

Inference at Scale: Opportunities and challenges

Vikram Sharma Mailthody
Sr. Research Scientist, NVIDIA Research

Disclaimer

“

The views expressed in the content belong to me and my collaborators, not NVIDIA, its affiliates, or employees

Vikram Sharma Mailthody

”

Lecture Overview

- Vacuums to tokens
 - Intro to modern day AI datacenters
 - AI factories
 - Datacenter architecture
 - System architecture
 - Open challenges
- AI inference
 - How it works
 - KV caching
- Inference at scale using Dynamo
 - Aggregated approach
 - Disaggregated Inference
 - KV aware router, Planner, and others
 - Open challenges
 - Fault tolerance
 - Load balancing

Credits: Wen-mei Hwu, Junda Chen

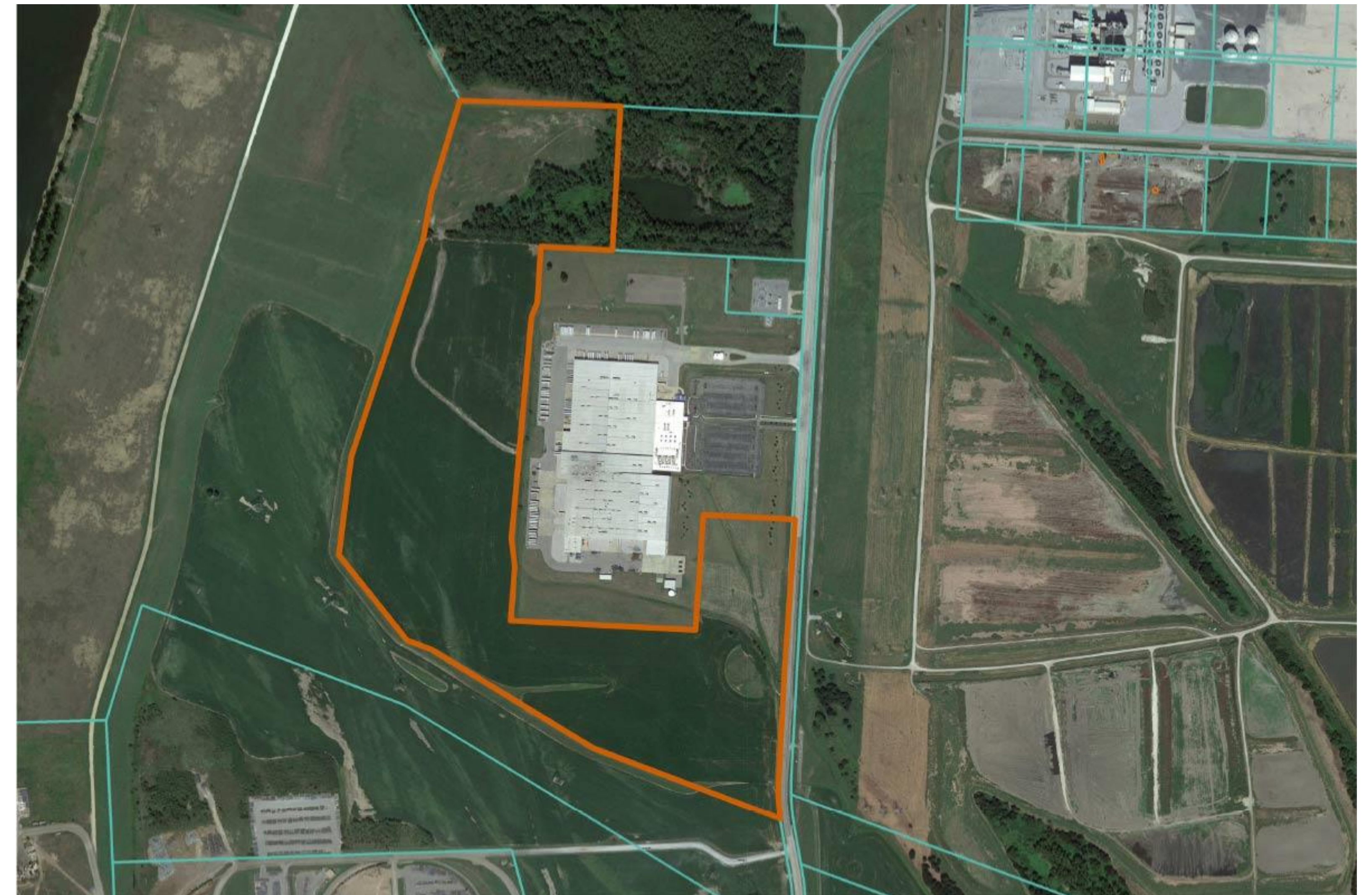
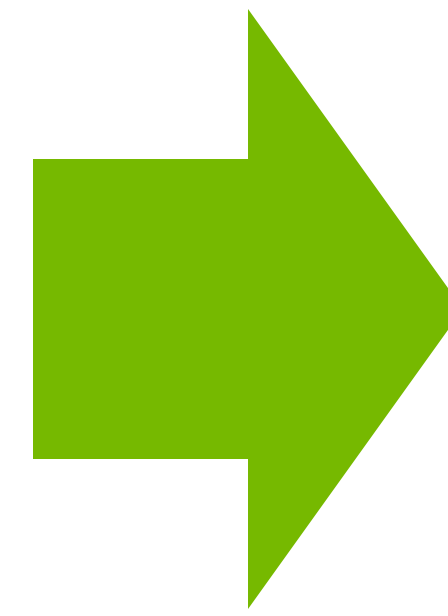


**AI WILL
REPLACE
PROGRAMMERS**

**PROGRAMMERS
WHO USE AI WILL
REPLACE
PROGRAMMERS WHO DON'T**

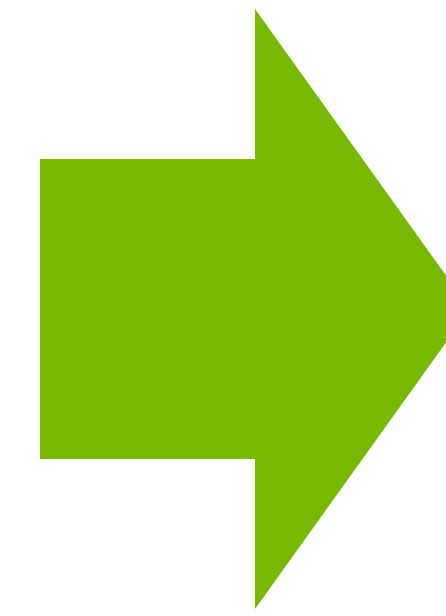
AI Today

Vacuums to tokens in <122 days 😊



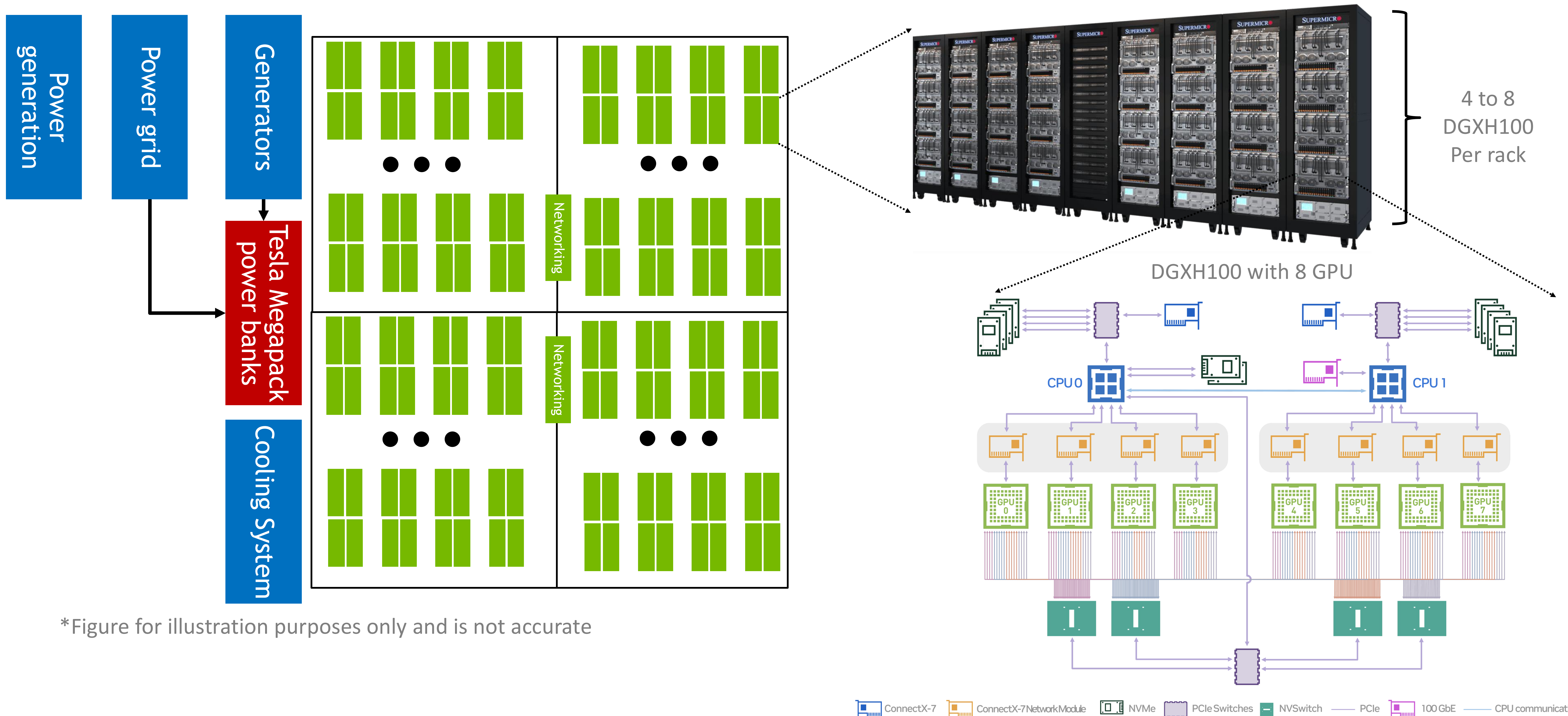
xAI 100K Cluster → 200K soon

>150MW Datacenter



AI factory - xAI 100K Cluster → 200K → ???

4 data halls, each data hall is 25K GPUs, liquid cooled racks, each rack with 64 GPUs



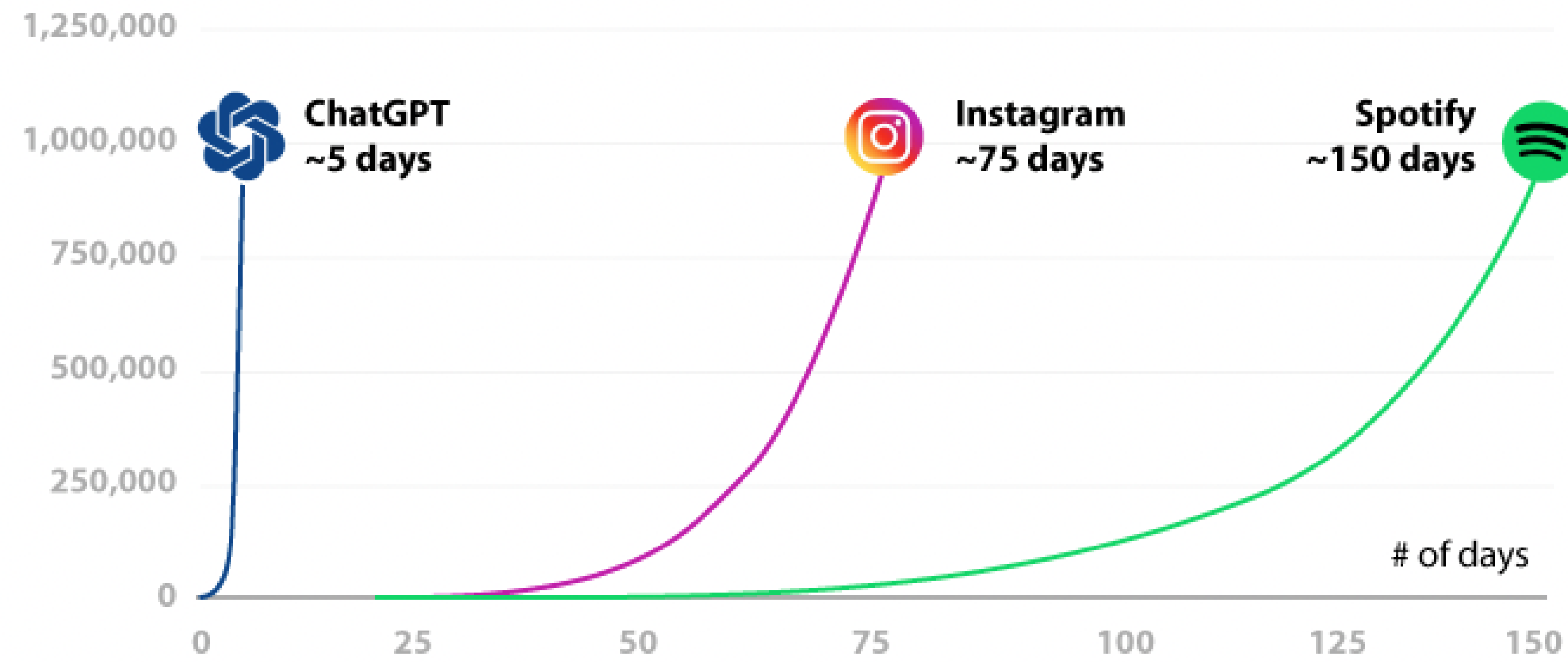
*Figure for illustration purposes only and is not accurate

AI Success

Industry transformation is happening

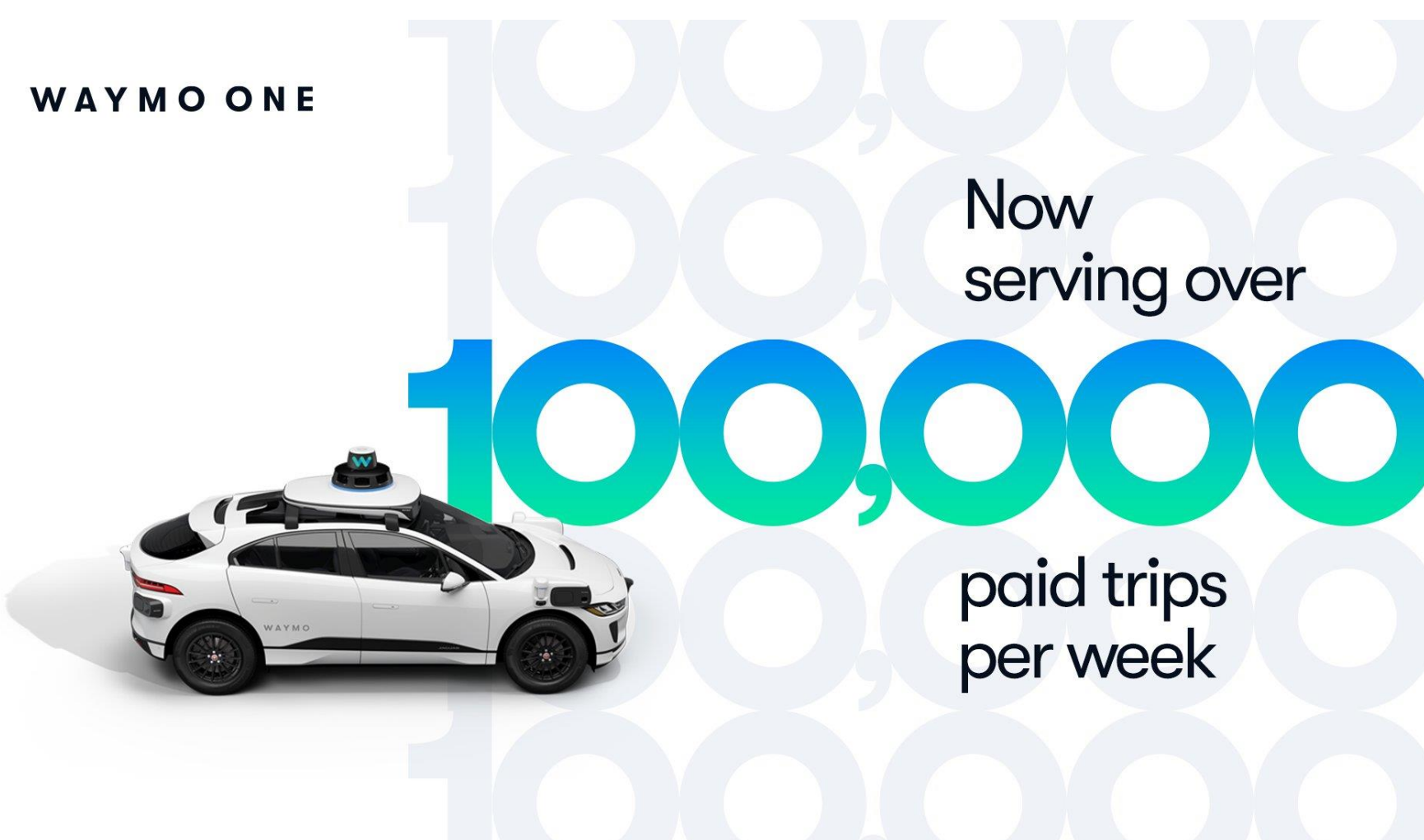
Chatbots

~ Path to 1 million users* (# of days from launch)



Sources: Google, Subredditstats, Media Reports

Self driving cars



Business functions in which respondents' organizations are regularly using AI, by industry,¹
% of respondents

	Total	Media and telecom	Insurance	Technology	Healthcare	Consumer goods and retail	Professional services	Travel and logistics	Energy and materials	Financial institutions	Advanced manufacturing	Engineering and construction	Pharma and medical products
Knowledge management	40	34	64	46	54	28	58	36	33	34	29	39	35
Marketing and sales	39	45	52	49	31	51	46	34	33	35	29	26	46
IT	34	38	55	56	32	32	21	32	39	32	40	25	29
Service operations	33	46	60	45	27	34	32	47	32	34	22	28	21
Product and/or service development	31	32	40	49	33	21	33	34	28	29	30	23	41
Software engineering	26	33	39	58	22	19	13	19	30	22	32	13	19
Human resources	21	28	16	28	22	22	20	9	22	19	18	15	29
Risk, legal, and compliance	17	17	46	18	15	11	15	19	17	47	7	13	9
Strategy and corporate finance	17	17	6	20	17	9	22	22	20	15	16	15	19
Supply chain/inventory management	12	6	4	10	11	22	4	19	19	3	25	12	34
Manufacturing	10	5	0	9	6	13	1	1	21	1	26	14	17
Use in at least 1 business function, %	88	96	95	95	92	91	91	90	89	86	86	84	83

AI Flops outpacing Moore's law

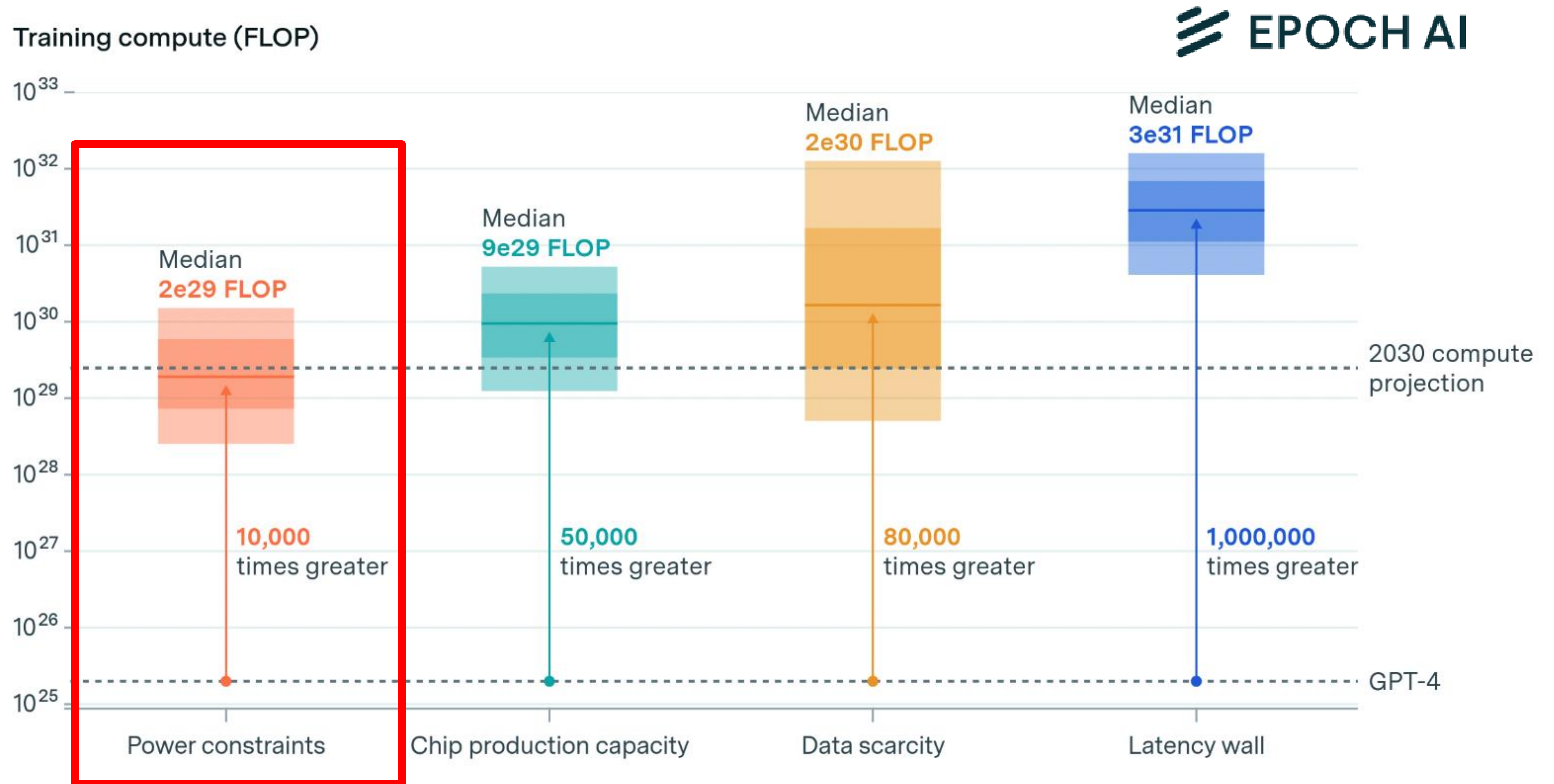
4.2X flops growth per year for LLM models. What about inference needs???



<https://epochai.org/blog/compute-trends>

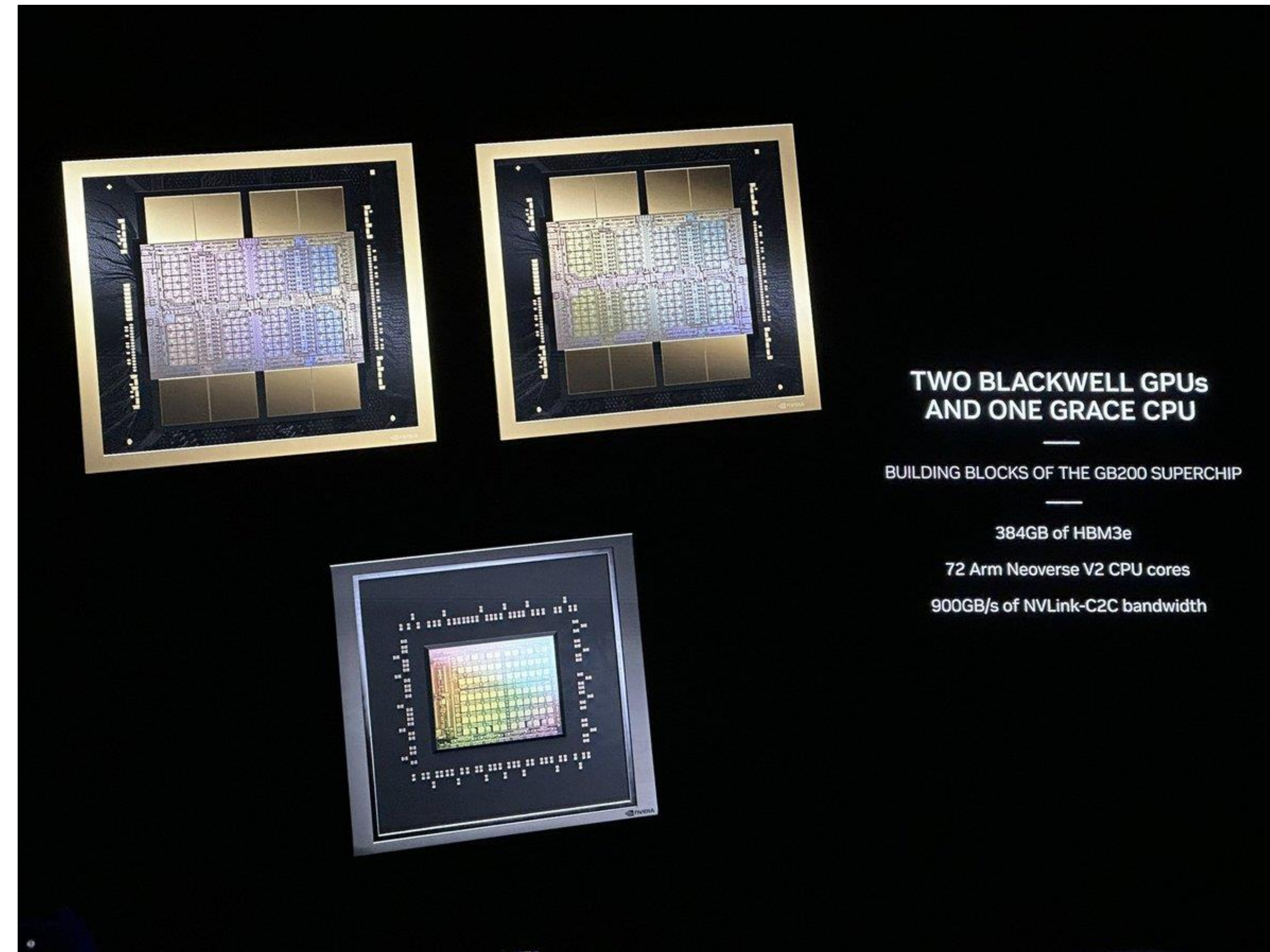
Power is THE primary limiter for AI training beyond 2030

Need to invest in power efficiency AND revamping power generation and delivery



NVL72 GB200 Pods

For power efficiency and higher performance



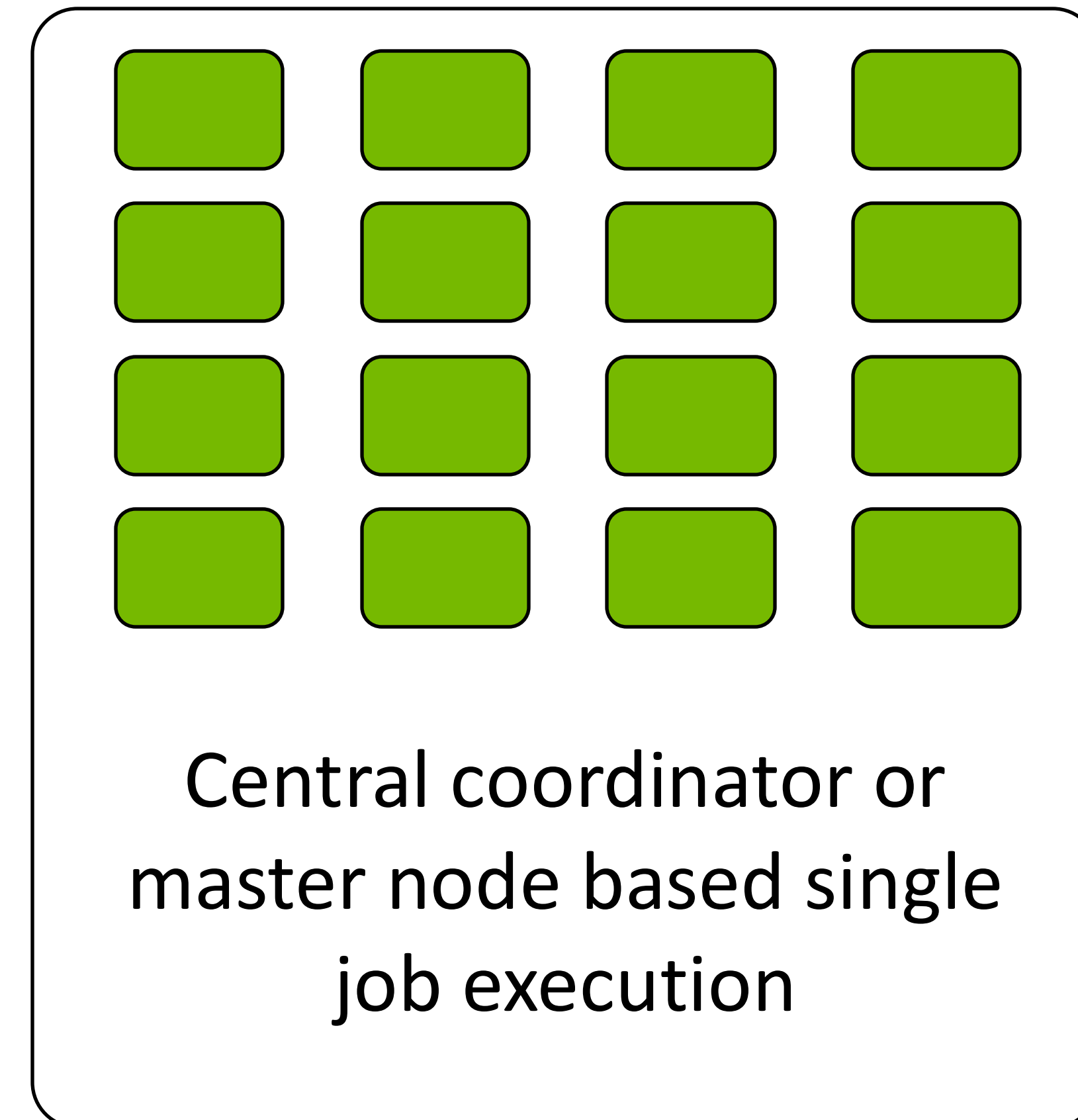
The background features a series of parallel, slightly curved lines in various shades of green, creating a sense of depth and movement. Overlaid on these lines are several overlapping, rounded rectangular shapes in different green tones, some appearing to be layered on top of others. The overall effect is a modern, abstract design.

AI Inference

AI Inference

Different beast compared to AI Training

Parallelism strategies
Failure and resiliency
Weak and strong scaling
Topology awareness
Preprocessing data
Optimizers
....



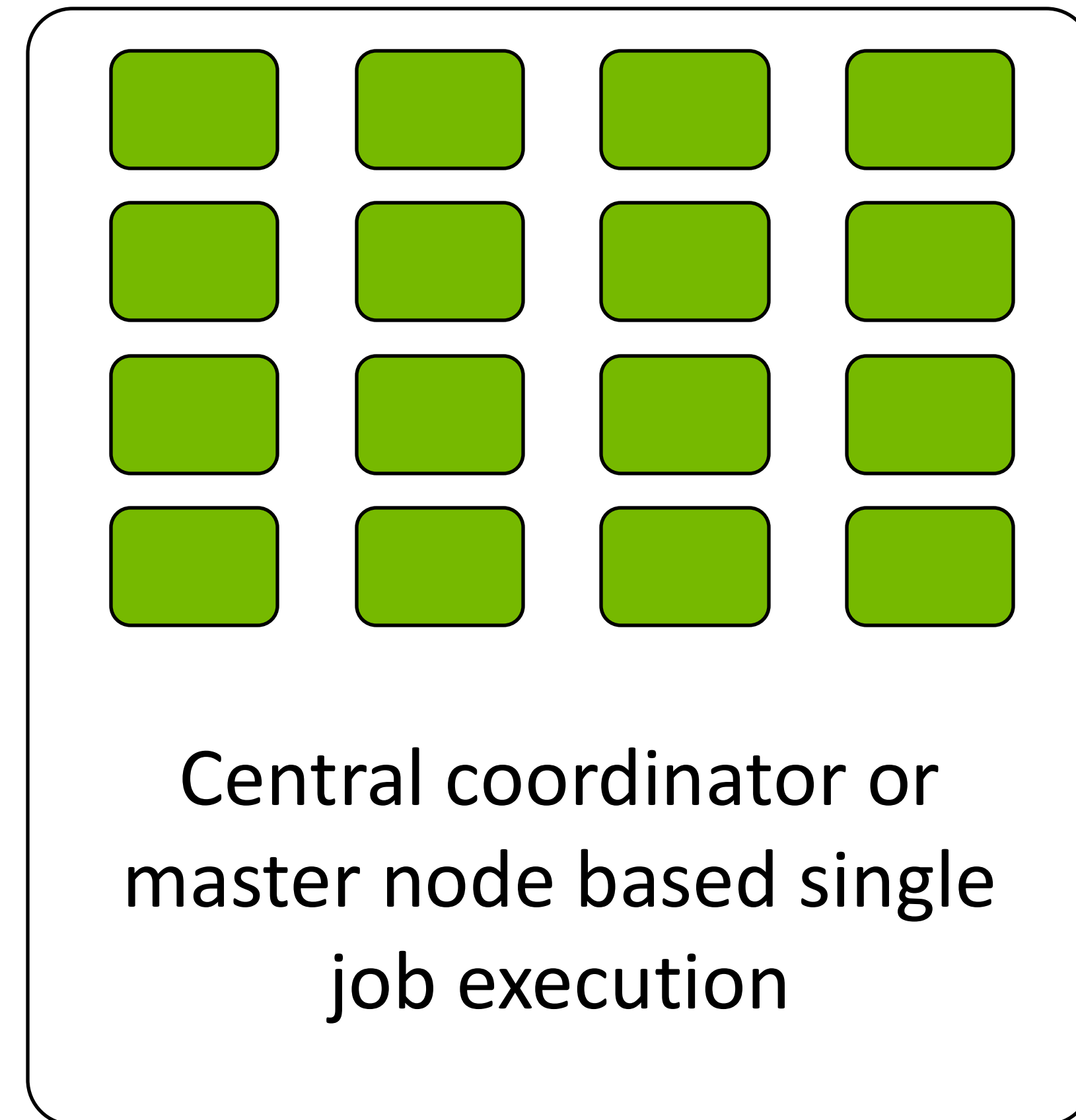
LLM Training

AI Inference

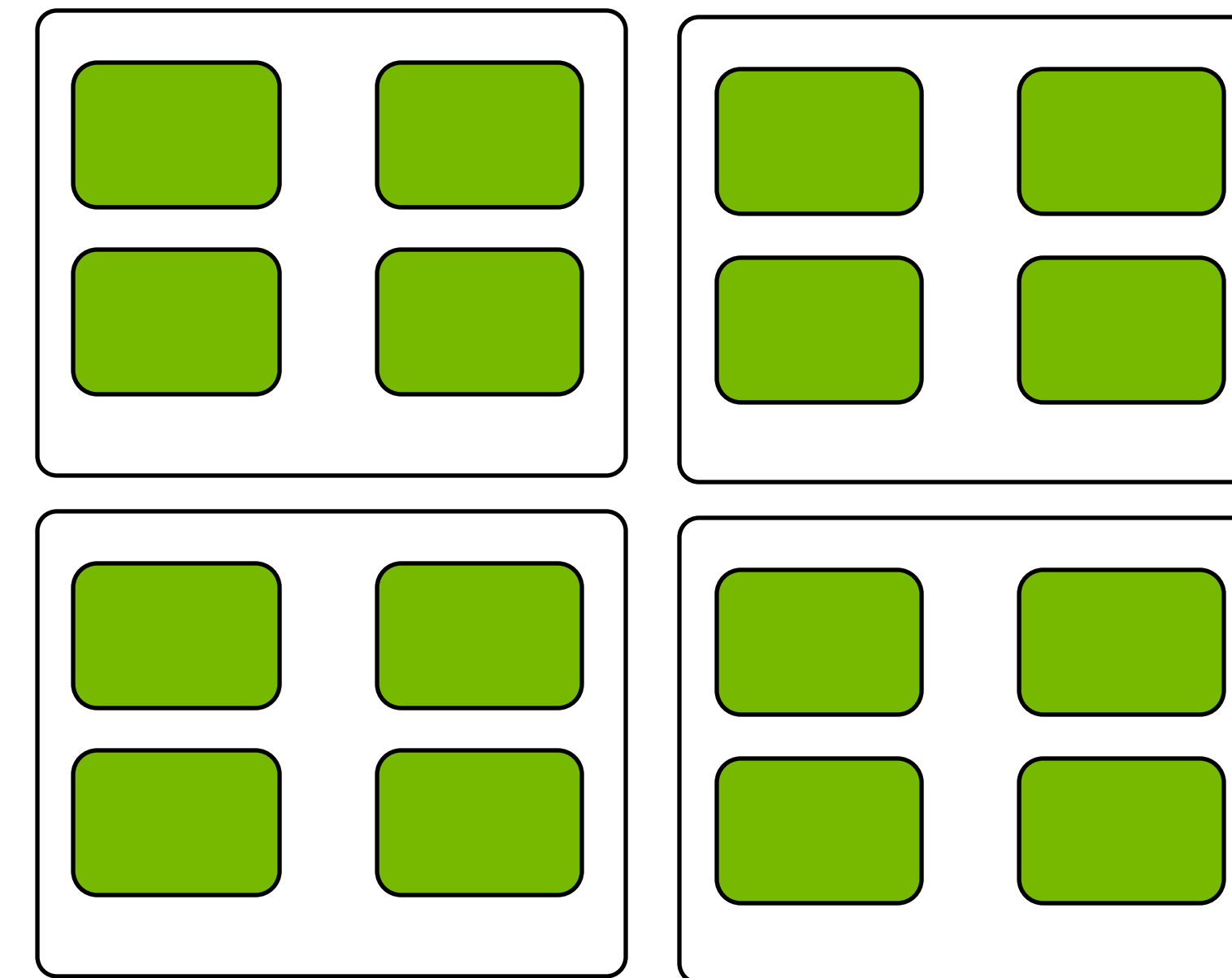
Different beast compared to AI Training

Parallelism strategies
Failure and resiliency
Weak and strong scaling
Topology awareness
Preprocessing data
Optimizers

....



LLM Training



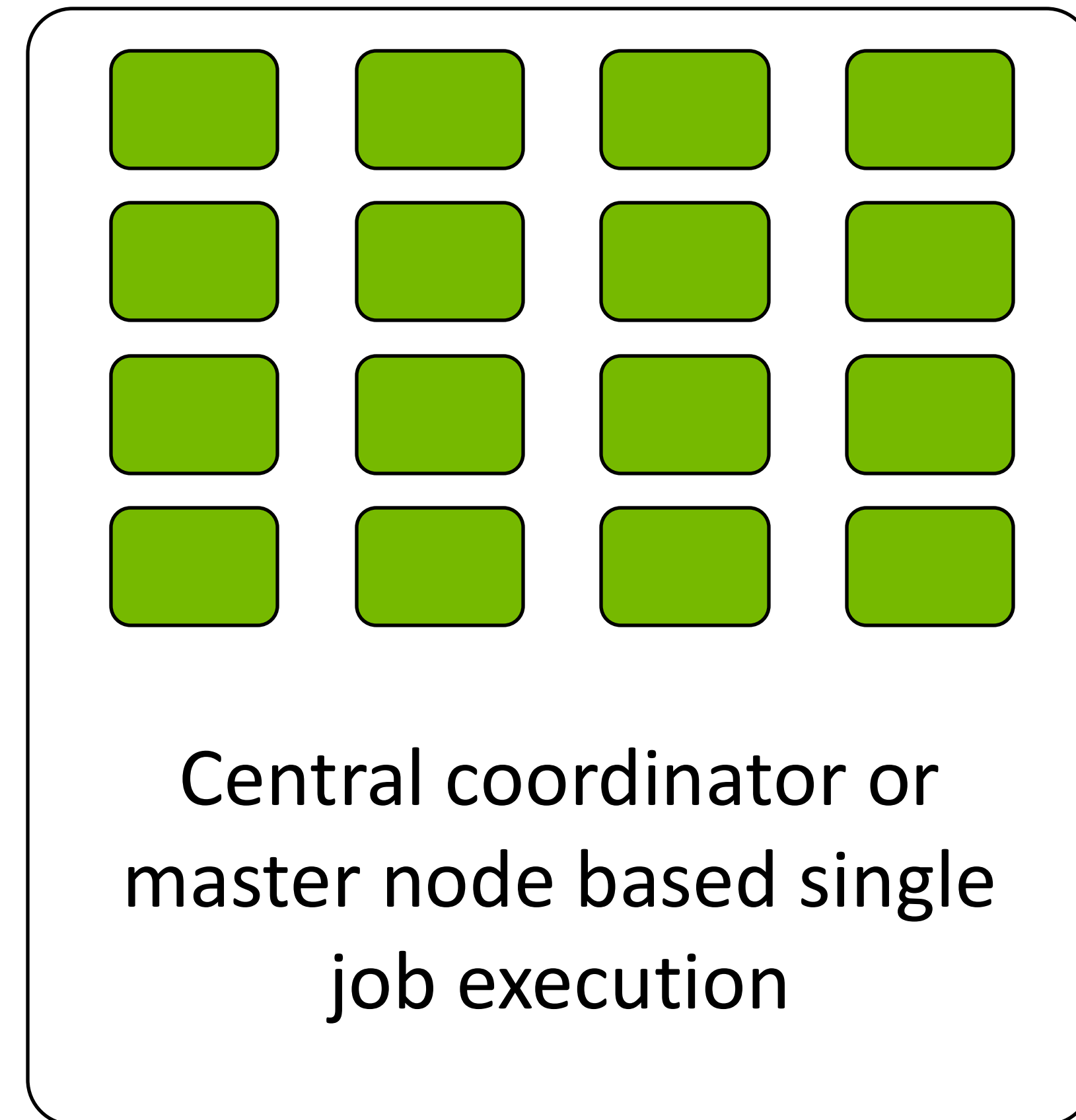
LLM Inference

AI Inference

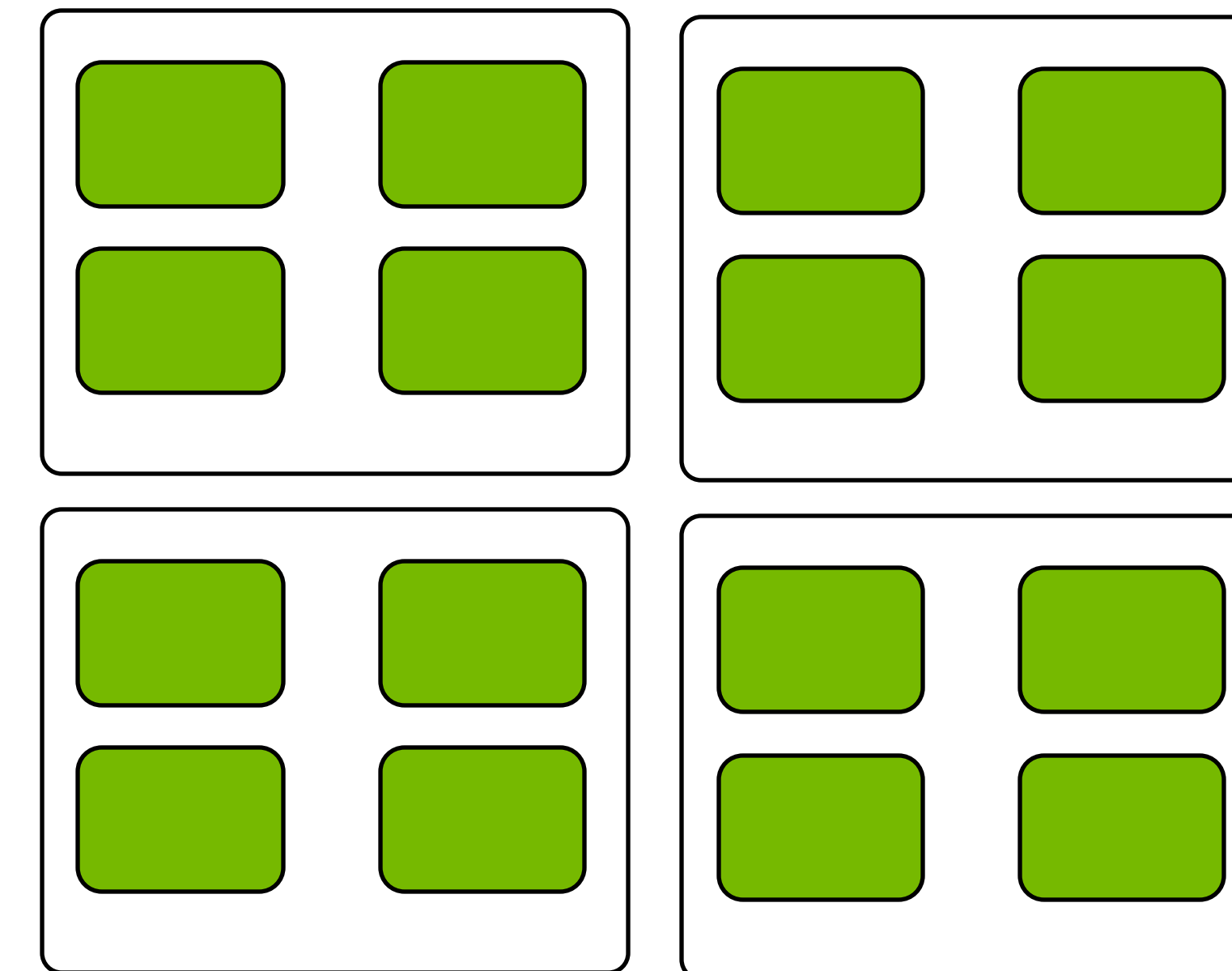
Different beast compared to AI Training

Parallelism strategies
Failure and resiliency
Weak and strong scaling
Topology awareness
Preprocessing data
Optimizers

....



LLM Training



LLM Inference

Inflight batching
KV cache optimization
Quantization
Parallelism strategies
User-defined scheduling

...

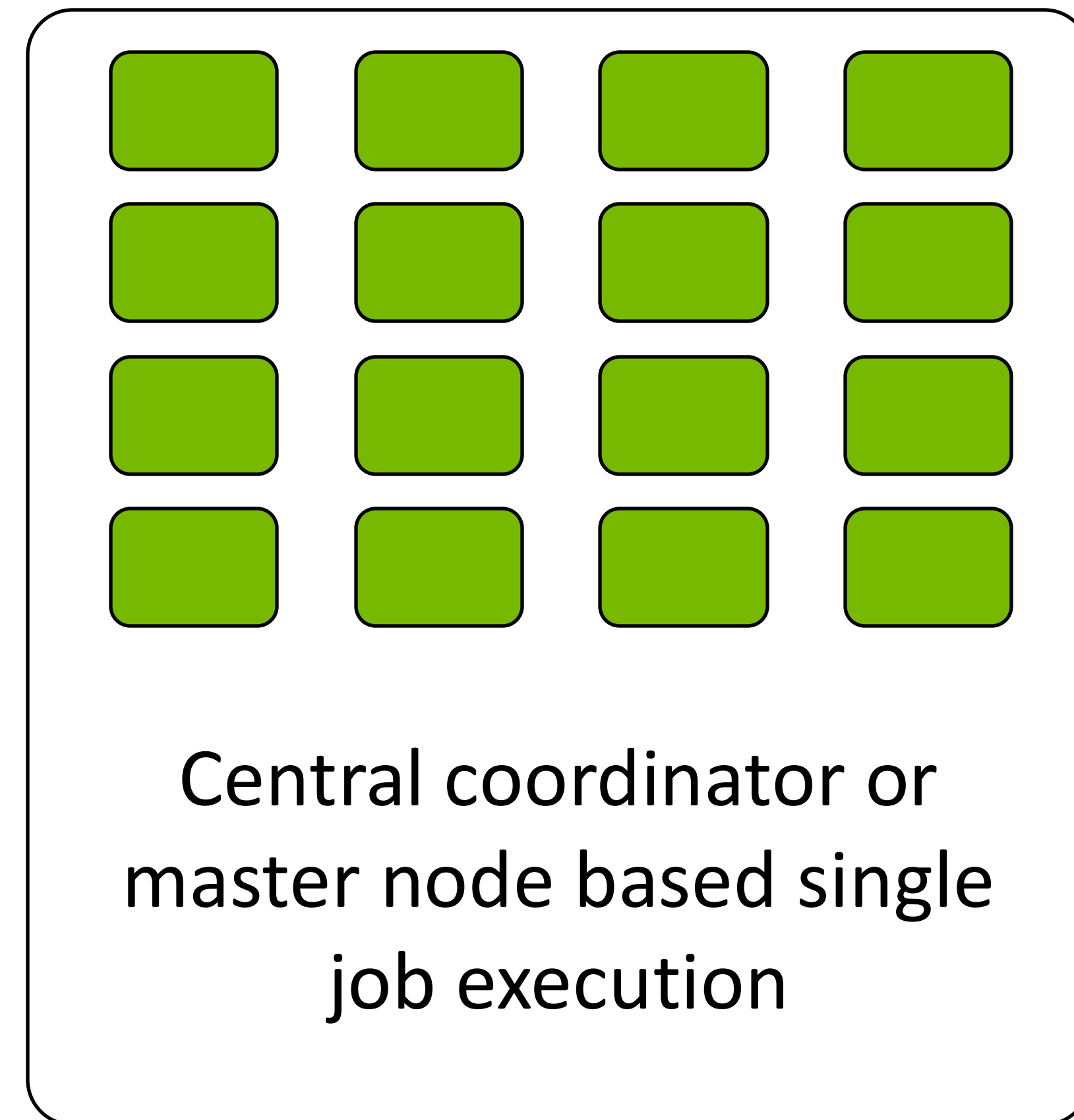
- Unique challenges and lot of pending optimizations in both stack

AI Inference

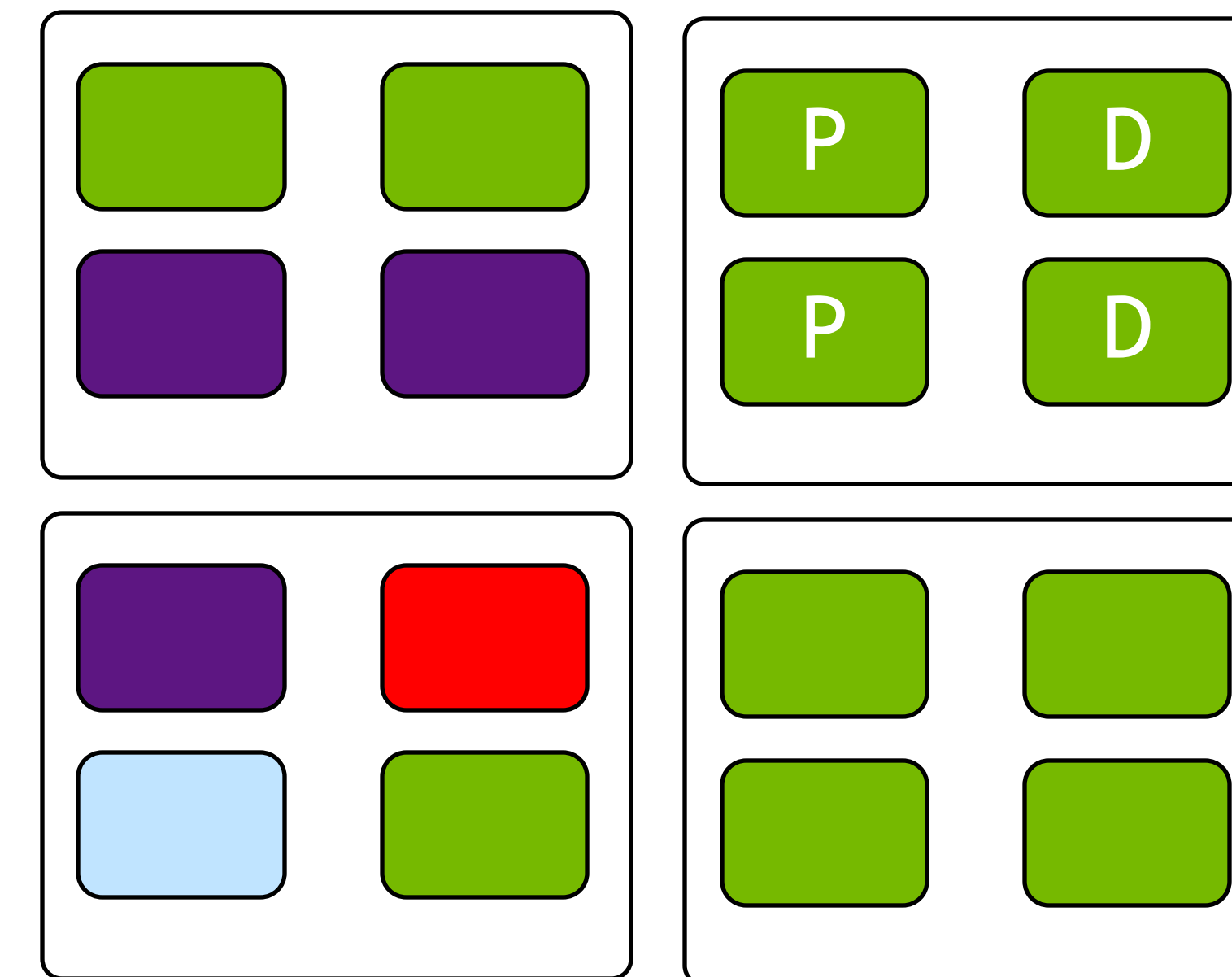
Different beast compared to AI Training

Parallelism strategies
Failure and resiliency
Weak and strong scaling
Topology awareness
Preprocessing data
Optimizers

....



LLM Training



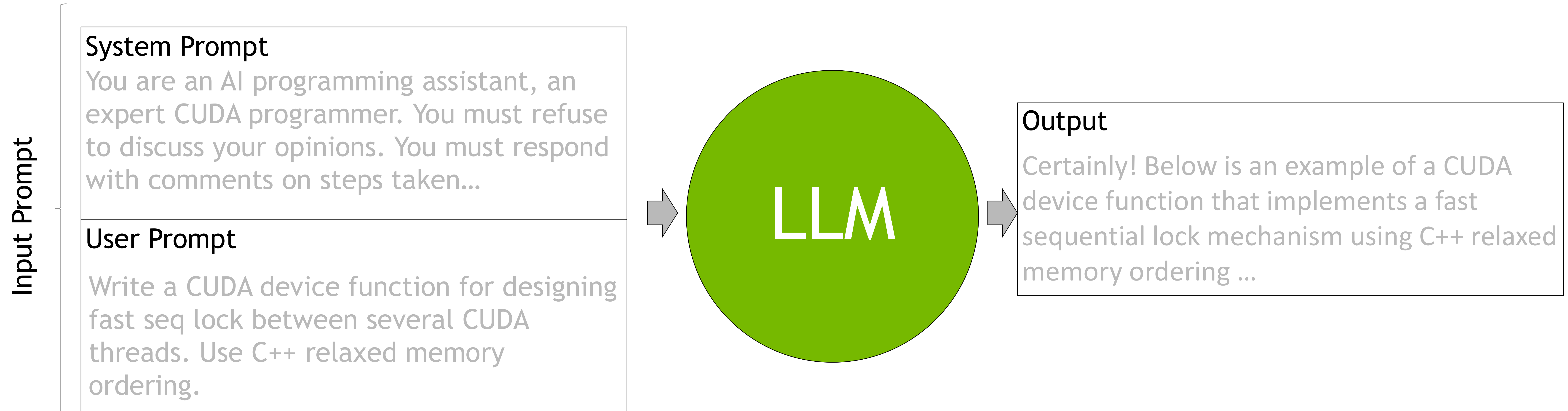
LLM Inference

Inflight batching
KV cache optimization
Quantization
Parallelism strategies
User-defined scheduling
...

- Unique challenges and lot of pending optimizations in both stack

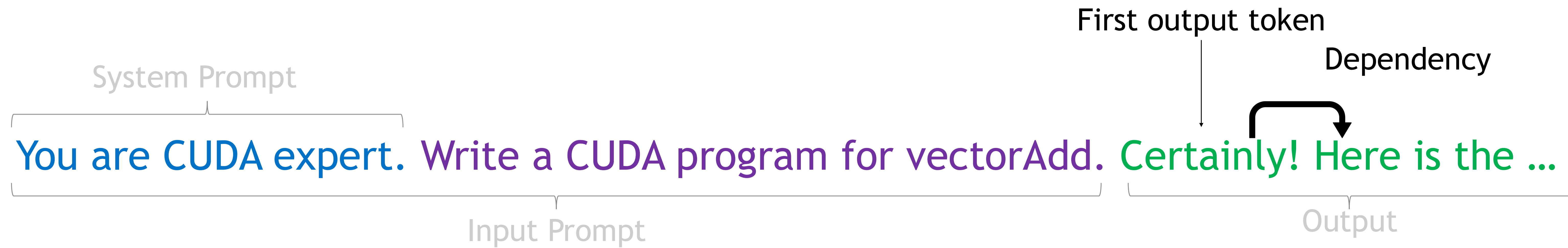
LLM Inference

Typical workflow of LLM inference



- System prompts are instructions to language model that constraints the output generation in a specified manner
- System prompts are common across several user prompts (from the same user or different set of users)
- Input prompt = system + user prompt and can consume even upto 32K tokens.
- Trend is towards longer prompts with "prompt rules" → leads to better quality response from the model.
- Input prompt processing accounts for a significant-to-dominant fraction of the overall processing time

LLM Inference



$$P(x) = P(x_1) \cdot P(x_2|x_1) \dots P(x_n|x_1, \dots x_{n-1})$$

Autoregressive decomposition

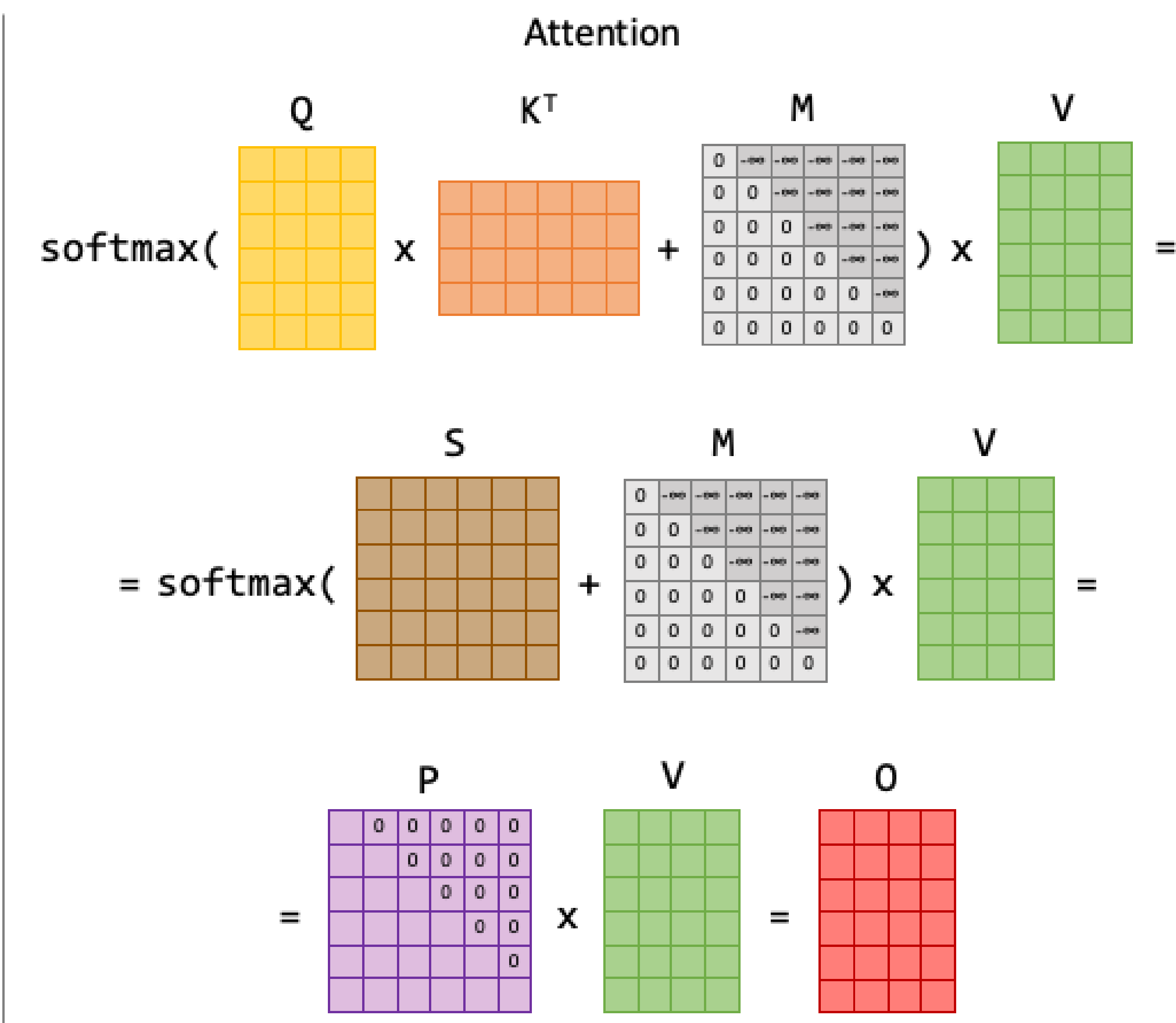
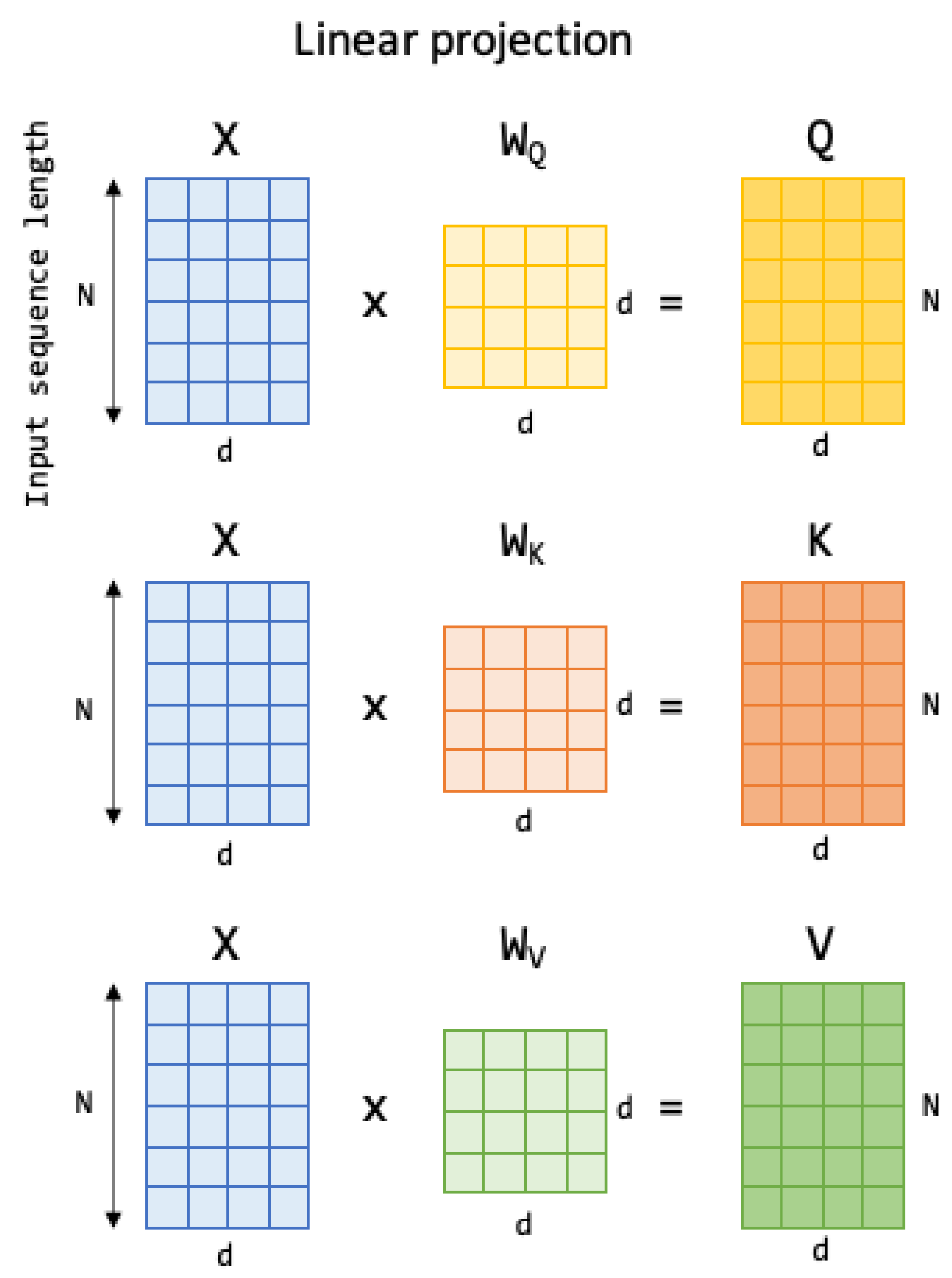
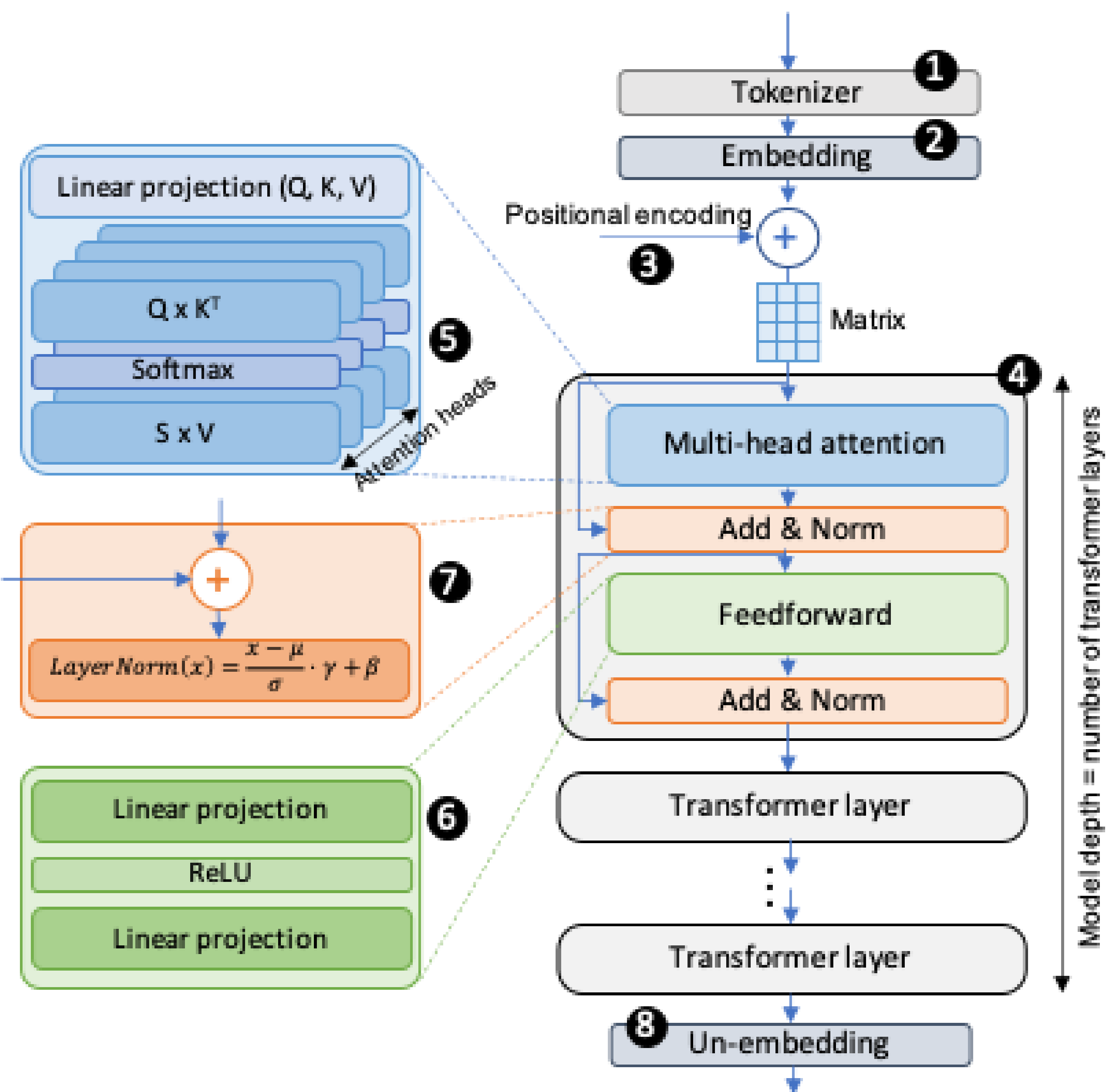
$$P(x_{n+t+1} || x_1, \dots x_{n+t})$$

Each output token depends on previous generated token

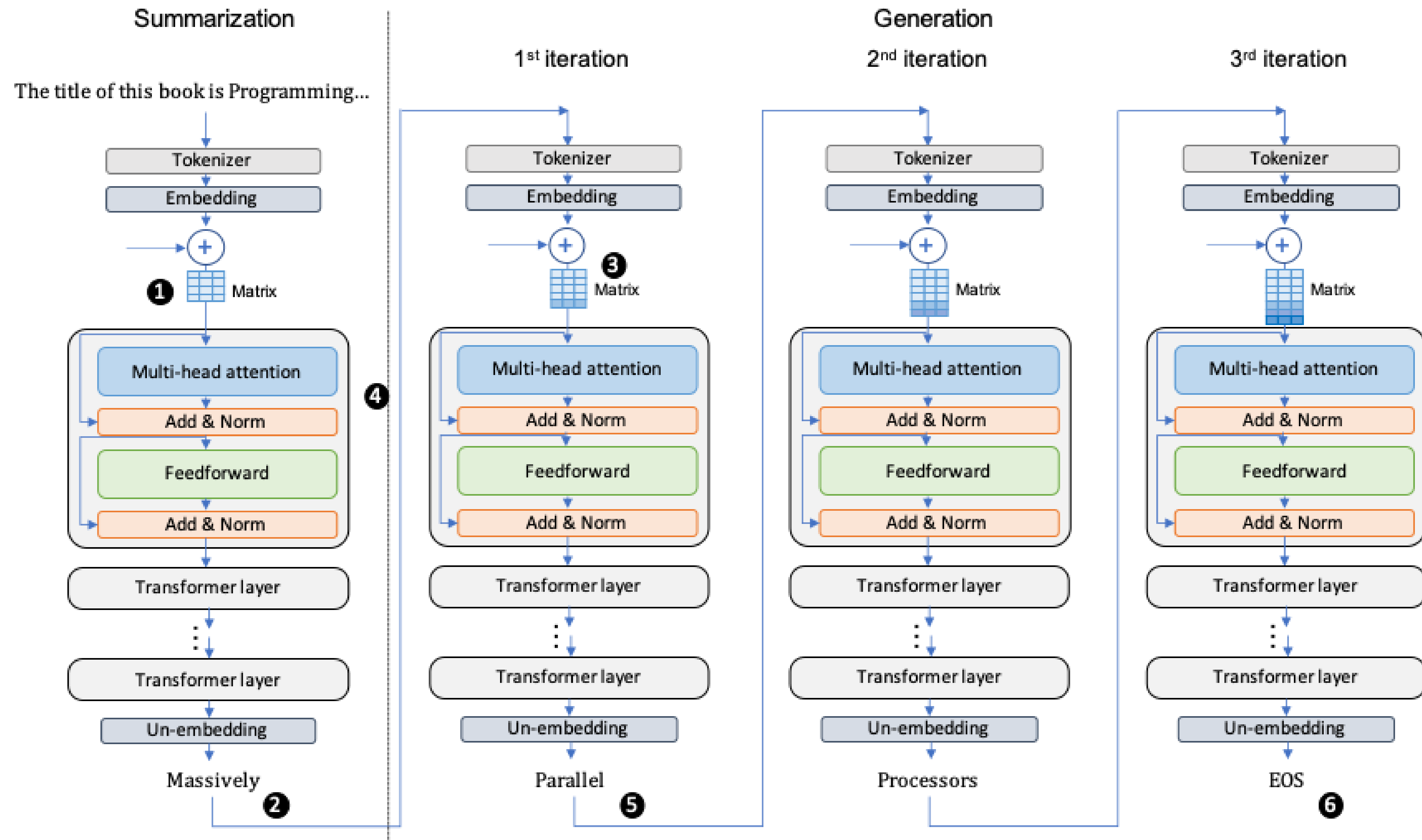
Prefill - Matrix multiplication

Decode - Matrix-vector multiplication

Transformer Computation

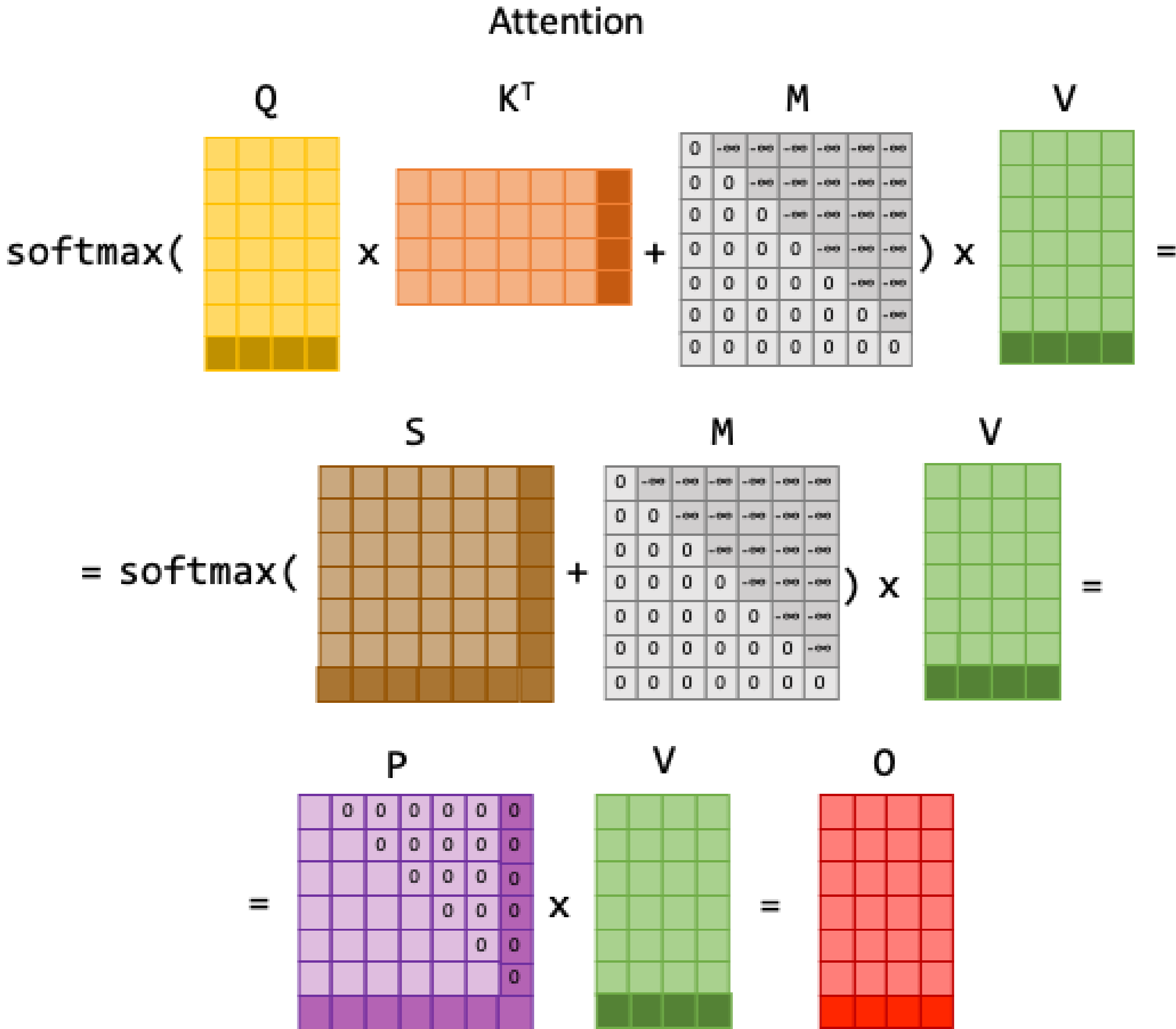
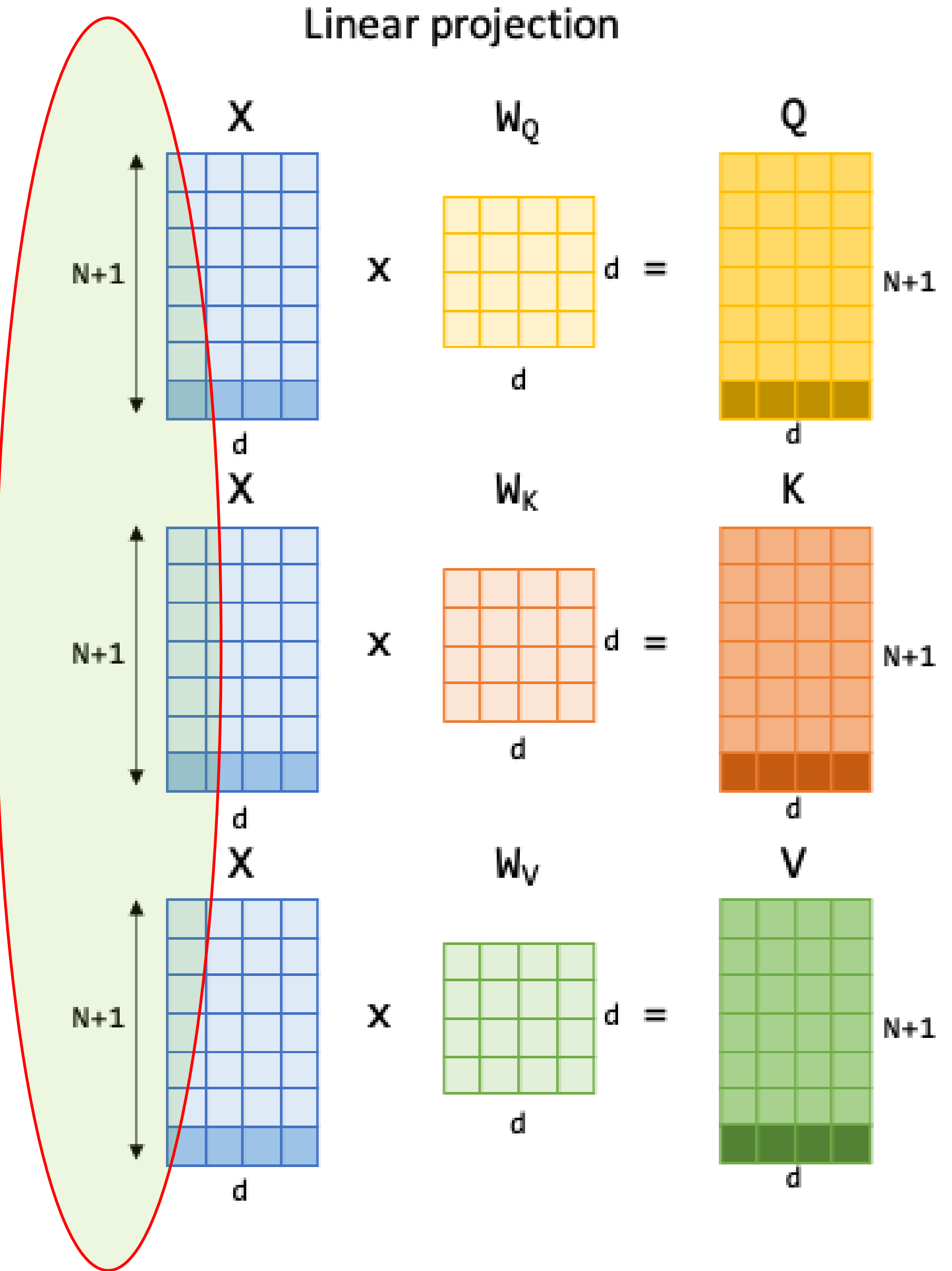


Next token prediction



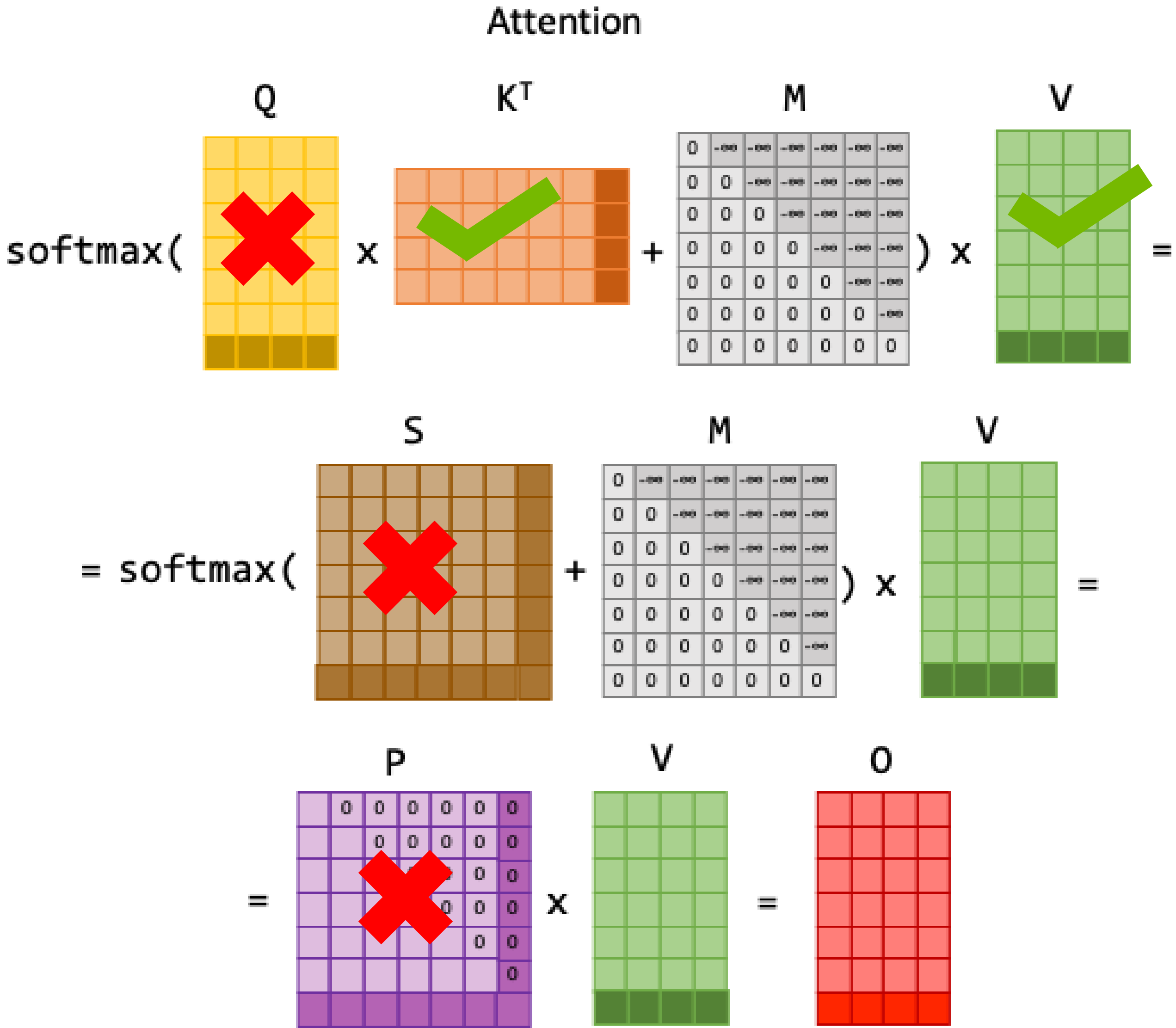
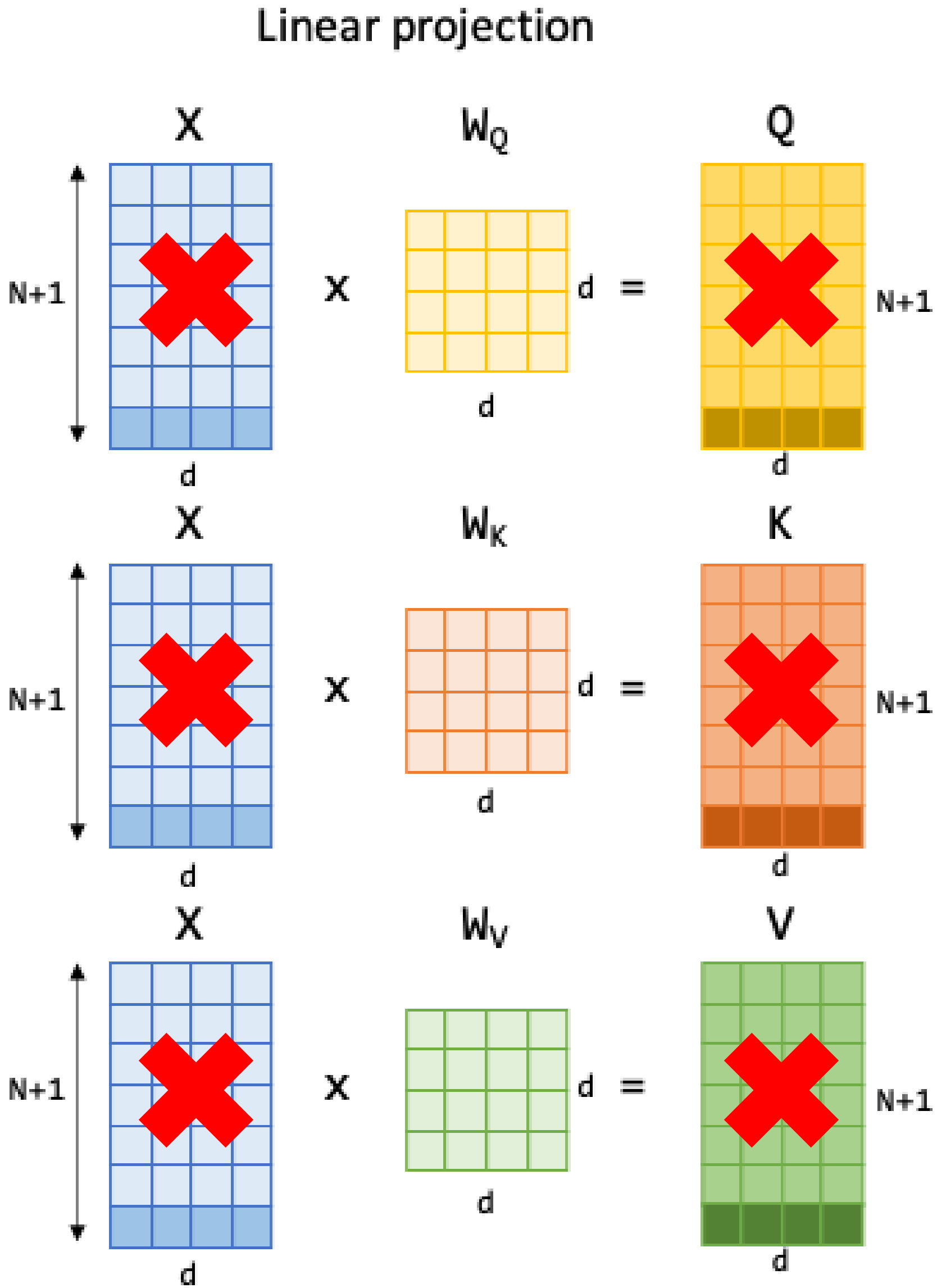
Transformer Computation

- Five matrix multiplications are performed in each transformer layer
 - Seq len (N) can grow to millions
 - GEMM becomes expensive
- Observation:
 - one additional row in the input matrix X from one iteration of the decoding phase to the next
 - Avoid repeating the computation exploiting GEMM property



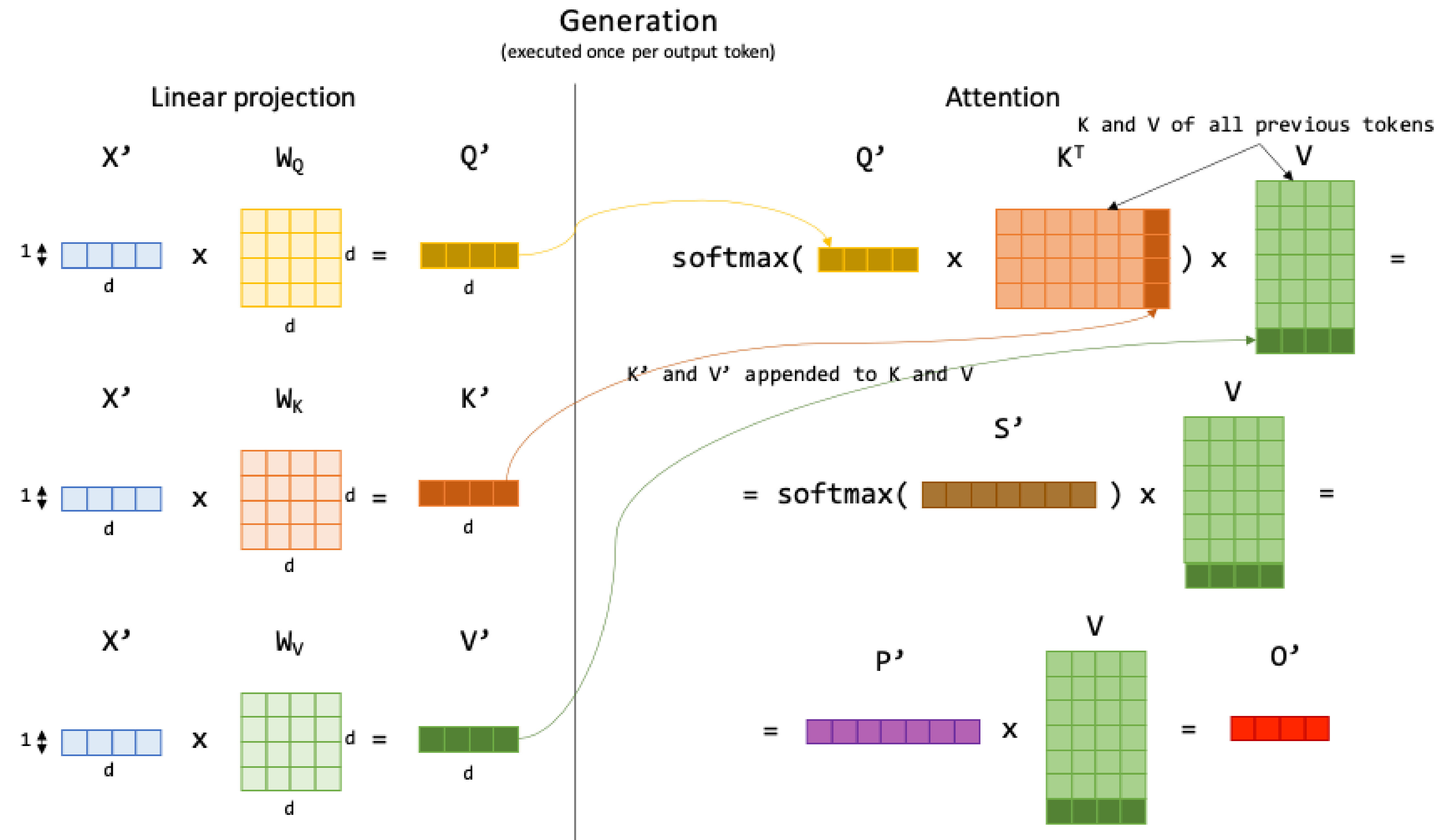
Transformer Computation

- Five matrix multiplications are performed in each transformer layer
- Seq len (N) can grow to millions
 - GEMM becomes expensive
- Observation:
 - one additional row in the input matrix X from one iteration of the decoding phase to the next
 - Avoid repeating the computation exploiting GEMM property



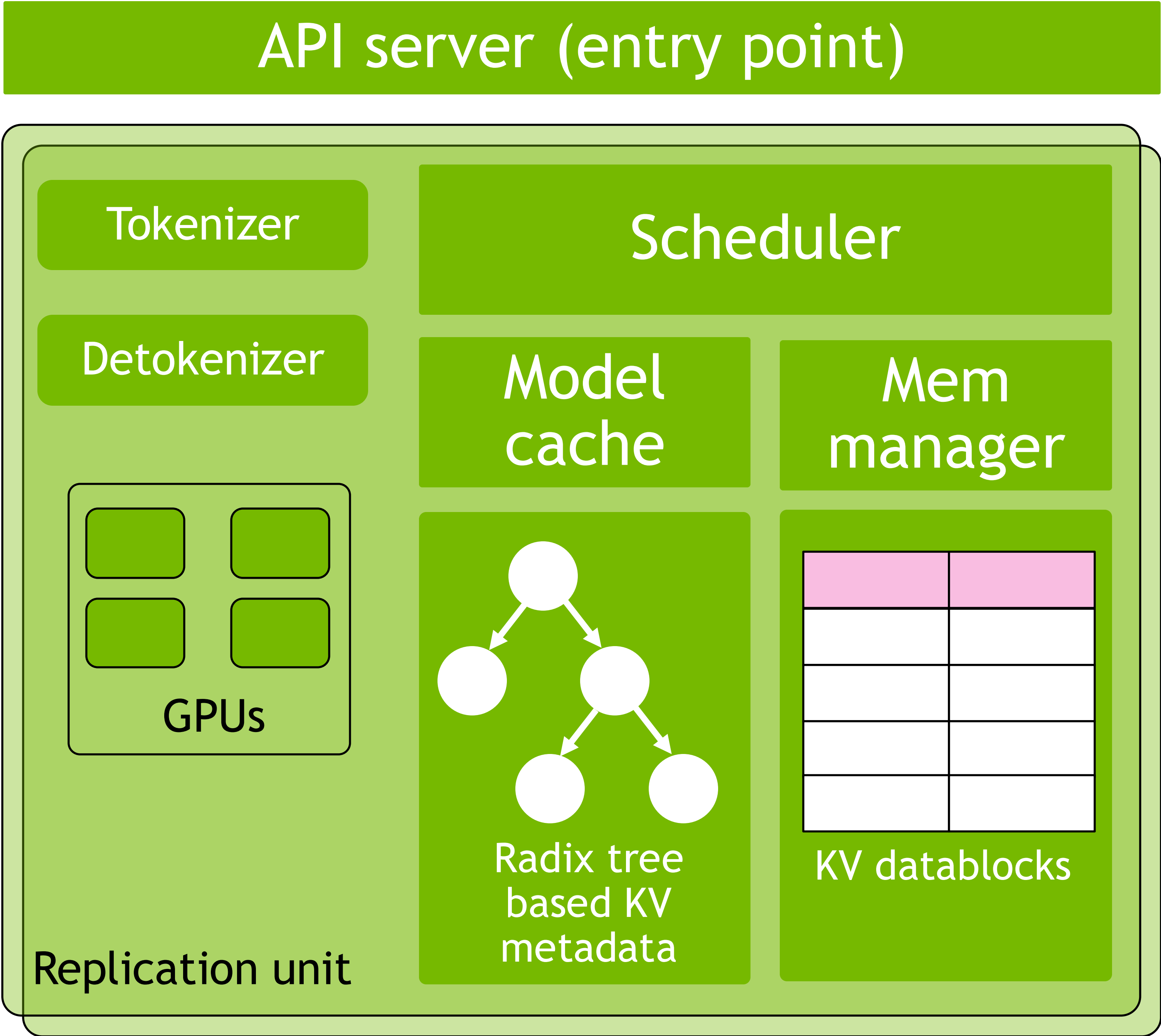
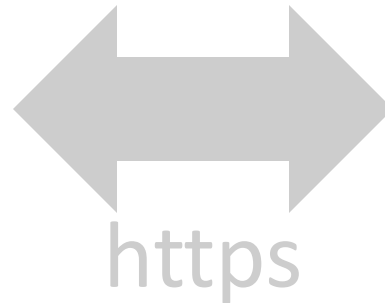
KV Caching

- Reusing the computation saves a lot of compute at the cost of memory
- We move from GEMM \rightarrow GEMV
 - Compute \rightarrow memory bound
- But memory needs to be allocated and managed efficiently
 - [PagedAttention](#) (vLLM)
 - [RadixAttention](#) (SGLang) across requests



LLM Inference Framework

An example using SGLang



Naïve K8s based
autoscaling approach

Good Model != Success

Good Model + Good System == Success

^
That hacks the inference
pareto frontier

Applications have Diverse SLO

TTFT = Time for First Token, TPOT = Time Per Output Token

Chat

Search

Program

TTFT

TPOT

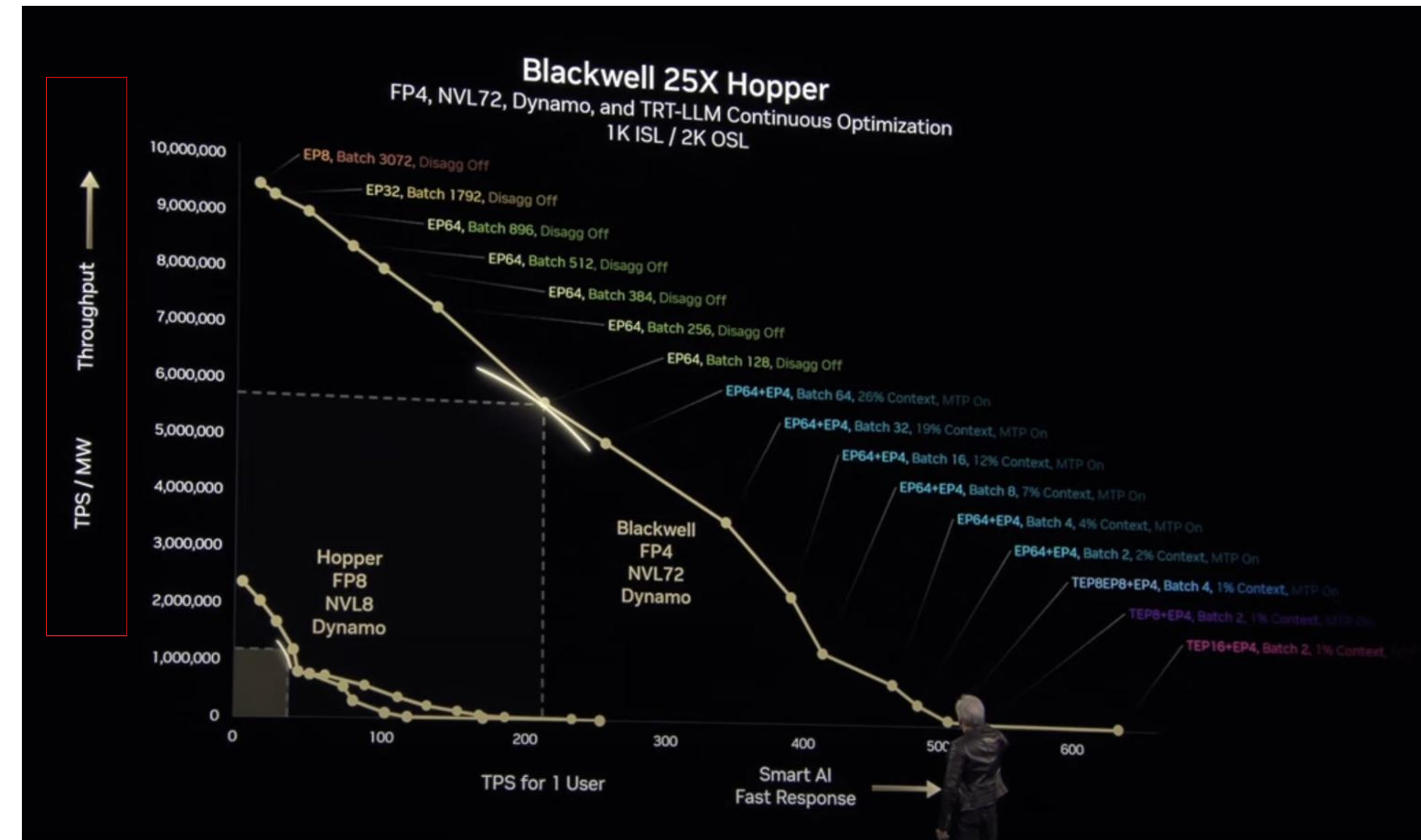
Low ($< 1s$)	Match read speed (~ 100ms)
Very Low (~ 200ms)	Match read speed (~ 100ms)
Very Low (~ 200ms)	Very Low (~ 50ms)

Measure Serving Cost

Throughput = # Requests / sec

≈ 💰 per request

Throughput as a proxy of dollar per request



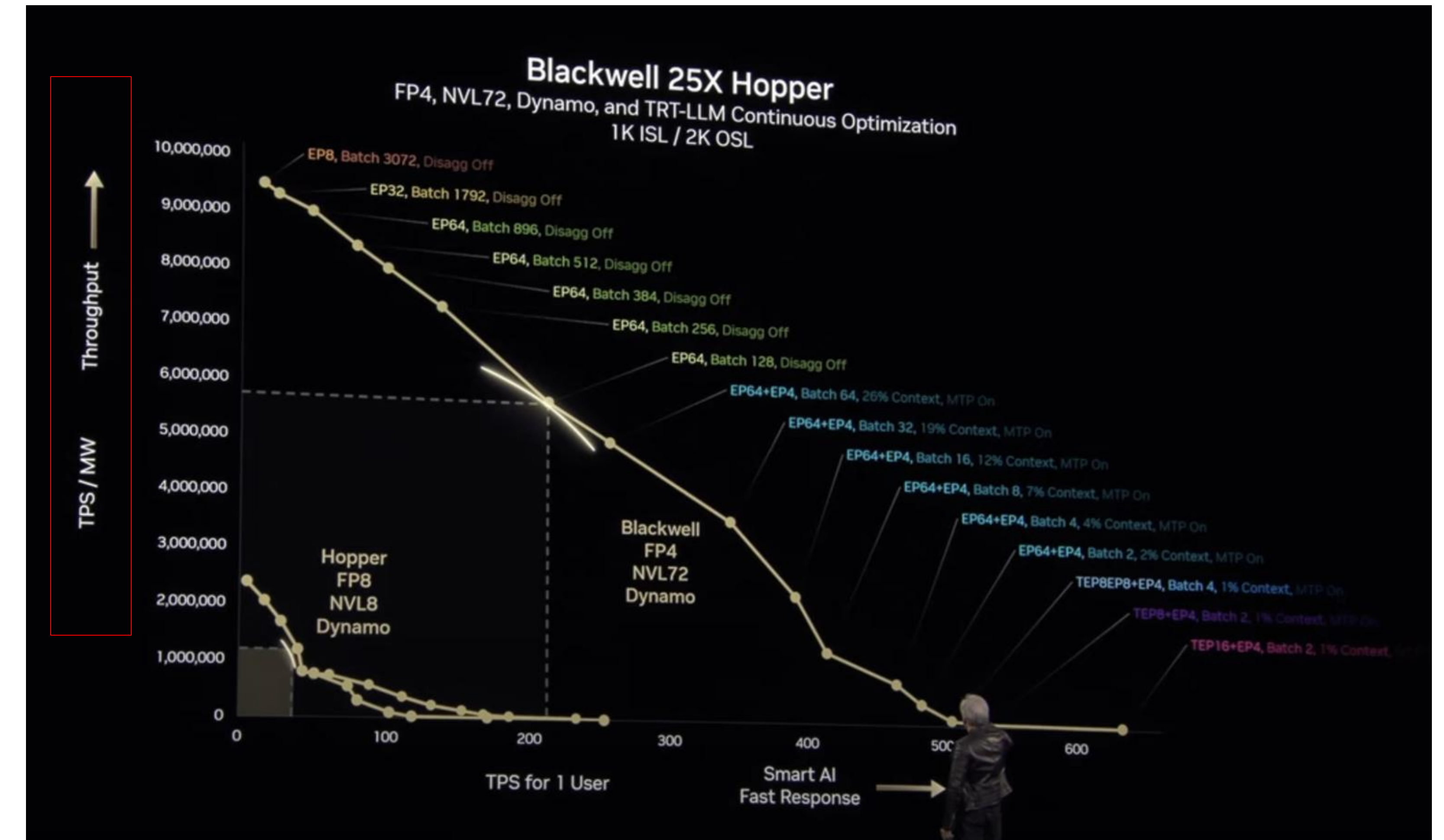
Jensen, GTC'25

Measure Serving Cost

Throughput = # Requests / sec

≈ 💰 per request

Throughput as a proxy of dollar per request

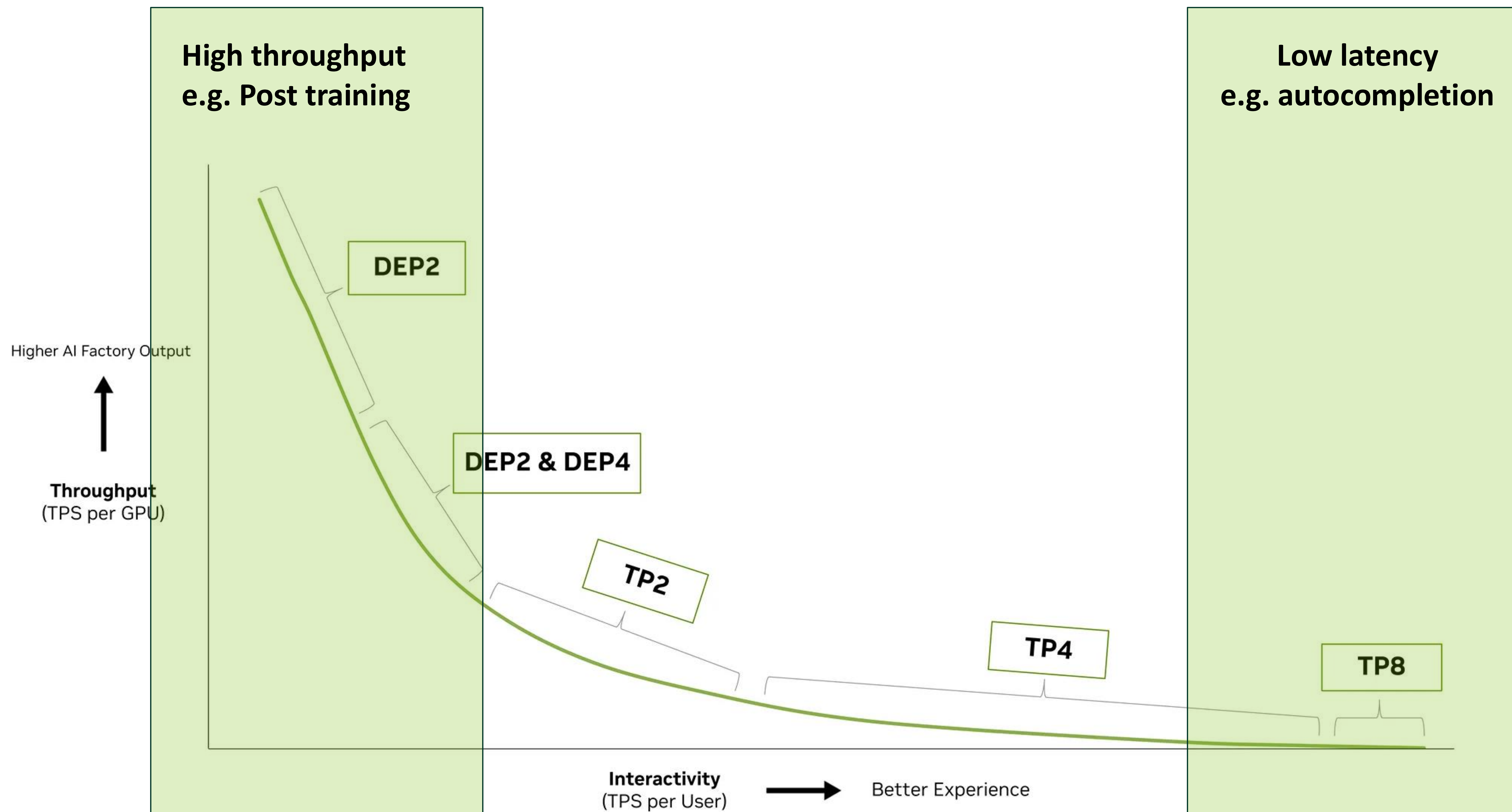


Jensen, GTC'25

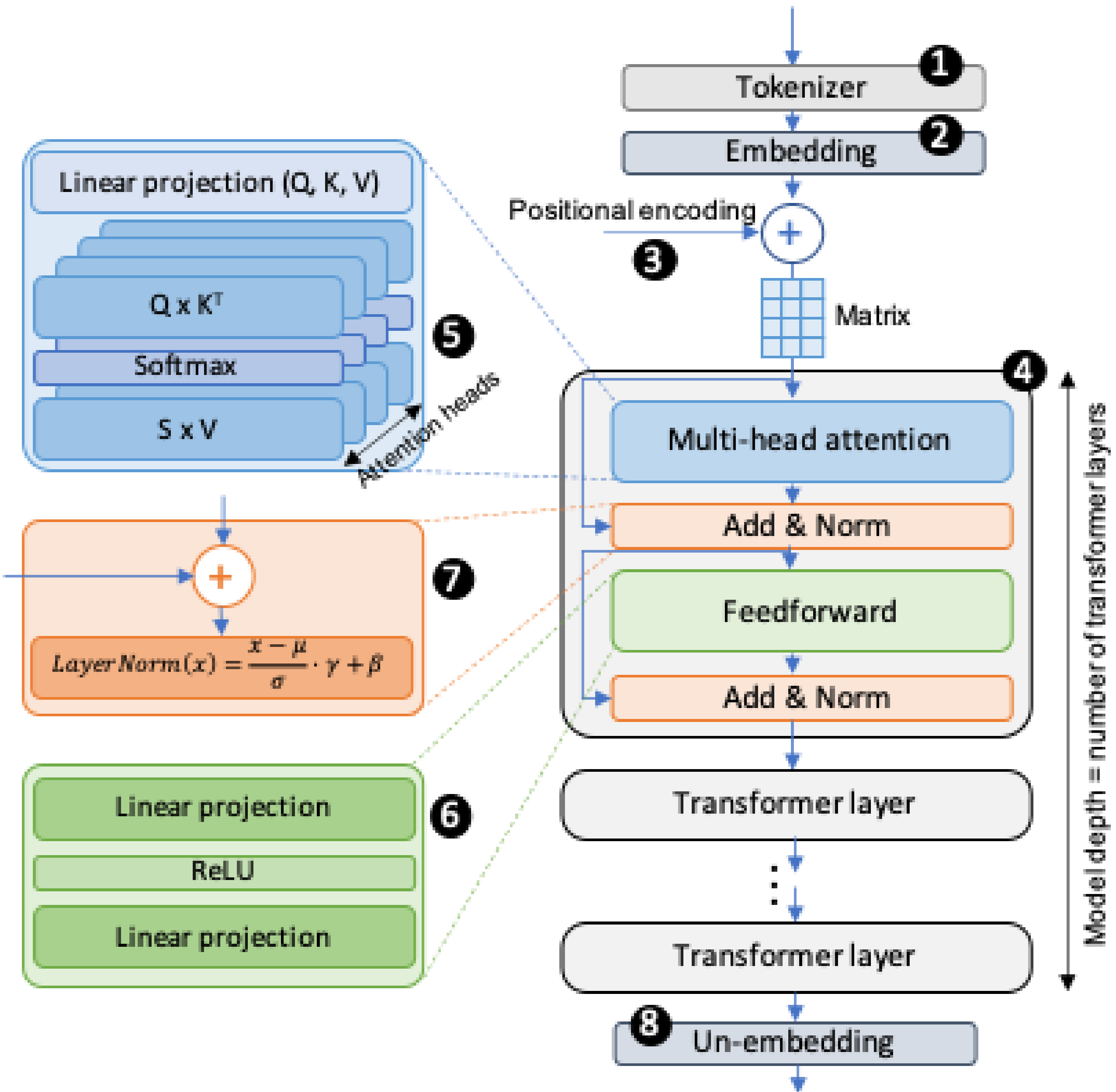
Goodput = # Requests / sec satisfying the SLO constraint

≈ 💰 per good request

Pareto Curve



How would you build an inference service at scale?



Challenges Serving AI Inference at scale

- Need to vary the parallelism strategy for hitting performance SLAs
- Multiple requests with varying degrees of ISL/OSLs
 - Some doing multiturn
 - Some from agents
- Heterogenous hardware
 - Different generations and across platforms
- Need for fault tolerance
- Multiple models to serve
- AI factory TCO is key at scale
 - Optimize for goodput
 - Optimize for power
 - Optimize for maximizing KV cache hit rate
- ...

Inference Demand Analysis – Broad Estimates

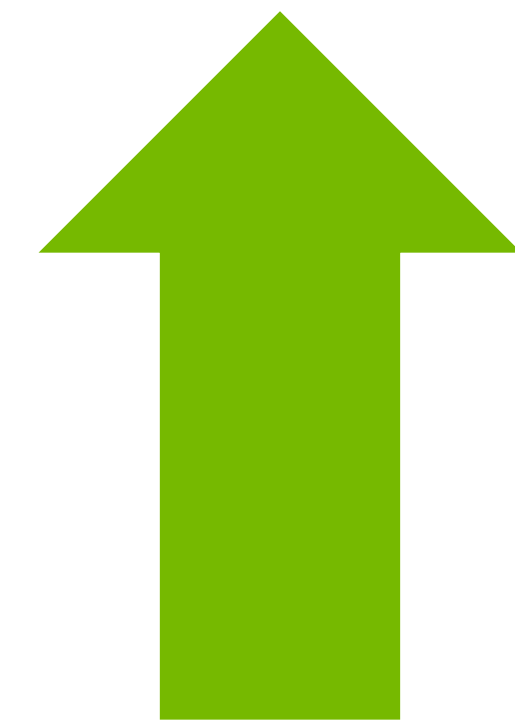
Inference Tokenomics - OAI Inference Economics						
Measure	Unit	2025	2028	2030	The Internet	Analysis Method
Global (Ex China) Population	M	6,770	6,852	6,906	6,770	OpenAI Data Disclosures
WAU Penetration Rate	%	12.9%	35.5%	44.3%	66.0%	
Year End Weekly Active Users (WAU)	M	875	2,435	3,061	4,468	
MAU to WAU Ratio	%	80.0%	80.0%	80.0%	80.0%	
Monthly Active Users (MAU)	M	1,094	3,044	3,826	5,585	AI Tokenomics Model
ChatGPT M Tokens per MAU	MTok/MAU/Month	1.95	2.99	3.74		
ChatGPT Tokens Per Day	Trillion Tok/Day	71	303	477		
API Tokens Per Day	Trillion Tok/Day	12	237	482		InferenceMAX™ Benchmarking
GPUs Required	GPUs	1,298,905	3,918,757	3,666,803		
AI Datacenter Capacity	MW	1,370	8,794	12,009		AI Datacenter Model

1. DeepSeek R1 FP8 on H200 using SGLang. Uses 8k input, 1k output tokens, 43 interactivity. For 2025 only, future periods are for different models and GPU assumptions.

Source: Semianalysis OCP'25

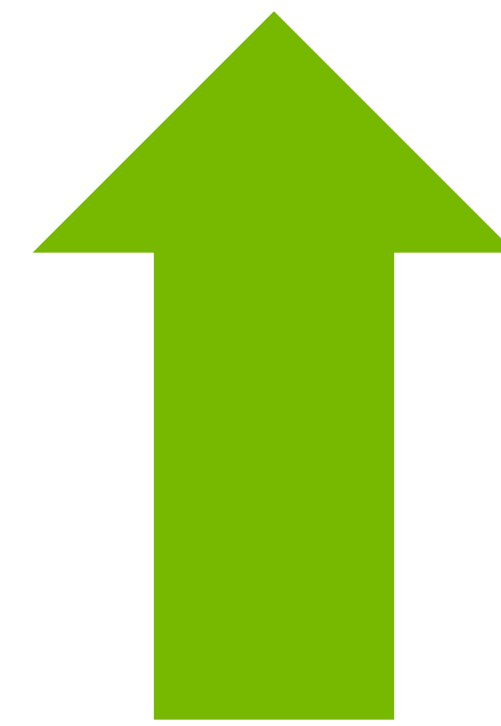
LLM Deployments in 2025

Scale Never Seen Before



Input Length

1M+ with RAG, Video understanding, Code...



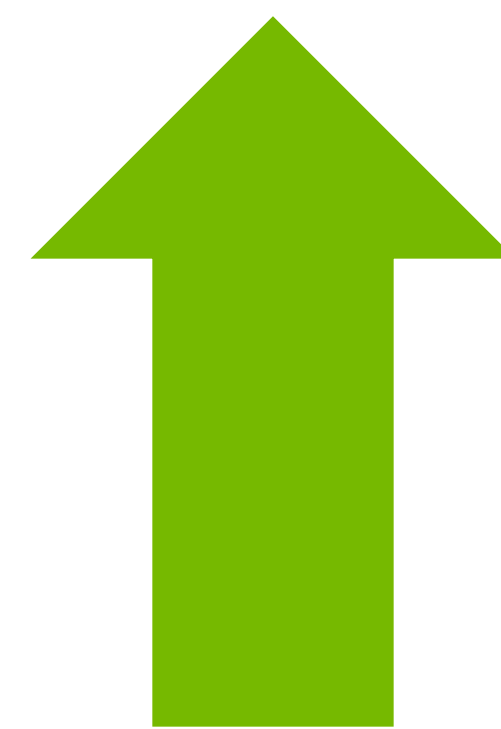
Output Length

Reasoning, Agentic Workflows...



Model Size

O(100B) – O(1T) parameters



Users + UX

Fast tokens served = Revenue generated

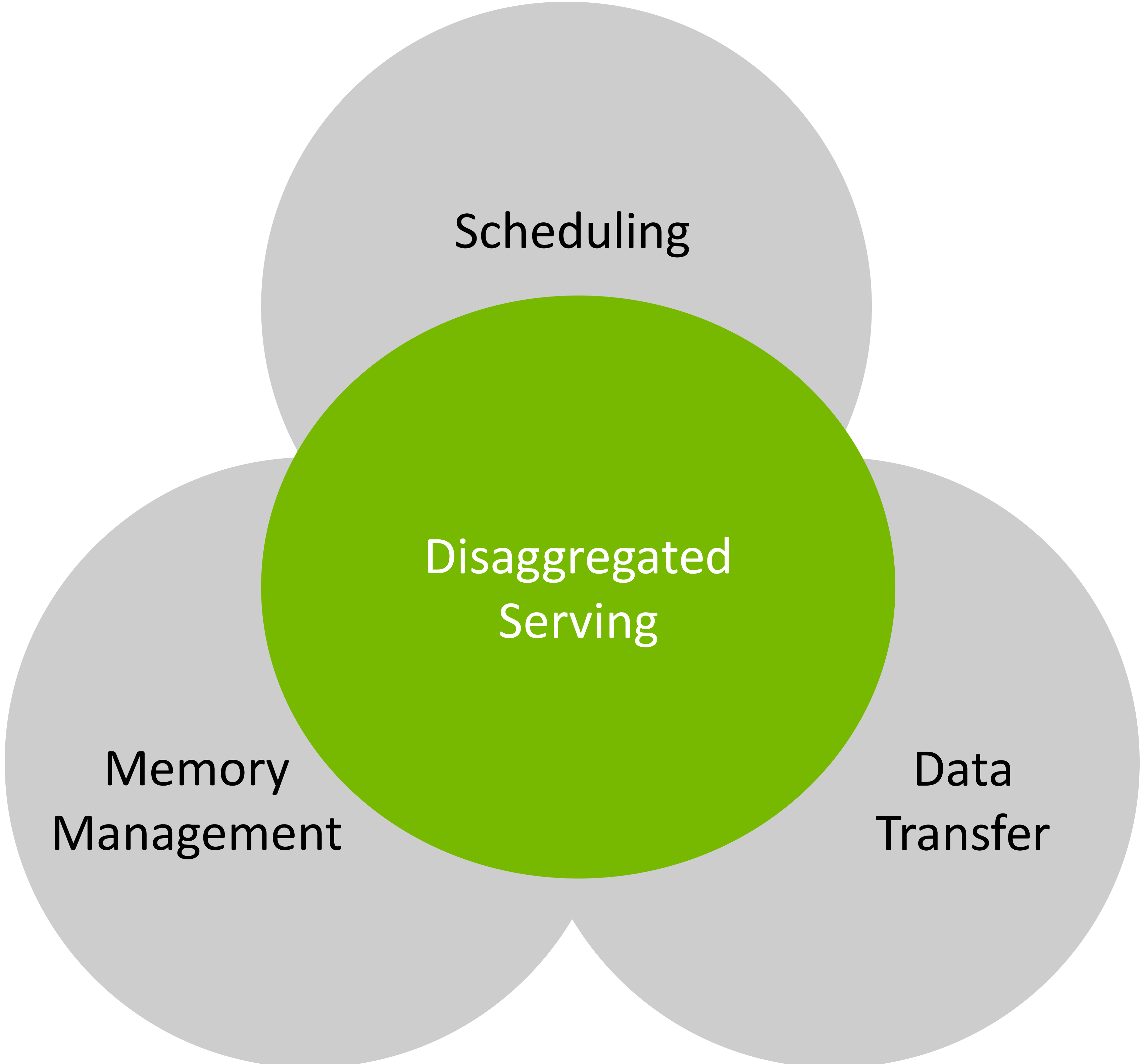


Deployment Constraints

HW availability, power, rack space, cooling...

NVIDIA Dynamo: Distributed Inference At Scale

Systematic approach to AI inference at scale



System level optimization

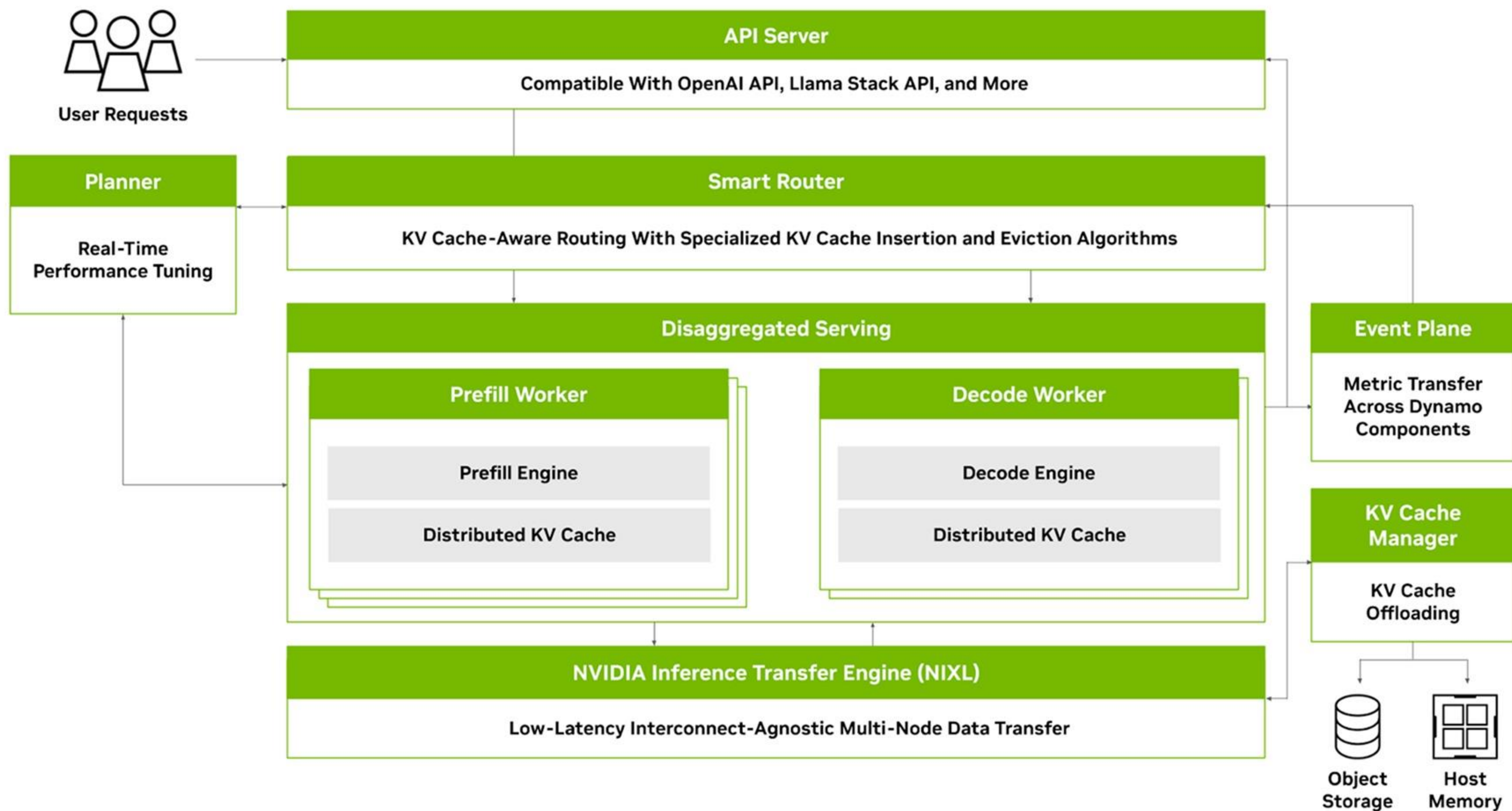
Supports all inference engines

Modular Components

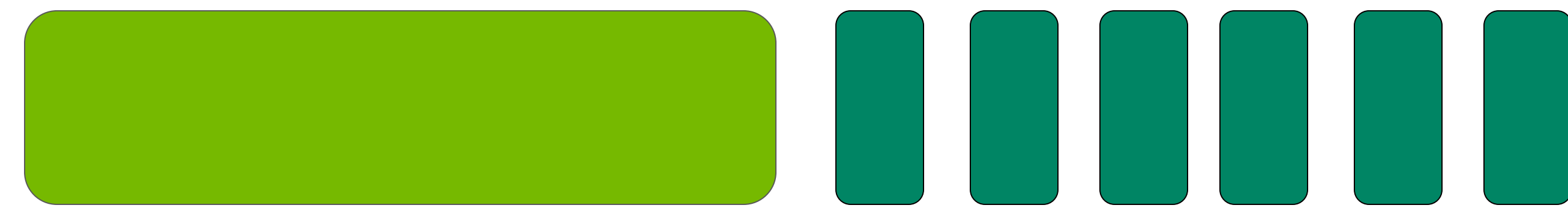
Production grade serving

NVIDIA Dynamo: Distributed Inference At Scale

Modular design, independently scalable yet highly efficient



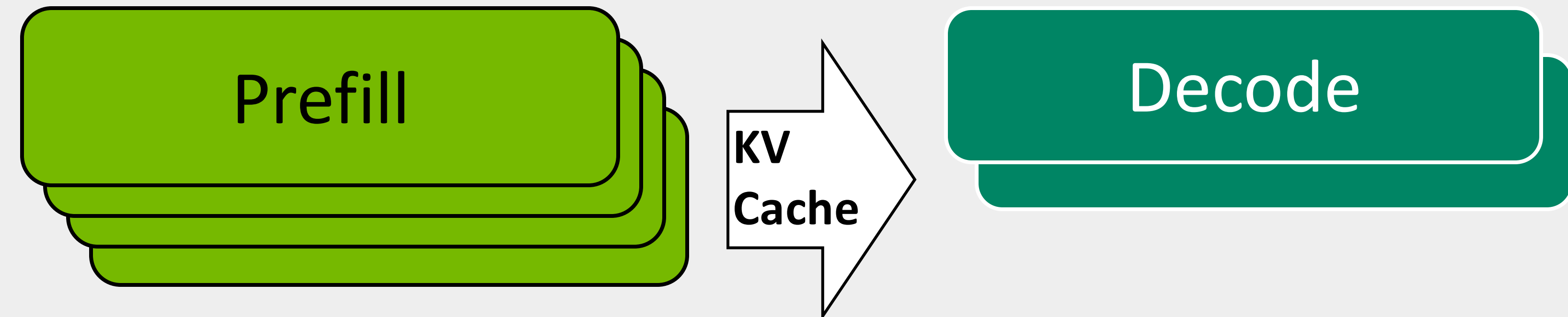
Disaggregated serving



Prefill
Compute
bound

Decode
≠ Memory bound

Scale prefill & decode independently
on separate GPUs



Apply suitable parallelism for prefill and decode
separately

Why do Disaggregated Serving?

Q1. Prefill and decode interference.

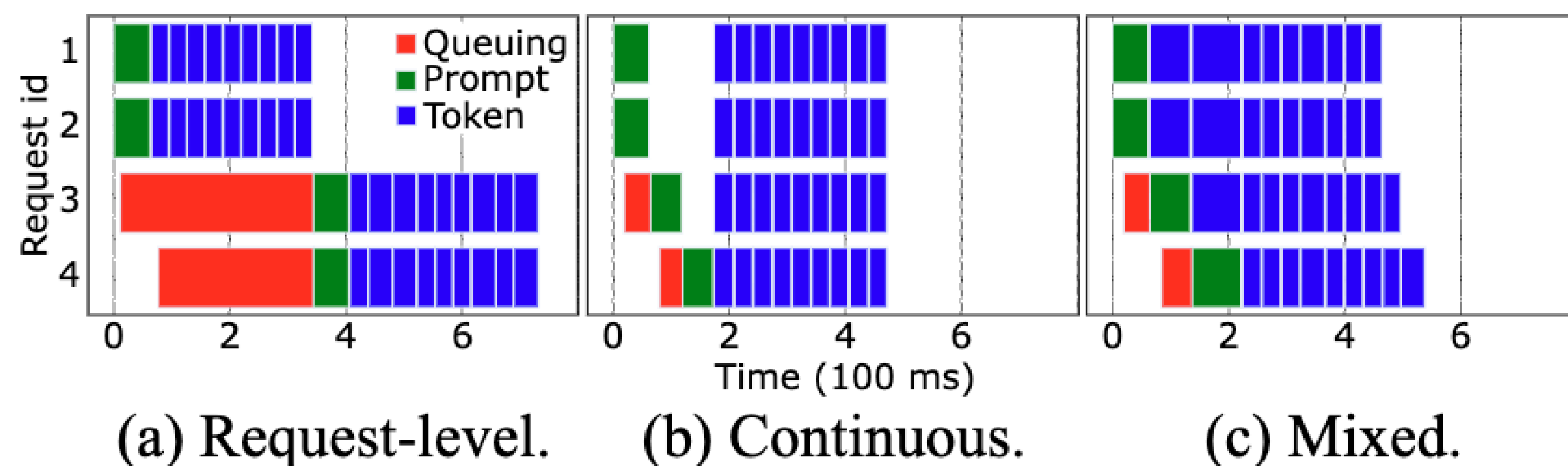
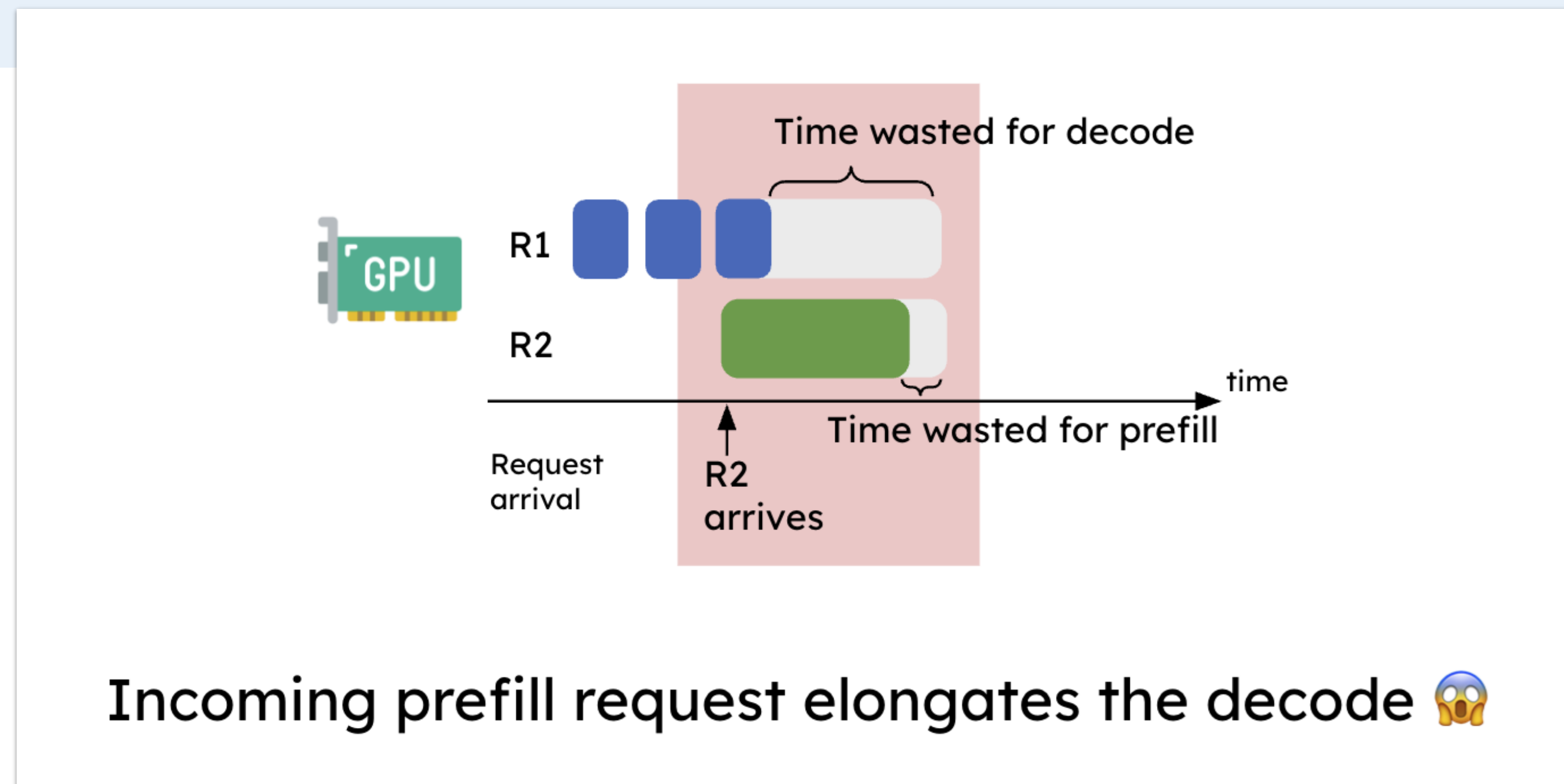


Fig. 2: Batching mechanisms and their latency impact on the prompt and token phases.

Q2. Resource and Parallelism Coupling

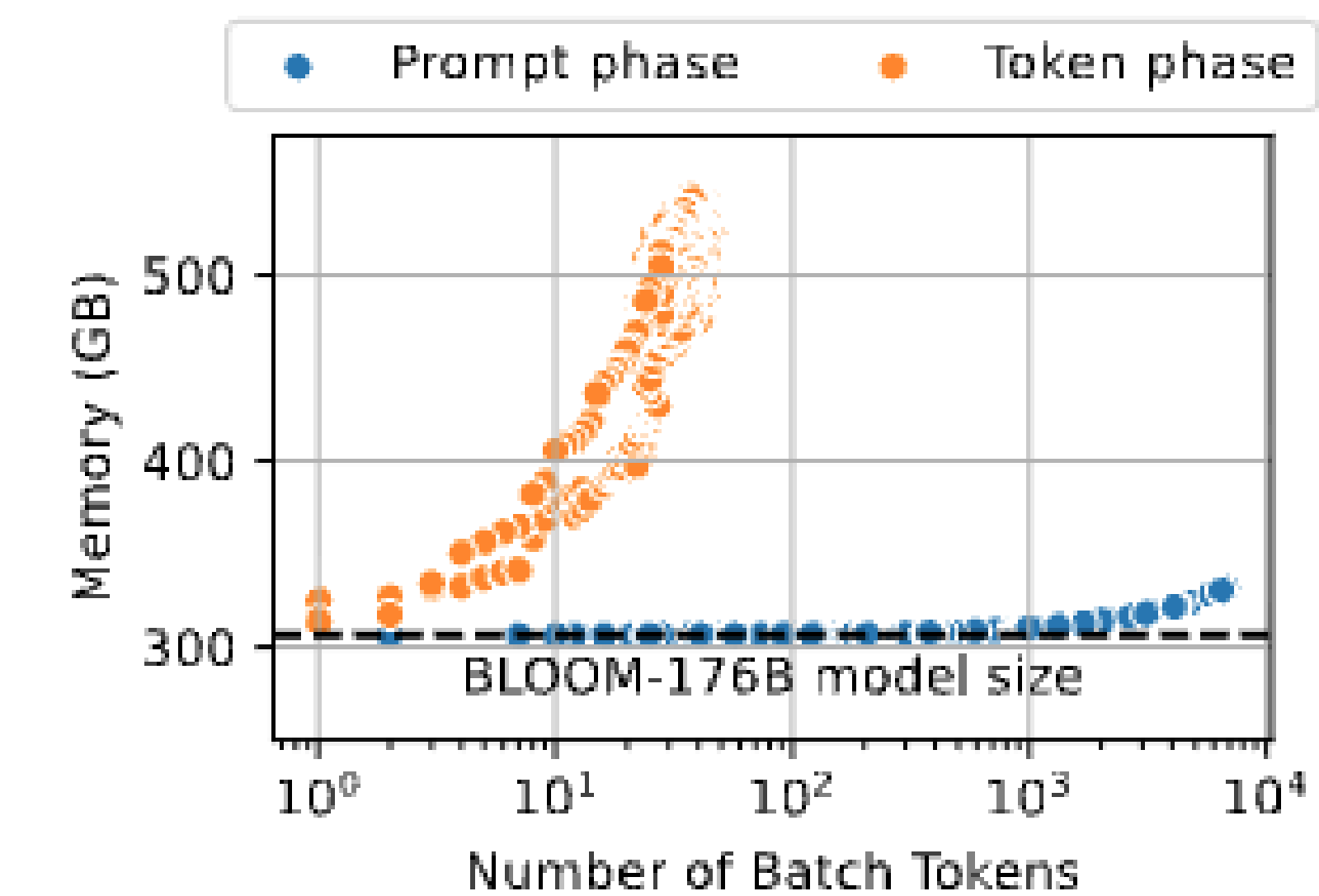
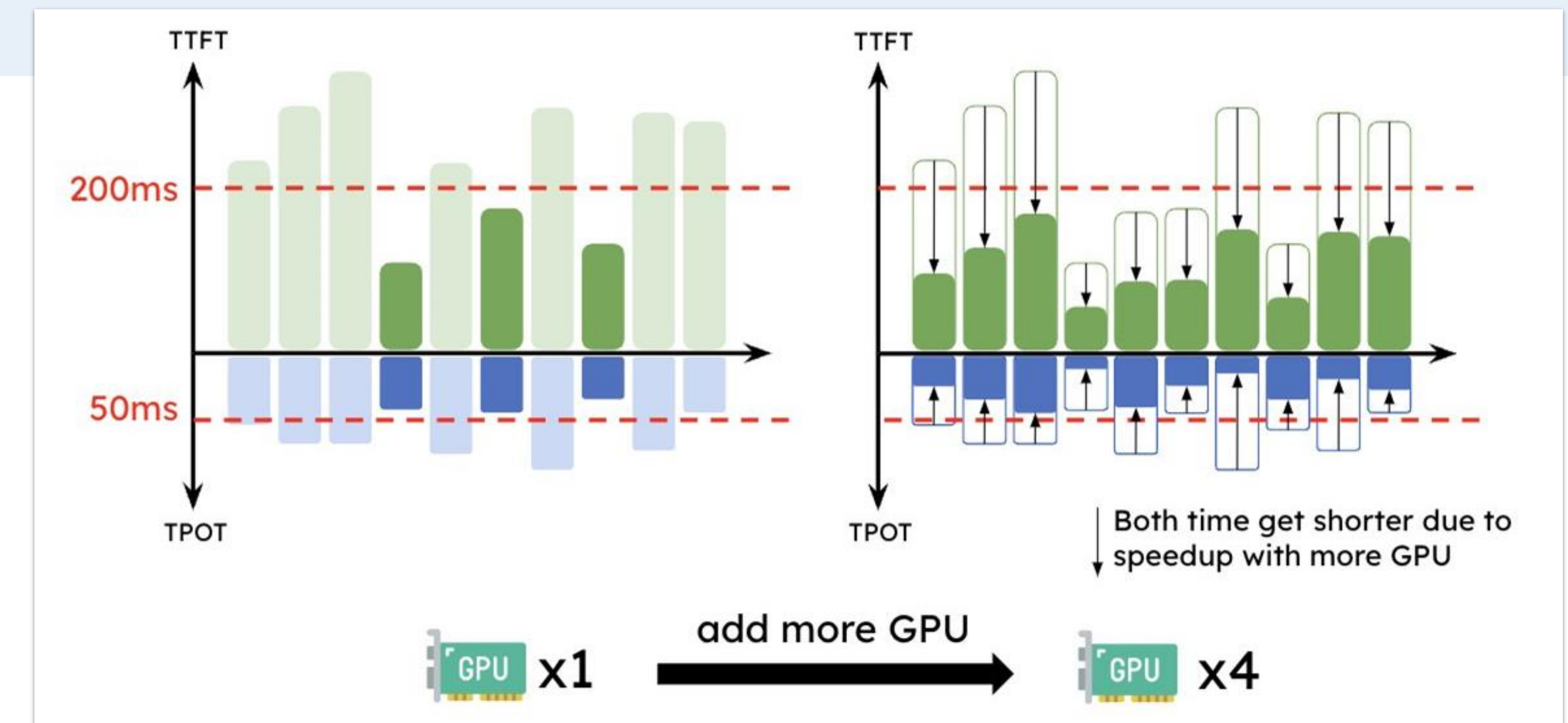
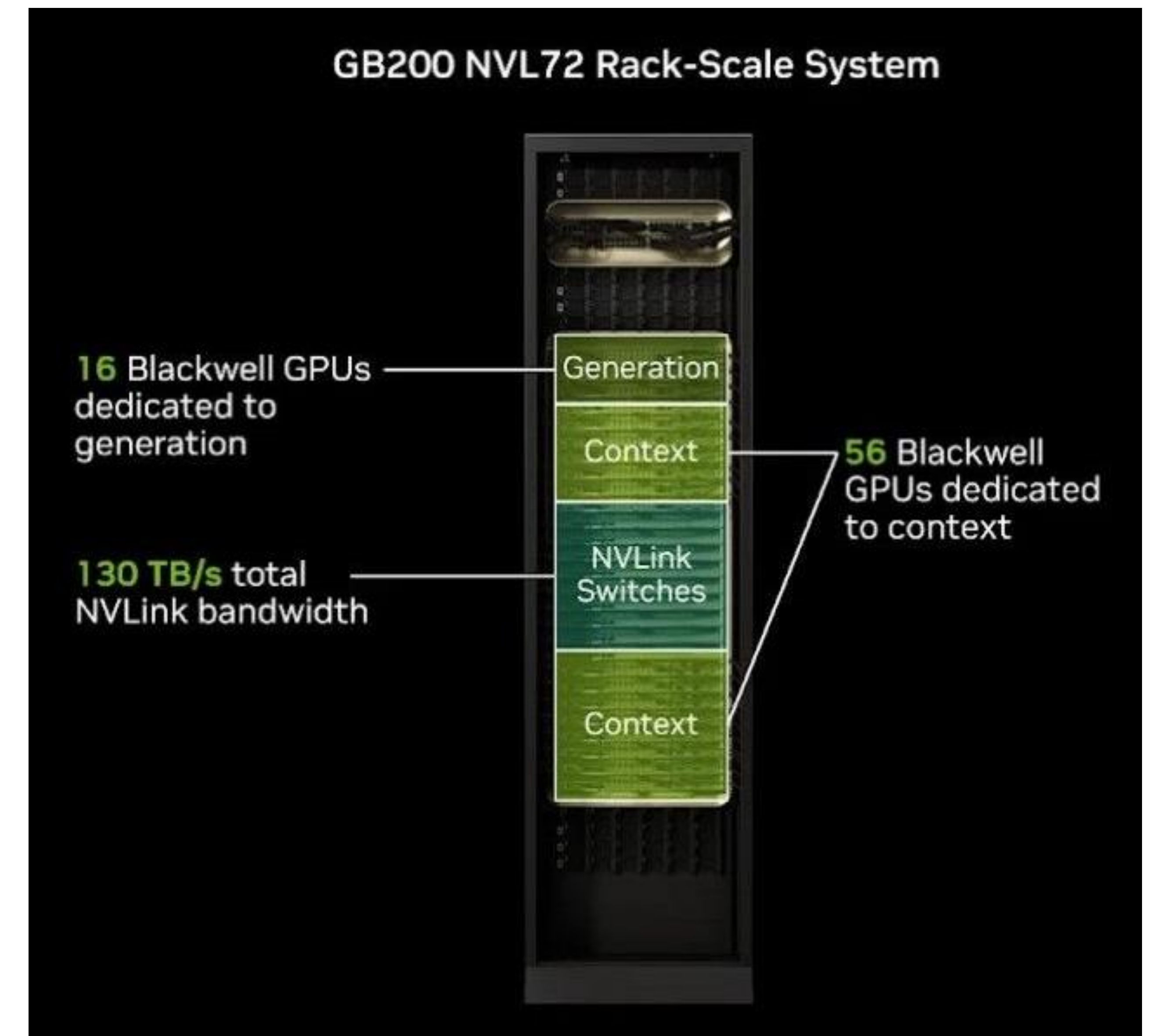


Fig. 7: Required memory with batching in prompt/token phases.

Implementing Disaggregated Serving

- Good proxy example for learning
 - https://docs.vllm.ai/en/latest/examples/online_serving/disaggregated_serving/
- Use for production
 - https://github.com/ai-dynamo/dynamo/blob/main/examples/backends/vllm/deploy/disagg_kvbm_2p2d.yaml

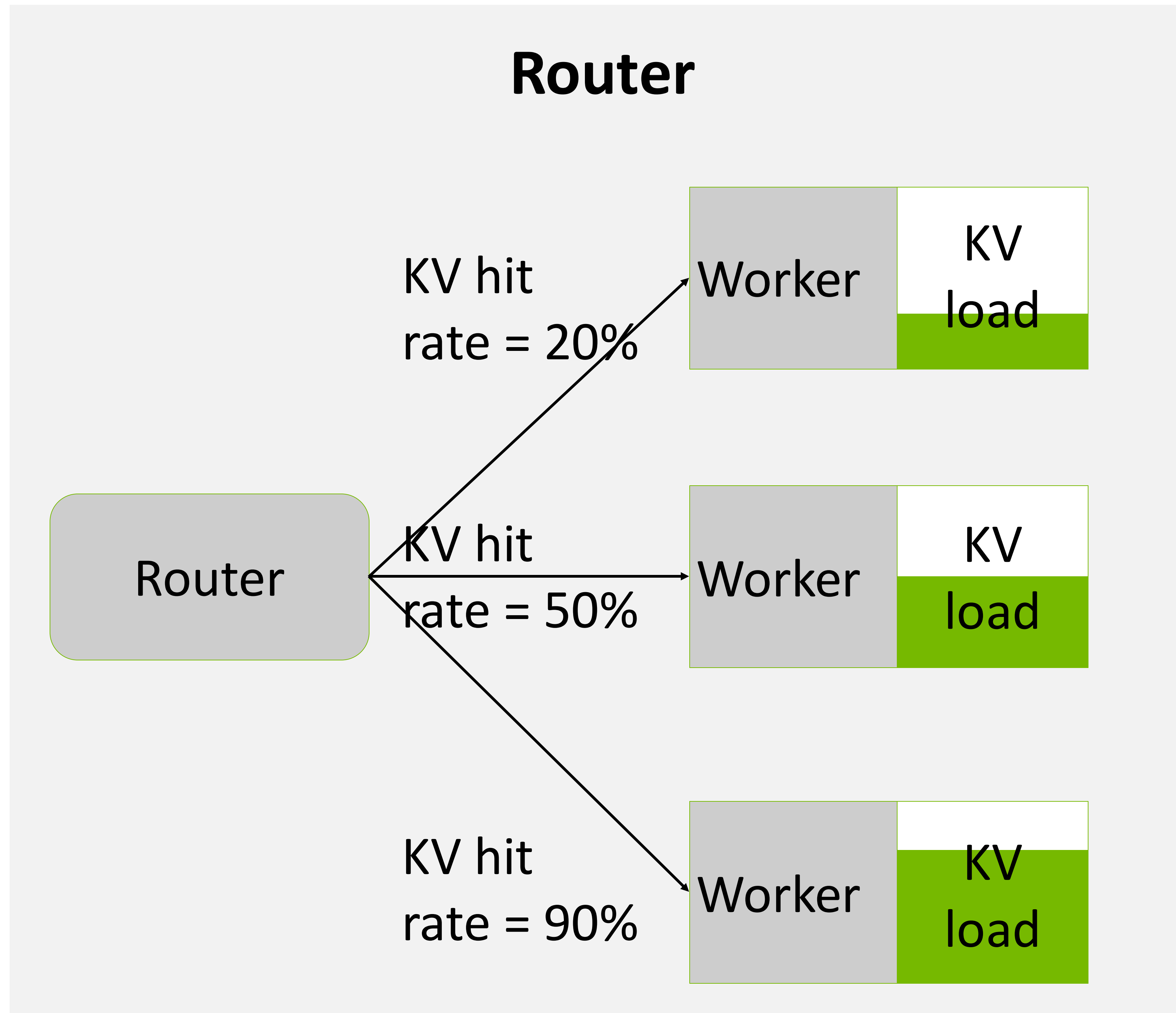


Disaggregated Inference is one feature of Dynamo

At scale inference needs solving the AI factory TCO optimization problem

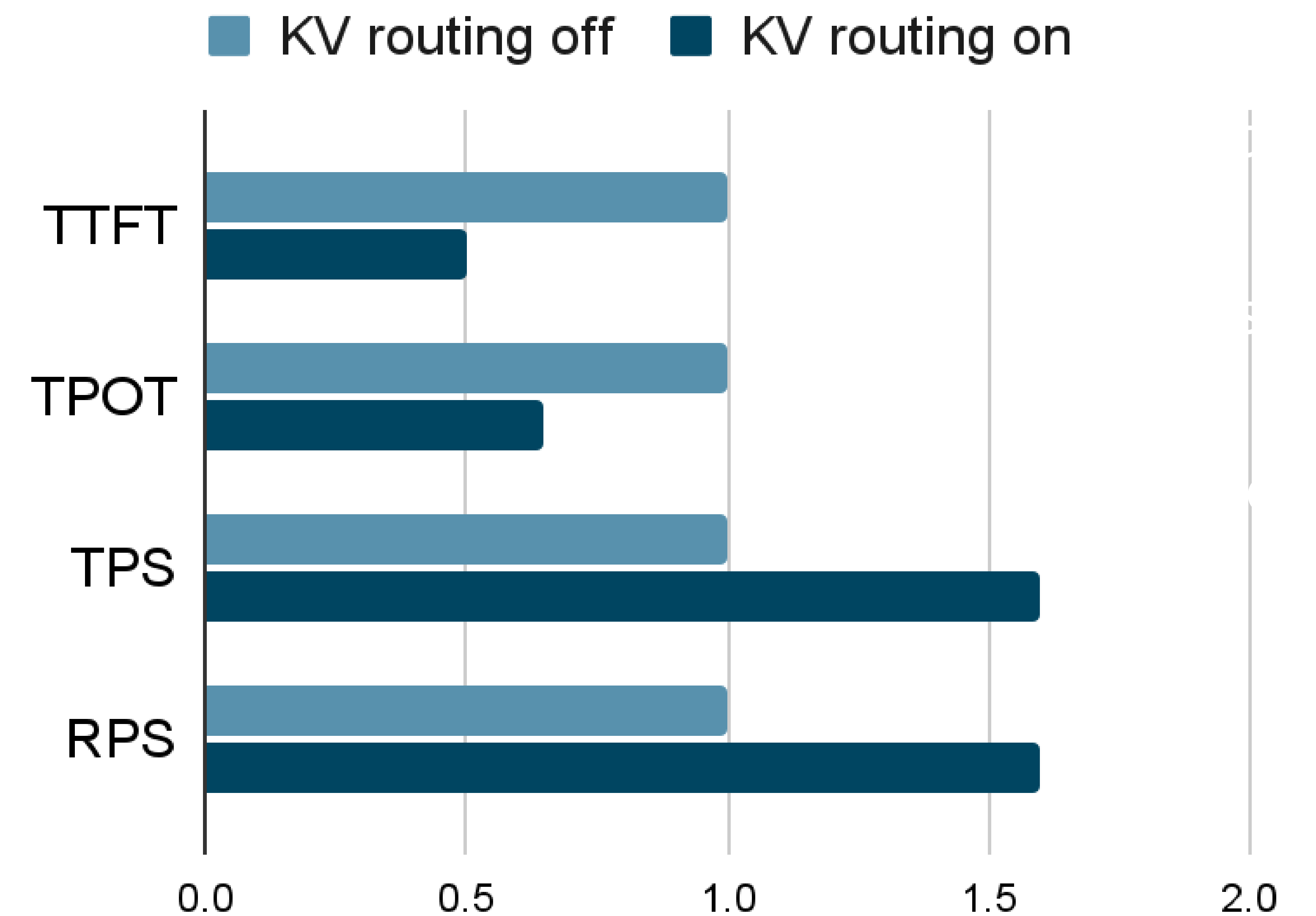
Scheduling

KV aware routing archives 2X TTFT, 1.5X TPOT and 1.6X TPS for Qwen3 480B at Baseten



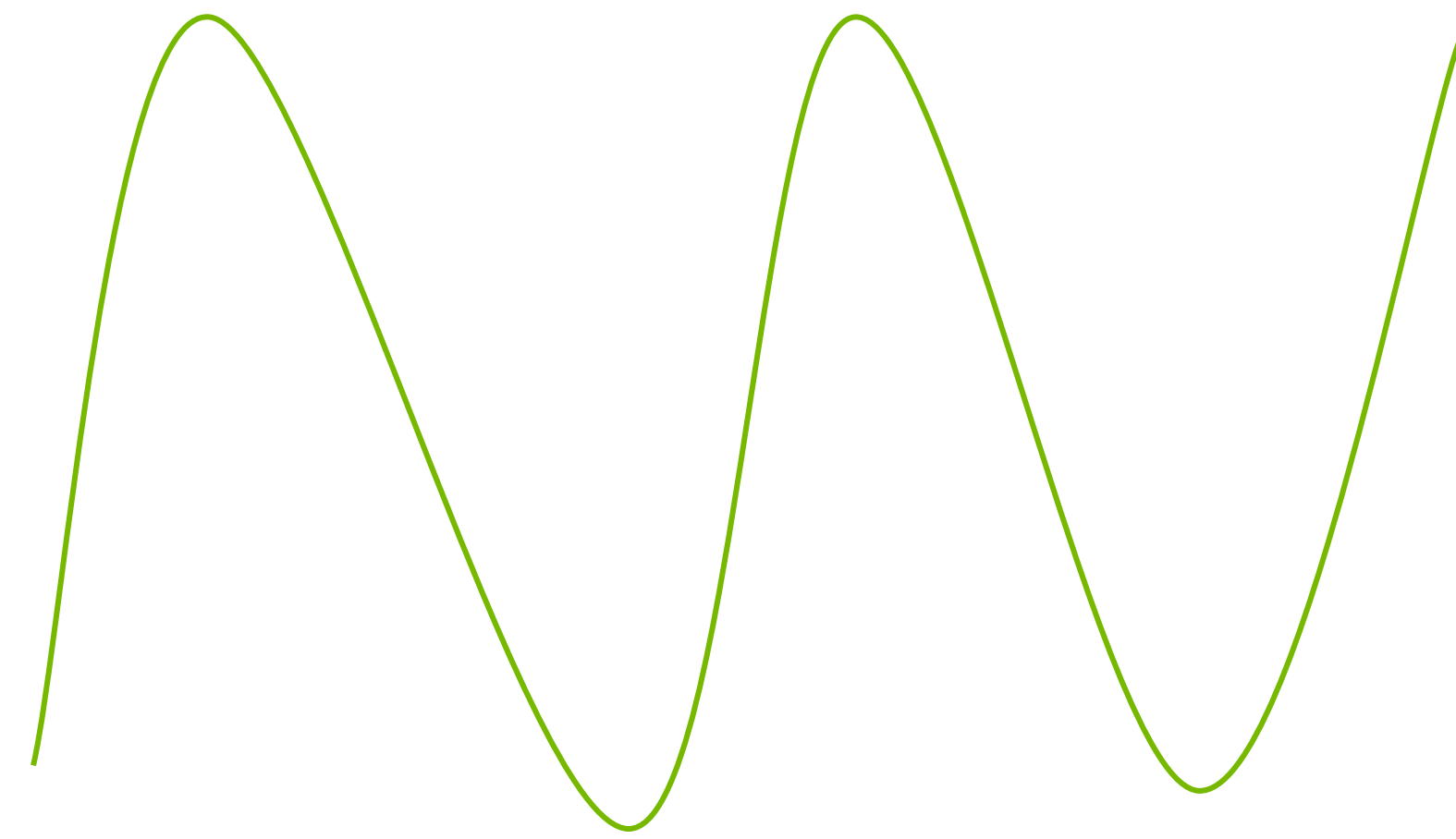
 baseten

Qwen3 480B Coder KV aware routing performance at Baseten

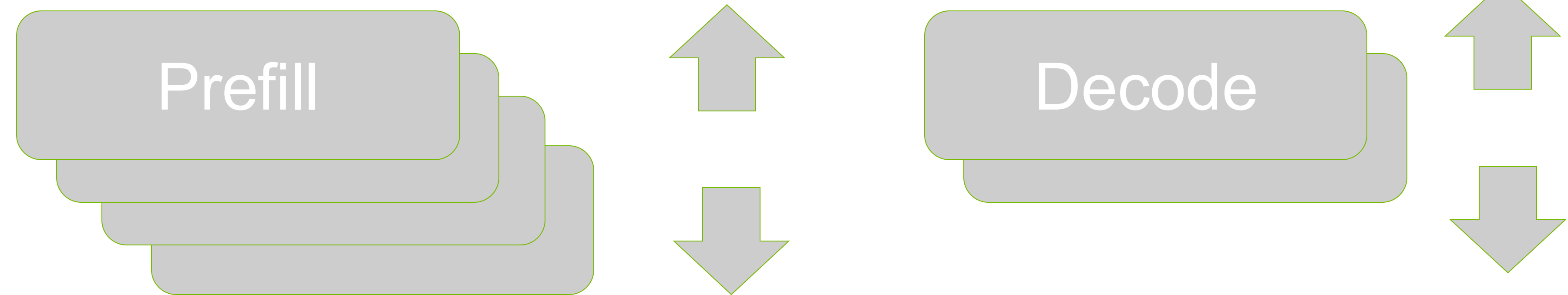


Online perf optimizer

LLM traffic can change rapidly



Planner autoscales in real time w.r.t SLA (TTFT & TPOT)



Initial profiling

Metrics Collection
TTFT/ITL

Load Predictor
ARIMA/Prophet

Interpolator
With profiler info

Scaling logic

Alibaba ACK

80% reduction in TTFT
With 5% less GPUs

Based on performance benchmarks presented by Alibaba at APSARA conference on 2025.09.24.

Source: <https://yunqi.aliyun.com/2025/session?agendaId=6062> at 2 hr 50 min mark

Offline Perf Configurator

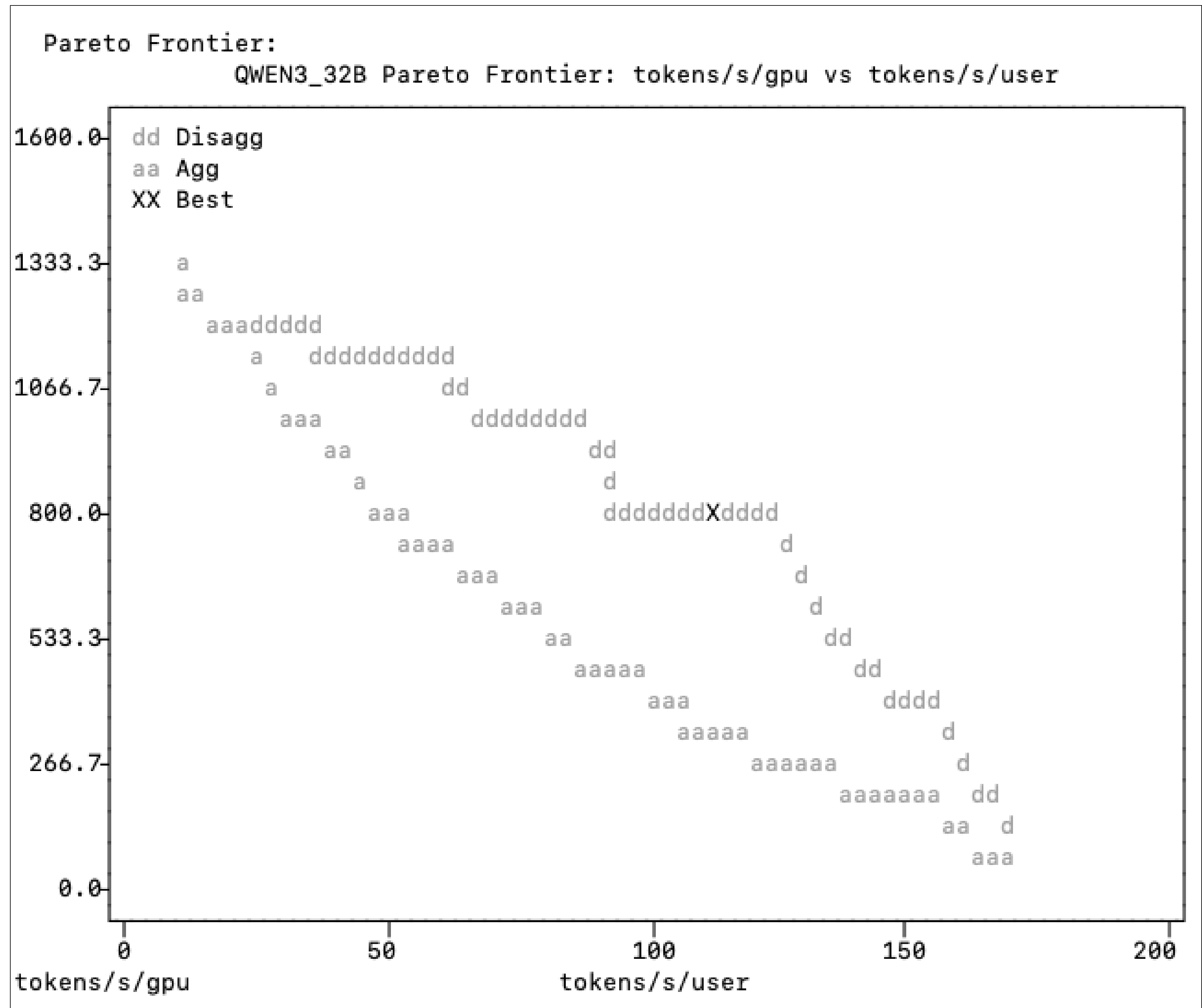
AIConfigurator

- . Find best config for disaggregated serving
- . Compare projected perf gain
- . <1 min using laptop vs 2 days using GPUs
- . Generate yaml to deploy with Dynamo
- . Modular and standalone

aiconfigurator cli

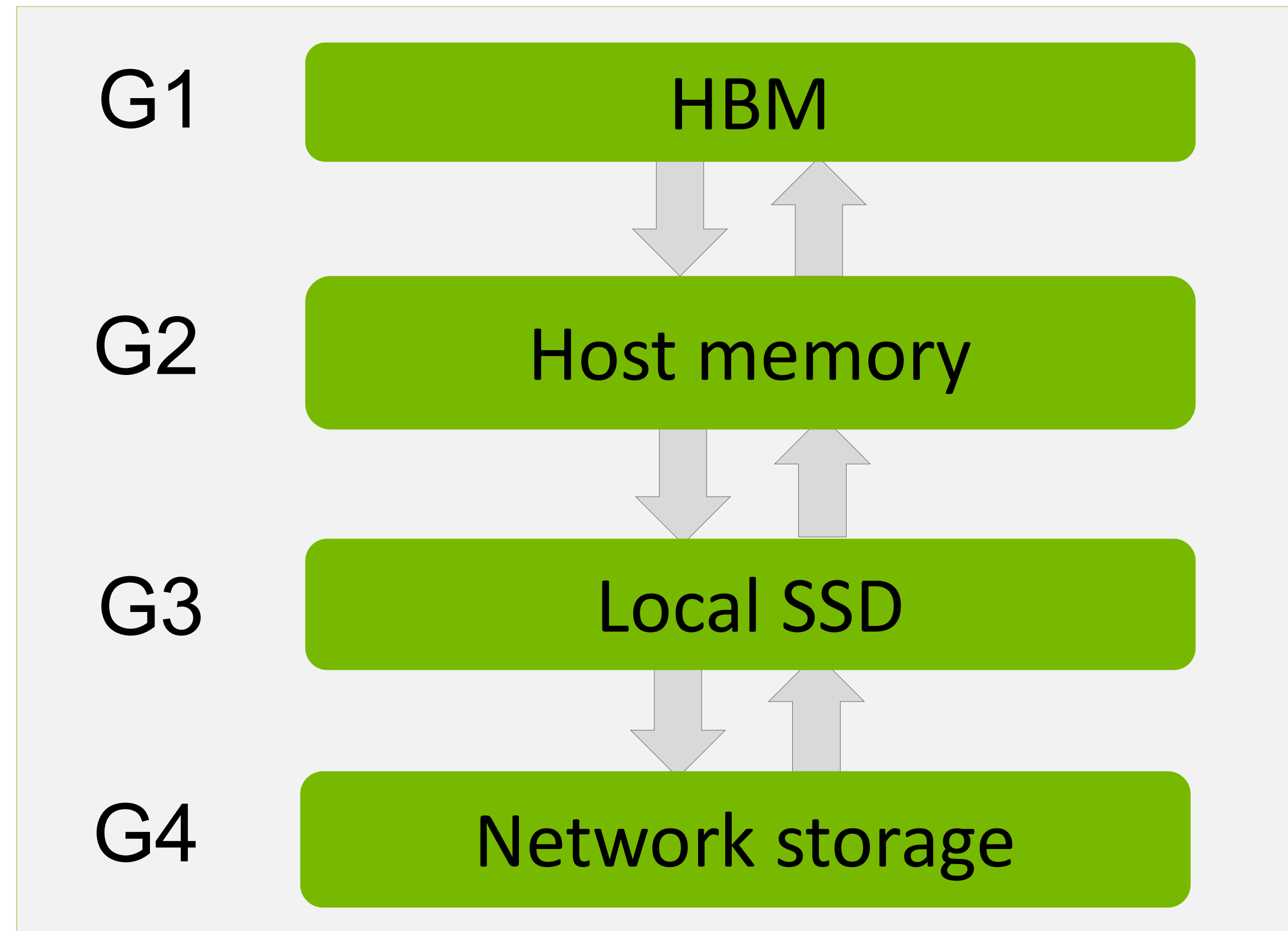
```
--model QWEN3_32B  
--total_gpus 32  
--system h200_sxm  
--isl 4000  
--osl 500  
--ttft 300  
--tpot 10
```

SLA requirements

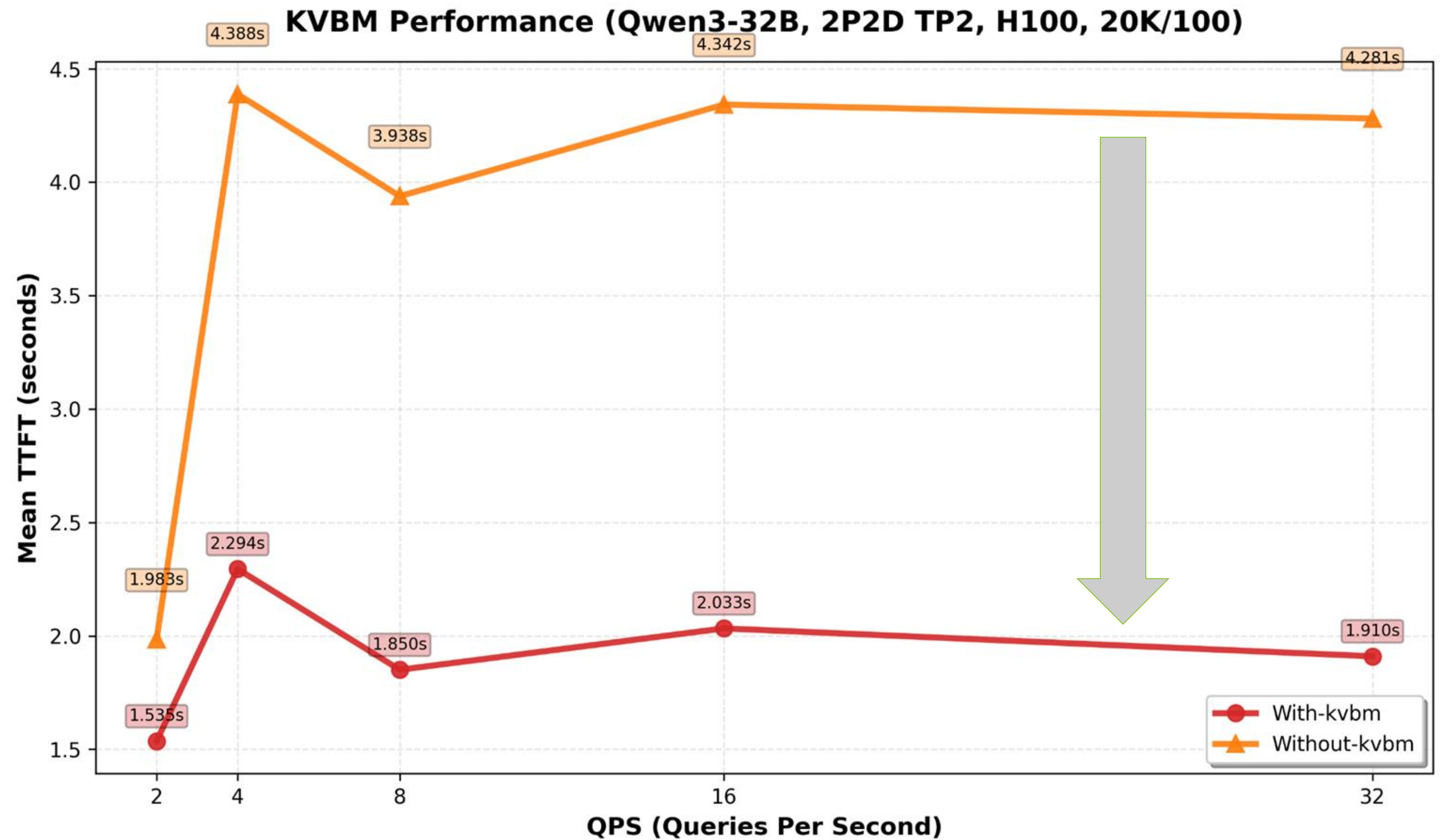


Memory Management

Leverage all memory and storage available in the data center



- ### Use cases
- Multi-turn conversations
 - Code generation
 - Agent calls
 - Generative recommender



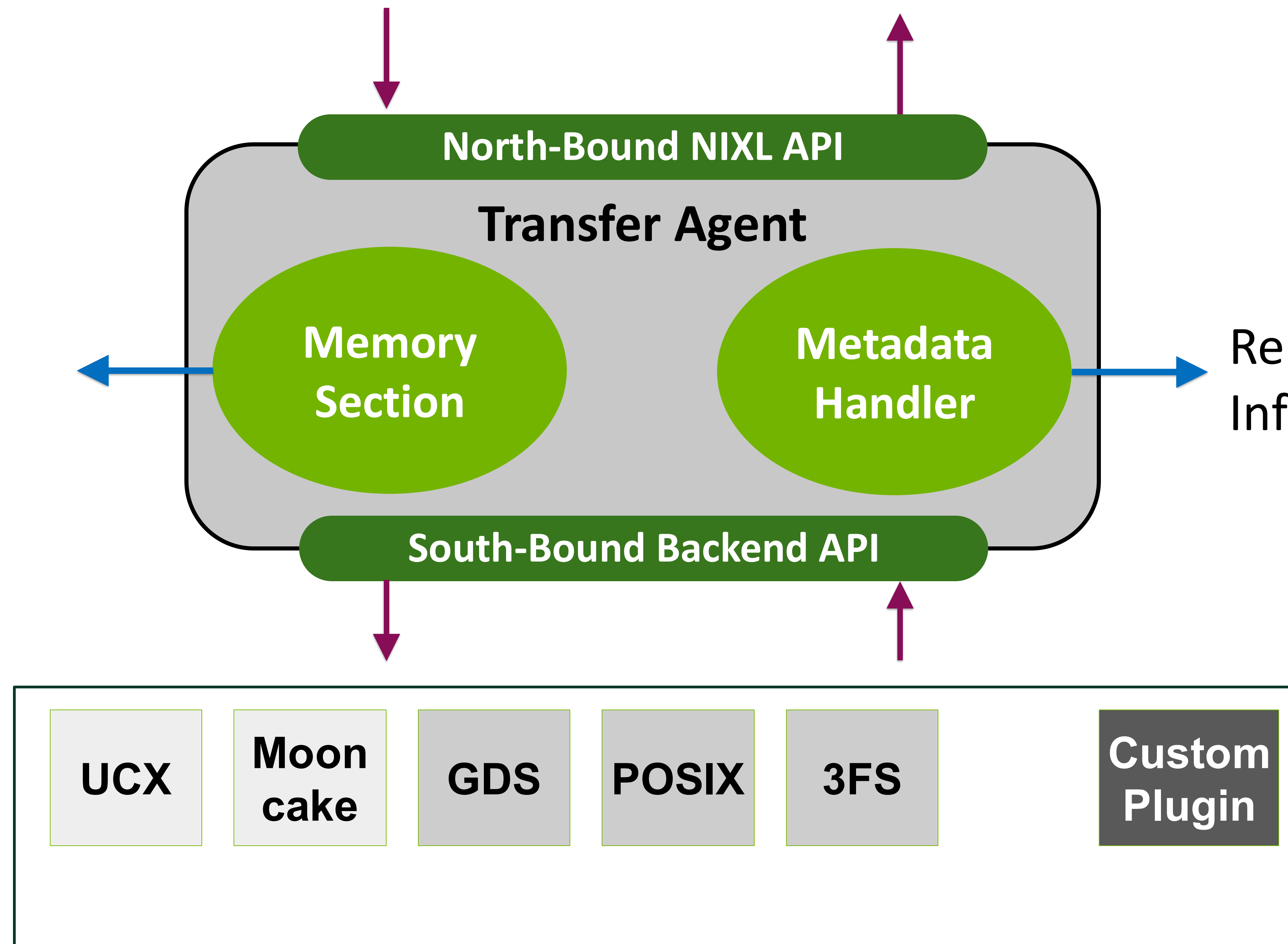
NIXL

NIXL cross-node
& cross-mem
buffer list primitive

NIXL Transfer
Handles with repost
capability

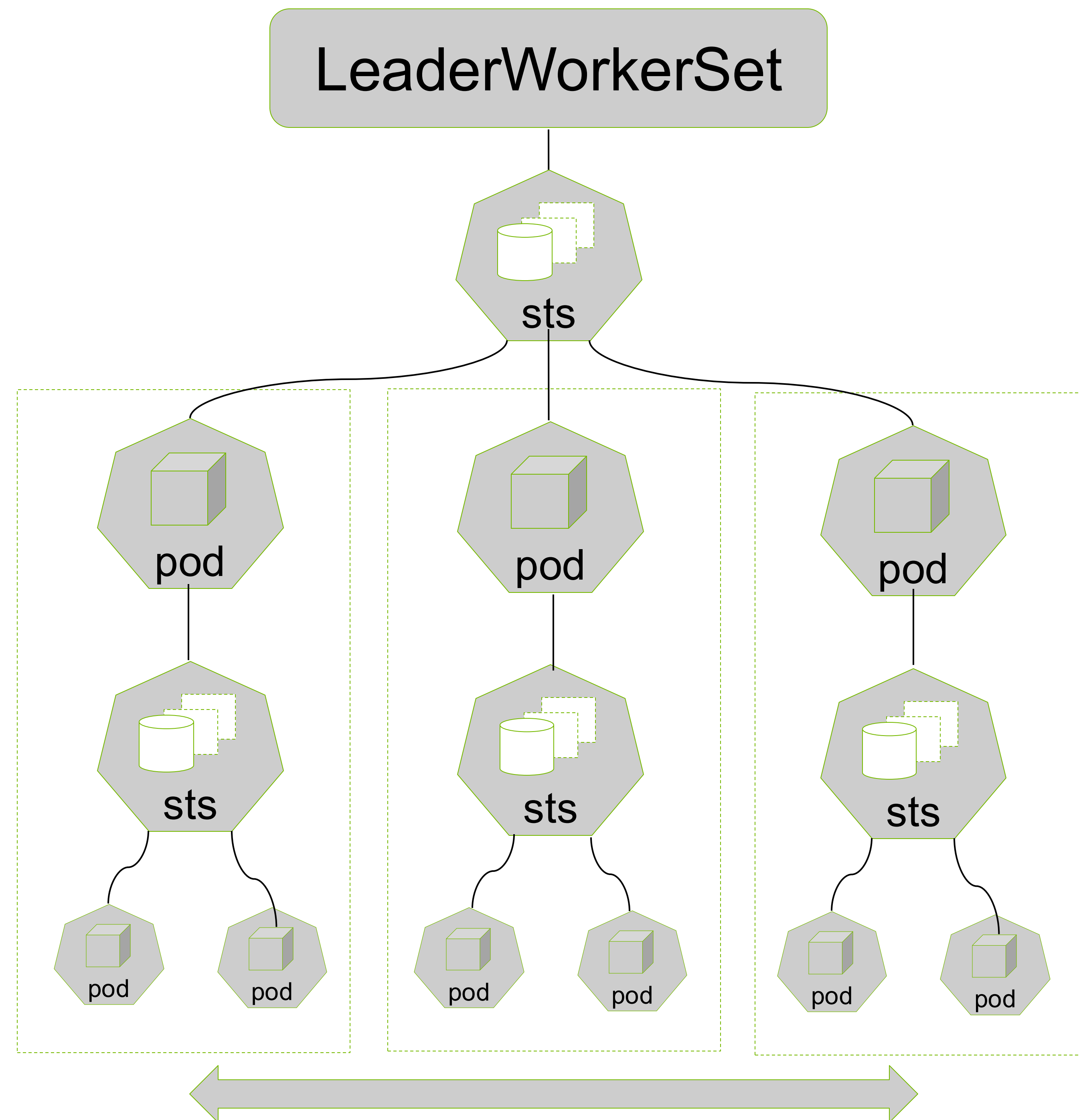
Local Registered
Memory Info

Remote Agent
Info

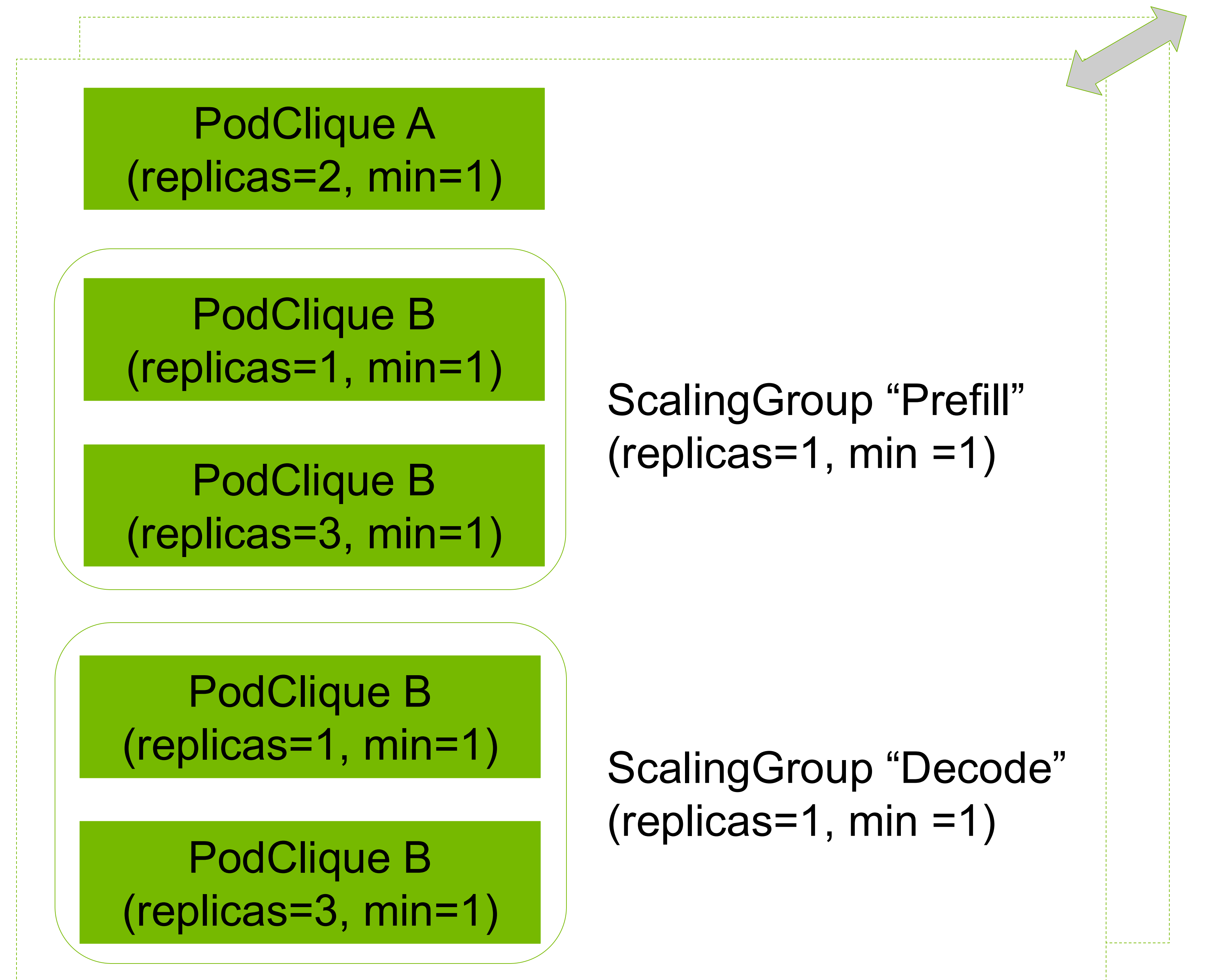


Topology Optimized Dynamic K8 Scheduling

Gap in scaling disaggregation



Grove



Fault tolerance

Ensure detection, resilience, and recovery

High availability

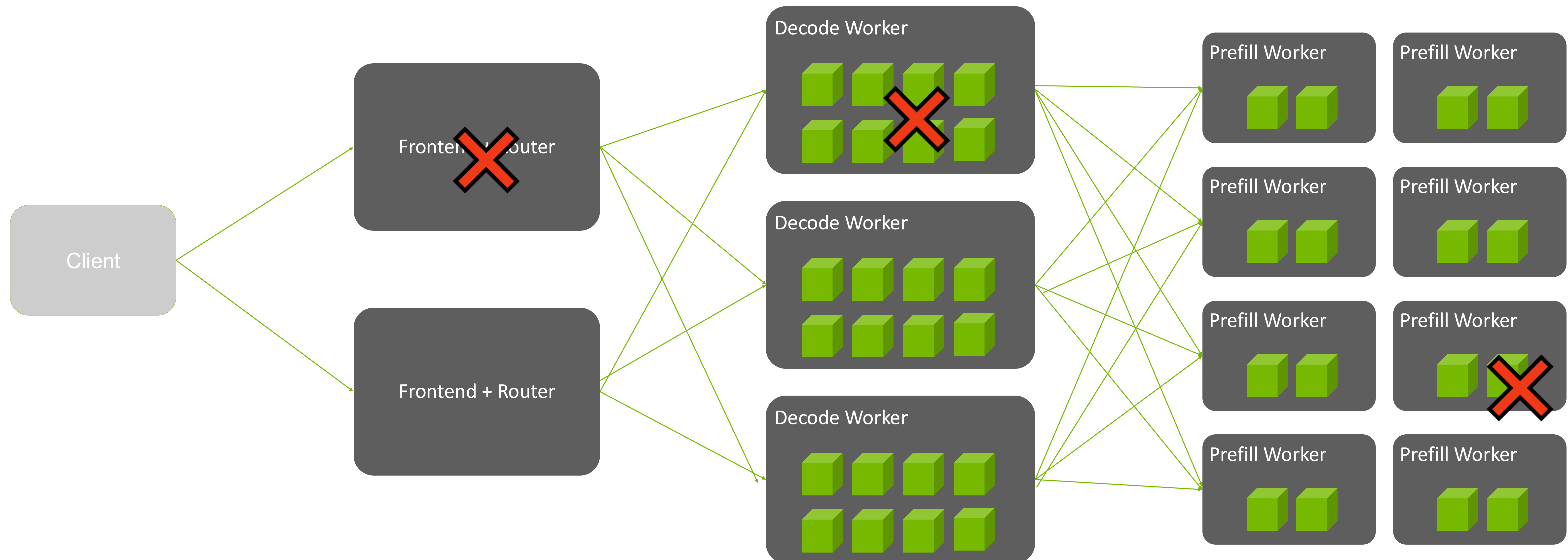
Automatic restart with shared router state

Request cancellation

Cancel request to avoid computation

Request migration

Migrate partial response at token level to avoid Decode



Challenges and opportunities

- We are barely at the start of large-scale AI inference deployment
- Heterogenous model deployment is least explored at scale
 - End to end agentic tasks need multiple models to closely exist together for performing the task
 - KV cache and other efficiency management
 - Fast data transfer mechanism
- Fault tolerance and performance measurement tools for inference at scale do not exist
 - AIPerf – a stab at performance measurement at scale from API standpoint of view
 - Agentic performance modelling?
 - Multimodal performance characterization?
 - Fault injection frameworks for inference?
 - Dynamo is working towards building something here but it is quite a task itself
- Speculative decoding, diffusion-based reasoning model, ... all open for exploration
- Enterprise RAG with CAG (cache augmented generation) at scale is unsolved issue
- How to evolve the stack fast with rapid innovations in model development?
- ...

Conclusions

- AI inference is unsolved problem, and you can contribute and learn with us!
- <https://github.com/ai-dynamo/dynamo> (★)
- Join the discord group: <https://discord.com/invite/D92uqZRjCZ>

```
# Create a virtual environment and activate it
uv venv venv
source venv/bin/activate

# Install Dynamo from PyPI (choose one backend extra)
uv pip install "ai-dynamo[sglang]==0.6.1" # or [vllm], [trtllm]
```

- Getting started with Dynamo: <https://www.youtube.com/watch?v=1bRmskFCnqY>
- Disaggregated inference demo: https://www.youtube.com/watch?v=UDJy_5_Czw
- KV Aware router demo: <https://www.youtube.com/watch?v=PRCZZKQirN8>
- AIConfigurator demo: <https://www.youtube.com/watch?v=KuZXeolOfKk>

