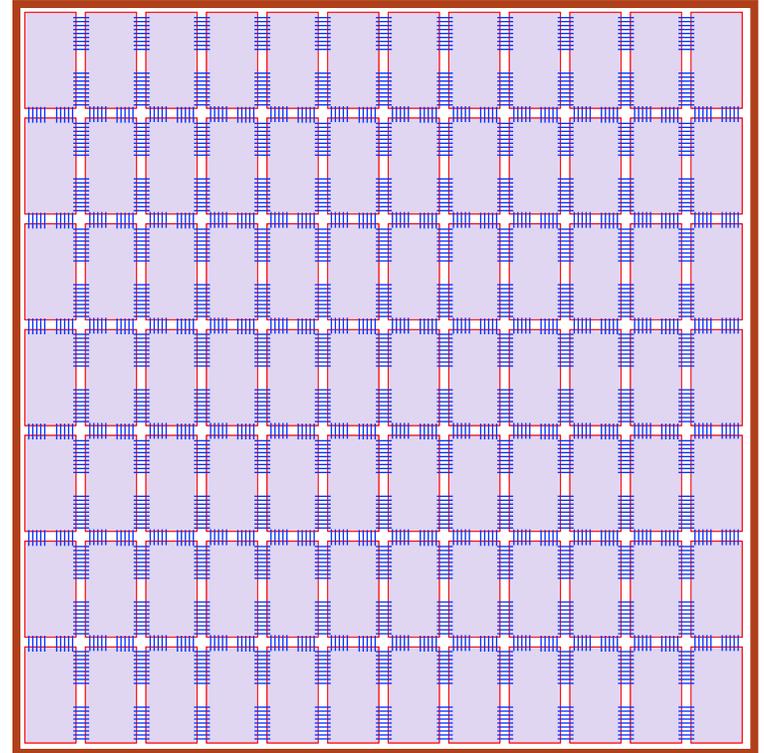
**84 Die**
**900k Cores**

# WSE-3 Interconnect

**Connection between dies crossing reticle boundaries**

- Invented process in first generation WSE

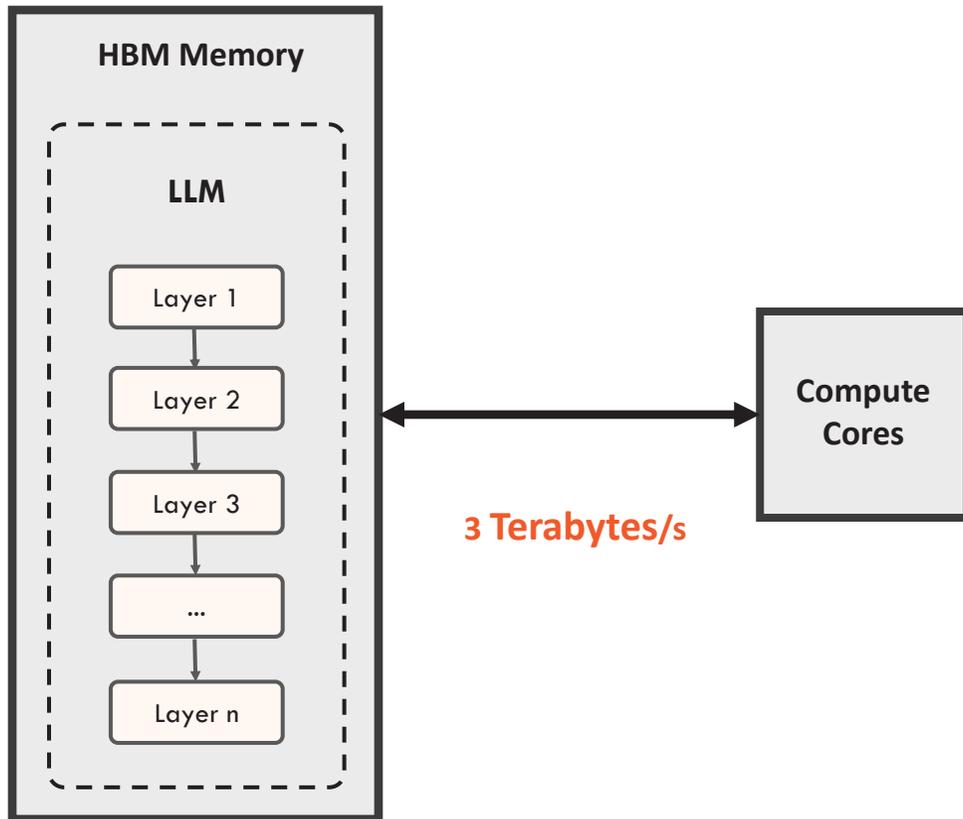- Extended to 5nm in collaboration with TSMC

**Co-designed with die-level fabric and system software**

- Each die has 2D mesh fabric connecting all cores

- Extend across die boundaries at full performance

- Uniform fabric at die level and wafer level

- Built-in redundancy to route around failures

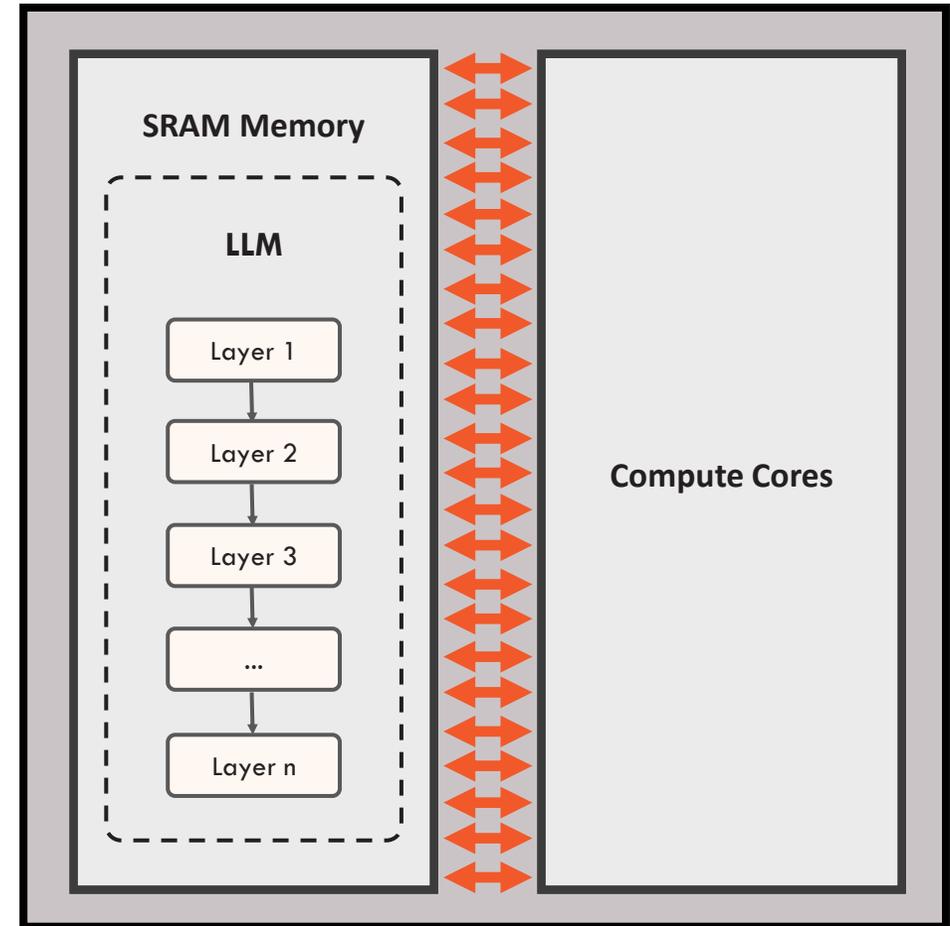- Software always sees a fully uniform 2D mesh

# Wafer Scale, SRAM based, Compute-in-Memory
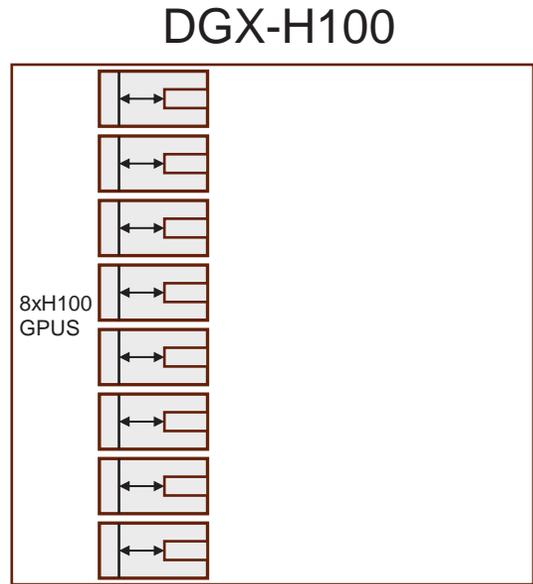# 7,000x more memory bandwidth than the GPU

**Nvidia H100**

**Cerebras Wafer Scale Engine 3**



**HBM Memory**

**LLM**

Layer 1

Layer 2

Layer 3

...

Layer n

**Compute Cores**

**3 Terabytes/s**

**SRAM Memory**

**LLM**

Layer 1

Layer 2

Layer 3

...

Layer n

**Compute Cores**

**21 Petabytes/s**

# Multi-GPU memory bandwidth aggregation
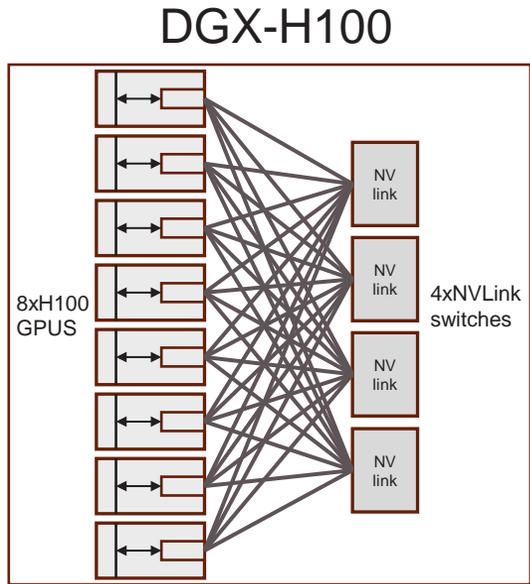
DGX-H100



**To achieve higher aggregate memory bandwidth...**

- Integrate 8 GPUs into a single server

- Accessed with tensor parallel execution

- 8x higher aggregate memory bandwidth

|  | Each H100 | 8xH100 |
|---|---|---|
| Mem Bandwidth | 3.35 TB/s | 26.8 TB/s |

# Multi-GPU memory bandwidth aggregation

### DGX-H100



| | Each H100 | 8xH100 |
|---|---|---|
| Mem Bandwidth | 3.35 TB/s | 26.8 TB/s |
| IO Bandwidth | 9000GB/s<br>36x 100Gb/s serial | 7.2 TB/s<br>288x 100Gb/s serial |
| Power | 36W<br>(5.0 pJ/bit*2) | 288W IO + 300 Switch<br>588W Total |

*GPU estimate based on 5nm 100G serdes power with Nvidia H100 NVLink bandwidth

**To achieve higher aggregate memory bandwidth…**

- Integrate 8 GPUs into a single server
- Accessed with tensor parallel execution
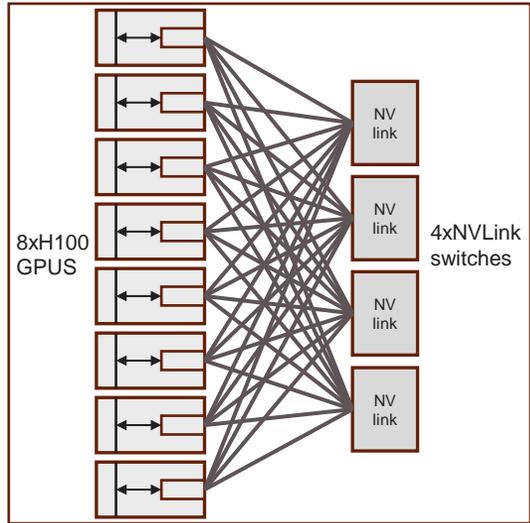- 8x higher aggregate memory bandwidth

**But it comes at a cost**

- Hundreds of high speed serial links
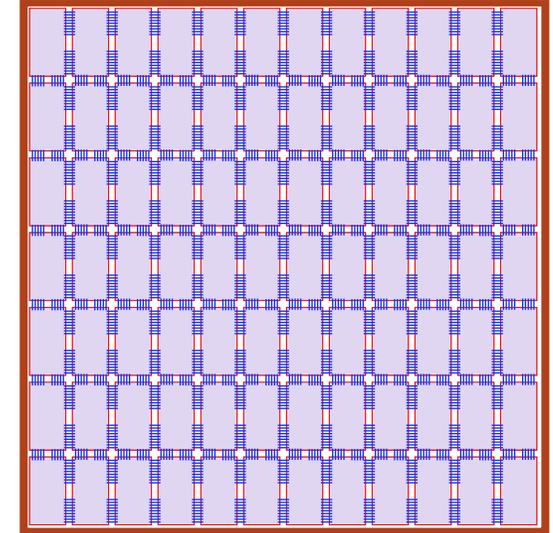- Multiple interconnect switch chips
- High cost
- High power – 0.5 kW

# Multi-GPU vs. Wafer-Scale Integration

DGX-H100



Wafer Scale Engine



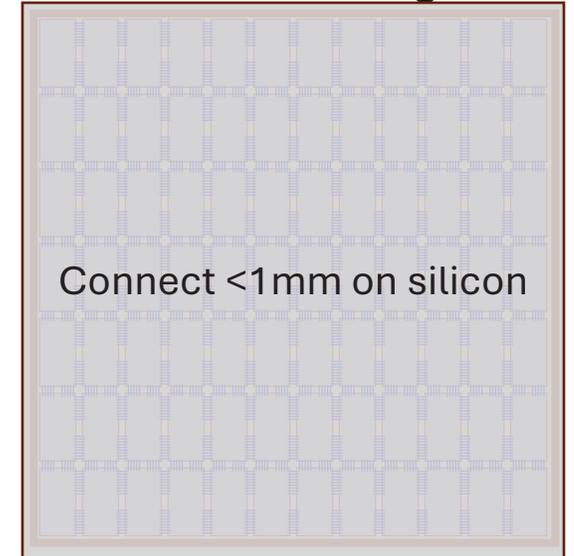| | Each H100 | 8xH100 | | Each Die | 84xDie |
|---|---|---|---|---|---|
| Mem Bandwidth | 3.35 TB/s | 26.8 TB/s | **800x more Mem bandwidth** | 255 TB/s | 21 PB/s |
| IO Bandwidth | 9000GB/s 36x 100Gb/s serial | 7.2 TB/s 288x 100Gb/s serial | **33x more Inter-Die IO** | 2880GB/s 480x 24Gb/s parallel | 242TB/s 40320x 24Gb/s parallel |
| Power | 36W (5.0 pJ/bit*2) | 288W IO + 300 Switch 588W Total | **6x Lower Power** | 1.1W (0.05 pJ/bit*2) | 97W |

*GPU estimate based on 5nm 100G serdes power with Nvidia H100 NVLink bandwidth

# Multi-GPU vs. Wafer-Scale Integration

DGX-H100

Hundreds of serdes serial links, multiple connectors, PCBs, switch chips

Wafer Scale Engine

Connect <1mm on silicon

|  | Each H100 | 8xH100 |  | Each Die | 84xDie |
|---|---|---|---|---|---|
| Mem Bandwidth | 3.35 TB/s | 26.8 TB/s | **800x more Mem bandwidth** | 255 TB/s | 21 PB/s |
| IO Bandwidth | 9000GB/s 36x 100Gb/s serial | 7.2 TB/s 288x 100Gb/s serial | **33x more Inter-Die IO** | 2880GB/s 480x 24Gb/s parallel | 242TB/s 40320x 24Gb/s parallel |
| Power | 36W (5.0 pJ/bit*2) | 288W IO + 300 Switch 588W Total | **6x Lower Power** | 1.1W (0.05 pJ/bit*2) | 97W |

*GPU estimate based on 5nm 100G serdes power with Nvidia H100 NVLink bandwidth

# Multi-GPU memory bandwidth does not scale well
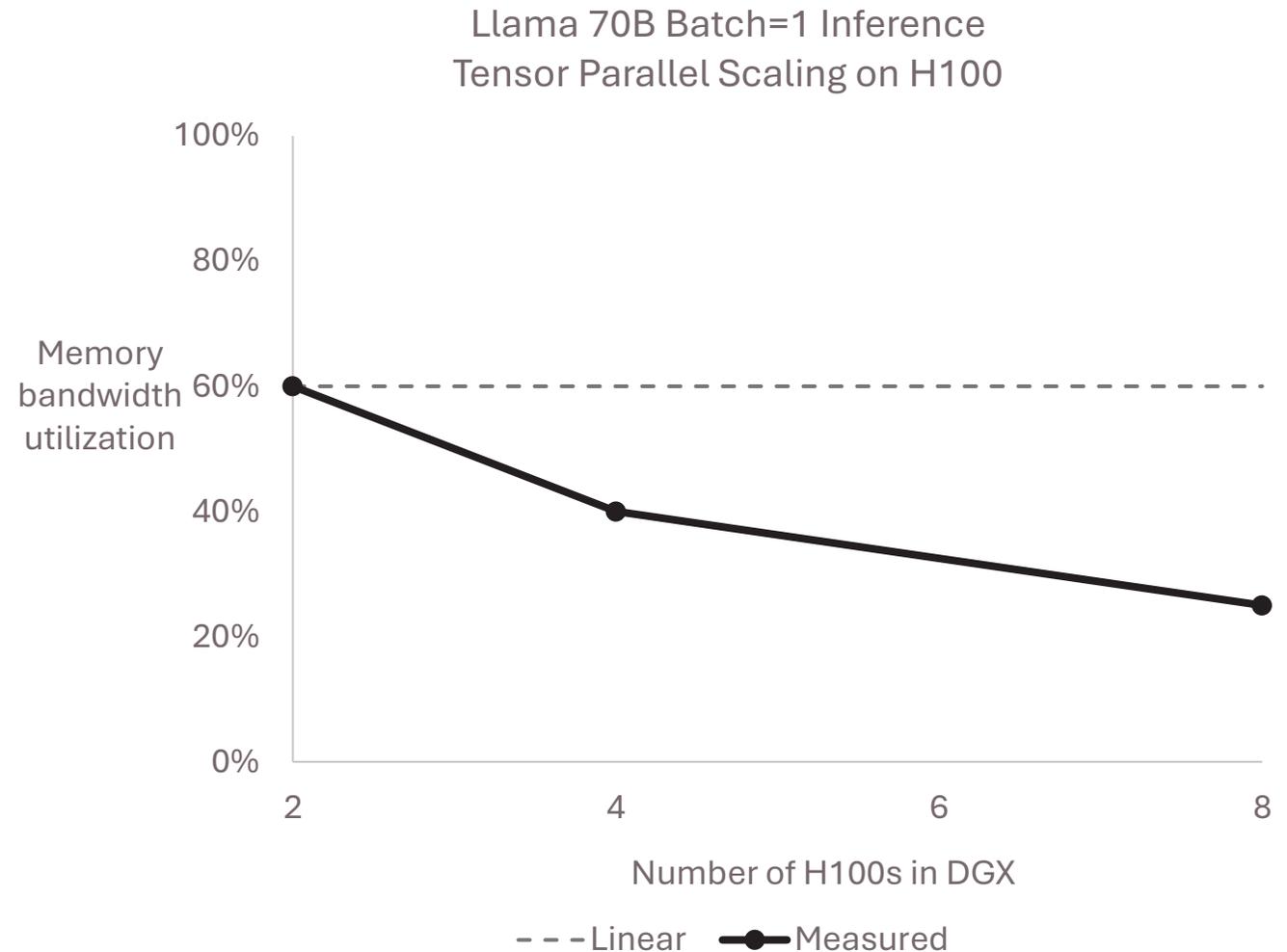
**Higher memory bandwidth on paper only**

- In reality, tensor parallel scales poorly due to interconnect bandwidth and latency overhead

- Measured memory bandwidth utilization drops from 60% on 2 H100s to 25% on 8 H100s

- That's even on the highest performance interconnect within a DGX
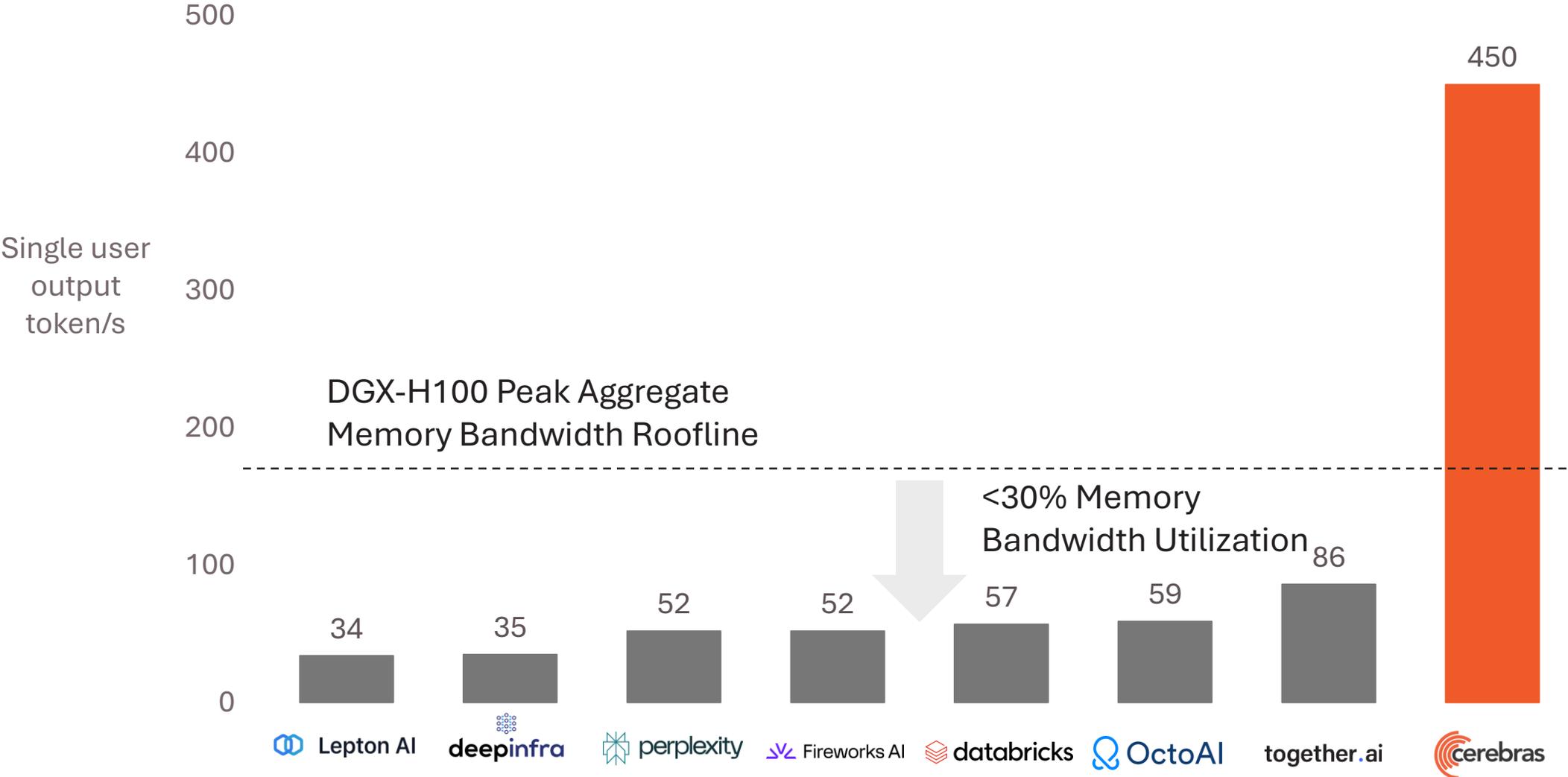
**Multi-GPU architecture has poor scaling**

- 4x theoretical higher memory bandwidth

- 1.7x performance only

- 42% scaling efficiency

**This is why GPU cannot scale inference generation performance beyond DGX**

**Lower IO limits performance scaling**

Llama 70B Batch=1 Inference
Tensor Parallel Scaling on H100



Number of H100s in DGX

- - - Linear   —●— Measured

# Real world low memory bandwidth utilization on Llama3.1-70B

Single user output token/s

500

450

400

300

DGX-H100 Peak Aggregate
Memory Bandwidth Roofline

200

<30% Memory
Bandwidth Utilization    86

100

34          35          52          52          57          59

0

Lepton AI    deepinfra    perplexity    Fireworks AI    databricks    OctoAI    together.ai    cerebras

Lowest Latency Generation

# Pipeline execution on a single chip

**Massive memory bandwidth enables opposite execution model**

- GPU: multiple chips to run a single layer

- Cerebras: fraction of a chip to run a single layer

# Pipeline execution on a single chip

**Massive memory bandwidth enables opposite execution model**

- GPU: multiple chips to run a single layer

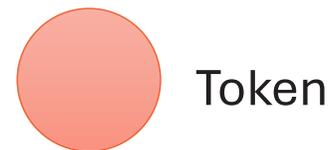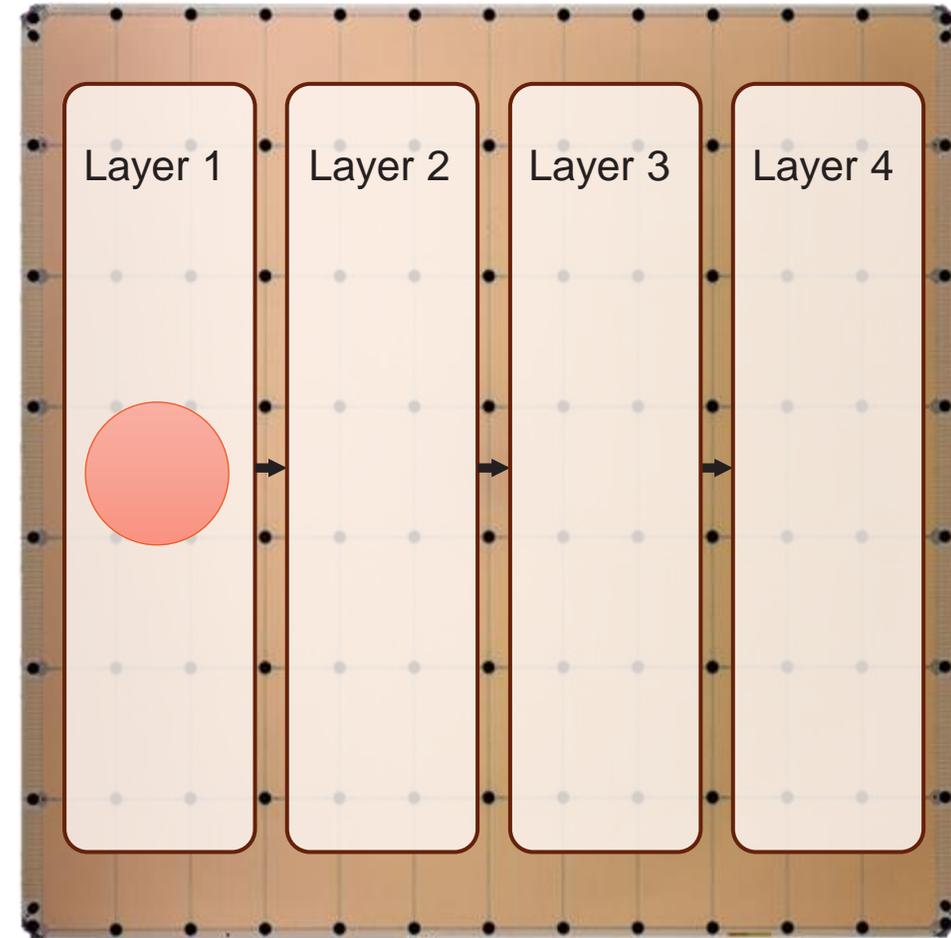- Cerebras: fraction of a chip to run a single layer

**Pipeline mapping of model**

- Model layers are mapped to wafer regions

- Size of wafer region determined based on memory and compute requirements

- Model weights and KV cache stored locally in the region memory close to the compute
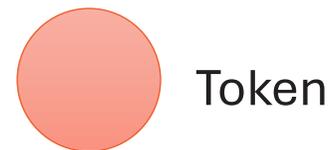
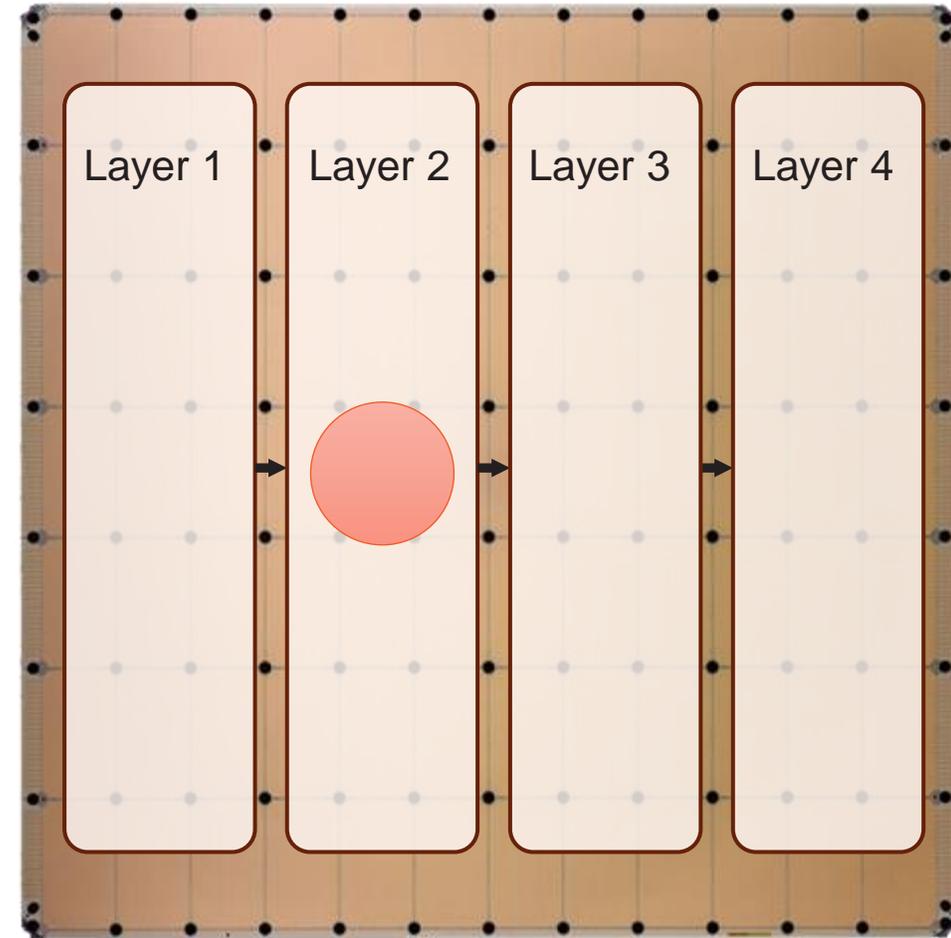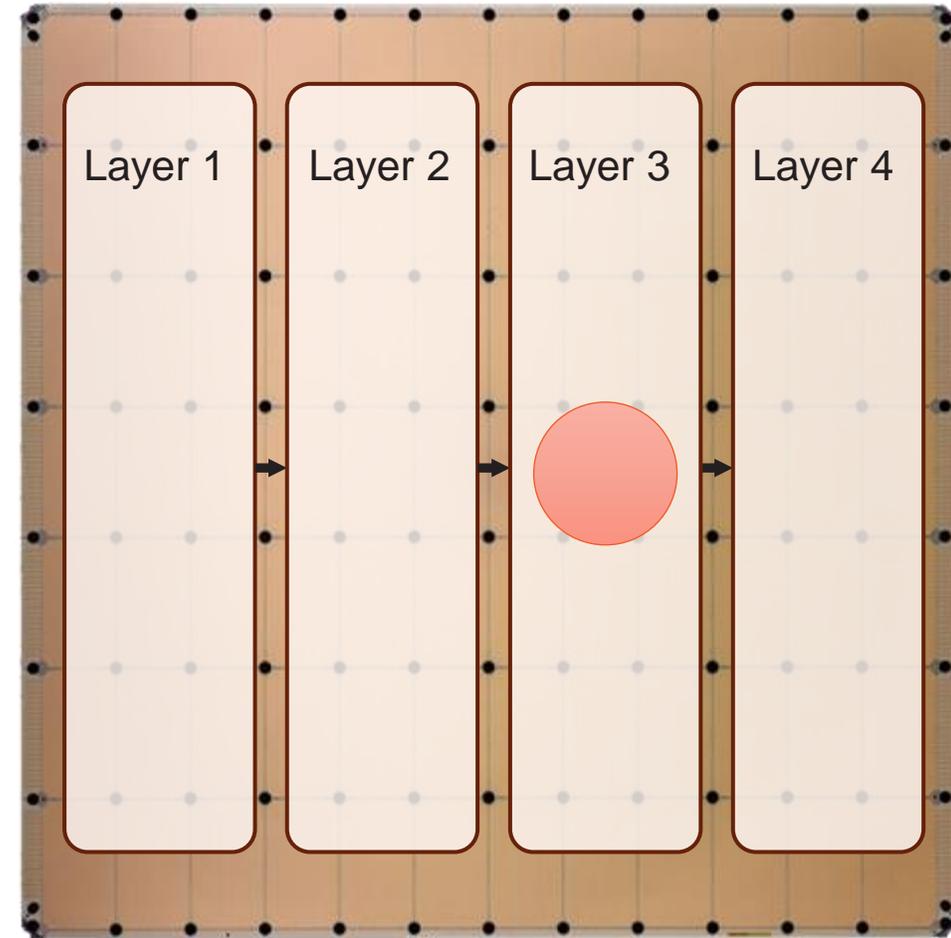# Pipeline execution

## A low latency pipeline

- Each wafer region processes 1 token
  - Enough memory bandwidth to run local batch size 1
  - Memory to feed compute datapath at full speed
  - Enabling optimized performance for matrix*vector
  - Local fabric interconnect to maintain low latency



Token

# Pipeline execution
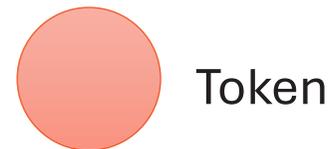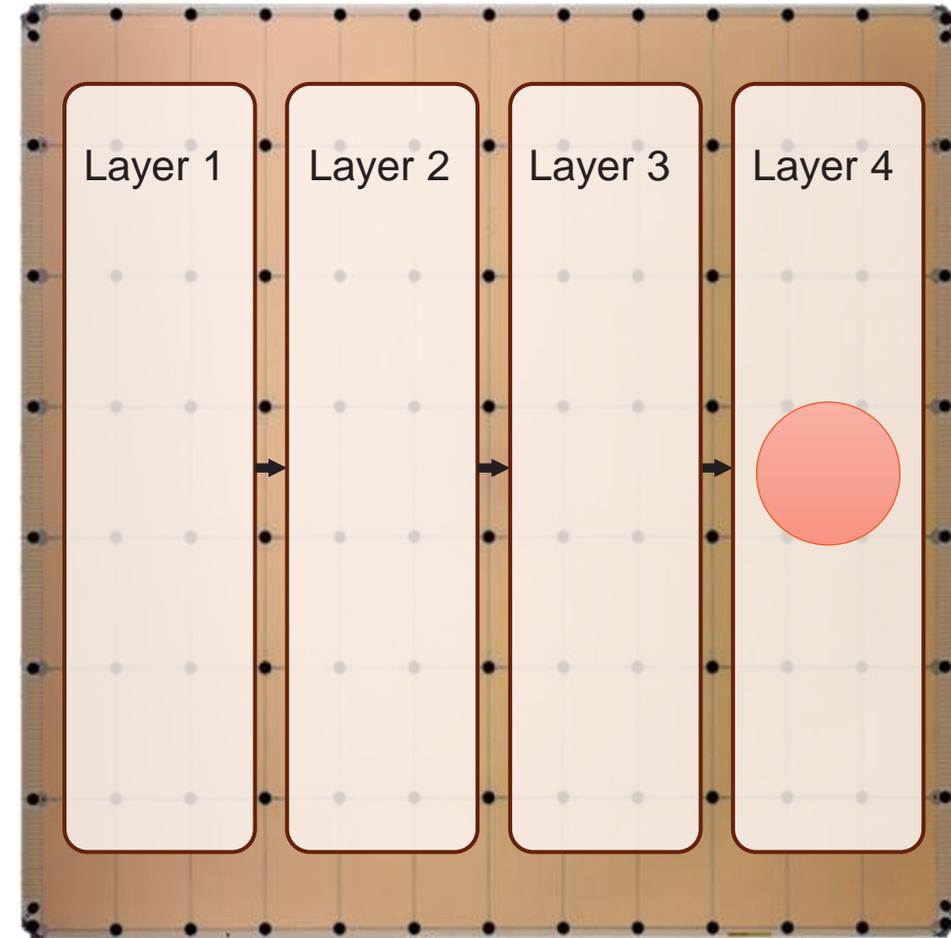
## A low latency pipeline

- Each wafer region processes 1 token
  - Enough memory bandwidth to run local batch size 1
  - Memory to feed compute datapath at full speed
  - Enabling optimized performance for matrix*vector
  - Local fabric interconnect to maintain low latency

- Regions are placed to reduce latency
  - The next region is physically adjacent
  - Virtually no latency between pipe stages
  - Possible because it's all done on-chip



Token

# Pipeline execution

## A low latency pipeline

- Each wafer region processes 1 token
  - Enough memory bandwidth to run local batch size 1
  - Memory to feed compute datapath at full speed
  - Enabling optimized performance for matrix*vector
  - Local fabric interconnect to maintain low latency

- Regions are placed to reduce latency
  - The next region is physically adjacent
  - Virtually no latency between pipe stages
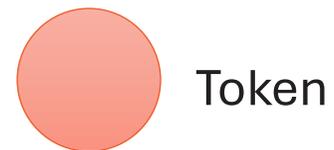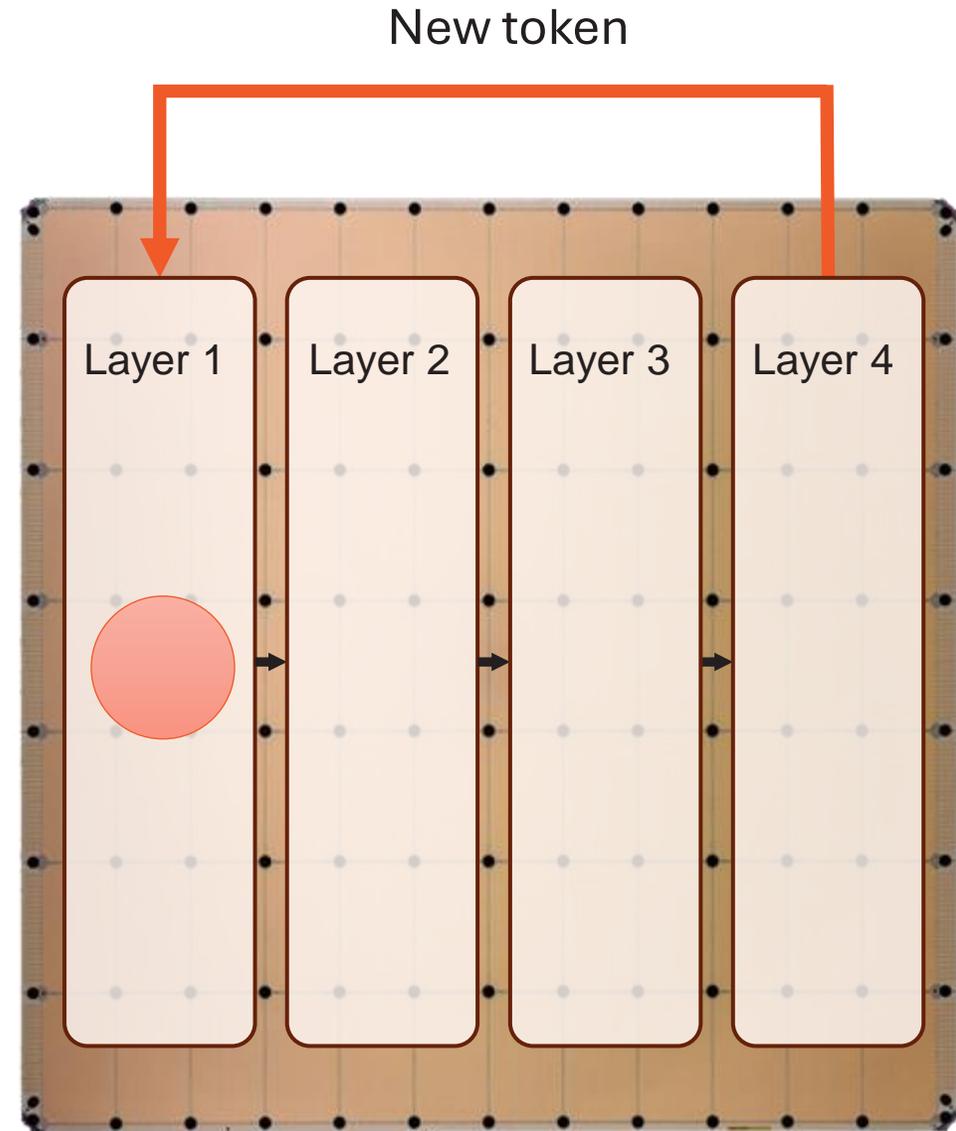  - Possible because it's all done on-chip



Layer 1    Layer 2    Layer 3    Layer 4

Token

# Pipeline execution

**A low latency pipeline**

- Each wafer region processes 1 token
    - Enough memory bandwidth to run local batch size 1
    - Memory to feed compute datapath at full speed
    - Enabling optimized performance for matrix*vector
    - Local fabric interconnect to maintain low latency

- Regions are placed to reduce latency
    - The next region is physically adjacent
    - Virtually no latency between pipe stages
    - Possible because it's all done on-chip



Layer 1  Layer 2  Layer 3  Layer 4

Token

# Pipeline execution

**A low latency pipeline**

- Each wafer region processes 1 token
  - Enough memory bandwidth to run local batch size 1
  - Memory to feed compute datapath at full speed
  - Enabling optimized performance for matrix*vector
  - Local fabric interconnect to maintain low latency

- Regions are placed to reduce latency
  - The next region is physically adjacent
  - Virtually no latency between pipe stages
  - Possible because it's all done on-chip

- Last region output cycled back to generate next token

**Pipelined execution enables super-fast token generation**



New token

Layer 1 | Layer 2 | Layer 3 | Layer 4
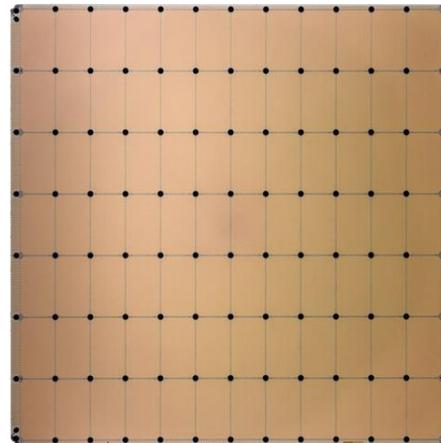
Token

# Scalability

# Large models fit on our large on-chip memory



**Llama3.1-8B**
8 billion parameters
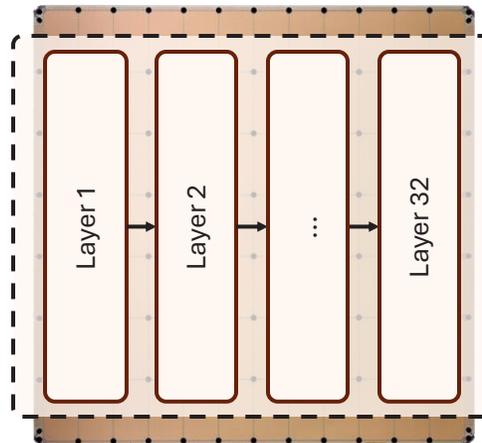16 GB of memory (FP16)

**WSE-3**
44GB of SRAM

# Large models fit on our large on-chip memory

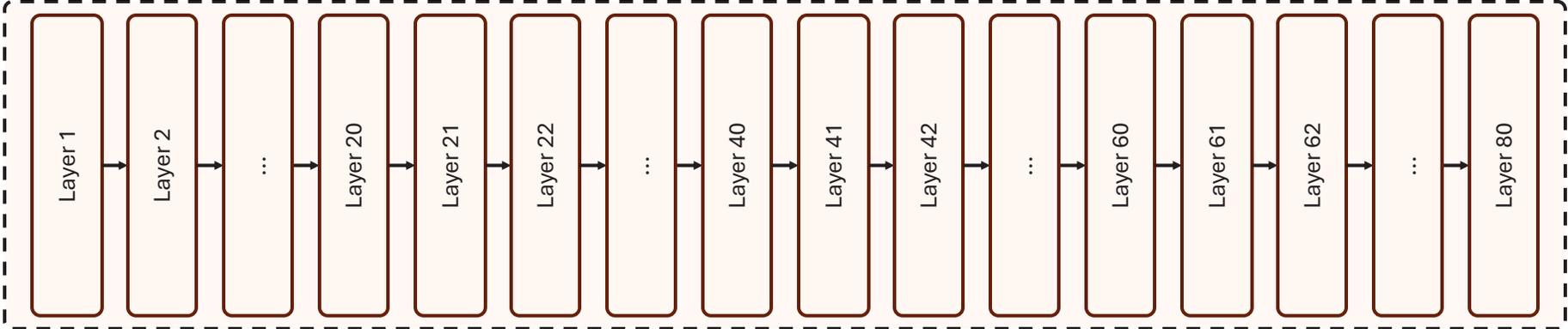- Models that fit entirely on the wafer are mapped directly



**Llama3.1-8B**
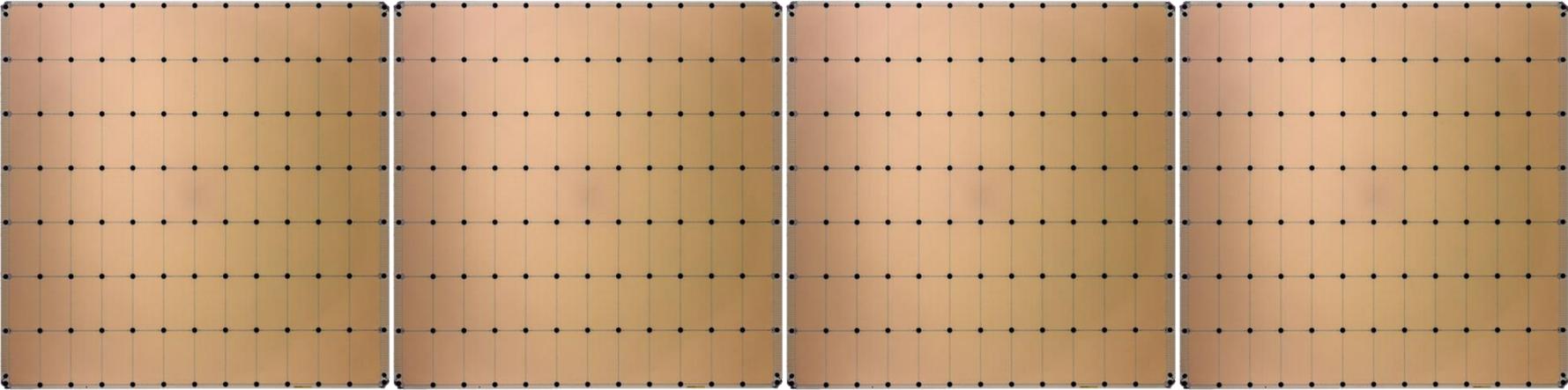8 billion parameters
16 GB of memory (FP16)

**WSE-3**
44GB of SRAM

# Naturally scales to multiple wafers



**Llama3.1-70B**
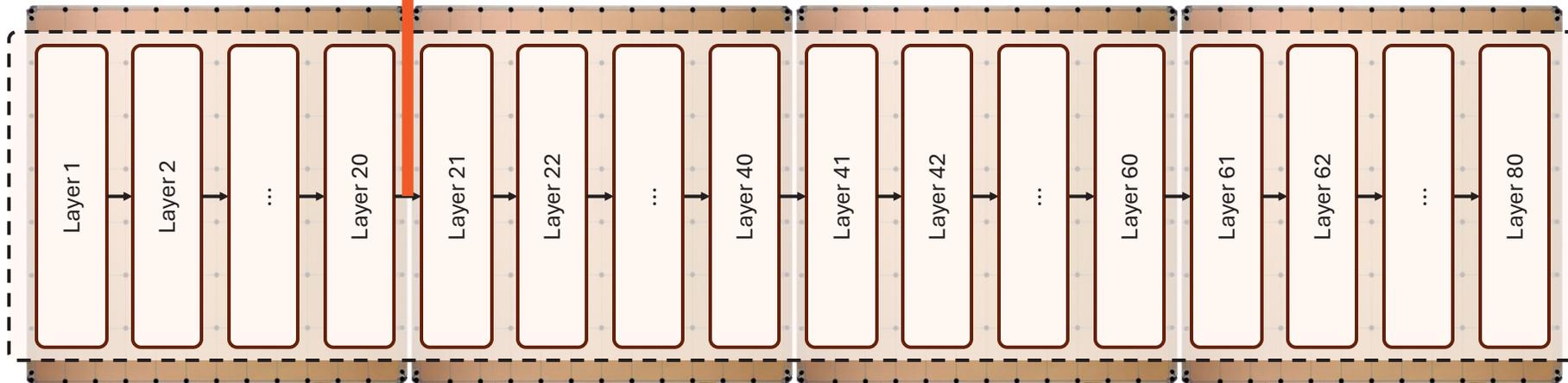70 billion parameters
140 GB of memory (FP16)

**4x WSE-3**
176GB of SRAM

# Naturally scales to multiple wafers

- Models that require more memory are mapped to multiple wafers
- Keep almost all high communication on wafer, using internal high bandwidth fabric
- Transfer only activations between wafers, requiring relatively lower bandwidth
- Using CS-3 low latency RDMA over Ethernet interconnect between systems

**Latency:**
CS-3 IO is <5us
4 hops required
<1% impact to latency

**Bandwidth:**
CS-3 IO is 1.2Tbps
100Gbps required
<10% of available



**Llama3.1-70B**
70 billion parameters
140 GB of memory (FP16)

**4x WSE-3**
176GB of SRAM

# High Throughput

# The nasty GPU latency vs. throughput tradeoff

**GPUs are designed for high throughput**

- High throughput requires high batch size

- Batch size = concurrent users

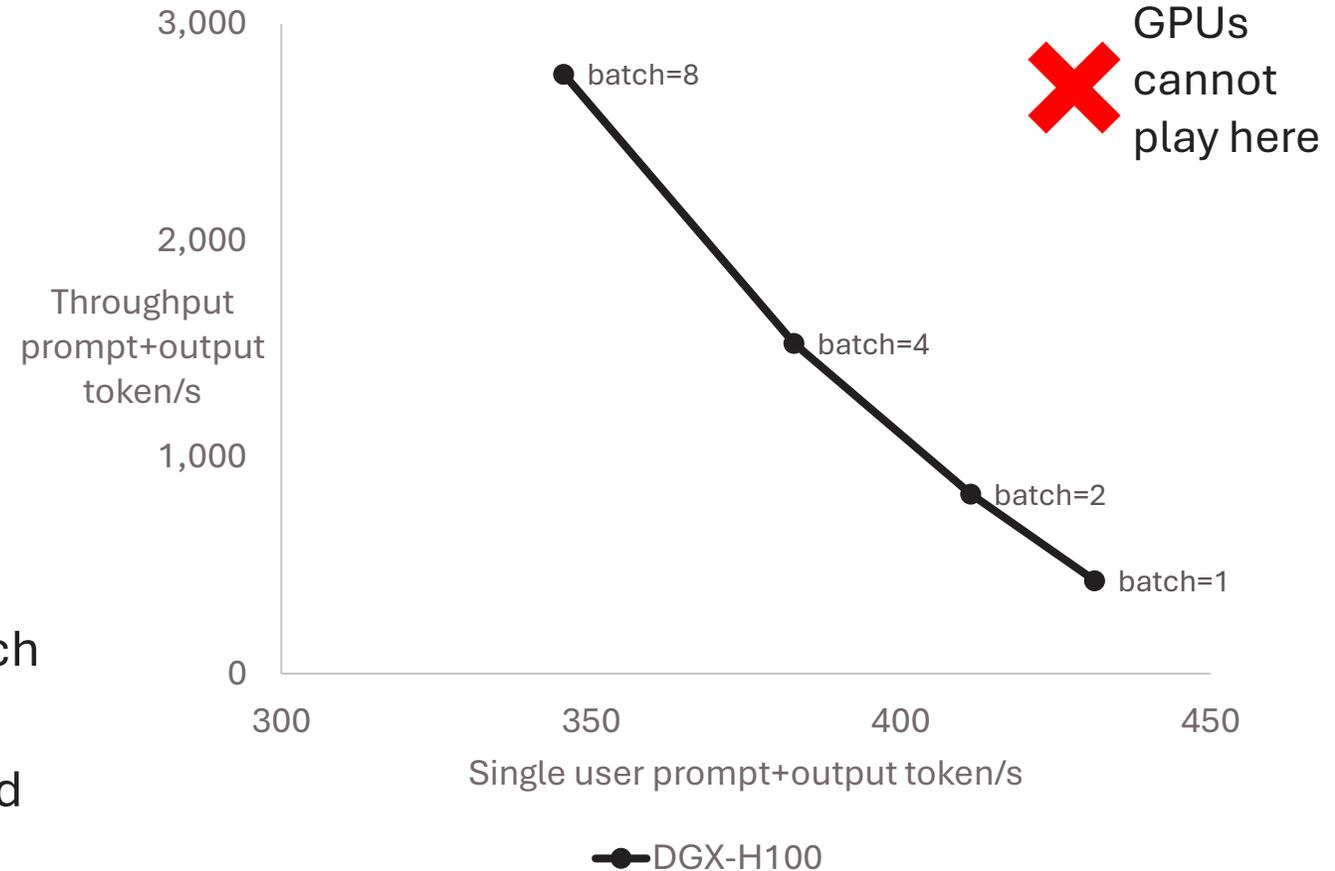- But high batch size results in worse latency

**Latency is critical in generative inference**

- Latency = single user speed

- Better latency enables better user experience

- Better latency enables agentic workflows

**The GPU latency vs. throughput tradeoff**

- Single user speed low to start, even at low batch

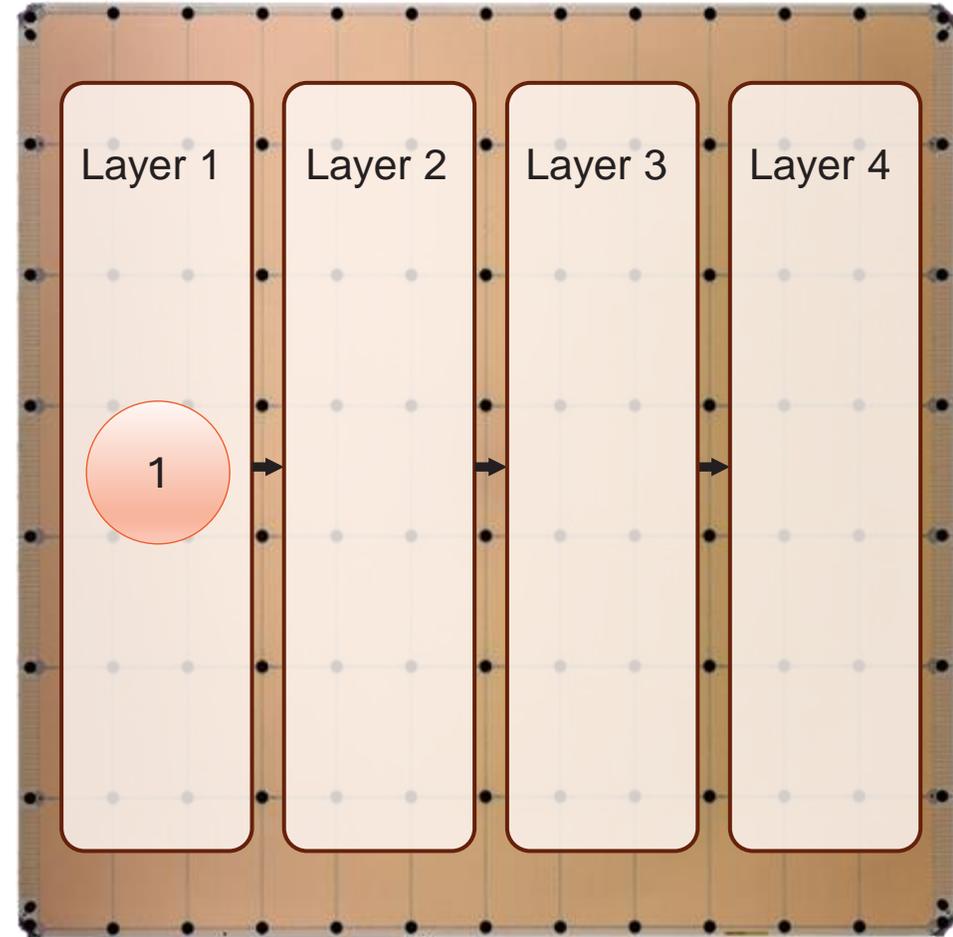- To achieve higher throughput requires higher batch, resulting in even lower single user speed

Llama 70B Throughput vs. Single User Speed
128 Prompt Tokens,  20 Output Tokens



GPUs cannot play here

Throughput prompt+output token/s

batch=8

batch=4

batch=2

batch=1

Single user prompt+output token/s

DGX-H100

# Cerebras enables **low latency and high throughput**

**More than enough memory bandwidth for single user enables high multi-user throughput**
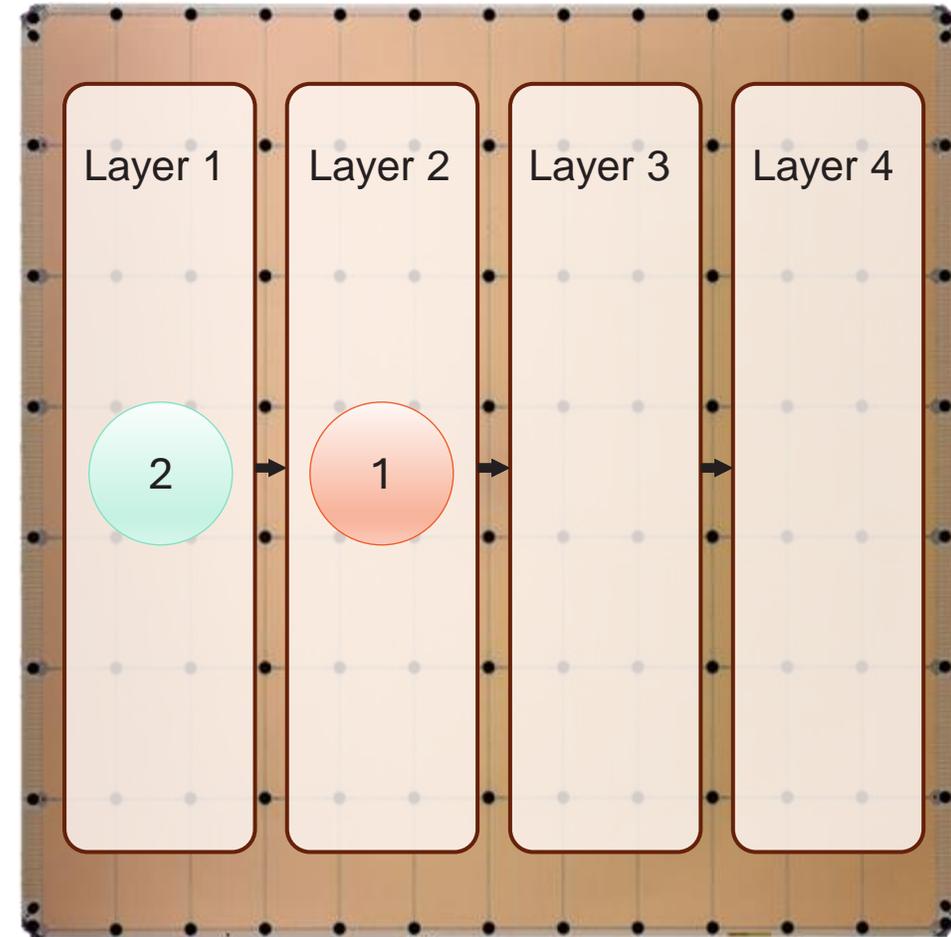
- Single user is using only a fraction of the total memory bandwidth

- We can use the extra bandwidth to support multiple users

# Cerebras enables **low latency and high throughput**

**More than enough memory bandwidth for single user enables high multi-user throughput**
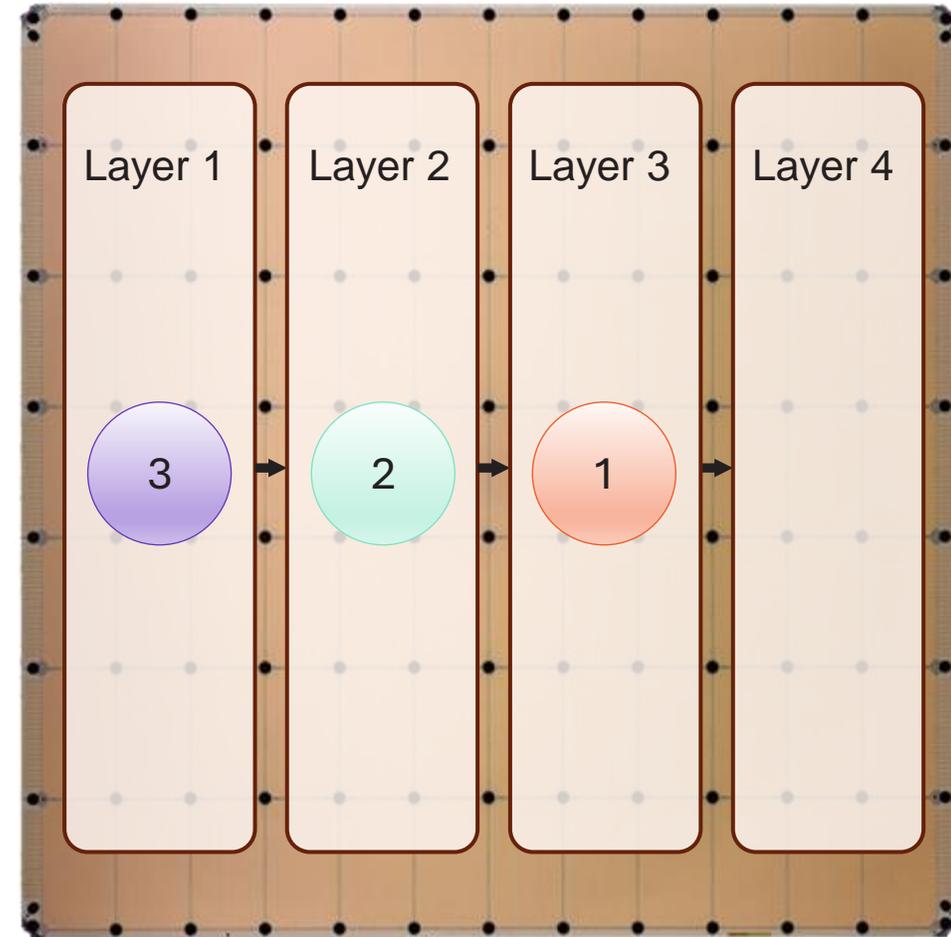
- Single user is using only a fraction of the total memory bandwidth

- We can use the extra bandwidth to support multiple users

- Additional users all run in parallel

- Each accessing the model simultaneously

# Cerebras enables **low latency and high throughput**

**More than enough memory bandwidth for single user enables high multi-user throughput**

- Single user is using only a fraction of the total memory bandwidth

- We can use the extra bandwidth to support multiple users

- Additional users all run in parallel

- Each accessing the model simultaneously

- Every user gets full performance
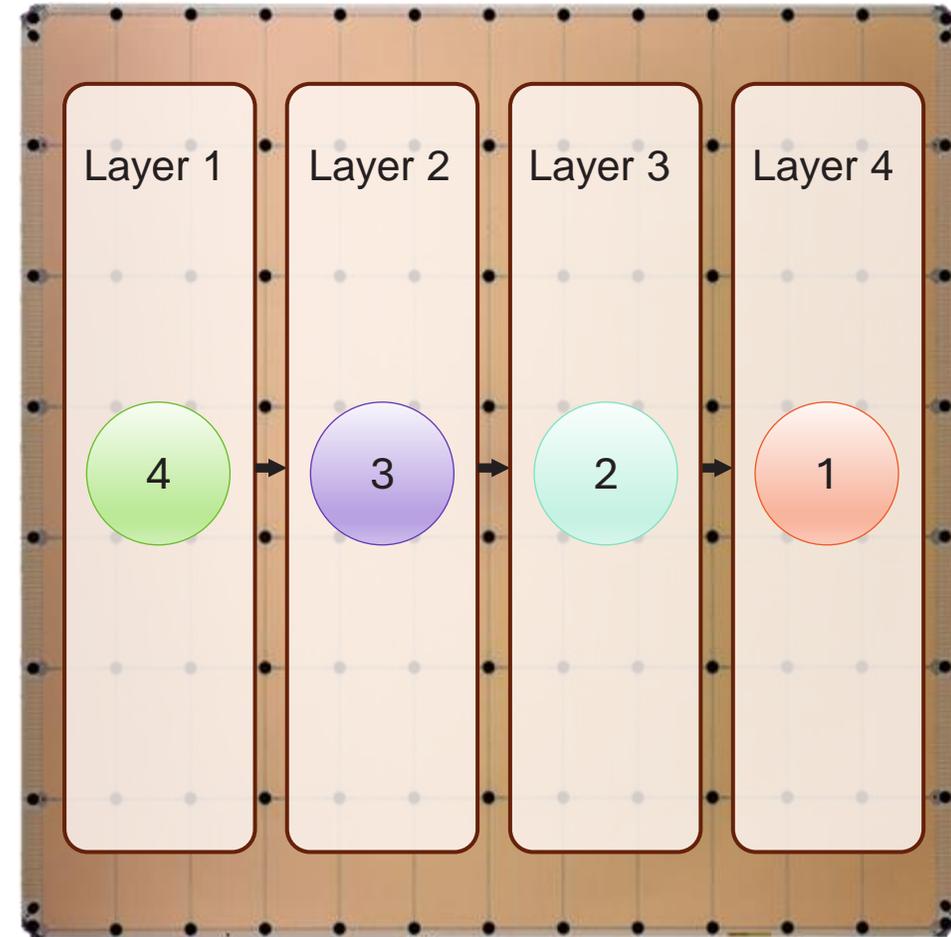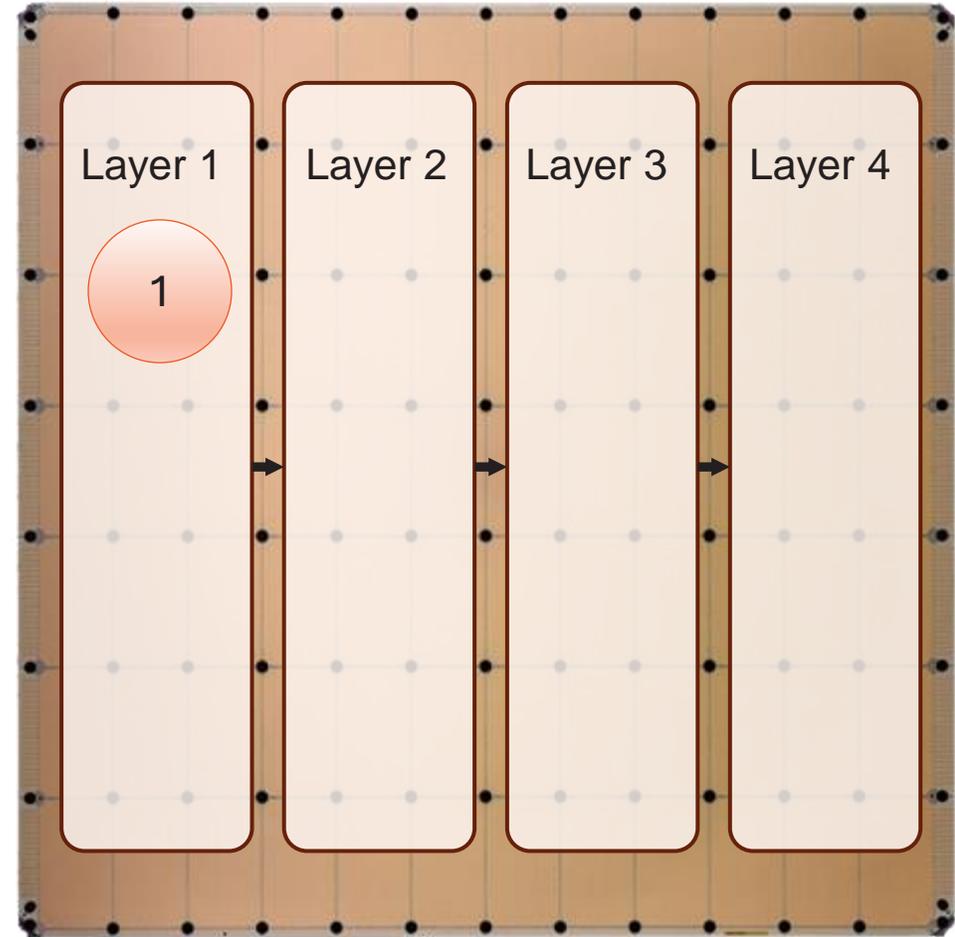
- All pipe stages to run at the same time

# Cerebras enables **low latency and high throughput**

**More than enough memory bandwidth for single user enables high multi-user throughput**

- Single user is using only a fraction of the total memory bandwidth

- We can use the extra bandwidth to support multiple users

- Additional users all run in parallel

- Each accessing the model simultaneously

- Every user gets full performance

- All pipe stages to run at the same time

**Full pipeline model parallelism on single chip**

# High performance prompt processing

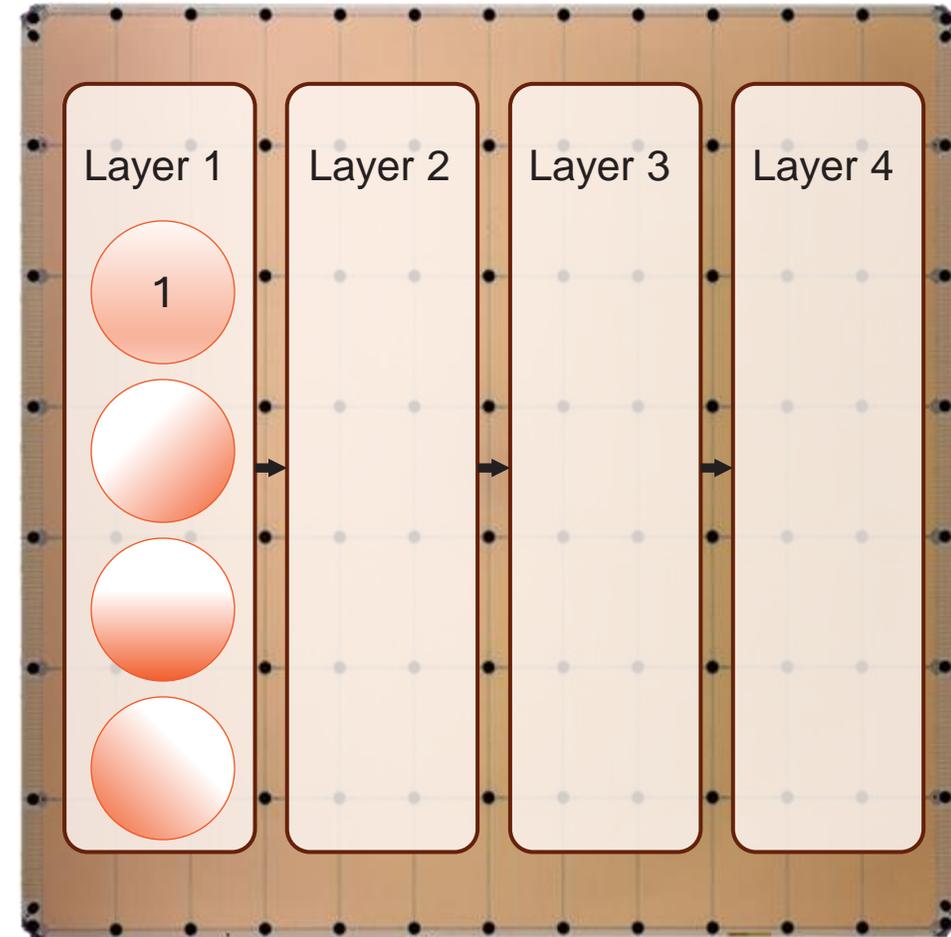**Prompt tokens can be processed in parallel**

- Prompt is known upfront so no need to wait for full pipeline execution for next token

# High performance prompt processing

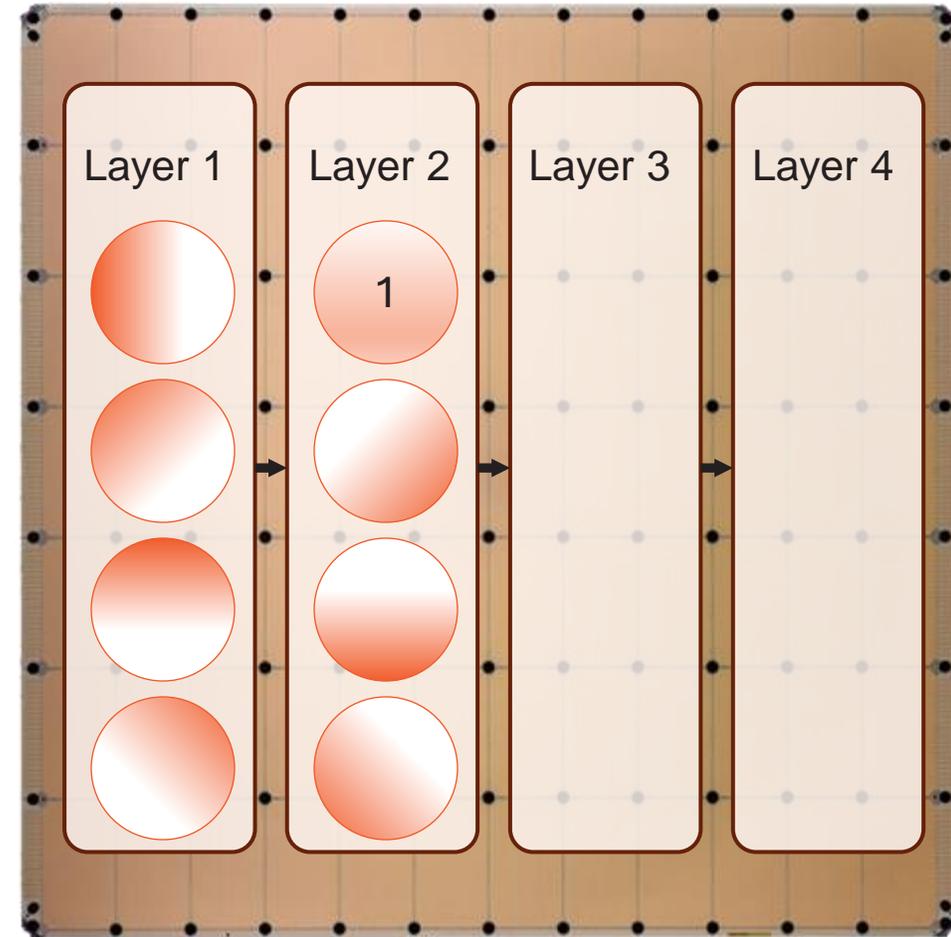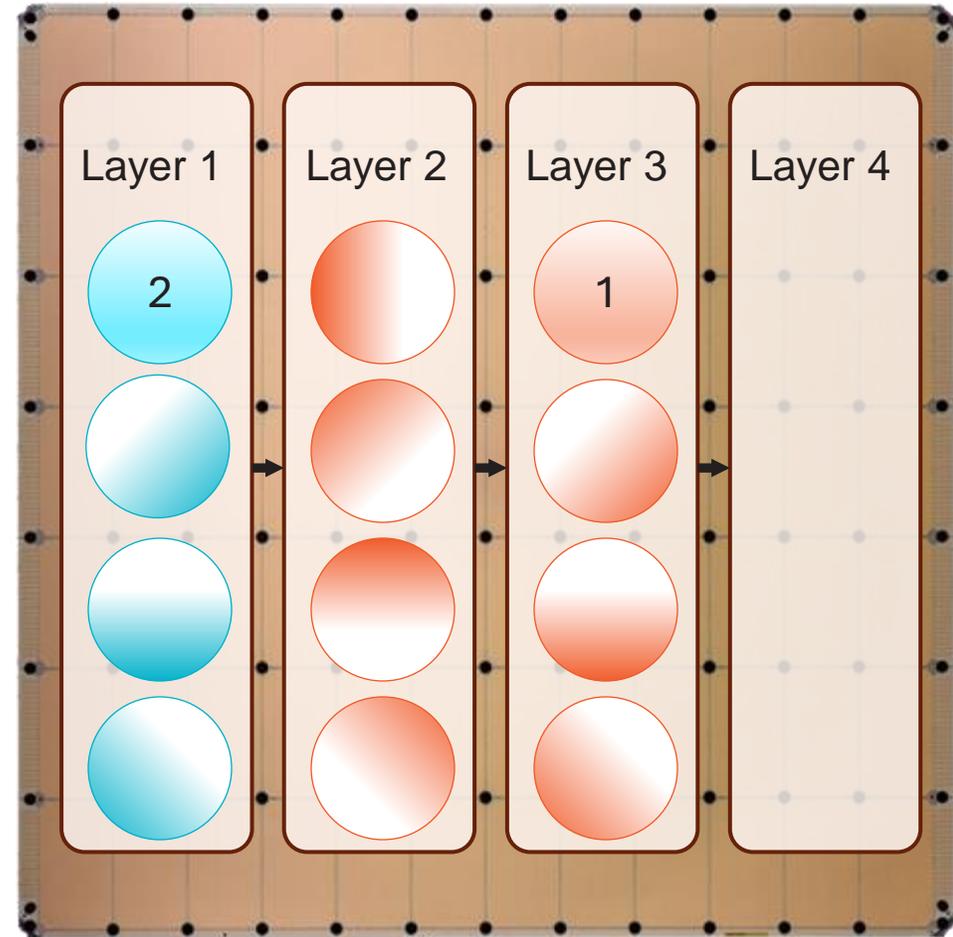**Prompt tokens can be processed in parallel**

- Prompt is known upfront so no need to wait for full pipeline execution for next token

- Run multiple single user prompt tokens in the same pipe stages

# High performance prompt processing

**Prompt tokens can be processed in parallel**

- Prompt is known upfront so no need to wait for full pipeline execution for next token

- Run multiple single user prompt tokens in the same pipe stages

- Run multiple single user prompt tokens in multiple pipe stages

# High performance prompt processing

**Prompt tokens can be processed in parallel**

- Prompt is known upfront so no need to wait for full pipeline execution for next token

- Run multiple single user prompt tokens in the same pipe stages

- Run multiple single user prompt tokens in multiple pipe stages

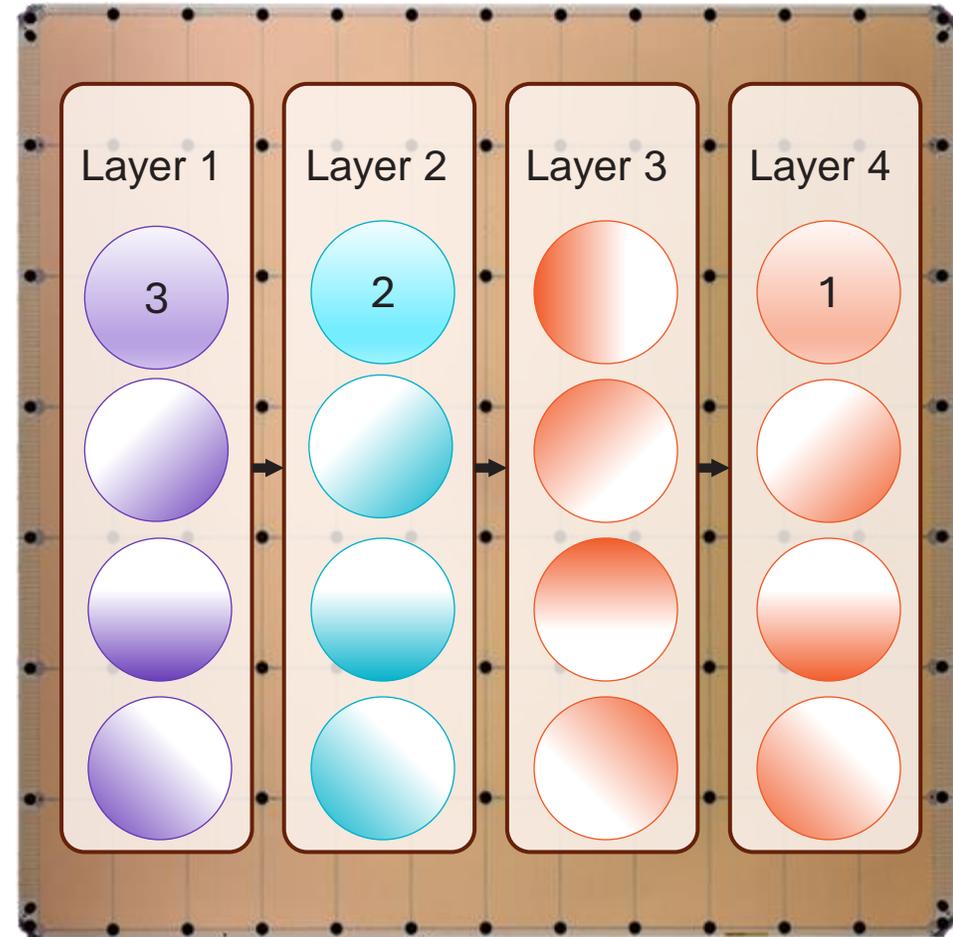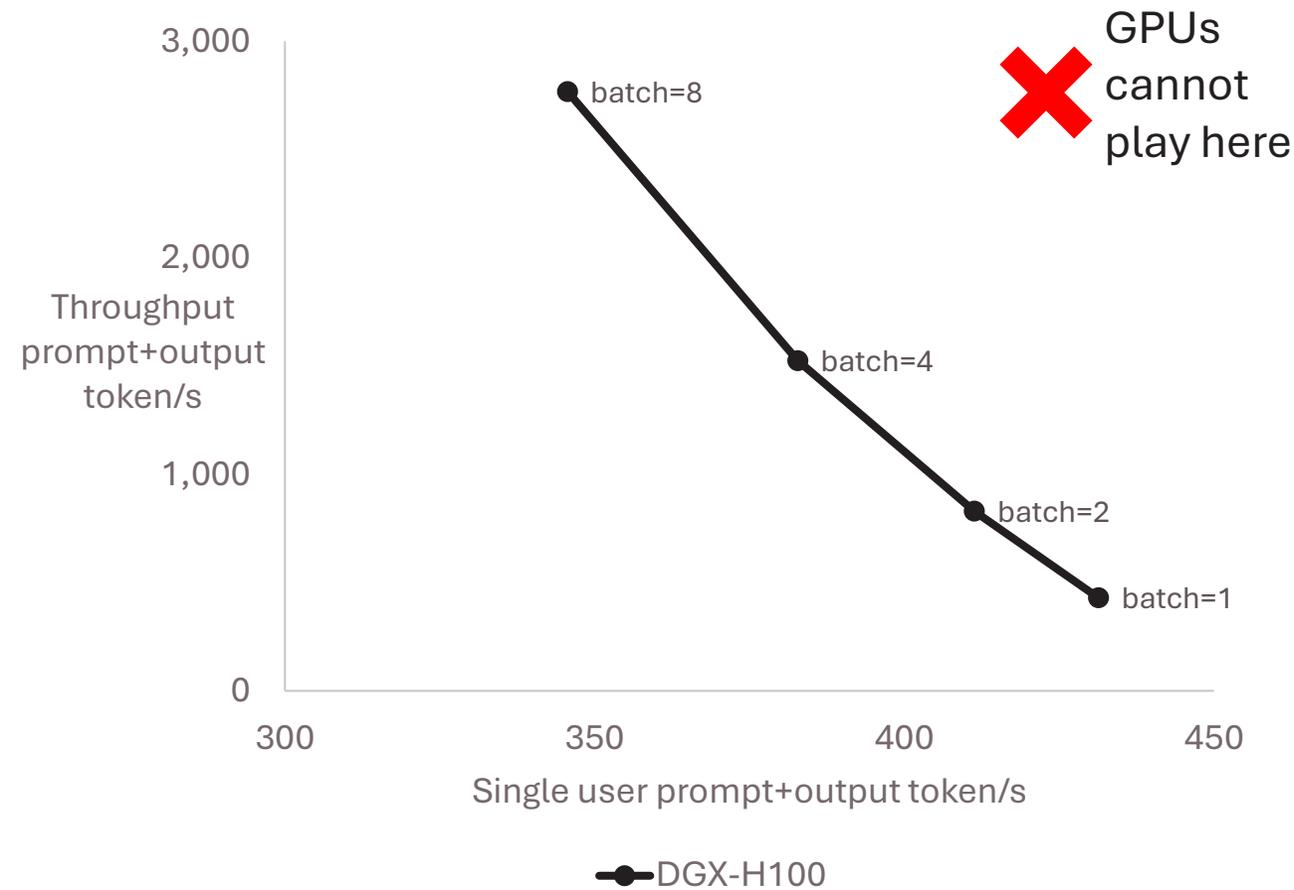- Can achieve higher single user prompt speed using empty pipe stages when fewer users

# High performance prompt processing

**Prompt tokens can be processed in parallel**

- Prompt is known upfront so no need to wait for full pipeline execution for next token

- Run multiple single user prompt tokens in the same pipe stages

- Run multiple single user prompt tokens in multiple pipe stages

- Can achieve higher single user prompt speed using empty pipe stages when fewer users

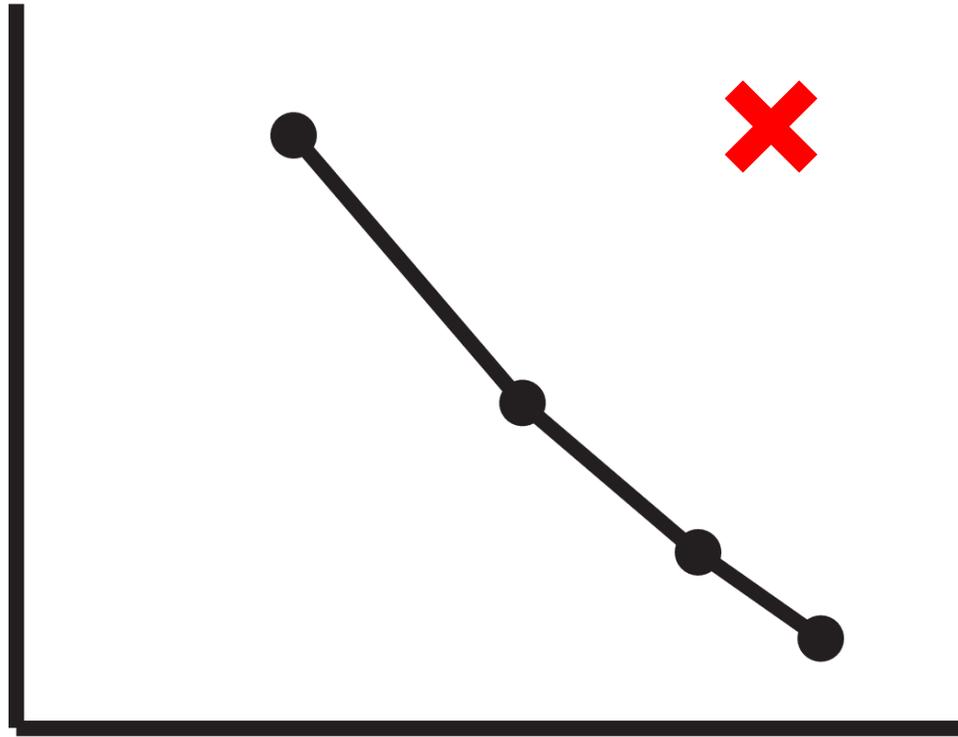**Flexible faster prompt processing driving maximum throughput**

# The latency vs. throughput tradeoff

### Llama 70B Throughput vs. Single User Speed
### 128 Prompt Tokens,  20 Output Tokens



**GPUs cannot play here**

**Where is Cerebras on this graph?**

# Zoom out 10x because Cerebras is so much faster...


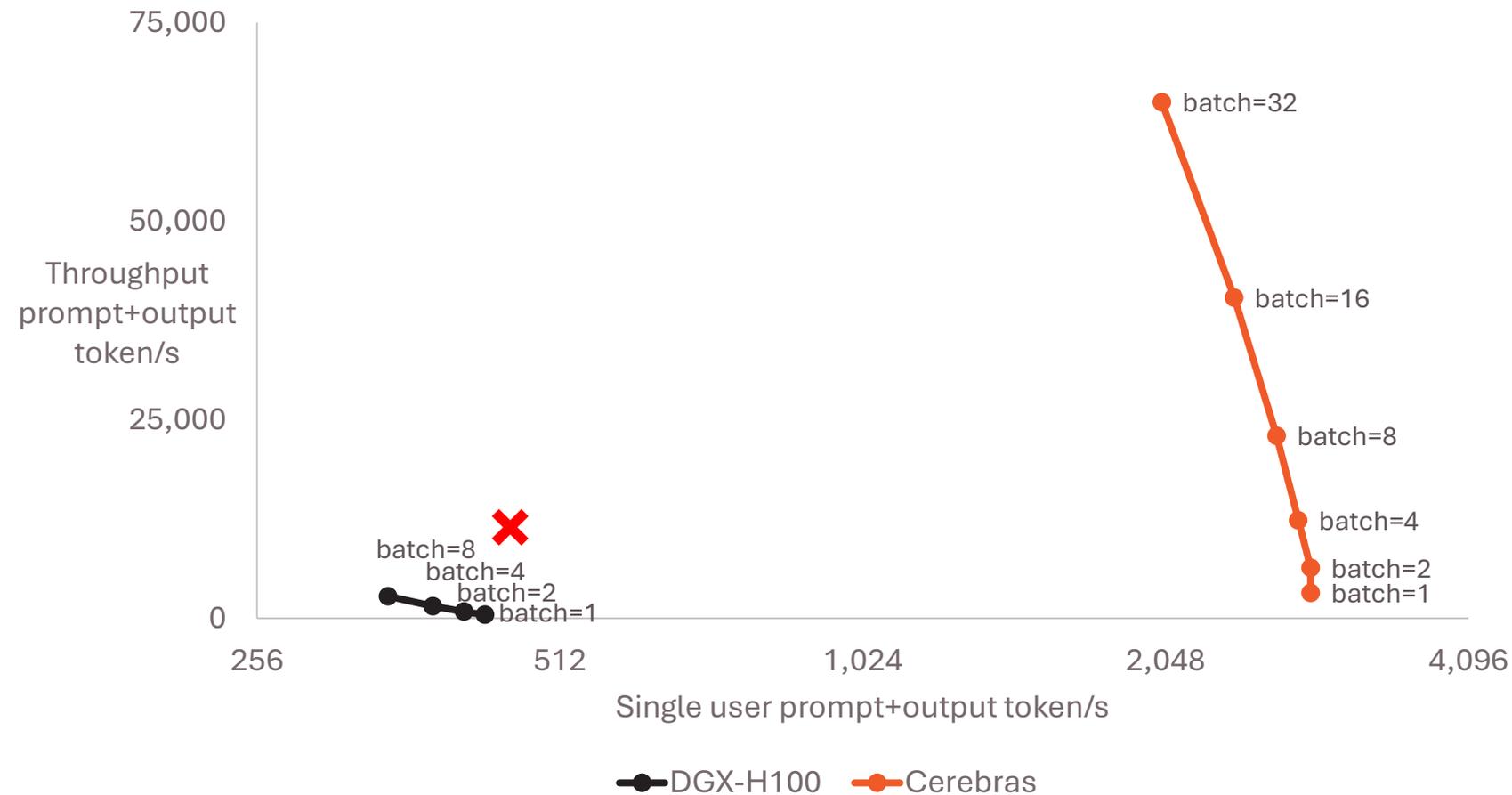
Where is Cerebras on this graph?

# Zoom out 10x because Cerebras is so much faster...

Where is Cerebras on this graph?

# Wafer-scale architecture enables highest single user speed **<u>and</u>** high throughput

Llama70B Throughput vs. Single User Speed
128 Prompt Tokens,  20 Output Tokens

With these throughput techniques, we expect

**20-40x** higher throughput

**5-20x** single user speed

**lower cost** per token

Throughput prompt+output token/s

75,000

50,000

25,000

0

batch=32

batch=16

batch=8

batch=4

batch=2
batch=1

batch=8
batch=4
batch=2
batch=1

256          512          1,024          2,048          4,096

Single user prompt+output token/s

—●— DGX-H100          —●— Cerebras

# And this is just the beginning...

We are already working on many techniques to improve further:

| Technique | Speed | Footprint & Throughput |
|---|:---:|:---:|
| Speculative decoding | ✓ | |
| KV cache optimizations | | ✓ |
| Quantization | ✓ | ✓ |
| Sparsity | ✓ | ✓ |
| More to come... | ✓ | ✓ |

**We are continually improving performance
and supporting larger models everyday**

Cerebras Inference Service

# Cerebras Inference Service Launching Today!

Go to **inference.cerebras.ai** to try it out!

# Cerebras Inference Service Launching Today!

## Llama3.1-8B

1,800 tokens/s

**Free tier**
30 requests/m
1M daily token limit

**Paid tier**
10¢ per M tokens

## Llama3.1-70B

450 tokens/s

**Free tier**
30 requests/m
1M daily token limit

**Paid tier**
60¢ per M tokens

## Coming Soon

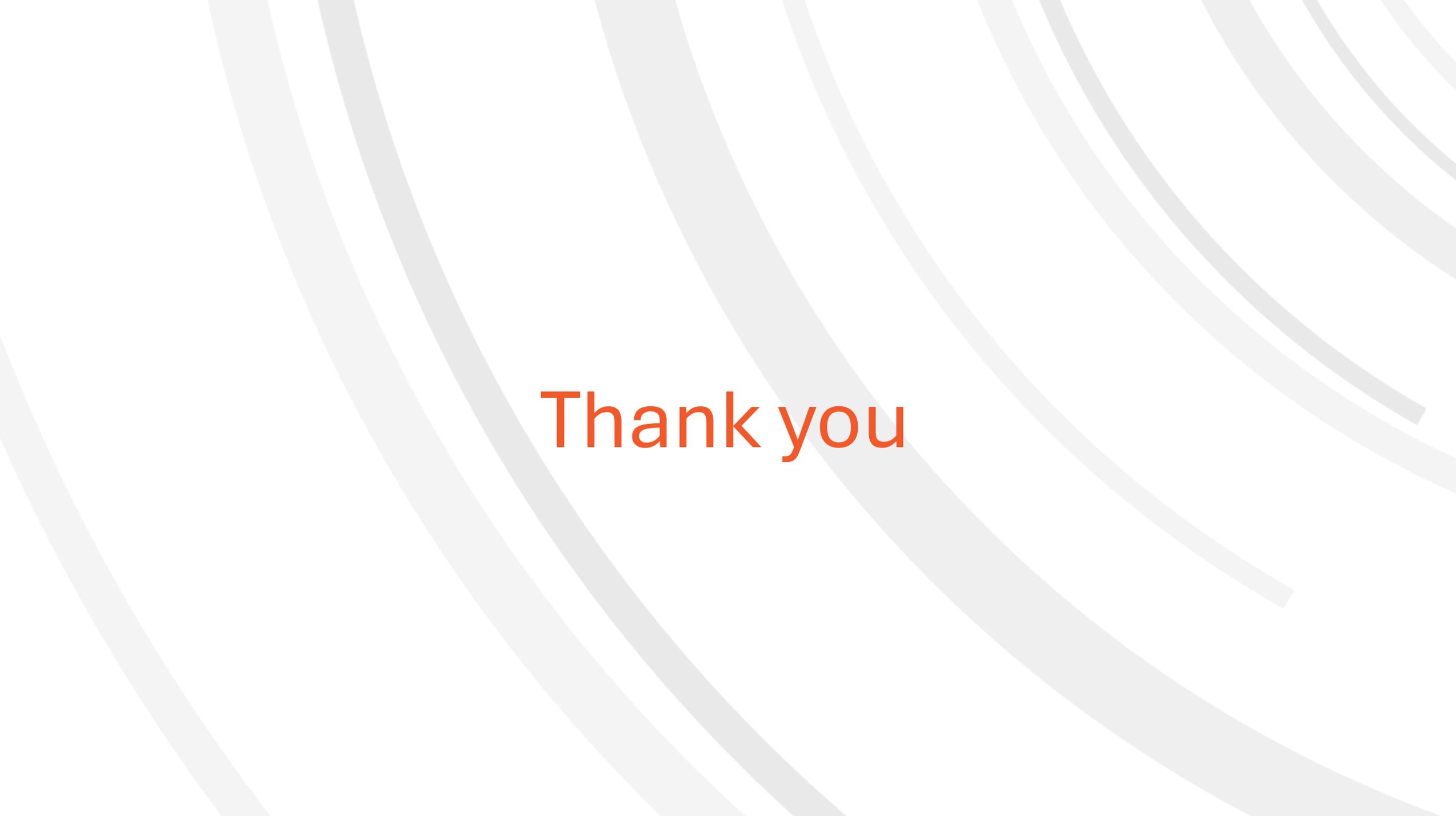Llama-405B

Mistral Large 2

Cohere Command R

Whisper

Custom fine-tunes

**inference.cerebras.ai**
to experience the GPU impossible

Thank you