



LLM Inference Webinar

Sagar Desai, Dmitry Mironov



Welcome & Introduction

Sagar Desai is a Senior Solution Architect at NVIDIA. He specializes in large-scale generative AI solutions, including LLMs, multimodal AI, and retrieval-augmented generation (RAG). He excels at designing and deploying scalable, reliable, and secure AI architectures that deliver business value.

He's adept at advanced LLM training techniques, such as FP8 precision, and optimizing model convergence and GPU utilization.

Drawing on his background in MLOps, Sagar designs comprehensive deployment pipelines that ensure seamless, high-throughput inference with concurrent usability. Focused on accuracy, scalability, and reliability, he empowers organizations to unlock the full potential of generative AI technologies.

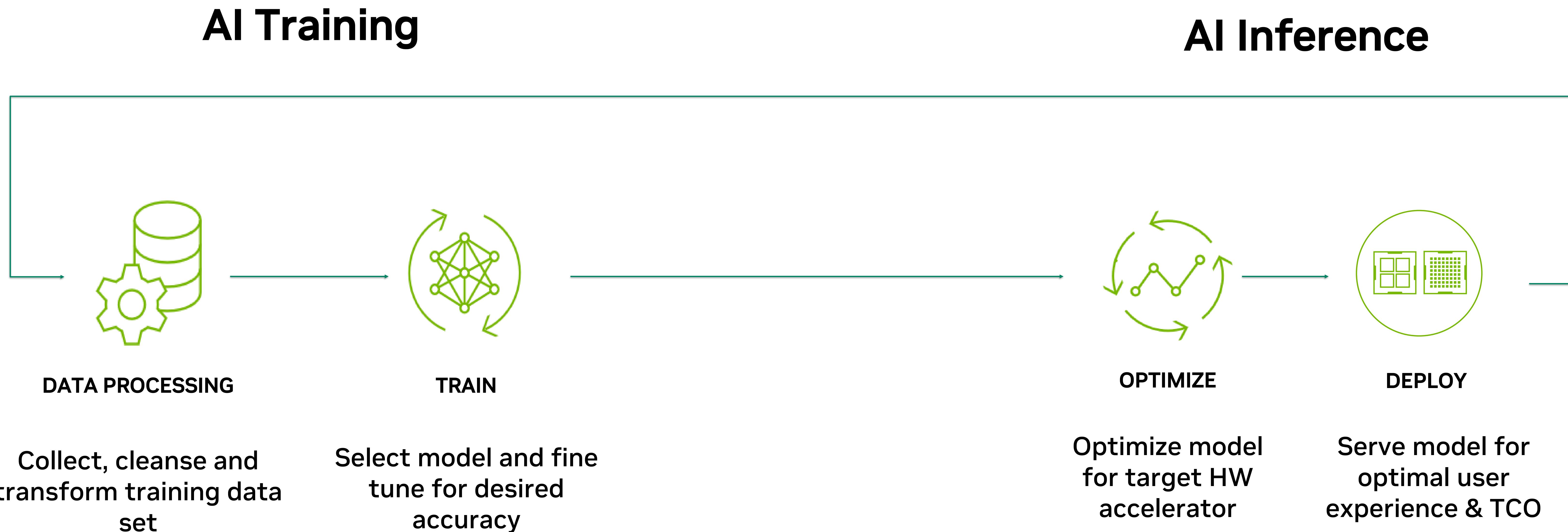


Agenda

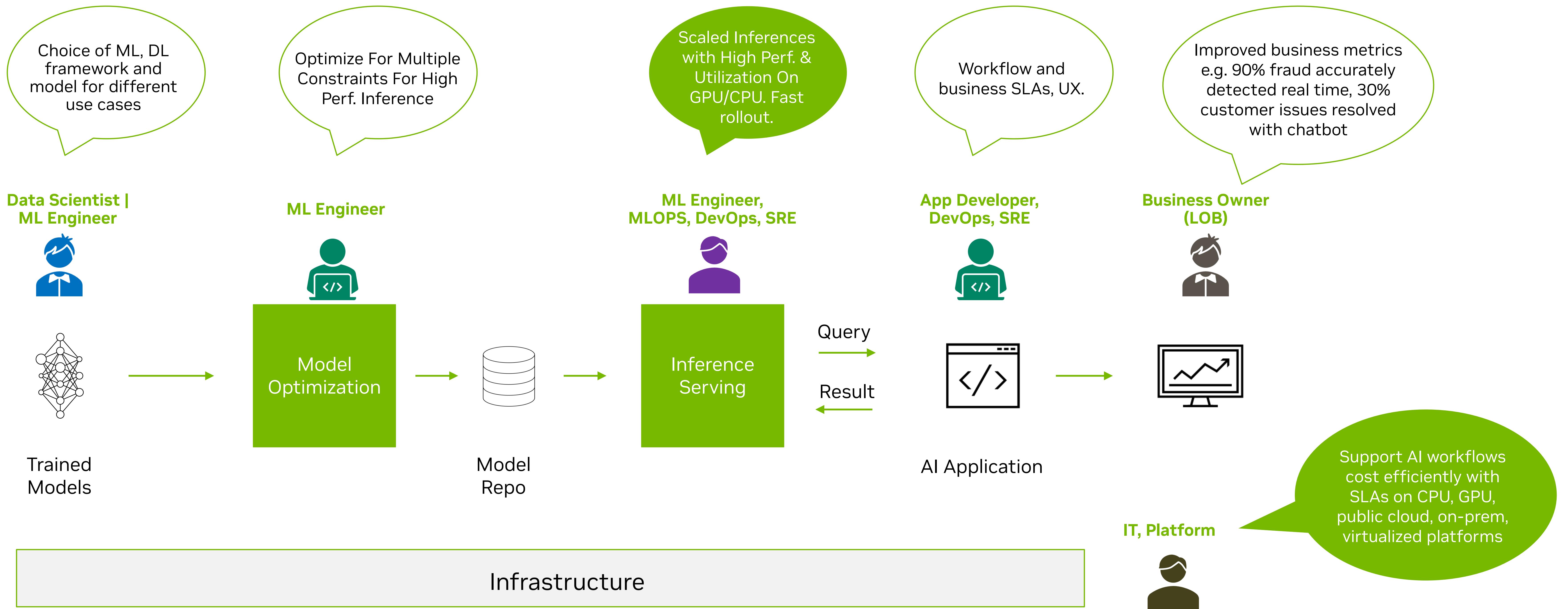
- Deep understanding of LLM inference theory and applications
- Efficient prompt processing and token generation techniques.
- NVIDIA Inference options
- NVIDIA NIM Microservices, TensorRT-LLM, Triton
- NVIDIA Dynamo
- Demo
- Case Studies

Challenges of LLM Inference in Production

Inferencing in the End-to-end AI Workstream



AI Inference Workflow



LLM Inference Fundamentals

Inference Latency (SLOs)

Tokens

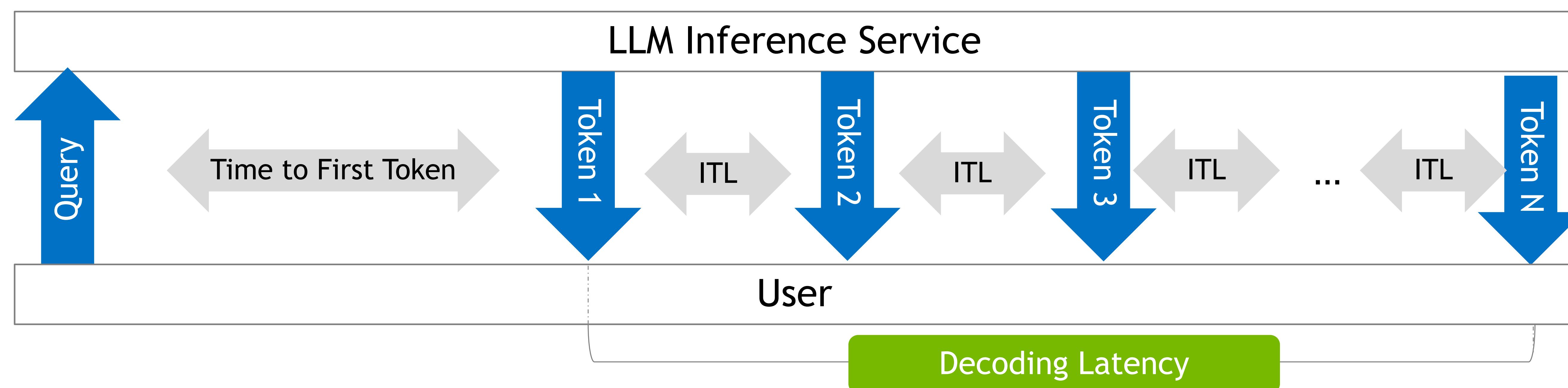
- Input sequence length of tokens (**ISL**) are fed into a model.
- LLM generates of output sequence length (**OSL**) of tokens one at a time

Time to First Token (TTFT)

- Time it takes for the model to generate the first token after receiving a request
- This part of execution is called **Prefill**
- Critical for chat-like applications

End-to-end Latency (E2E)

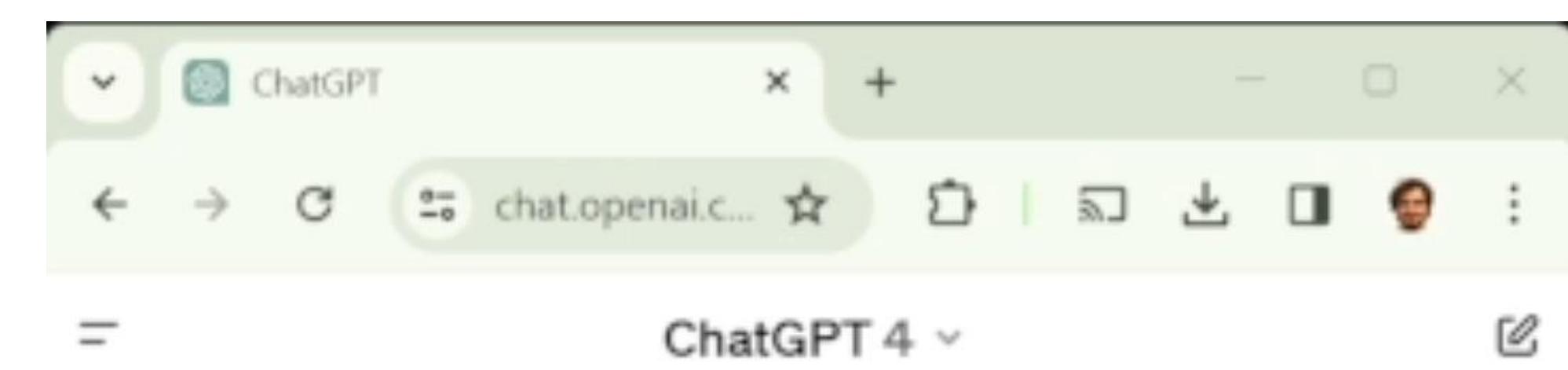
- Time it takes for a model to generate a complete response to a user's prompt.
- $E2E\ Latency = TTFT + \text{Decoding}\ Latency$
- Inter-token Latency (**ITL**) is an average time between output tokens
- Critical, when full response has to be processed further: guardrails, tool calling



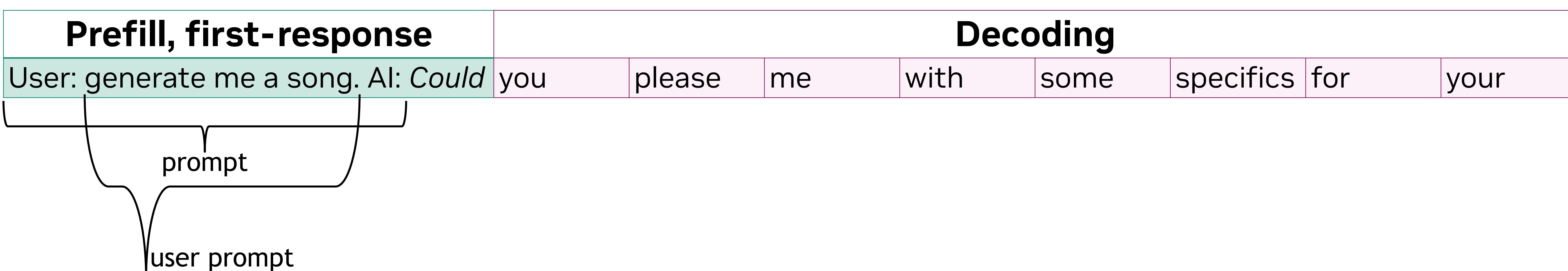
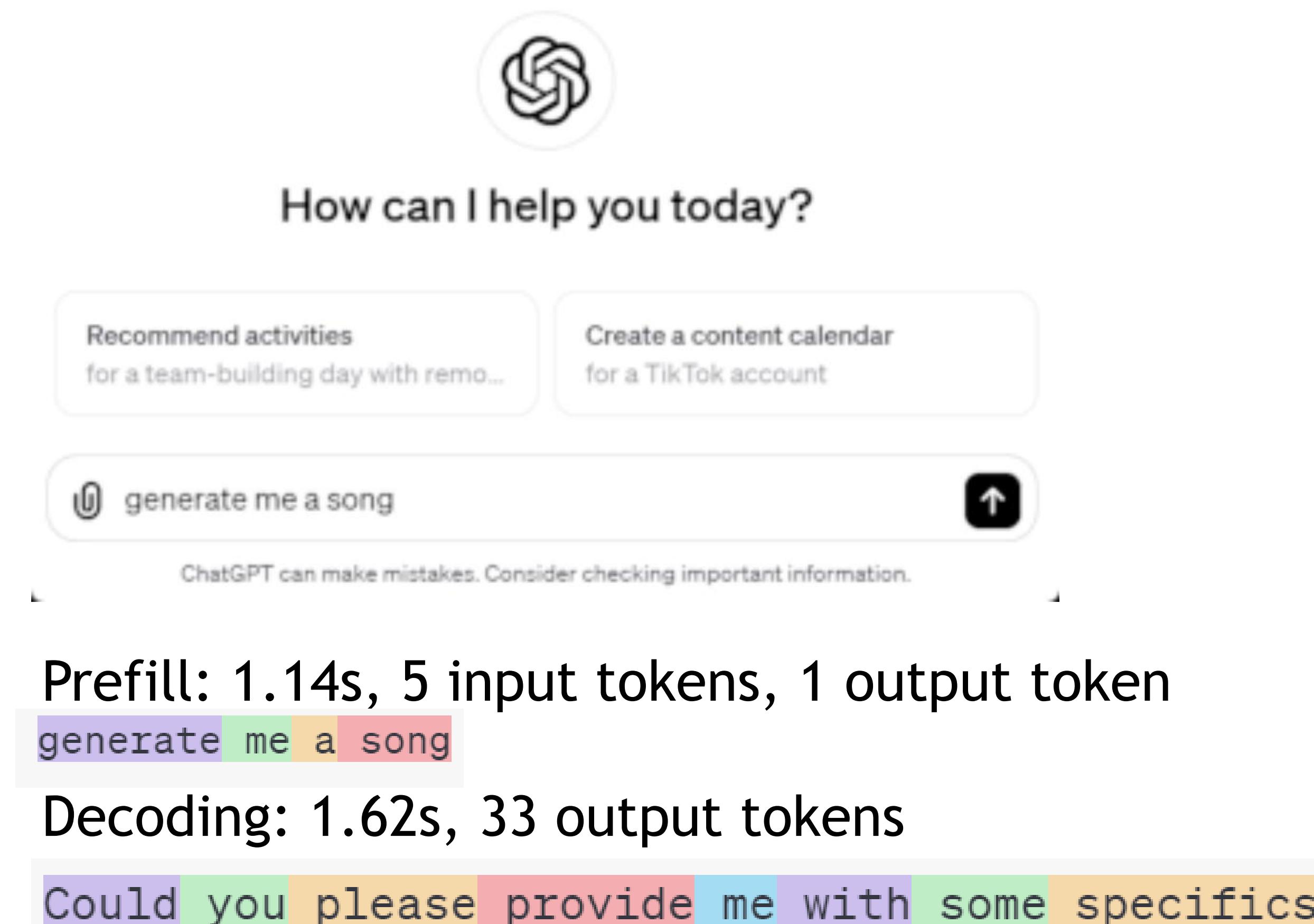
Two Stages of LLM Execution

Prefill vs Decoding

- **Prefill** = time to first token (~word)
 - Loading the user prompt into the system
 - From the request reception to the first token
 - Depends only on the number of input tokens
 - Populate KV-cache for all the tokens from the prompt.
 - Compute-bound for most of the reasonable prompt lengths
- **Decoding** = inter-token latency
 - Generating the response token by token, word by word
 - Inter-token latency depends on the total token number, both input and output tokens.
 - Usually memory-bound

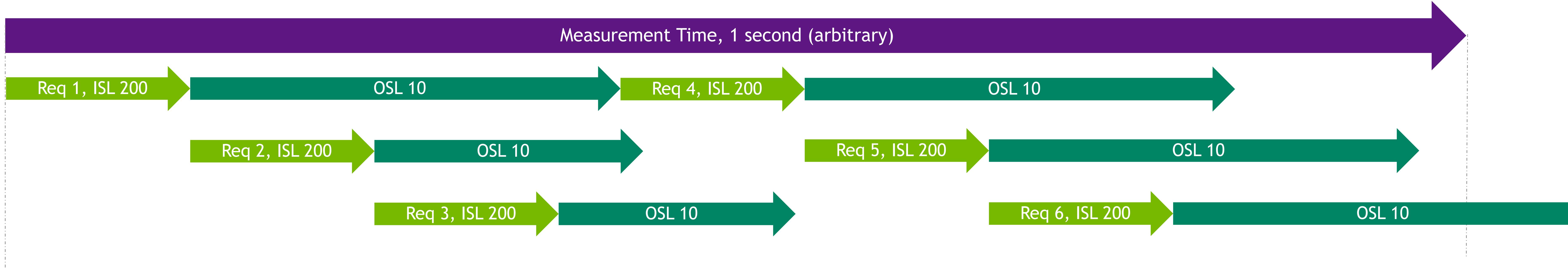


Realtime



Inference Throughput Example (Costs)

Assuming 1 GPUs per model, 1 instance



- Requests per second per deployments (**RPS**) = 5 started and completed requests / 1 second = 5 (200 tokens in, 10 tokens out)
- Output tokens per second per deployments (**TPS**) = 57 output tokens / 1 second = 57 (200 tokens in, 10 tokens out)
- **Requests per second per GPU** = RPS / #GPU = 5 / 1 = 5 (200 tokens in, 10 tokens out)
 - If 2X RPS needed keeping latency the same, need 2X servers.
 - **Best definition of throughput. Direct relationship to costs. Always try to calculate it.**

NVIDIA AI Inference Solutions

Challenges with AI inference in Production

IT leaders chose AI inference solutions that meet enterprise grade requirements



Fast

Meets use case and user experience latency constraints



Interoperable

Supports multiple development frameworks
Adapt to SOTA architectures, e.g., reasoning agents



Scalable

Easily add new models and support spikes in user demand

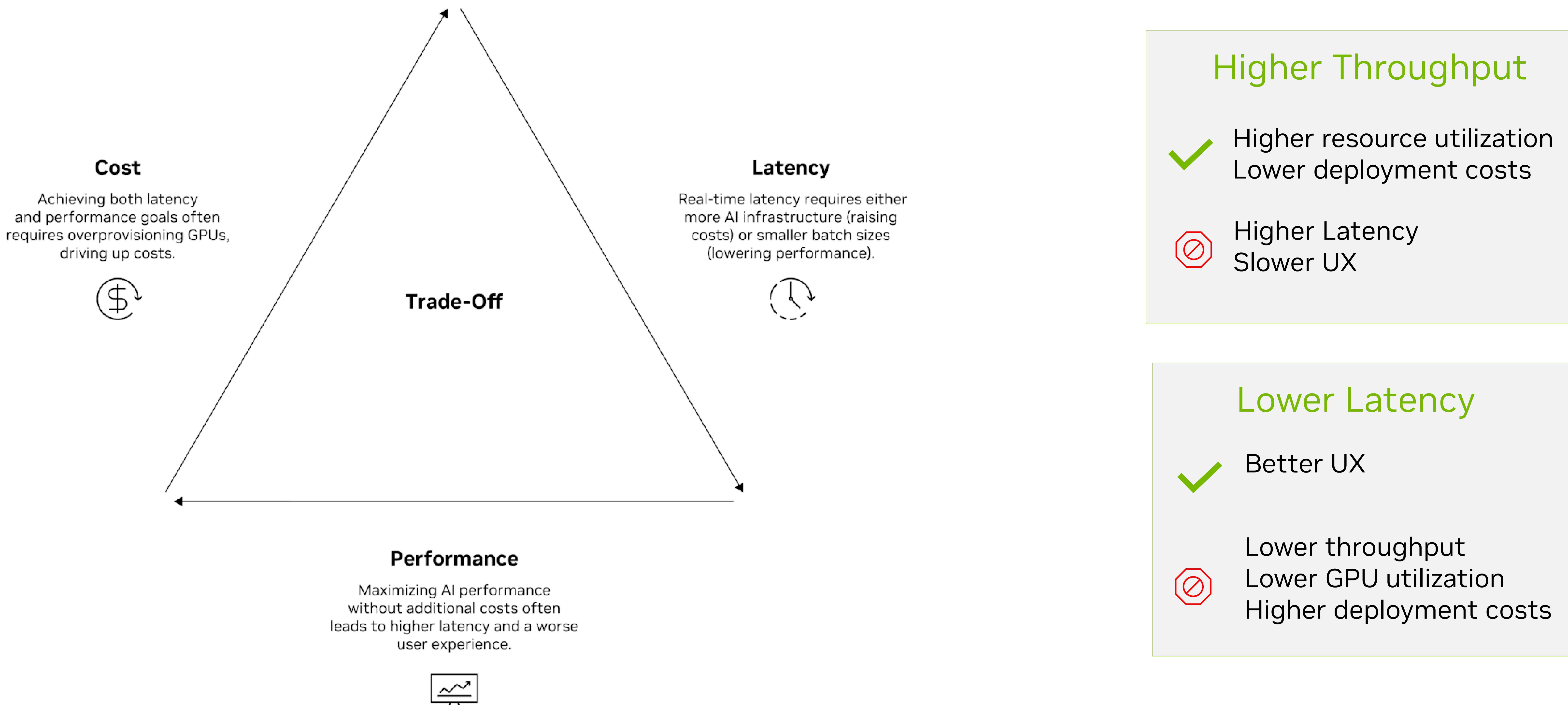


Cost

Meets budget constraints and delivers max ROI

Deploying LLM Inference: Throughput vs Latency

Balance between throughput and latency is crucial for optimal performance and user experience



Inferencing in the End-to-end AI Workstream



TensorRT-LLM

TensorRT-LLM in the *DL Compiler Ecosystem*

TensorRT-LLM builds on TensorRT Compilation

- **TensorRT-LLM**

- Inference runtime & compiler specifically designed for LLMs
- LLM specific optimizations:
 - KV Caching & Custom MHA Kernels
 - Inflight batching, Paged KV Cache (Attention)
 - Multi-GPU, Multi-Node
 - Grammar support
 - & more
- *ONLY for LLMs*

TensorRT-LLM

- LLM specific optimizations:
- KV Caching
 - Multi-GPU, Muti-Node
 - Custom MHA optimizations
 - Paged KV Cache (Attention)
 - etc...

- **TensorRT**

- General purpose Deep Learning Inference Compiler
 - Graph rewriting, constant folding, kernel fusion
 - Optimized GEMMs & pointwise kernels
 - Kernel Auto-Tuning
 - Memory Optimizations
 - & more
- *All AI Workloads*

TensorRT

- General Purpose Compiler
- Optimized GEMMs & general kernels
 - Kernel Fusion
 - Auto Tuning
 - Memory Optimizations
 - Multi-stream execution

TensorRT-LLM GPU Support



GB200 NVL72



Blackwell



Hopper



Ada Lovelace

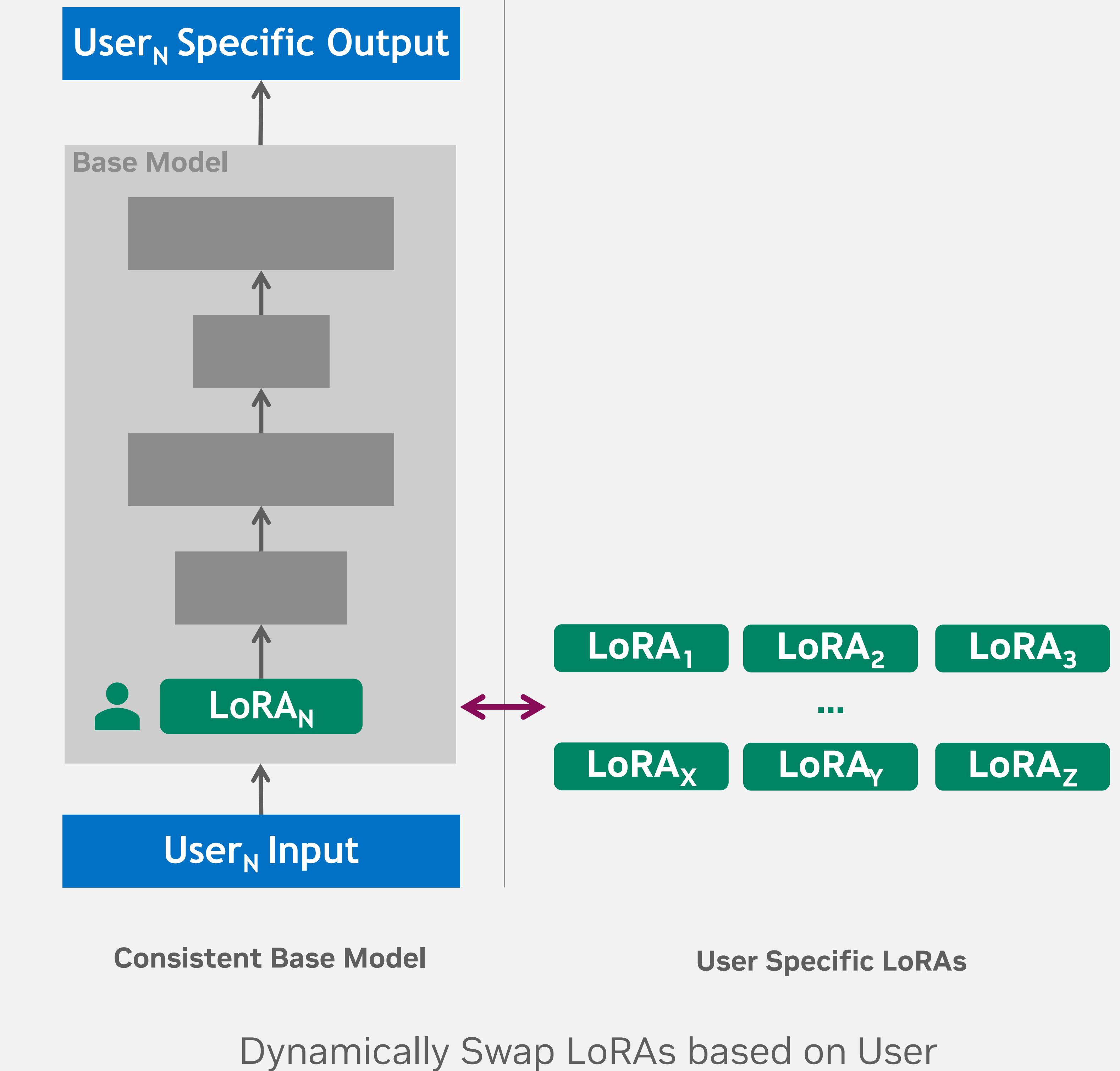


Ampere

LoRA & Customization

Efficiently Supporting Customer User Experience

- LoRA & Prompt tuned models are support in TRT-LLM and NVIDIA NIM
- Support multiple customers with a single model
- Dynamically swap LoRA's at runtime
- SLoRA / LoRAX caching adapters on device
- Base model can be quantized for memory savings
 - QLoRA in progress

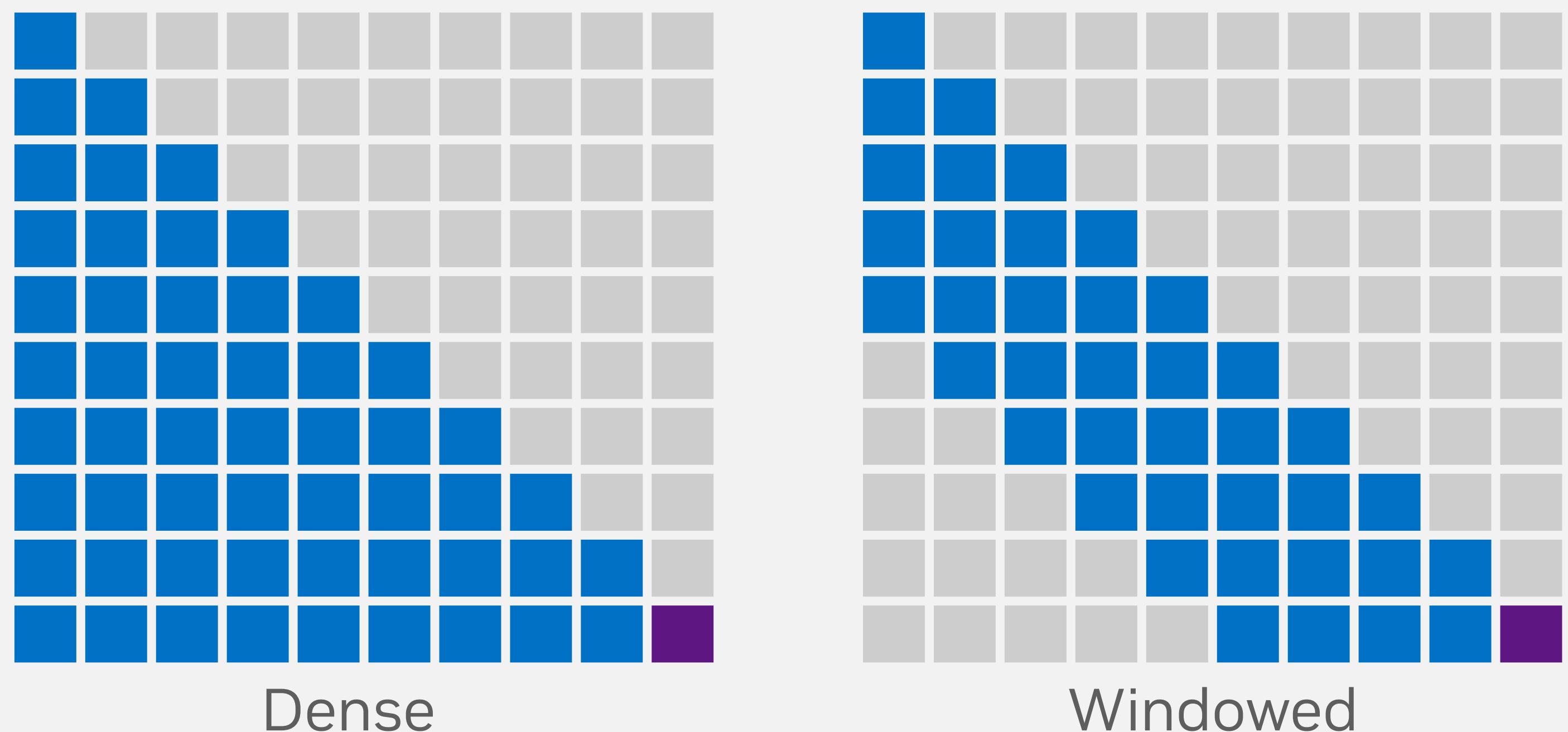


KV Cache & Attention Techniques

(Sliding) Window Attention, & Streaming LLM

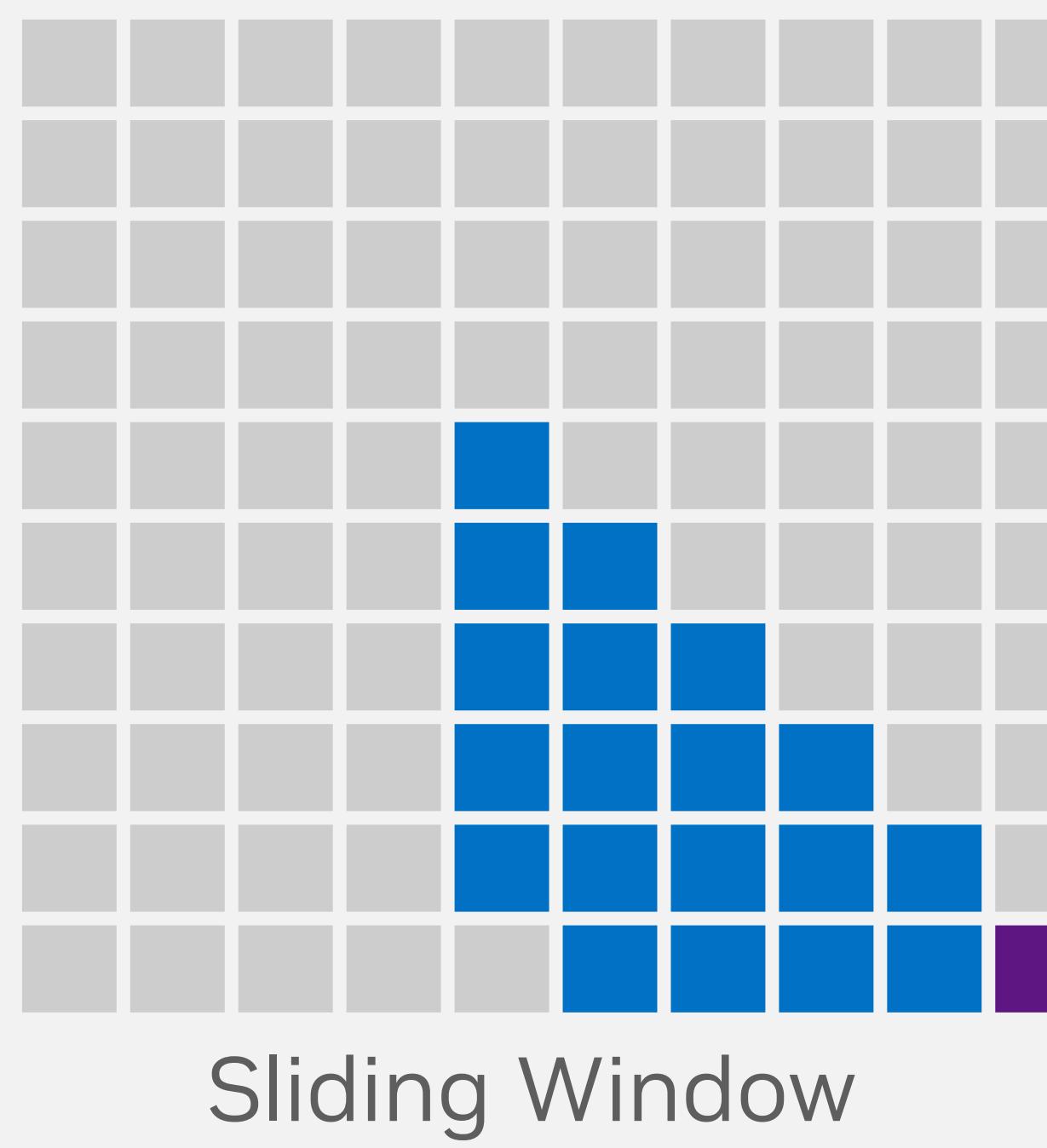
- Allow for longer (sometimes unlimited) sequence length
 - Reduces KV Cache Memory usage
 - Avoids OOM Errors
- (Sliding) Windowed Attention evict tokens based on arrival
 - Significantly reduces memory usage
 - Can negatively impact accuracy or require recomputing KV
- Streaming-LLM allows for unlimited sequence length
 - Does not evict Attention Sinks (important elements)
 - KV Cache stays constant size
 - Does not require recompute & does not impact accuracy
 - Particulary beneficial for multi-turn (ie. chat) usecases

Attention KV Cache Usage (*Less is Better*)

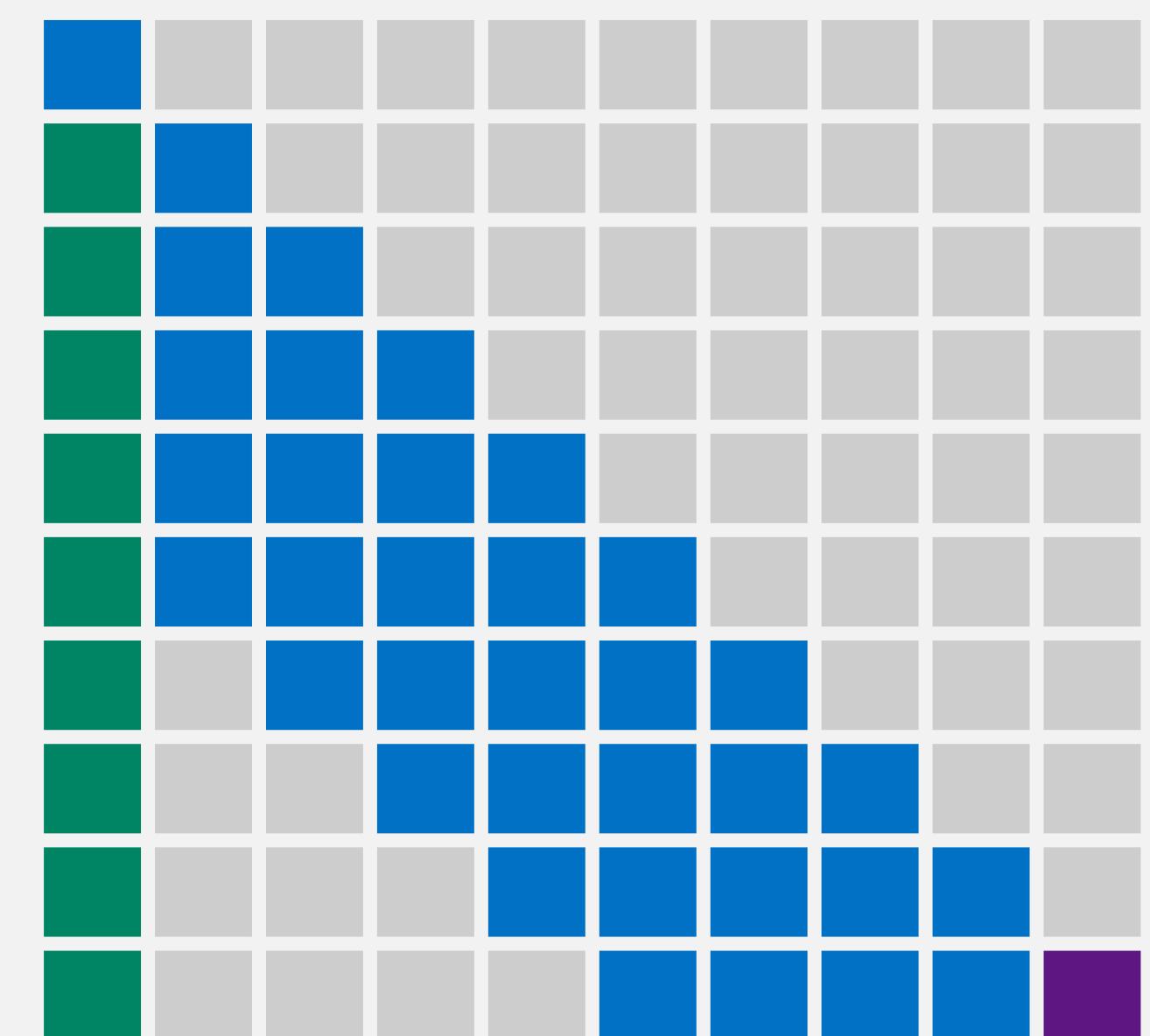


Dense

Windowed



Sliding Window



StreamingLLM

■ Free ■ Prev. Tokens ■ Curr. Token ■ Attn. Sync

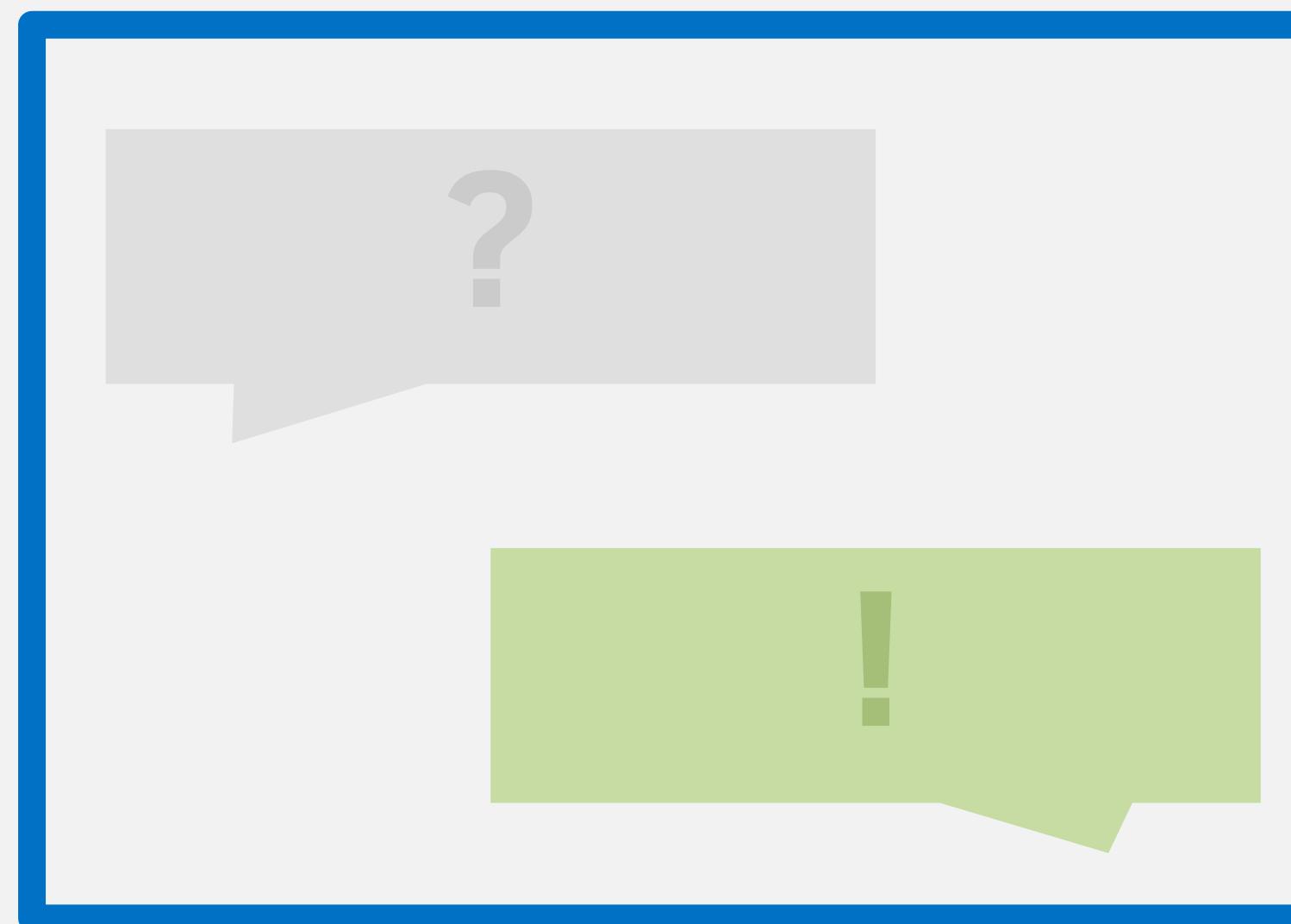
KV Cache Reusage

System Prompt Caching & Block reusage

Allows for interactive/ turn based systems & System Prompts

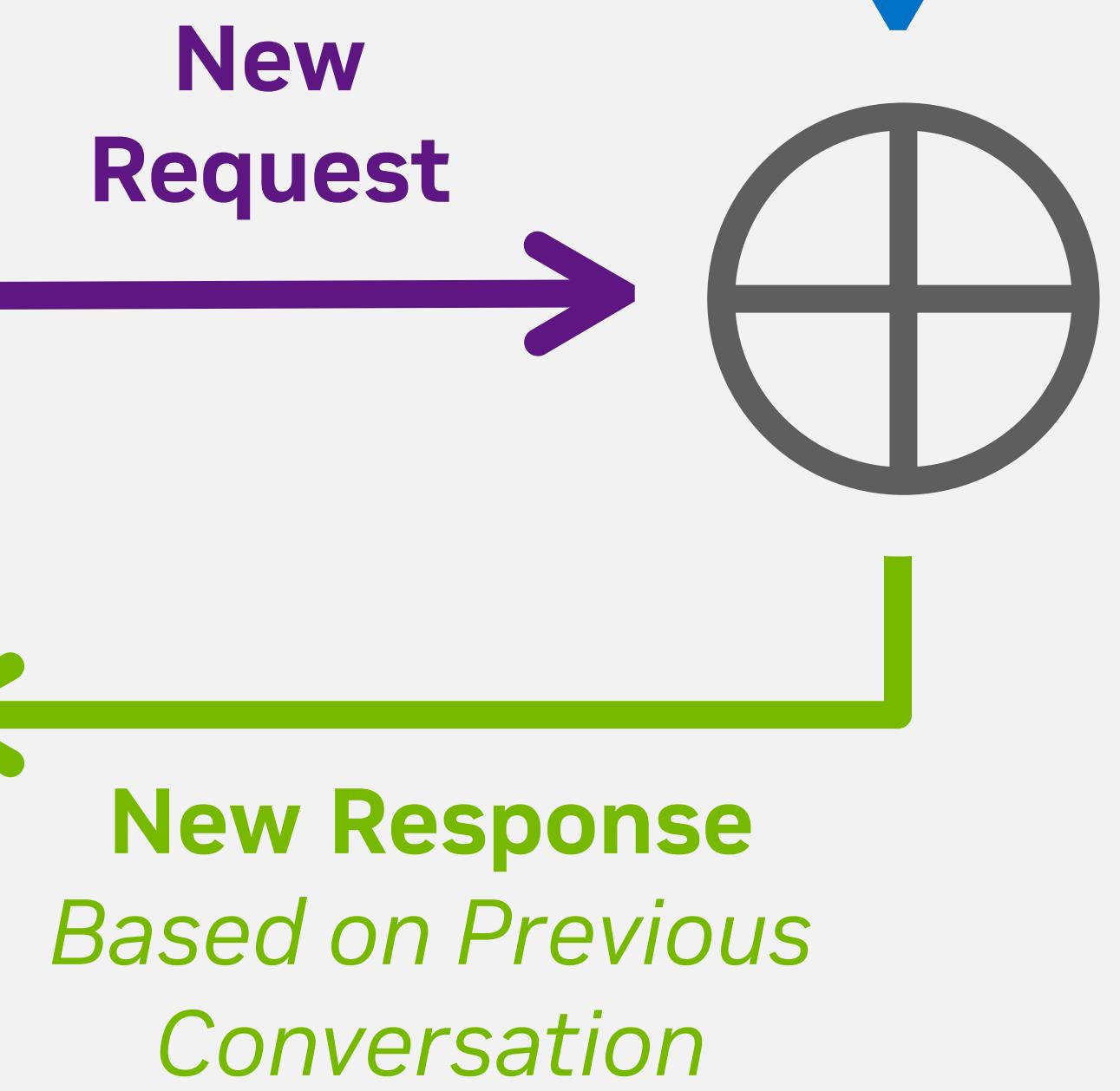
- Load prior KV cache blocks to avoid recomputation
 - Saves significant compute
 - Reduces Start-up time
- Block reusage allows for turn-based (chat) applications
 - Allows for additional options for intelligently reusing blocks
- System prompts allows for a preset KV cache for the LLM
 - E.g. to give rules, personality, or prior knowledge

Prior Session



Previous
KV Cache

New Session



Inflight Batching

Maximizing GPU Utilization during LLM Serving

TensorRT-LLM provides custom Inflight Batching to optimize GPU utilization during LLM Serving

- Replaces completed requests in the batch
 - Evicts requests after EoS & inserts a new request
- Improves throughput, time to first token, & GPU utilization
- Integrated directly into the TensorRT-LLM Triton backend
- Accessible through the TensorRT-LLM Batch Manager

Batch Elements	Iteration									...
	1	2	3	4	5	6	7	8	9	
R ₁						END				R ₅
R ₂			END							R ₆
R ₃				END						R ₇
R ₄					END					R ₈

Static Batching

Batch Elements	Iteration									...
	1	2	3	4	5	6	7	8	9	
R ₁						END	R ₇			...
R ₂			END	R ₅						...
R ₃				END	R ₆		END	R ₈		...
R ₄					END			END	R ₉	...

Inflight Batching

Context | Gen | EoS | NoOp

KV Cache Optimizations

Paged & Quantized KV Cache

Paged KV Cache improves memory consumption & utilization

- Stores keys & values in non-contiguous memory space
- Allows for reduced memory consumption of KV cache
- Allocates memory on demand

Quantized KV Cache improves memory consumption & perf

- Reduces KV Cache elements from 16b to 8b (or less!)
- Reduces memory transfer improving performance
- Supports INT8 / FP8 KV Caches

Both allow for increased peak performance

KV Cache Contents:

TensorRT-LLM optimizes inference on
NVIDIA GPUs ...

Block 0	TensorRT	LLM	is	...
Block 1				
Block 2	Hello	World		
Block 3				

Traditional KV Caching

B ₀	TensorRT	LLM	is	...
B ₁				
B ₂	Hello	World		
B ₃				

Paged KV Cache

B ₀	TRT	LLM	is	...				
B ₁								
B ₂	Hello	World						
B ₃								

Quantized Paged KV Cache

Request 1 | Request 2 | Wasted | Free

Quantization

Supported Precisions & Models

- Utilizes Transformer Engine on Hopper and Blackwell for FP8 and FP4
- Support many 8bit & 4bit methods
 - FP8, FP4, INT8/INT4 Weight only, INT8 Smooth Quant, AWQ, GPTQ
 - Support varies by model
- Reduced model size, memory bandwidth, & compute
 - Improves performance & allows for larger models per GPU
- Model optimization toolkit to quantize pre-trained models
 - Allows for per layer quantization strategies
- Currently requires all weights to be in same precision
 - Would like to relax this constraint going forward
- [Precision documentation](#)

	FP32	FP16	BF16	FP8	INT8	INT4
Volta (SM70)	Y	Y	N	N	Y	Y
Turing (SM75)	Y	Y	N	N	Y	Y
Ampere (SM80, SM86)	Y	Y	Y	N	Y	Y
Ada-Lovelace (SM89)	Y	Y	Y	Y	Y	Y
Hopper (SM90)	Y	Y	Y	Y	Y	Y

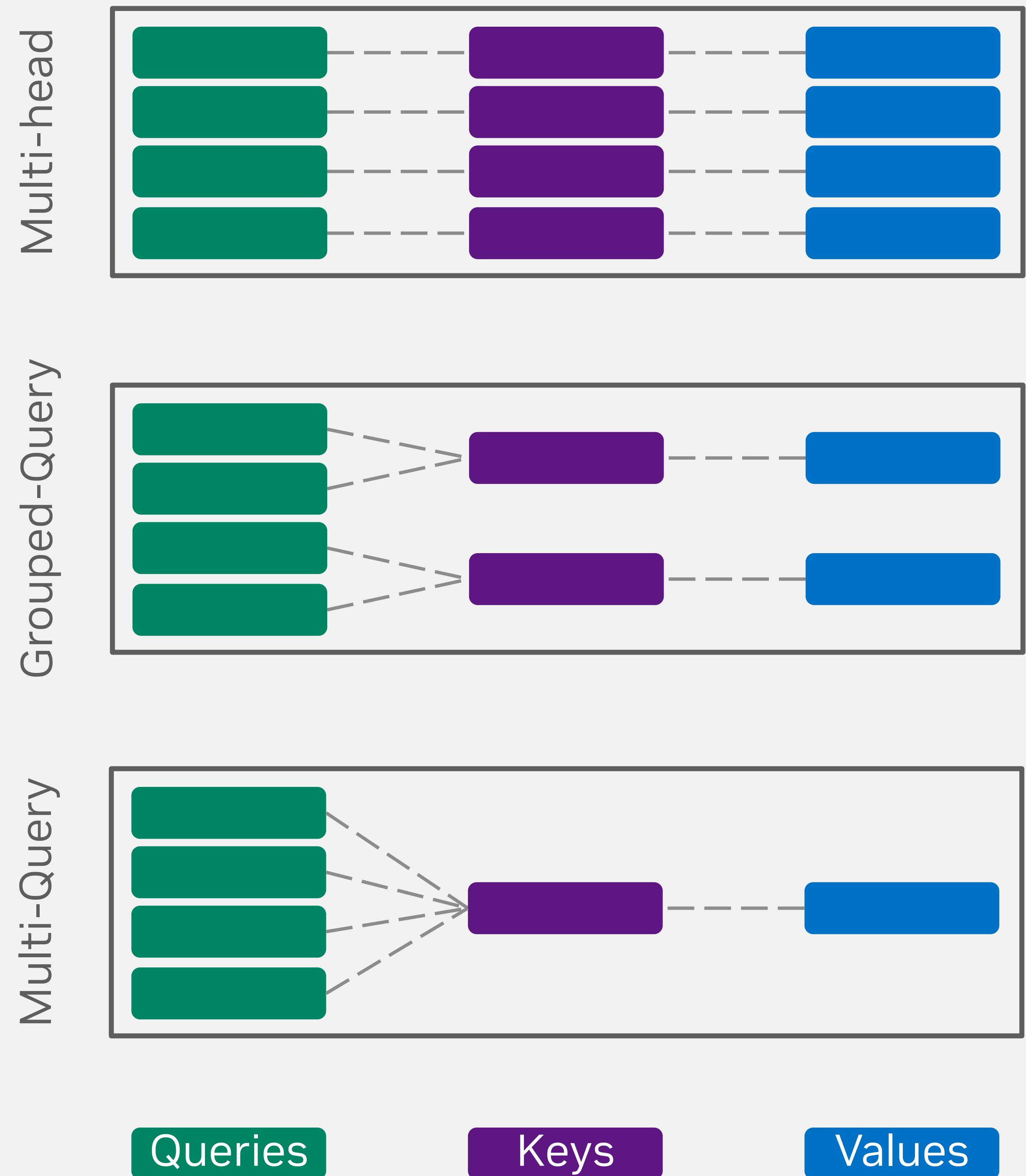
Model	FP32	FP16	BF16	FP8	W8A8 SQ	W8A16	W4A16	W4A16 AWQ	W4A16 GPTQ
Baichuan	Y	Y	Y	Y	Y	Y	Y	Y	Y
BERT	Y	Y	Y
BLIP-2	Y	Y	Y
BLOOM	Y	Y	Y	.	Y	Y	Y	.	.
ChatGLM	Y	Y	Y
ChatGLM-v2	Y	Y	Y
ChatGLM-v3
Flan-T5	Y	Y	Y
GPT	Y	Y	Y	Y	Y	Y	Y	Y	.
GPT-J	Y	Y	Y	Y	.	Y	Y	Y	.
GPT-NeMo	Y	Y	Y
GPT-NeoX	Y	Y	Y	Y
InternLM	Y	Y	Y	.	Y	Y	Y	.	.
LLaMA	Y	Y	Y	Y	Y	Y	Y	Y	Y
LLaMA-v2	Y	Y	Y	Y	Y	Y	Y	Y	Y
Mistral	Y	Y	Y	Y	Y	Y	Y	Y	.
MPT	Y	Y	Y	Y	Y	Y	Y	Y	.
OPT	Y	Y	Y
Phi	Y	Y	Y
Replit Code	Y	Y	Y	.	Y	Y	Y	.	.
SantaCoder	Y	Y	Y	.	.	Y	Y	.	.
StarCoder	Y	Y	Y	.	.	Y	Y	.	.
T5	Y	Y	Y

Quantization Examples Supported Models

Optimized Attention

Custom Implementations for Attention

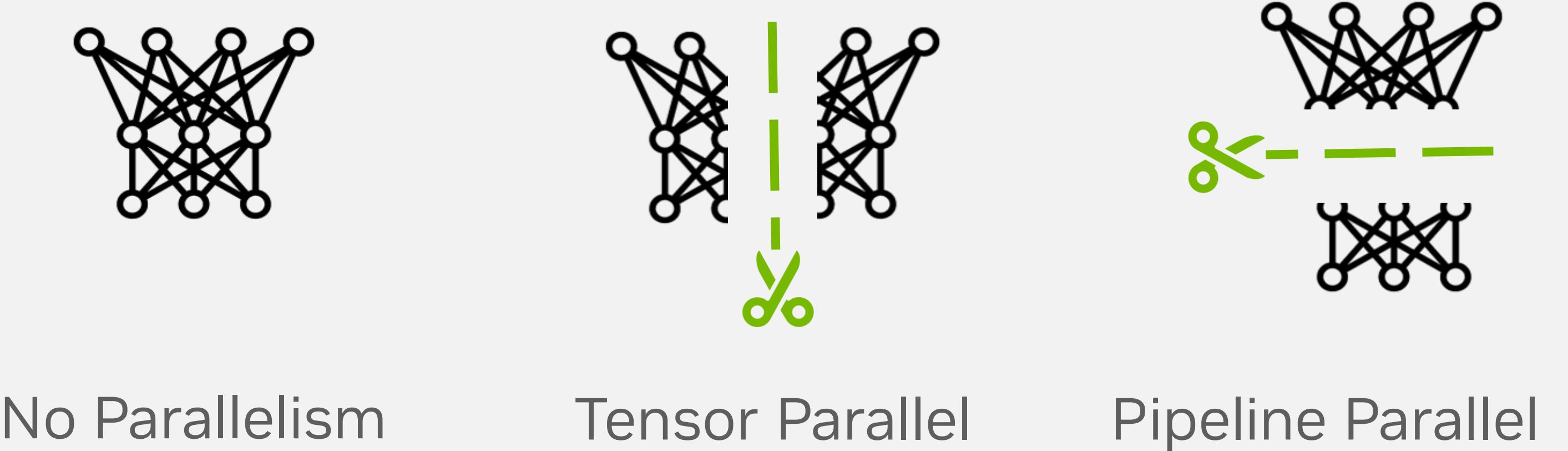
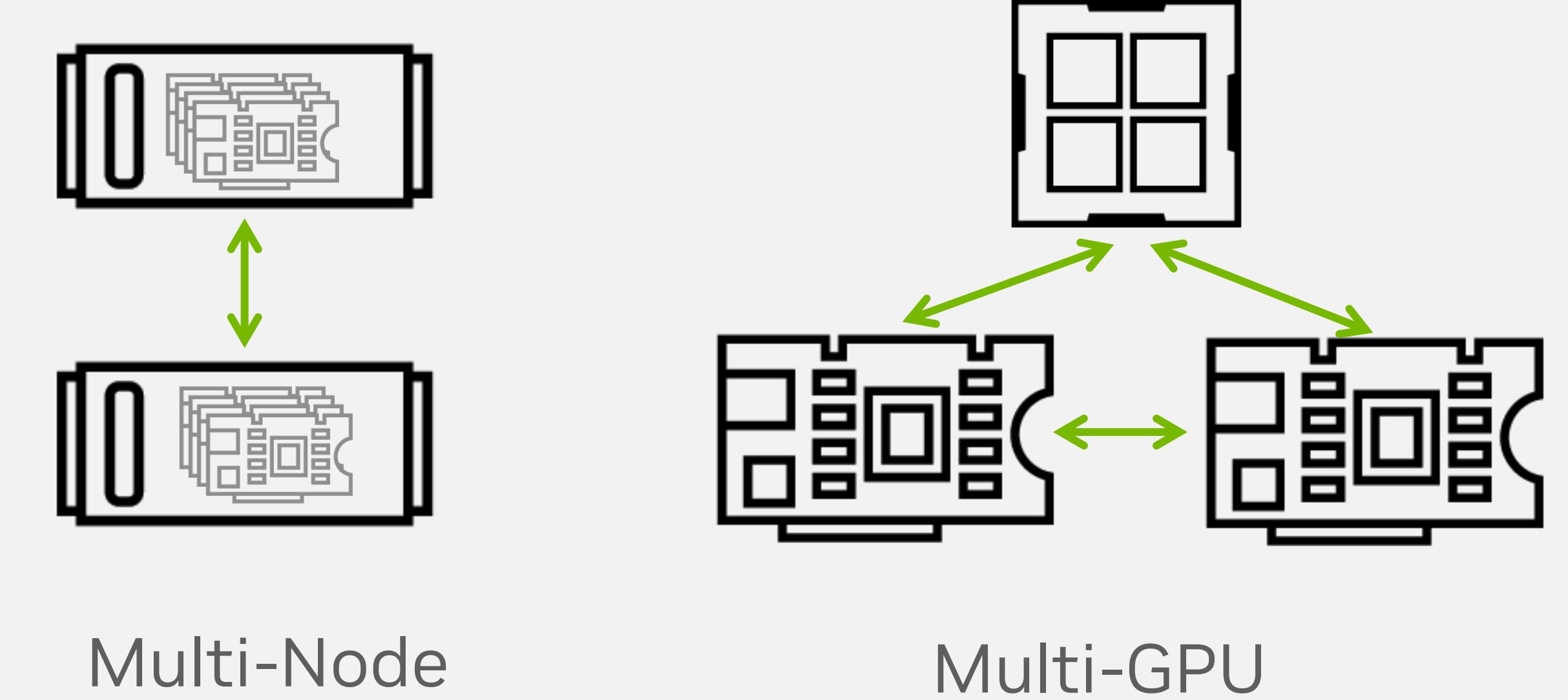
- Custom optimized CUDA kernels for Attention
 - Similar to FlashAttentionV2
- Optimized for A100 & H100 & B200
- Kernels for Encoder & Decoder, as well as context & prefill
- Supports MHA, MQA, GQA



Multi-GPU Multi-Node

Sharding Models across GPUs

- Supports Tensor & Pipeline parallelism
- Allows for running very large models (tested up to 530B)
- Supports multi-GPU (single node) & multi-node
- TensorRT-LLM handles communication between GPUs
- Examples are parametrized for sharding across GPUs



TensorRT-LLM Usage

NEW High-Level API

Experimental High-Level API

- Simplifies LLM workflow for common models & configs
- Few lines of Python Code to get started with TRT-LLM
- Static configs, and support for inflight-batching

Check out [examples/high-level-api](#) and provide feedback!

```
from tensorrt_llm import LLM, ModelConfig

# Instantiate the model
config = ModelConfig(model_dir=[...])
llm = LLM(config)

# Run Inference!
prompts = ["How do I get the fastest LLM Inference?"]
for output in llm.generate(prompts):
    print(output) # >> "Use TensorRT-LLM!"
```

Standard Usage

```
# Instantiate the model
config = ModelConfig(model_dir=[...])
llm = LLM(config, async=True)

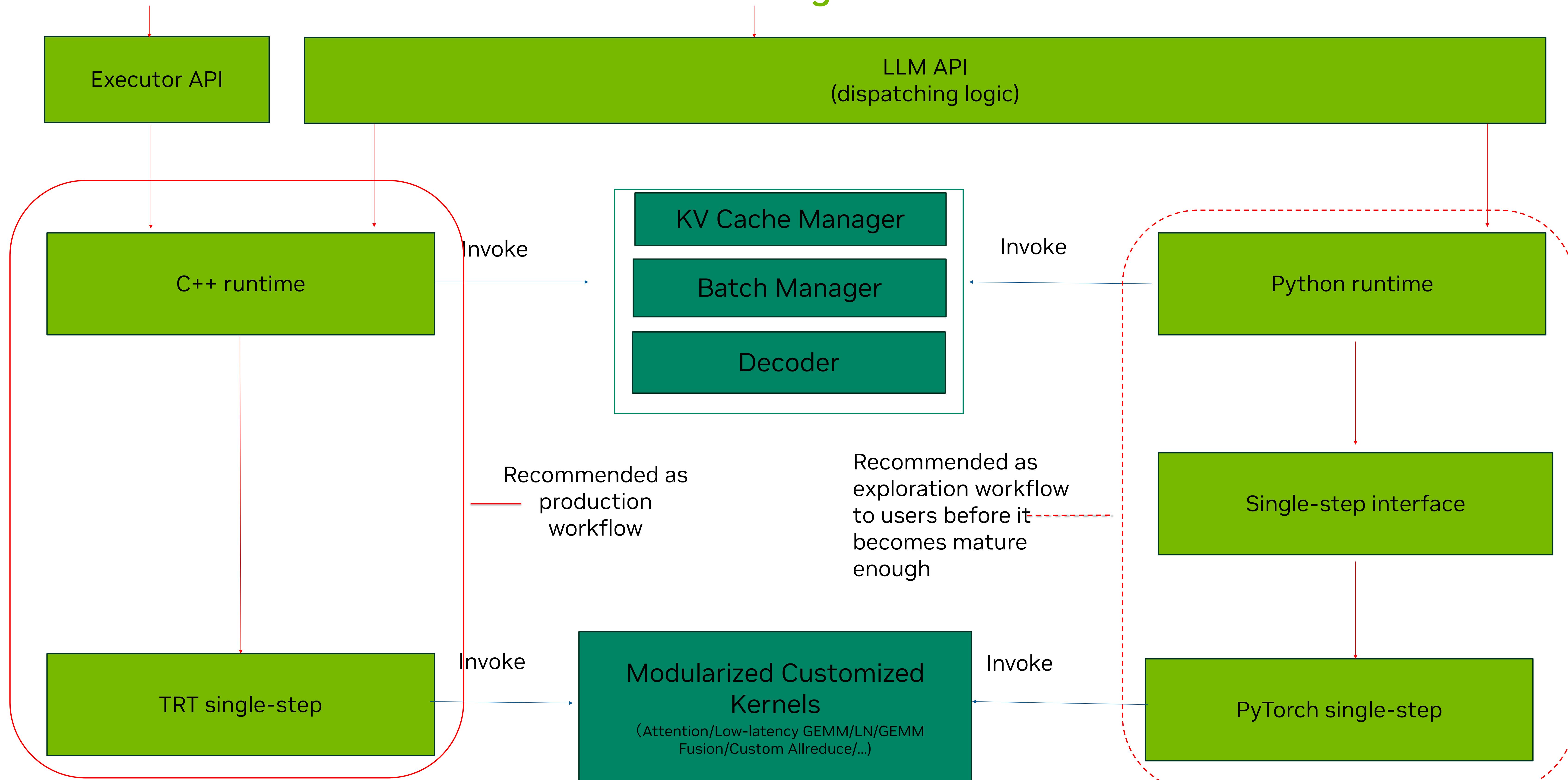
# Run Inference!
prompts = ["How do I get the fastest LLM Inference?"]

async for output in llm.generate_async(prompts, streaming=True):
    print(output)
```

Inflight Batching Usage

PyTorch Workflow release

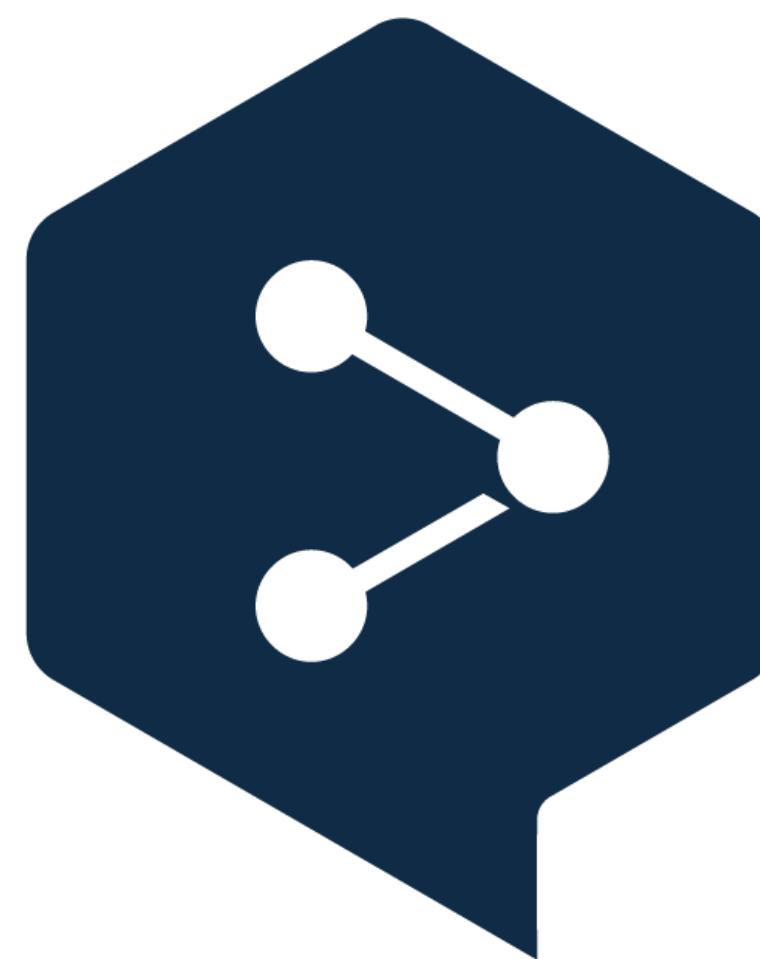
Starting v0.17



DeepL Use Case

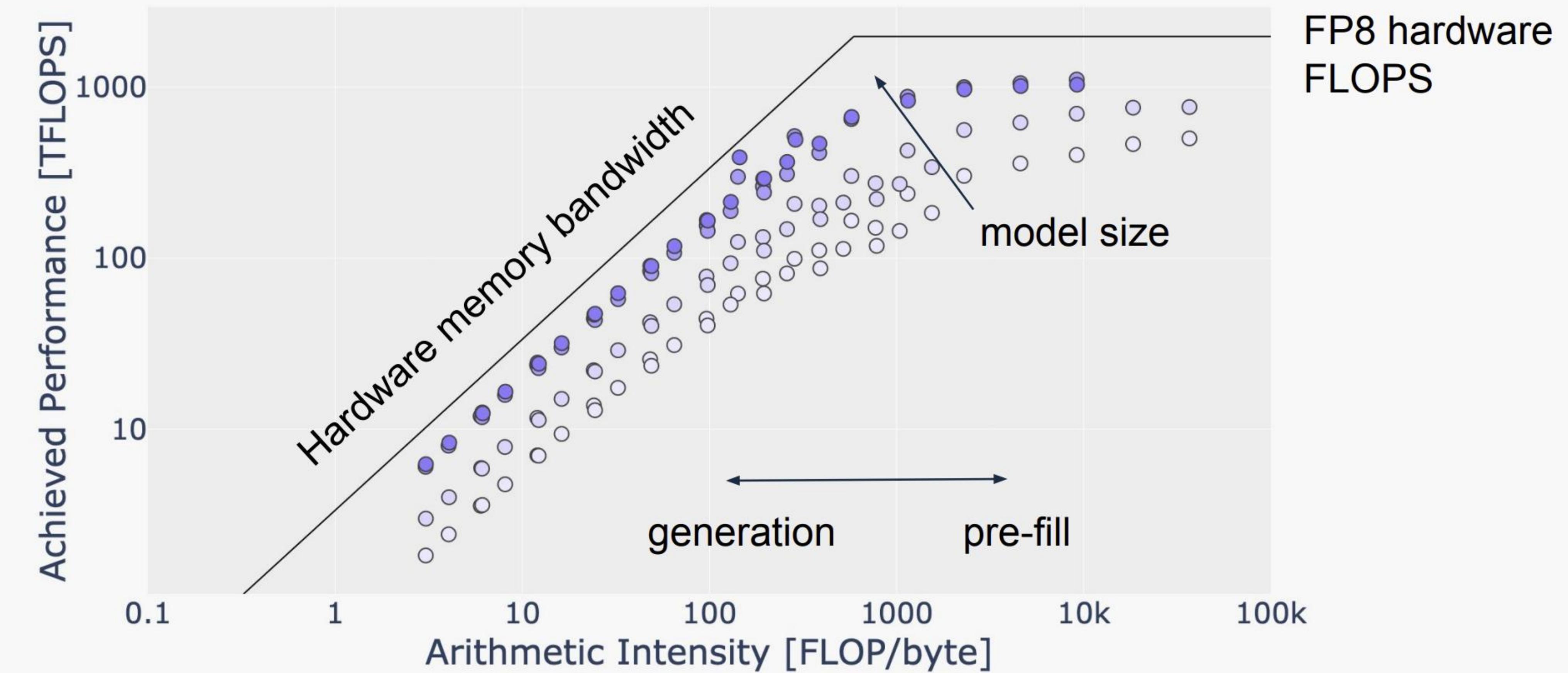
From FP8 LLM Training to Inference: Language AI at Scale

<https://www.nvidia.com/en-us/on-demand/session/gtc25-s72799/>



DeepL

TensorRT-LLM roofline model



Sweep over different model sizes (1B - 40B dense LLMs), batch sizes, generation and pre-fill lengths on NVIDIA H100

NVIDIA Dynamo Platform

The Operating System for AI Factories

NVIDIA Dynamo Platform

NVIDIA Dynamo

Distributed and Disaggregated
Generative AI Serving

NVIDIA Dynamo Triton (formerly Triton Inference Server)

Standardized Model Deployment
Across Every AI Workload

NVIDIA NIM

Fastest and Easiest Way to Deploy NVIDIA Dynamo Platform in Production

NVIDIA NIM Microservices

NVIDIA NIM Optimized Inference Microservices

Rapidly deploy reliable building blocks for accelerated generative AI anywhere



Portable Run cloud-native microservices anywhere, maintaining security and control of data and apps

Easy to Use Move fast with the latest agentic AI building blocks for reasoning, retrieval, images and more, deployed in minutes with standard APIs

Enterprise Supported Gain confidence with stable APIs, quality assurance, continuous updates, security patching, and support

Performance Optimize accuracy, latency and throughput to meet requirements with lowest TCO



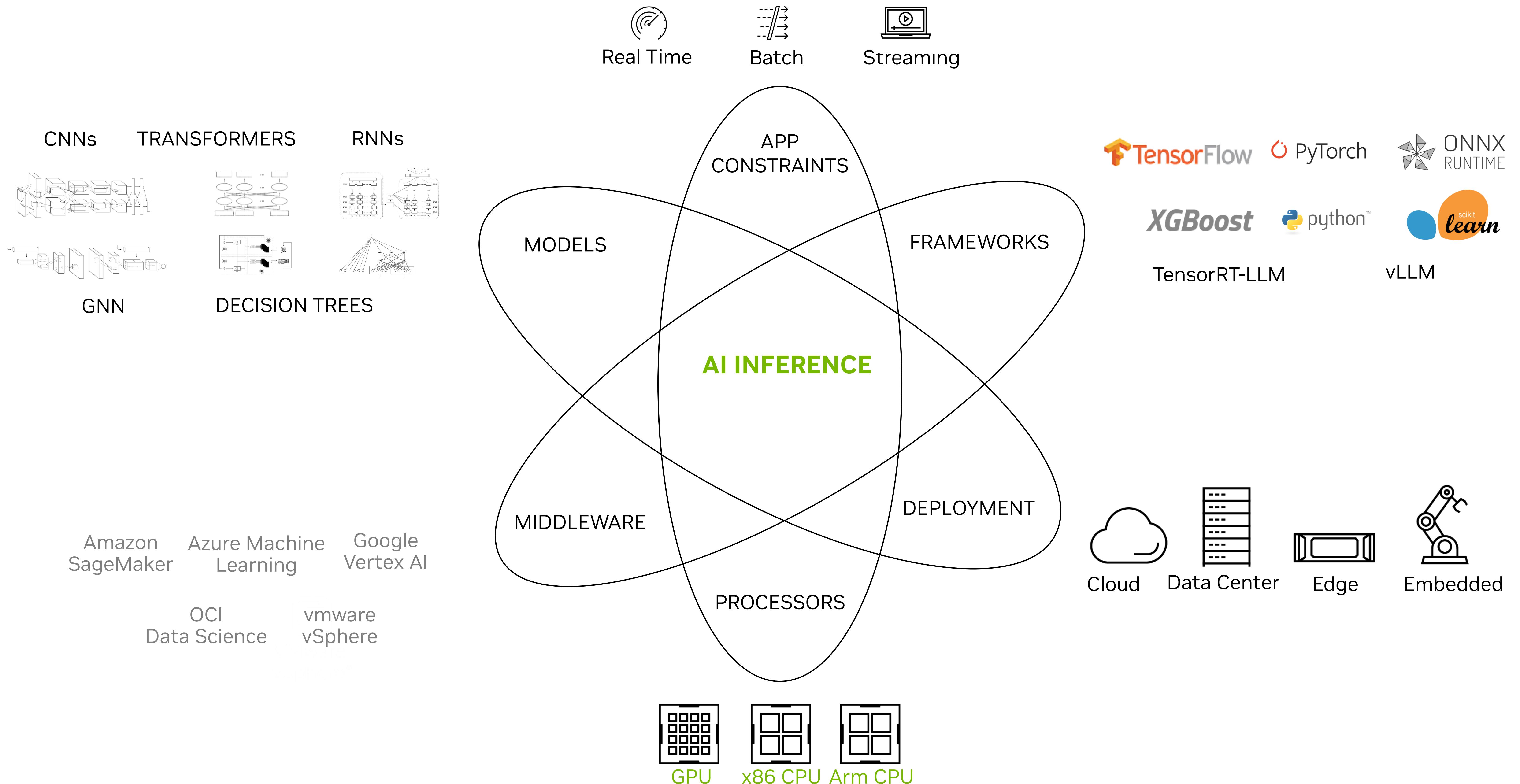
DGX &
DGX Cloud



NVIDIA Dynamo Triton

(Formerly Triton Inference Server)

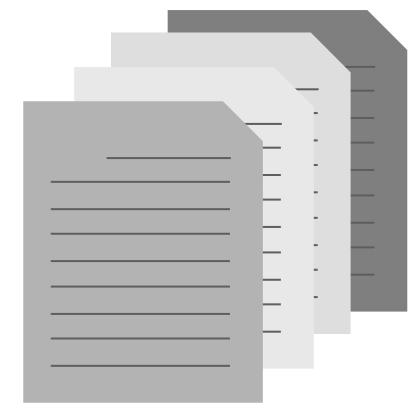
Challenges of AI Inference



NVIDIA Dynamo Triton (formerly Triton Inference Server)

Deploy models from all popular frameworks across GPUs and CPUs

Any Framework



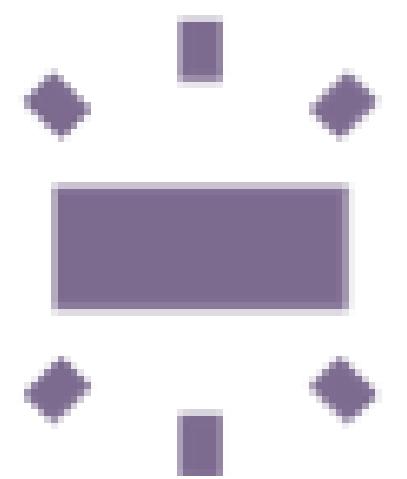
Supports Multiple Framework Backends
Natively e.g., TensorFlow, PyTorch, TensorRT, XGBoost, ONNX, Python, TensorRT-LLM, vLLM & More

Any Query Type



Optimized for Real Time, Batch, Streaming, Ensemble Inferencing

Any Platform



X86 CPU | Arm CPU | NVIDIA GPUs | MIG
Linux | Windows | Virtualization
Public Cloud, Data Center and Edge/Embedded (Jetson)

DevOps & MLOps



Integration With Kubernetes, KServe, Prometheus & Grafana
Available Across All Major Cloud AI Platforms

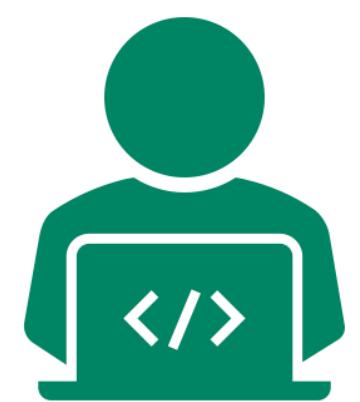
Performance & Utilization



Model Analyzer for Optimal Configuration
Optimized for High GPU/CPU Utilization, High Throughput & Low Latency

Benefits of NVIDIA Dynamo Triton

Standardize AI model deployment and execution across every workload

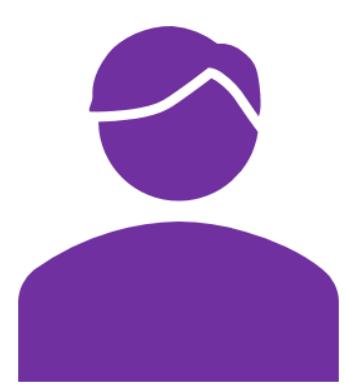


Standardized serving
across AI models and
frameworks

High performance
inference

Consistent deployment
across public cloud, on-
prem and edge

App developer | ML
engineer | MLOPS



IT, DevOps,
Platform team

Scale effortlessly with
application demand

CPUs and GPUs supported
as 1st class citizens

Easy to integrate in existing
production infrastructure

Start with CPU, evolve to
GPU

Faster model rollout

GPU accelerated inference
leads to \$ savings
(compute, floor space,
power)

Open source
Constant release of new
features

NVIDIA AI Enterprise
Enterprise grade & support
w/ SLA

Certified on major software
platforms

Welcome & Introduction



Dmitry Mironov is a Senior Deep Learning Solutions Architect at NVIDIA. His focus is on large-scale efficient deployment of LLMs and their inference.

He has co-authored NVIDIA Deep Learning Institute (DLI) workshops on the topic, including the one on sizing LLM inference systems. Prior to NVIDIA, Dmitry served as a co-founder and CTO of a startup.

He had been integrating computer vision into gold mining, transportation, energy, and other industries.

Demo 1: TRT-LLM + Triton

https://nvidia-my.sharepoint.com/:v/r/personal/dmitrym_nvidia_com/Documents/Videos/Cli pchamp/Triton-TRT-LLM/Exports/Triton-TRT-LLM.mp4?csf=1&web=1&e=pXDf5c

NVIDIA Dynamo

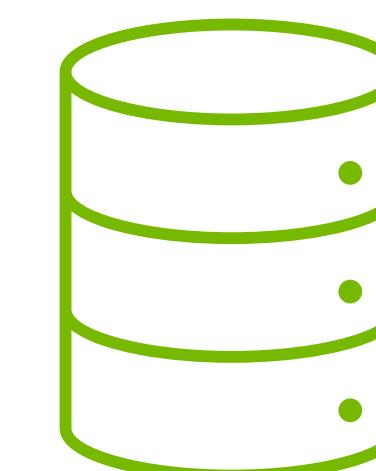
Pre-Training: Teaching AI Models Knowledge

What is NVIDIA?

NVIDIA is a multinational technology company that specializes in designing and manufacturing graphics processing units (GPUs), high-performance computing hardware, and artificial intelligence (AI) technologies. The company was founded in 1993 by Jensen Huang, Chris Malachowsky, and Curtis Priem.

NVIDIA is headquartered in Santa Clara, California, and has become one of the leading companies in the field of computer graphics, gaming, and AI.

Actual Llama 8B Response



400 TB Web Data

100 Trillion Tokens

Knowledge is not the same as Thinking

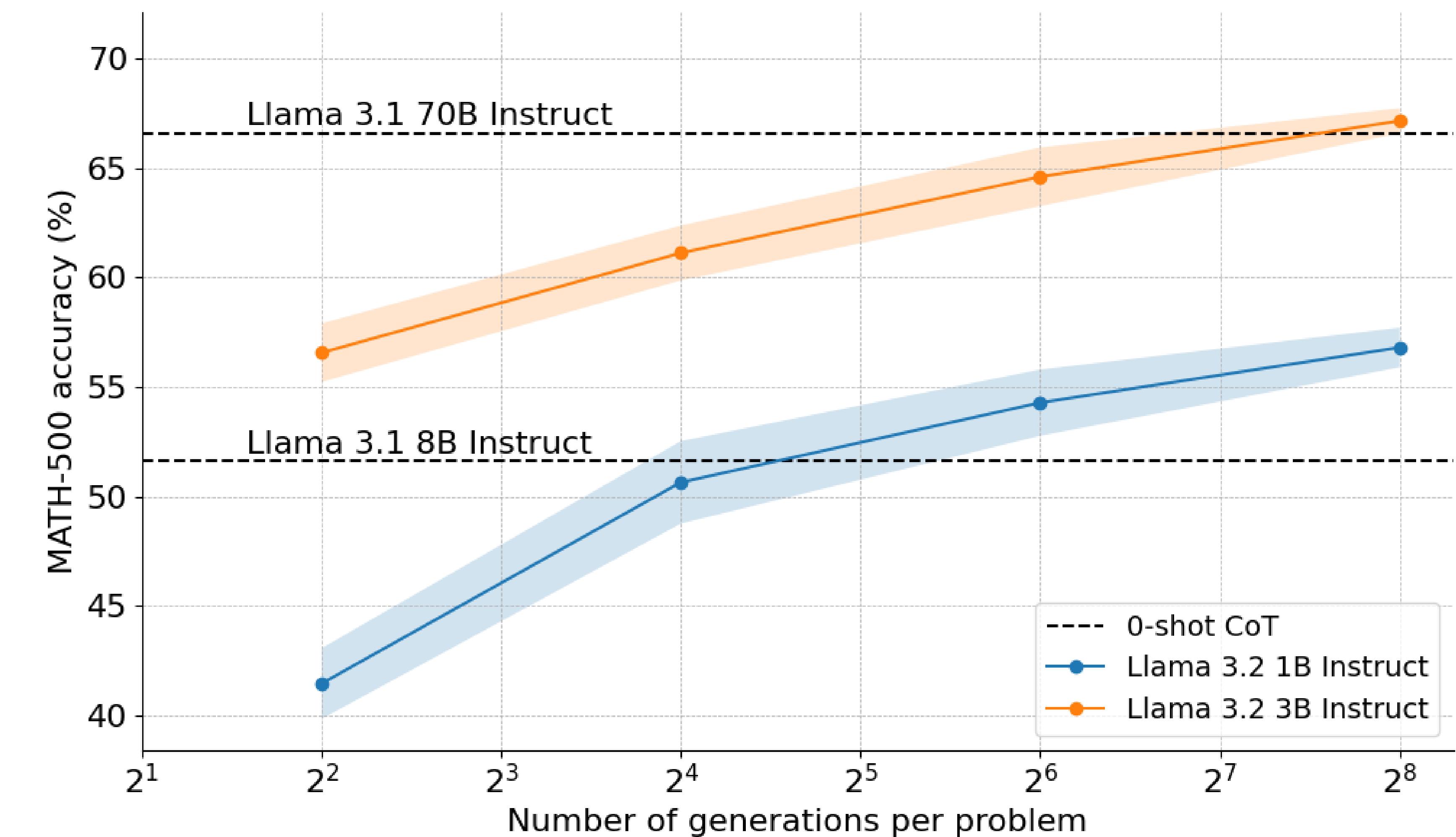
I need to figure out where to seat my family of 8 for dinner at a round table.
My in-laws don't get along with my parents.
My sister needs to sit next my little brother to help him eat.
My wife really doesn't want to sit next to my mom.
Where should each person be seated?



Actual Llama 8B Response

What Is Different for Reasoning?

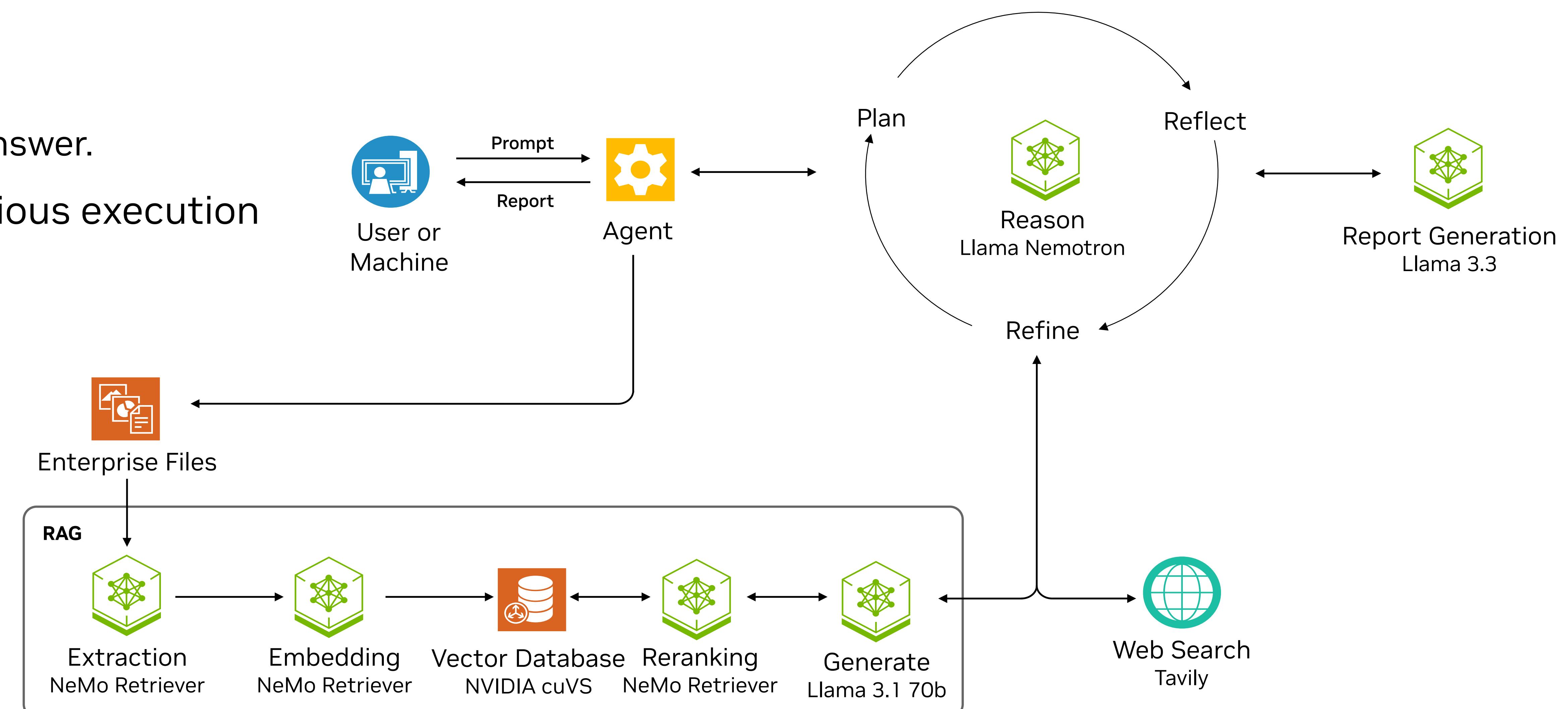
- Reasoning = Inference Time Compute
 - Spending more tokens to get higher accuracy
- Reasoning models have longer outputs.
 - More GPU time / request.
- Distributed Inference is coming.
 - System is optimized as a whole.
 - 2 GPUs are now better than 2 instances 1 GPU each
- Each query is more expensive – therefore input guardrails are common.



<https://huggingface.co/spaces/HuggingFaceH4/blogpost-scaling-test-time-compute>

How Do Multi-Agentic Workflows Impact Inference

- Complex workflows, where each step may be an LLM call, a tool call, guardrails, RAG pipeline, and other action.
- Lots of steps require E2E latency, because next steps depend on their full output.
- Thus, latencies are long (minutes to hours to days). More of batch-based, offline workflows.
- Complex Inference Time Compute
 - May require reward model to select the answer.
- One user request may include tens of various execution steps.



Inference Compute Requirements Scaling Exponentially

Fueled by reasoning models and AI agents



Larger Models

Hundreds of billions of parameters



Long Thinking Time

100x more thinking tokens



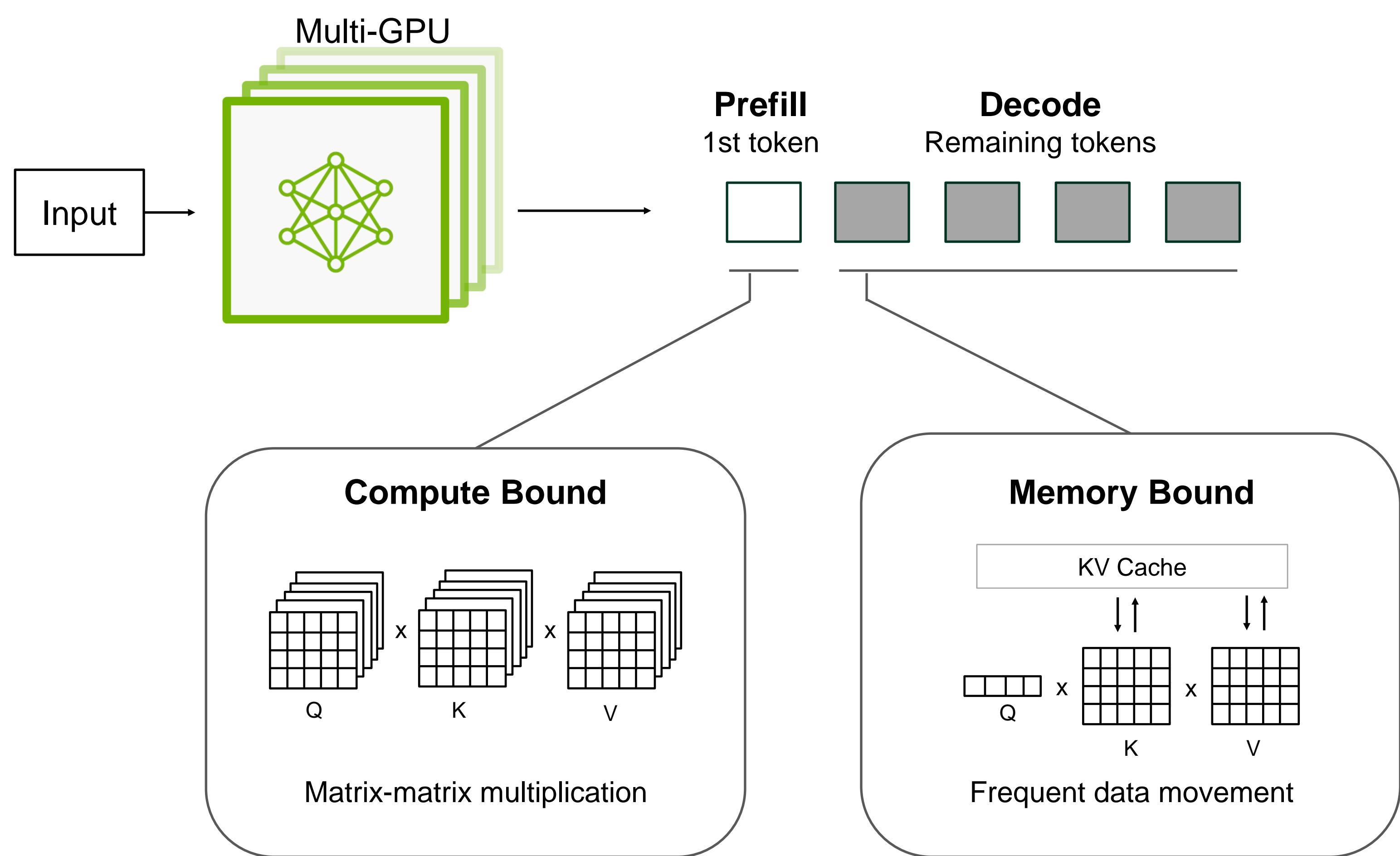
Larger Context

Millions of input tokens

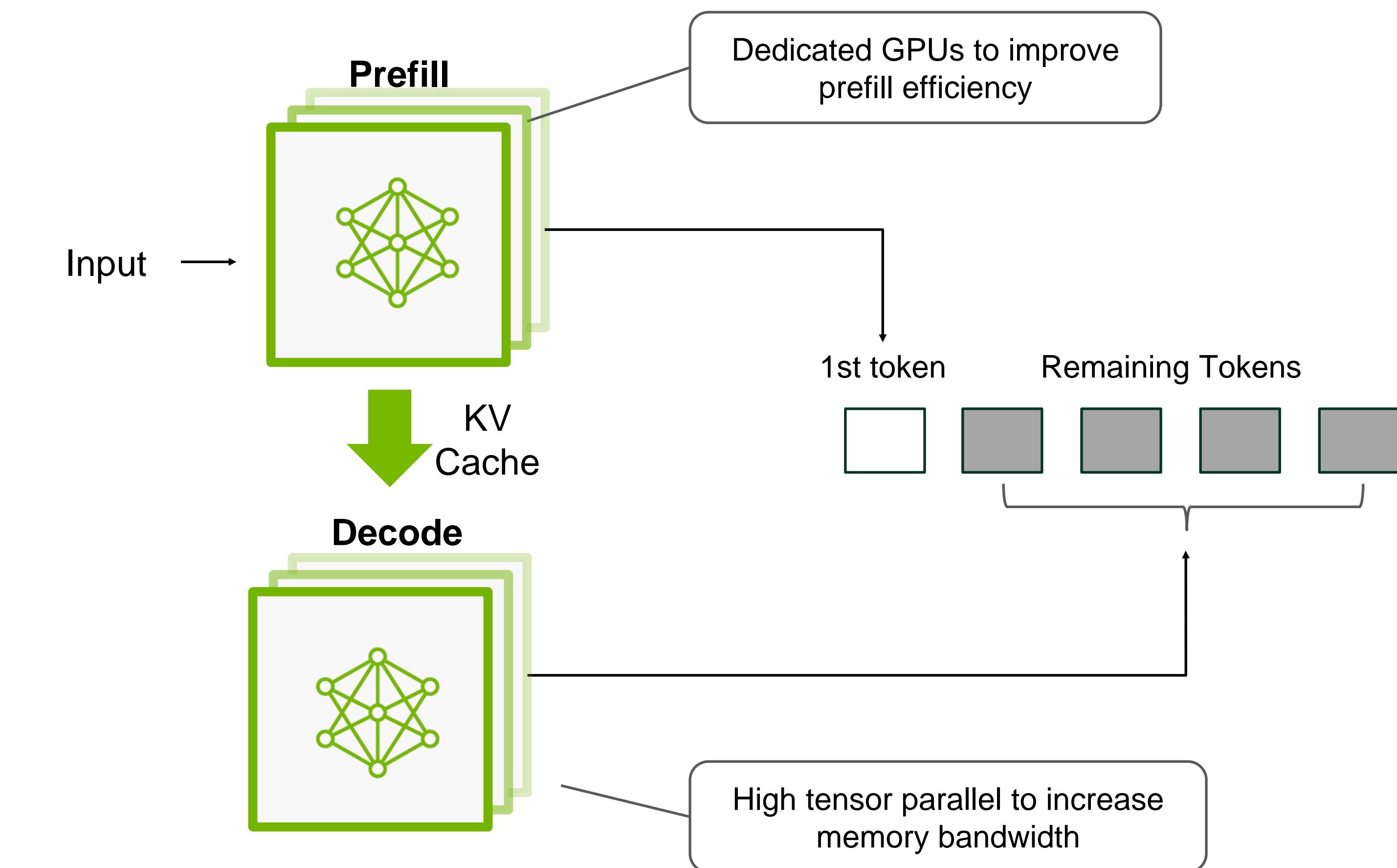
New Inference Optimization Techniques to Boost Inference

Disaggregated serving separates prefill and decode allowing each to be optimized independently

Traditional Serving



Disaggregated Serving



More flexibility to optimize cost and user experience

Announcing NVIDIA Dynamo

AI Inference Software for Reasoning Inference at Scale

30X

AI Factory Throughput
& Revenue
DeepSeek models
on Blackwell

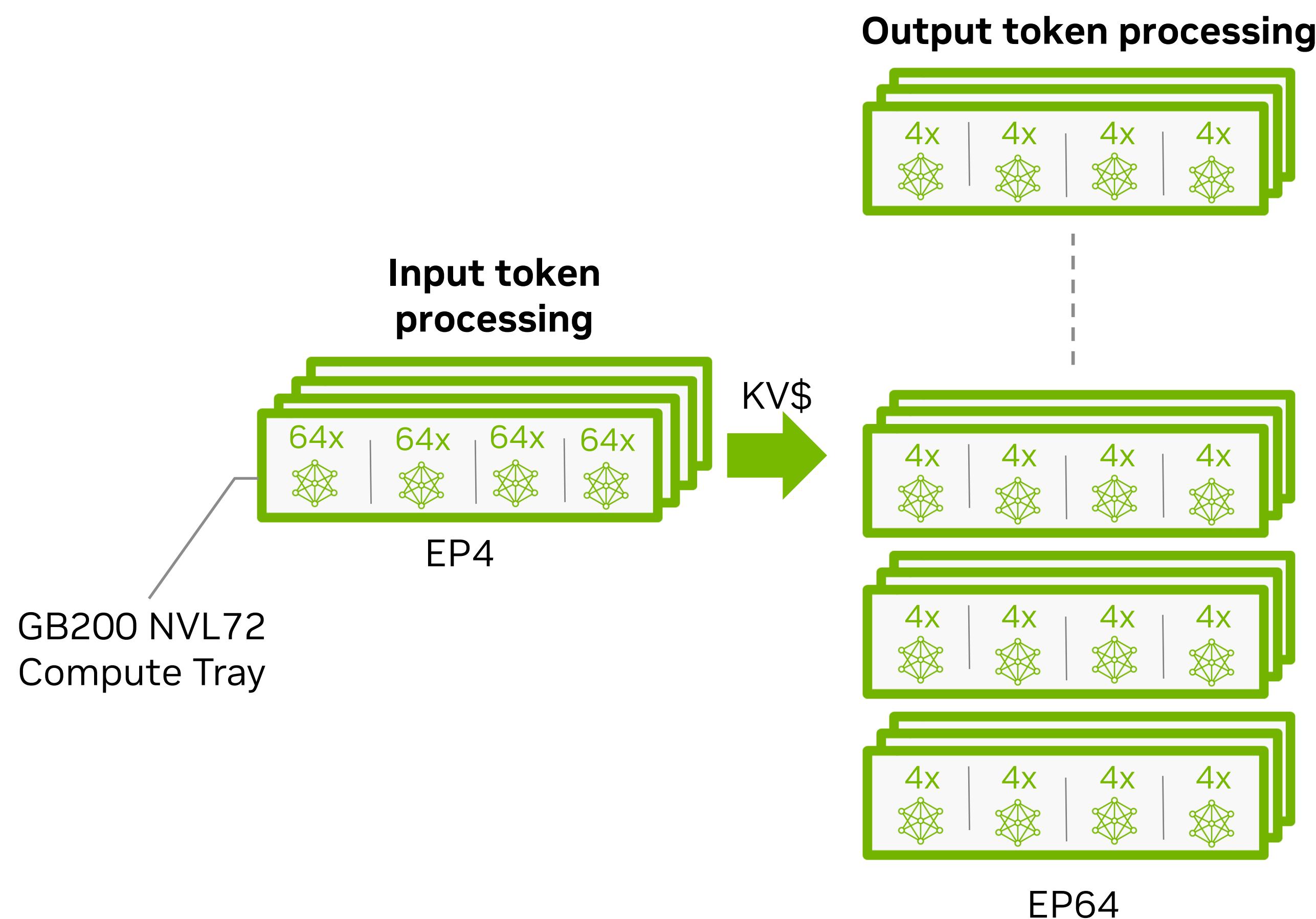
1000+

GPU Scale for
a single query

2X

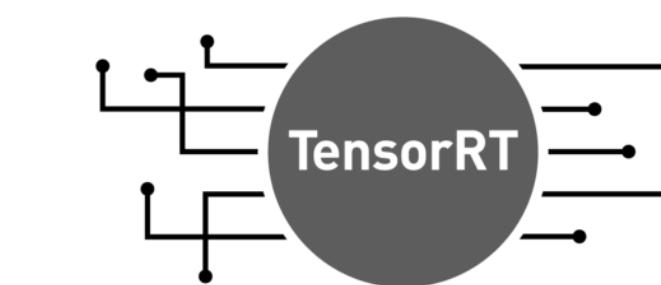
Throughput &
Revenue
Llama Models
On Hopper

Distributed and Disaggregated Serving DeepSeek R1

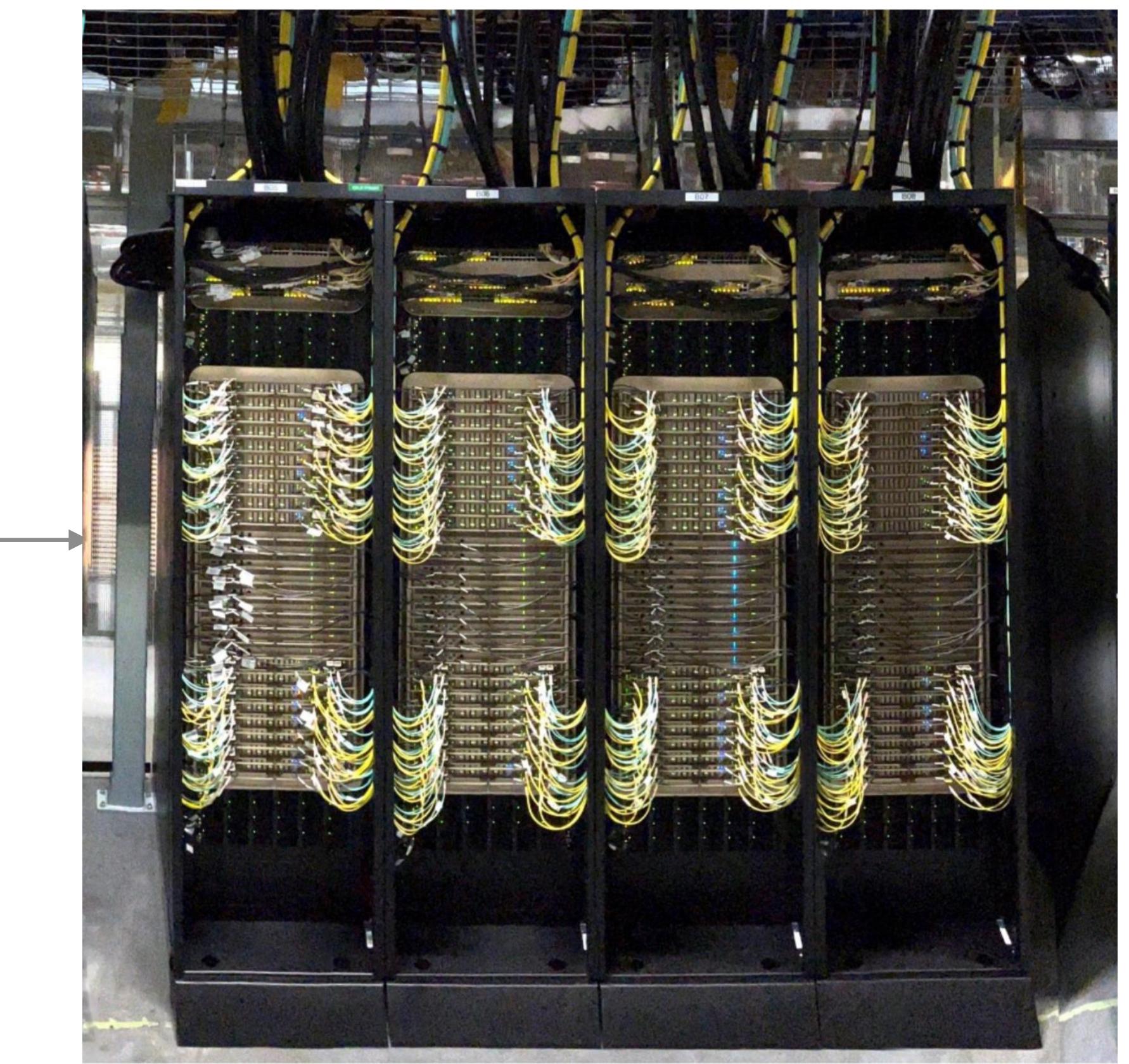
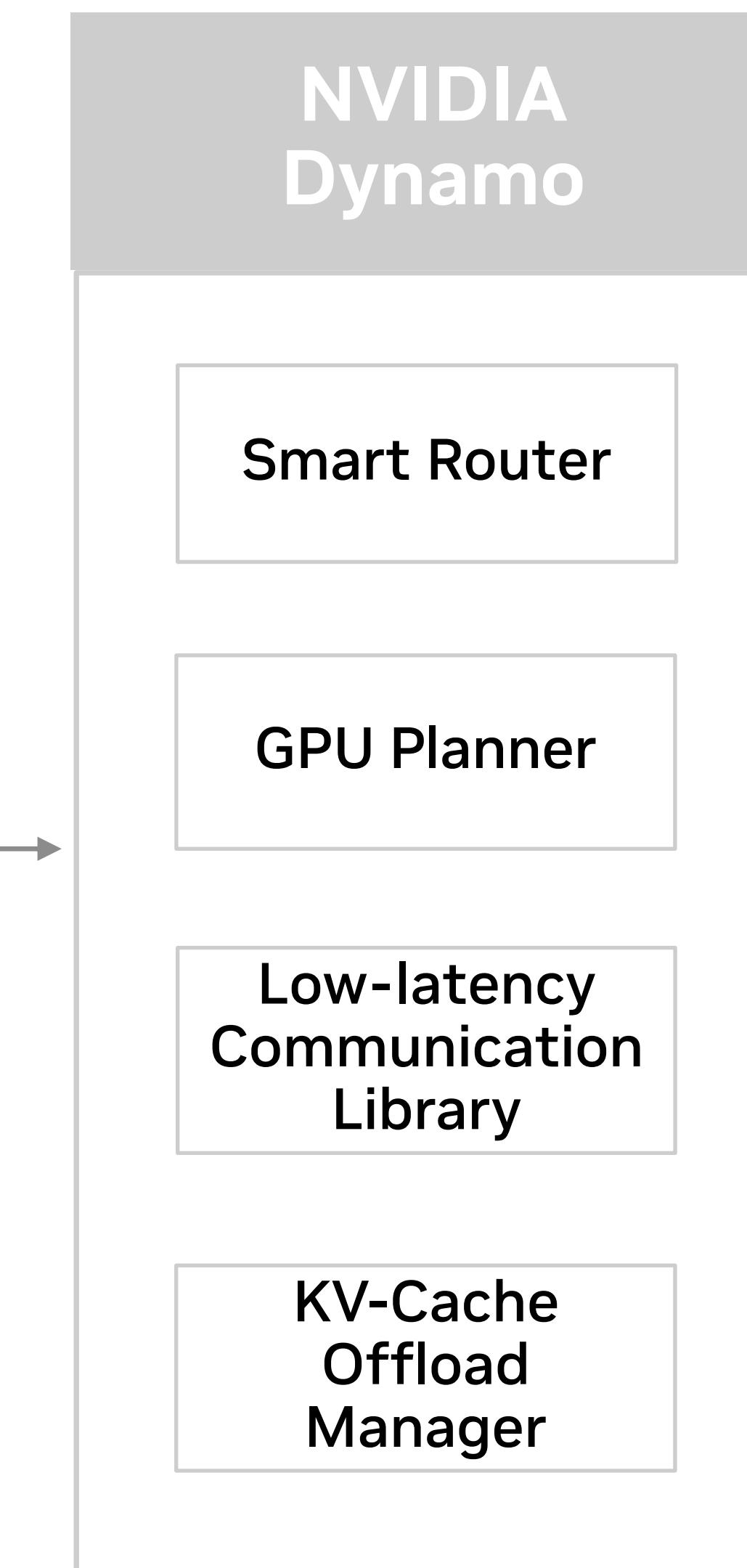


PyTorch

SGL

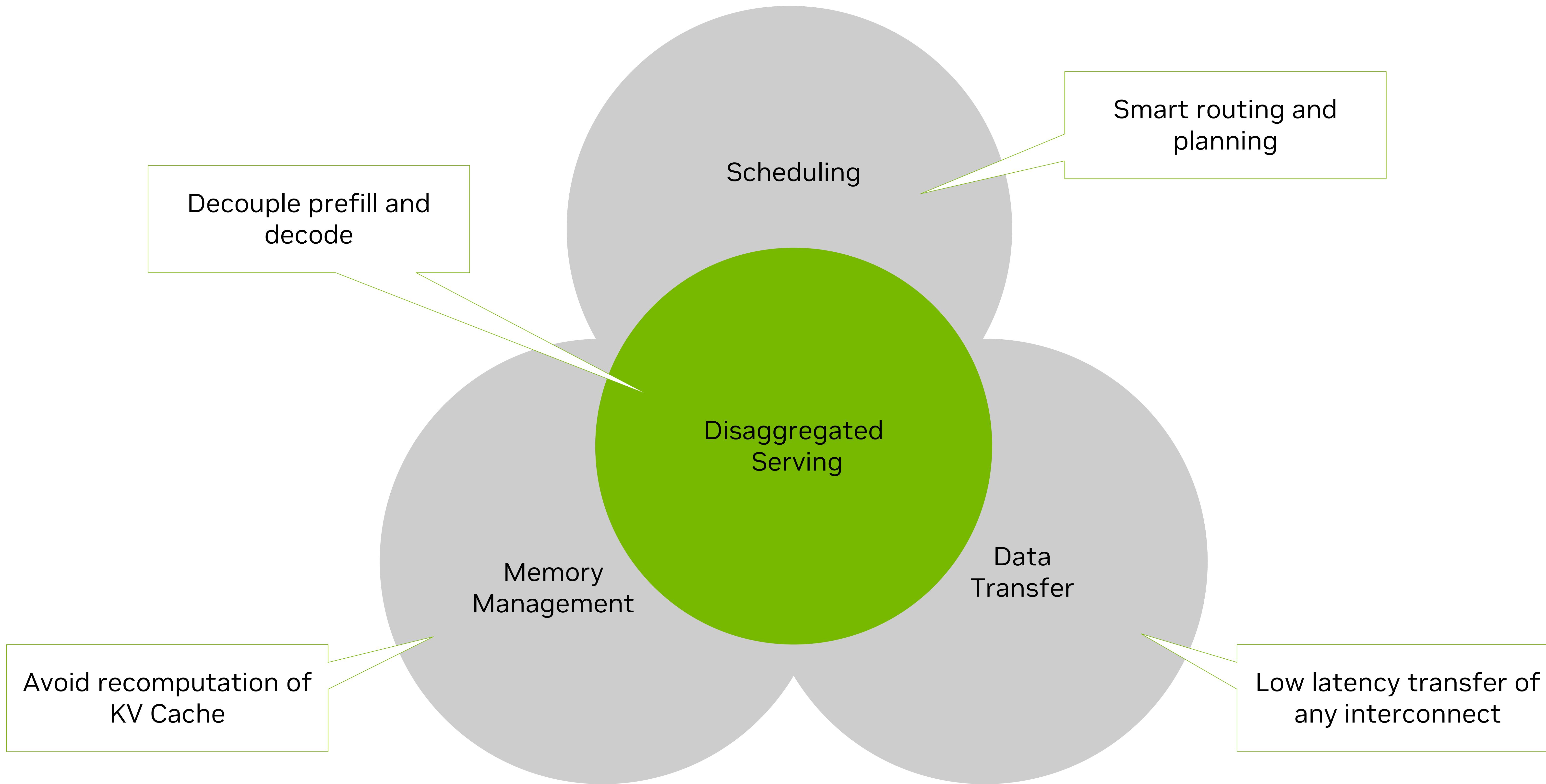


vLLM



NVIDIA Dynamo

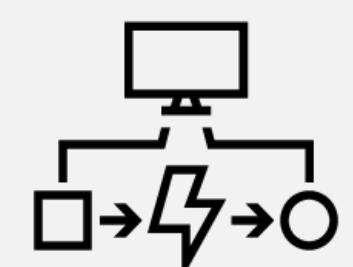
AI Inference Software for Reasoning Inference at Scale



NVIDIA Dynamo

A modular generative AI inference server designed for reasoning models and AI agents

NVIDIA Dynamo



Distributed Inference Serving

Seamlessly scale LLMs from a single GPU to thousands of GPUs



GPU Planning & Scheduling

Meet changing demand patterns w/o over or under provisioning of resources



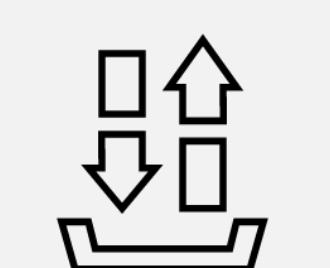
Smart Request Router

Free up GPU resources by reducing re-computations for similar requests



LLM Engine Plugin

Support for any TensorRT-LLM, vLLM, or SGLang modes

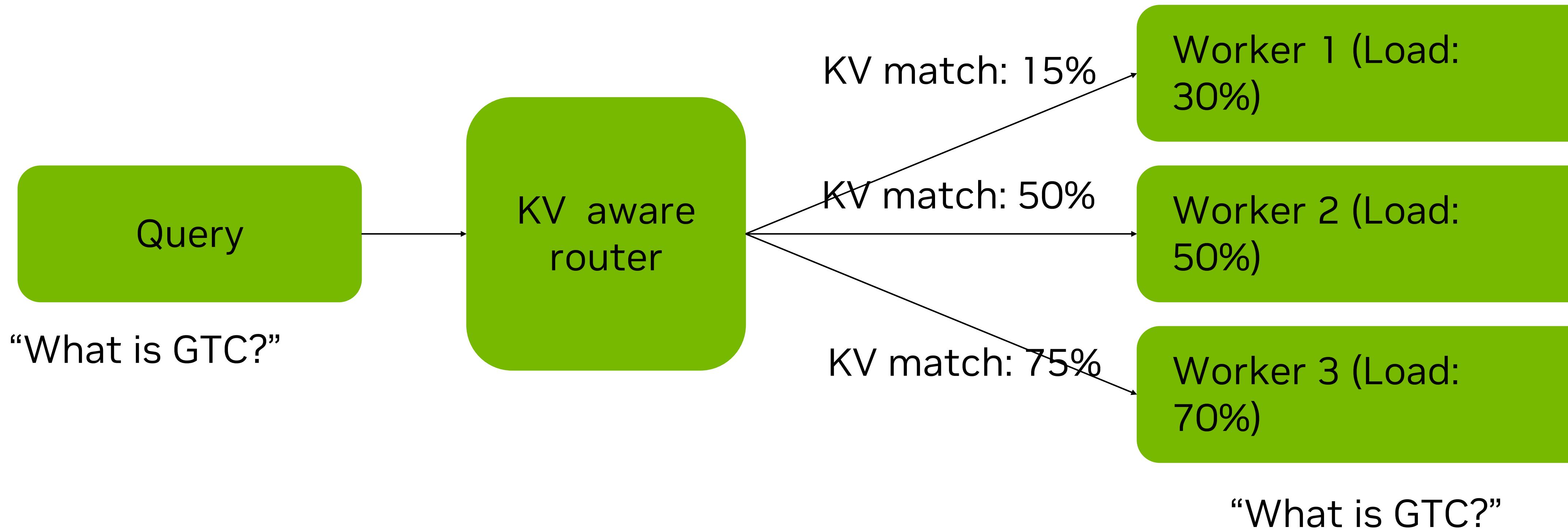
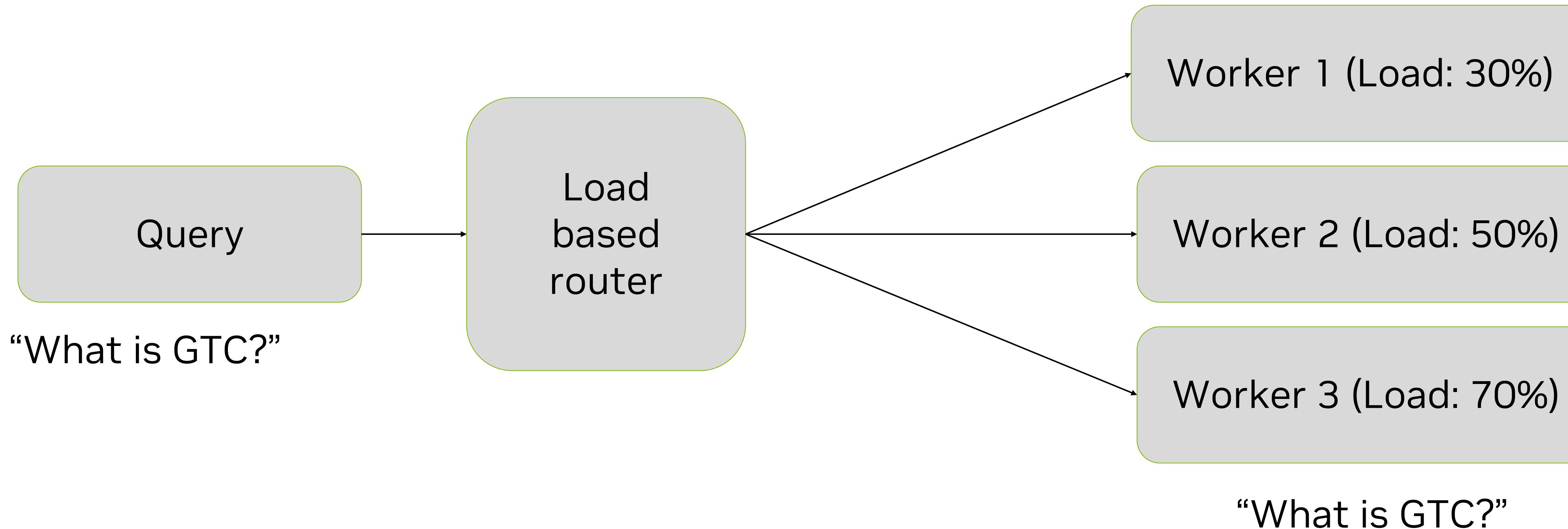


KV Cache Manager

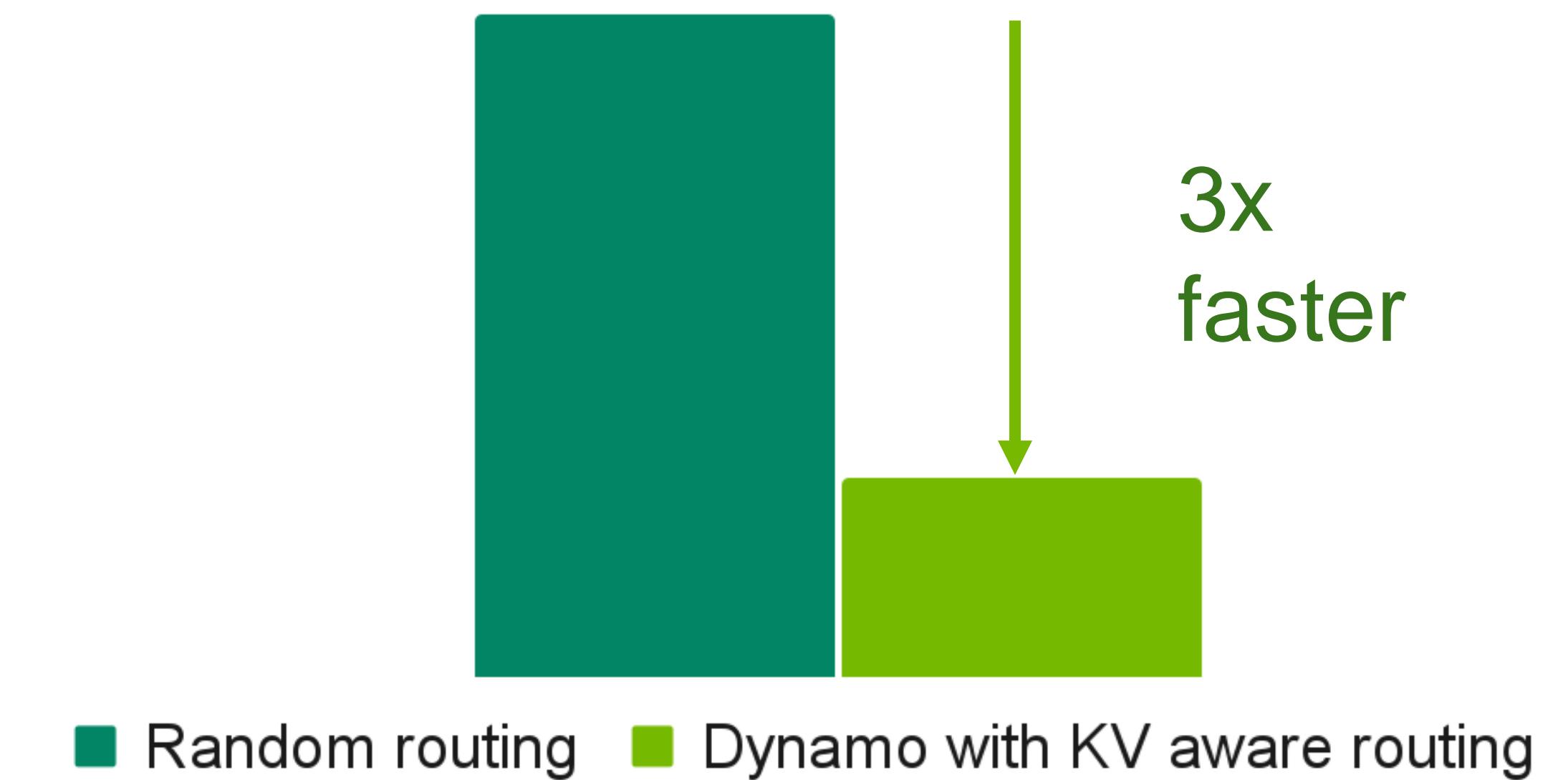
Preserve GPU memory by offloading context (KV\$) to cheaper storage

KV Cache Aware Routing

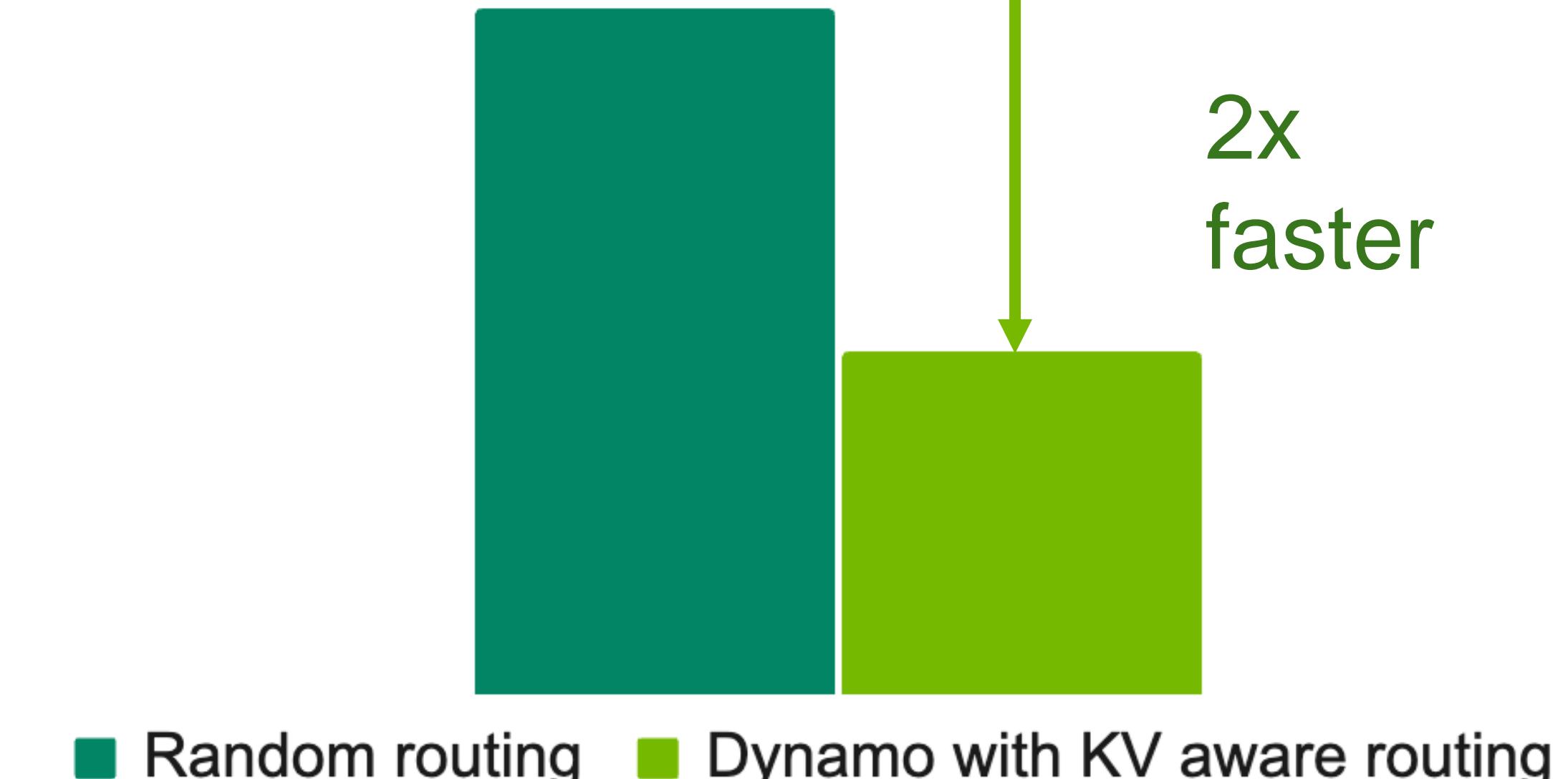
Significant boost in TTFT and End to End Latency with real data (100K requests with R1)



Time To First Token



Avg request latency

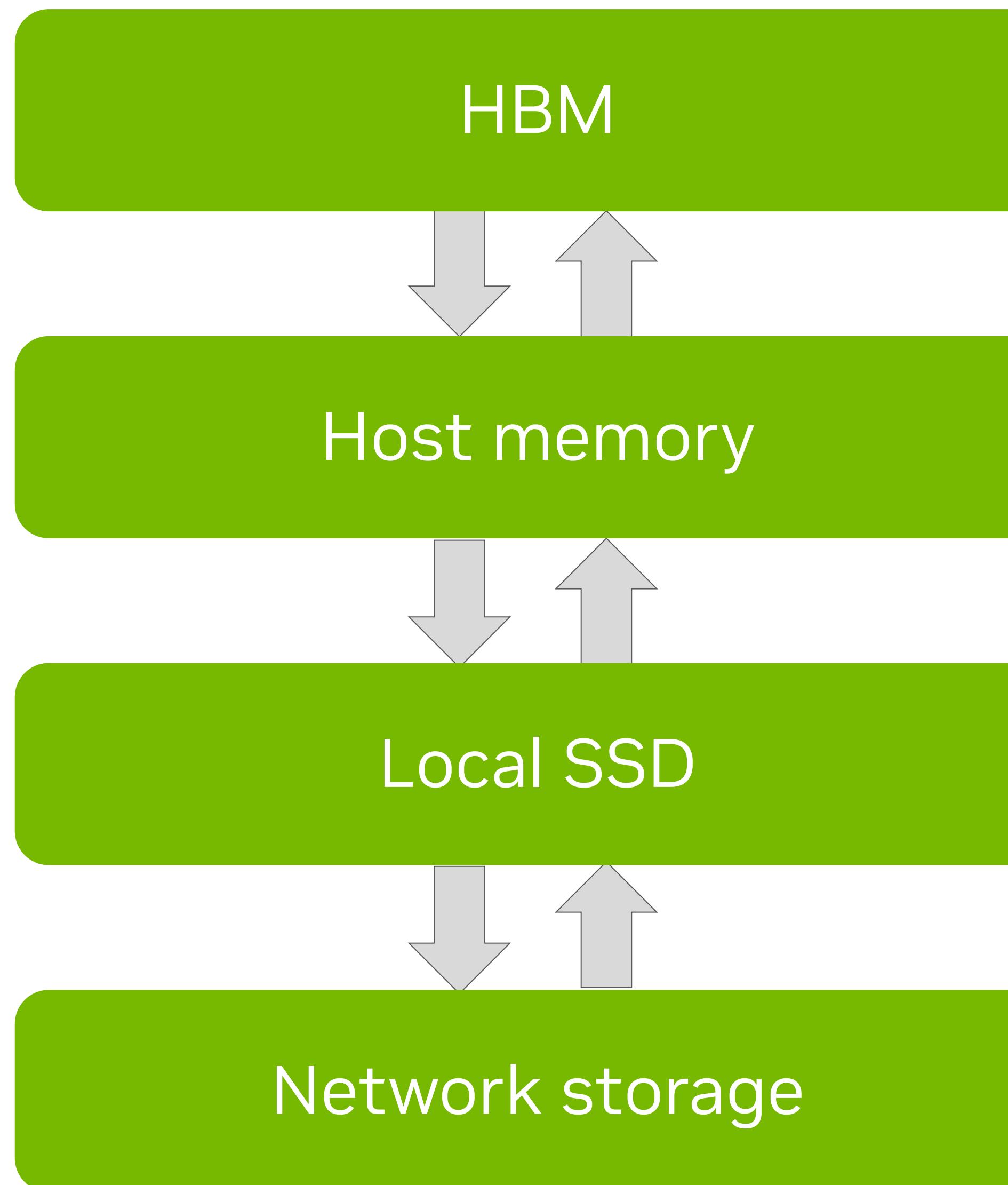


Tested with R1 Distilled Llama 70B over 2 nodes of 8 x H100s with vLLM 0.7.3

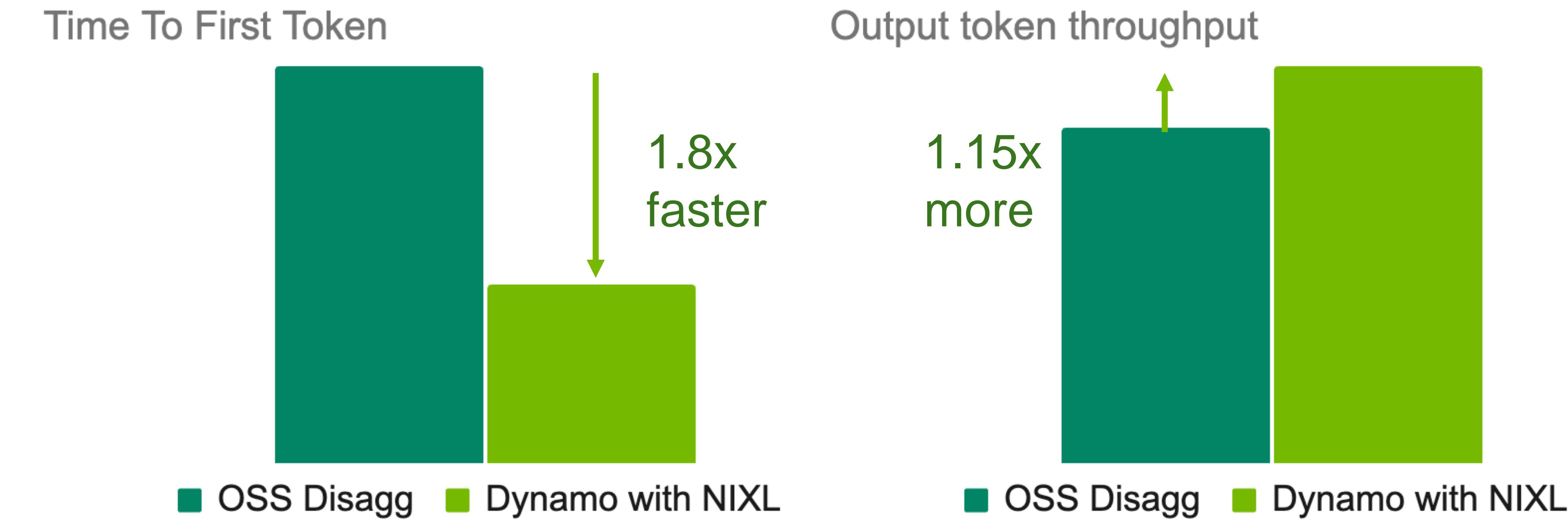


Data Transfer with NIXL (Nvidia Inference TranXfer Library)

Boosts both TTFT and throughput



- Low latency transfer between multiple nodes
- Dynamicity required for scaling and transfer to storage
- Common layer for inference specific network optimization.



Tested with vLLM 0.7.3

Dynamo as of Today

ai-dynamo

Overview Repositories 6 Projects Packages People

Dynamo

3 followers United States of America

Github: [nvidia.com/dynamo](https://github.com/nvidia/dynamo)

Find a repository...

dynamo Internal

A Datacenter Scale Distributed Inference Serving Framework

Rust 12 Apache-2.0 2 1 30 Upda

nixl Internal

NVIDIA Inference Xfer Library (NIXL)

- Apache 2 license and public CI
- Pip wheels on PyPi
- Rust for perf and Python for extensibility
- Discord for developer community
- Dynamo CLI
 - `dynamo run`: Quick start with model, input and output
 - `dynamo serve`: Construct graph of workers and serve
 - (EA) `dynamo build`: Containerize
 - (EA) `dynamo deploy`: Deploy to K8
- Three backends: TRT-LLM, vLLM, & SGLang
- Disaggregated serving with TRT-LLM and vLLM
- KV aware routing with TRT-LLM and vLLM
- (EA) KV manager with vLLM
- NIXL for RDMA and TCP (fallback for AWS EFA)

Customer Testimonial Highlights



Cohere

"Scaling advanced AI models requires sophisticated multi-GPU scheduling, seamless coordination and low-latency communication libraries that transfer reasoning contexts seamlessly across memory and storage. We expect NVIDIA Dynamo will help us deliver a premier user experience to our enterprise customers." **Saurabh Baji, Senior Vice President of Engineering at Cohere**



Perplexity AI

"Handling hundreds of millions of requests monthly, we rely on NVIDIA's GPUs and inference software to deliver the performance, reliability, and scale our business and users demand, "We'll look forward to leveraging NVIDIA Dynamo with its enhanced distributed serving capabilities to drive even more inference serving efficiencies and meet the compute demands of new AI reasoning models." **Denis Yarats, CTO of Perplexity AI.**



Together AI

"Scaling reasoning models cost-effectively requires new advanced inference techniques, including disaggregated serving and context-aware routing. Together AI provides industry leading performance using our proprietary inference engine. The openness and modularity of NVIDIA Dynamo will allow us to seamlessly plug its components into our engine to serve more requests while optimizing resource utilization—maximizing our accelerated computing investment." **Ce Zhang, CTO of Together AI.**

Demo 2: Dynamo

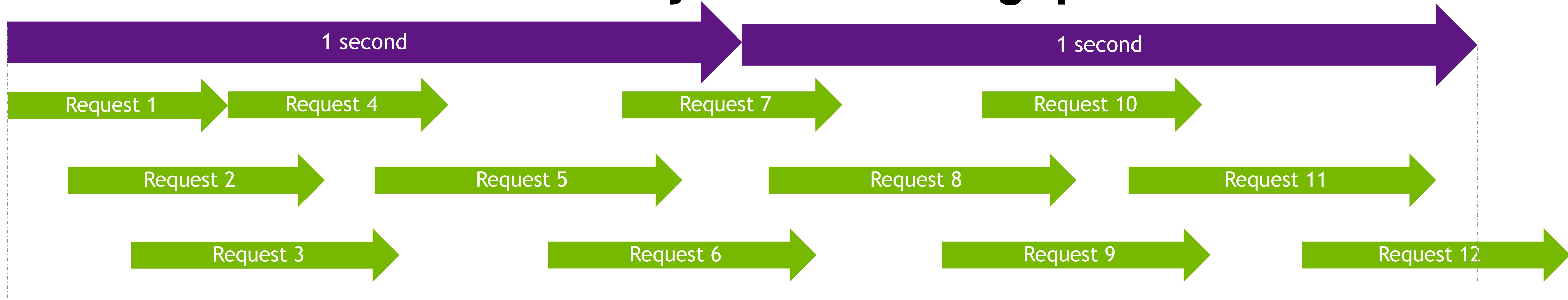
https://nvidia-my.sharepoint.com/:v/r/personal/dmitrym_nvidia_com/Documents/Videos/Cli pchamp/Dynamo%20Single%20GPU/Exports/Dynamo%20Single%20GPU.mp4?csf=1&web=1&e=adLKjk

Choosing the Optimal Solution for Your Needs

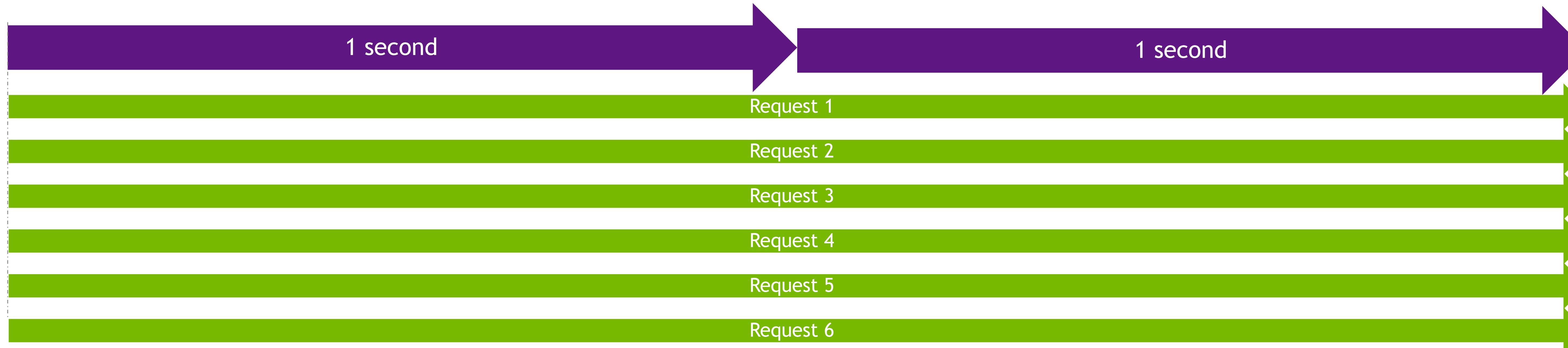
	NIM	Dynamo Triton (formerly Triton Inference Server)	Dynamo
Key Benefits	Fastest path for deployment	Standardizes AI model deployment and execution across every workload	Serves generative AI models in large scale distributed environments
Reasons to select	Out of the box performance with desired model.	Extensible deployment with any model and desired backend	Taking advantage of system level optimization such as disaggregated serving, smart routing with KV awareness
Supported via NVAIE	Yes	Yes	Planned
Support for VLMs	Yes	Yes	Planned
Support K8s	Yes	Yes	Yes
Separating Prefill and Decoding Nodes	Planned	No	Yes
Supports for any LLM	Coming soon	Yes	Yes

GenAI-Perf

Concurrency is NOT throughput



- Top and bottom: two different settings of a serving system, same model, same requests, same ISL and OSL.
- Top system: **concurrency** 3 (serving 3 concurrent requests), average E2E Latency 250 ms, RPS = 5.5
- Bottom system: **concurrency** 6, average E2E Latency 2000 ms (slower!), RPS = 3 (more GPUs needed!)



Concurrency is NOT Number of Concurrent Users

Example

- Application developers may call it 100 concurrent users:
 - some are typing
 - some are listening to the generated response
 - some are waiting for the interface to load
- **Convert Concurrent Users to Target Throughput**
- **Given**
 - 100 **concurrent users** at peak
 - Each session is 5 minutes
 - 10 requests to the LLM per session
- **Then**
 - Each user submits $10 \text{ requests} / 5 \text{ min} = 2 \text{ req / min / user}$
 - 100 concurrent users submit $100 * 2 = 200 \text{ req / min} = 3.3 \text{ req / s at peak (target RPS)}$

GenAI-Perf

Open-Source Benchmarking Client



1. Simulate LLM requests
 2. Measure time to get response
to calculate latency and throughput
- A green arrow points from the user icon towards the server icon. A second green arrow points back from the server icon towards the user icon.

GenAI-Perf



Any OpenAI-like API

Open Source Benchmarking tool for GenAI Performance measurement

Enables apples to apples comparison for evaluation and perf improvement.

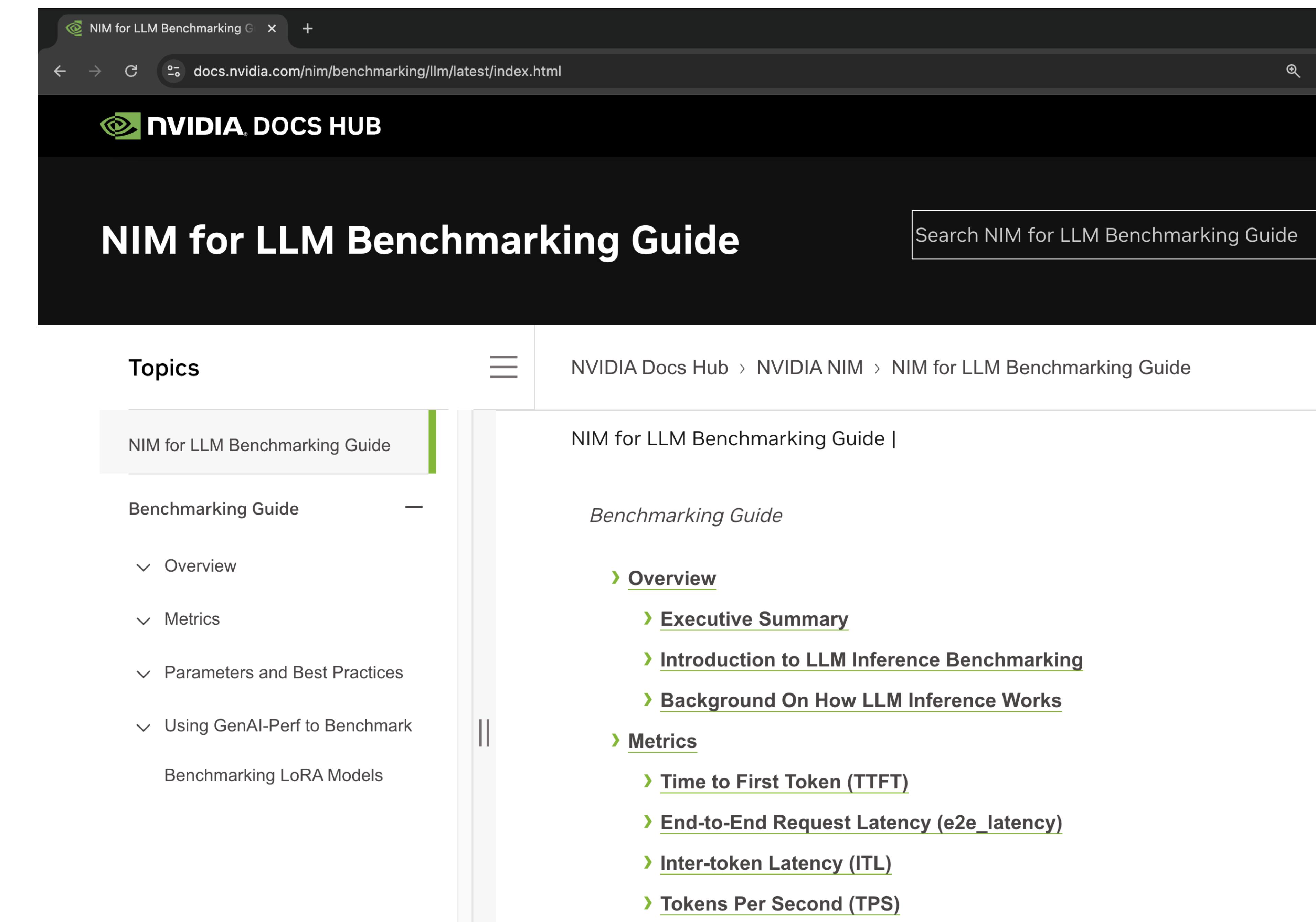
Measured via KServe & OpenAI API

Agnostic of underlying software and hardware

Provides refined implementation of metrics to guide action

Inter Token Latency (ITL): GenAI Perf's ITL metric focuses solely on decoding, excluding the initial token. This adjustment ensures a more targeted assessment of processing efficiency during decoding.

- Documented in NIM benchmarking guide
- <https://docs.nvidia.com/nim/benchmarking/llm/latest/index.html>



GenAI-Perf to Benchmark

The screenshot shows a web browser displaying the [NIM for LLM Benchmarking Guide](https://docs.nvidia.com/nim/benchmarking/llm/latest/step-by-step.html). The page title is "NIM for LLM Benchmarking Guide". A search bar at the top right contains the placeholder "Search NIM for LLM Benchmarking Guide" with a magnifying glass icon. The main content area features a green header "Using GenAI-Perf to Benchmark" which is highlighted with a green border. Below this, the text describes NVIDIA GenAI-Perf as a client-side LLM-focused benchmarking tool. The sidebar on the left is titled "Topics" and lists several steps: "Parameters and Best Practices", "Using GenAI-Perf to Benchmark" (which is expanded), "Step 1. Setting Up an OpenAI-Compatible LLama-3 Inference Service with NVIDIA NIM", "Step 2. Setting Up GenAI-Perf and Warming Up: Benchmarking a Single Use Case", "Step 3. Sweeping through a Number of Use Cases", "Step 4. Analyzing the Output", and "Step 5. Interpreting the Results". The "Using GenAI-Perf to Benchmark" section is the current focus.

NVIDIA Docs Hub > NVIDIA NIM > NIM for LLM Benchmarking Guide > Using GenAI-Perf to Benchmark

NIM for LLM Benchmarking Guide |

Using GenAI-Perf to Benchmark

NVIDIA [GenAI-Perf](#) is a client-side LLM-focused benchmarking tool, providing key metrics such as TTFT, ITL, TPS, RPS and more. It supports any LLM inference service conforming to the OpenAI API [specification](#), a widely accepted de facto standard in the industry. This section includes a step-by-step walkthrough, using GenAI-Perf to benchmark a Llama-3 model inference engine, powered by NVIDIA NIM.

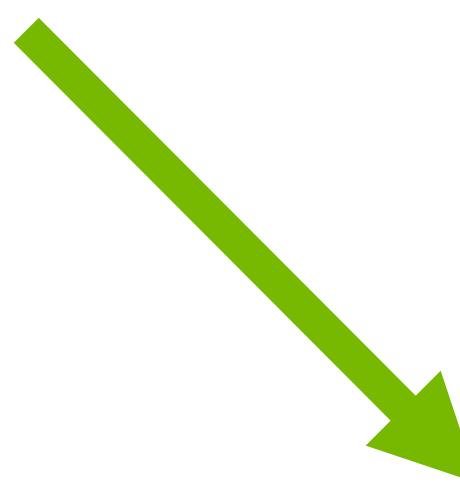
Step 1. Setting Up an OpenAI-Compatible LLama-3 Inference Service with NVIDIA NIM

GenAI-Perf command

Sample output generated by GenAI-Perf

```
export INPUT_SEQUENCE_LENGTH=200
export INPUT_SEQUENCE_STD=10
export OUTPUT_SEQUENCE_LENGTH=200
export CONCURRENCY=10
export MODEL=meta/llama3-8b-instruct

genai-perf profile \
-m $MODEL \
--endpoint-type chat \
--service-kind openai \
--streaming \
-u "http://meta-llama3-1-8b-instruct:8000" \
--synthetic-input-tokens-mean $INPUT_SEQUENCE_LENGTH \
--synthetic-input-tokens-stddev 0 \
--concurrency $CONCURRENCY \
--output-tokens-mean $OUTPUT_SEQUENCE_LENGTH \
--extra-inputs max_tokens:$OUTPUT_SEQUENCE_LENGTH \
--extra-inputs min_tokens:$OUTPUT_SEQUENCE_LENGTH \
--extra-inputs ignore_eos:true \
--measurement-interval 5000 \
--artifact-dir /genai-perf-results \
--profile-export-file ${INPUT_SEQUENCE_LENGTH}_${OUTPUT_SEQUENCE_LENGTH}_${CONCURRENCY}.json
-- \
-v \
--request-count $((10 * CONCURRENCY)) \
--max-threads=256
```



LLM Metrics

Statistic	avg	min	max	p99
Time to first token (ns)	85,485,242	27,402,273	152,621,817	130,194,943
Inter token latency (ns)	8,847,758	2,113,030	74,794,303	9,477,464
Request latency (ns)	1,848,822,497	1,844,511,394	1,924,017,143	1,905,132,459
Num output token	184	177	190	189
Num input token	200	198	201	200

Output token throughput (per sec): 995.61

Request throughput (per sec): 5.41

Sweeping across concurrencies

Running multiple GenAI-Perf calls

```
for concurrency in 1 2 5 10 50 100 250; do  
  
    local INPUT_SEQUENCE_LENGTH=$inputLength  
    local INPUT_SEQUENCE_STD=0  
    local OUTPUT_SEQUENCE_LENGTH=$outputLength  
    local CONCURRENCY=$concurrency  
    local MODEL=meta/llama3-8b-instruct  
  
    genai-perf \  
        -m $MODEL \  
        --endpoint-type chat \  
        --service-kind openai \  
        --streaming \  
        -u localhost:8000 \  
        --synthetic-input-tokens-mean $INPUT_SEQUENCE_LENGTH \  
        --synthetic-input-tokens-stddev $INPUT_SEQUENCE_STD \  
        --concurrency $CONCURRENCY \  
        --output-tokens-mean $OUTPUT_SEQUENCE_LENGTH \  
        --extra-inputs max_tokens:$OUTPUT_SEQUENCE_LENGTH \  
        --extra-inputs min_tokens:$OUTPUT_SEQUENCE_LENGTH \  
        --extra-inputs ignore_eos:true \  
        --tokenizer meta-llama/Meta-Llama-3-8B-Instruct \  
        --measurement-interval 10000 \  
done
```



Demo 3: GenAI-Perf

https://nvidia-my.sharepoint.com/:v/r/personal/dmitrym_nvidia_com/Documents/Videos/Cli-pchamp/GenAI-Perf%20Demo/Exports/GenAI-Perf%20Demo.mp4?csf=1&web=1&e=aEOnyA

Call to Action

- [Intro to LLMs](#)
- [Discover NVIDIA AI Enterprise Software](#)
- [Discover AI Inference Solution](#)
- [Explore NVIDIA AI Inference Performance](#)
- [Blog: Mastering LLM Techniques: Inference Optimization Technical Blog](#)
- [Explore NVIDIA TensorRT SDK](#)
- [NGC NIM Containers: Phind-CodeLlama-34B-v2-Instruct](#)
 - [NGC Containers: Llama-3.1-Nemotron-70B-Instruct](#)
 - [NGC Containers: Llama-3-Taiwan-70B-Instruct](#)
- [LLM Inference Sizing: Benchmarking End-to-End Inference Systems](#)
- [Transforming Medical Workflows With AI: A Deep Dive into CLLMs](#)

Claim your Free Self-Paced Course

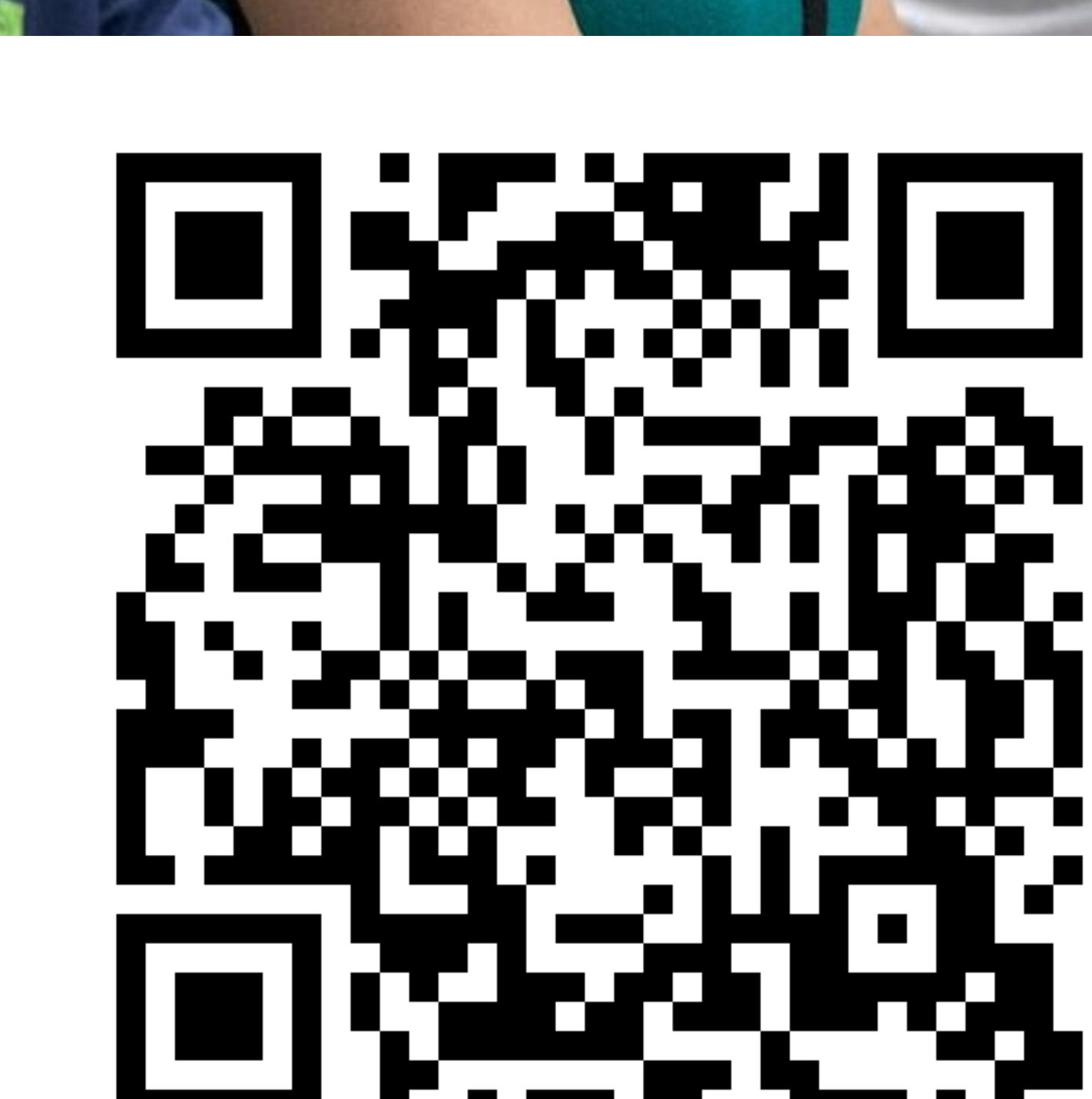
Access essential technical training

Sharpen your skills or learn a new technology. In partnership with NVIDIA Deep Learning Institute, we are offering a free self-paced course (worth up to \$90).

Courses on offer include:

- Accelerating End-to-End Data Science Workflows
- Building LLM Applications With Prompt Engineering
- Building RAG Agents with LLMs
- Building Real-Time Video AI Applications
- Deploying a Model for Inference at Production Scale
- Fundamentals of Accelerated Computing with CUDA Python
- Generative AI with Diffusion Models
- Get Started with Highly Accurate Custom ASR for Speech AI
- Sizing LLM Inference Systems
- Introduction to Deploying RAG Pipelines for Production at Scale
- Introduction to NVIDIA NIM™ Microservices

Scan the QR code to access the full course list and redeem your free training. (Redemption rule applies)





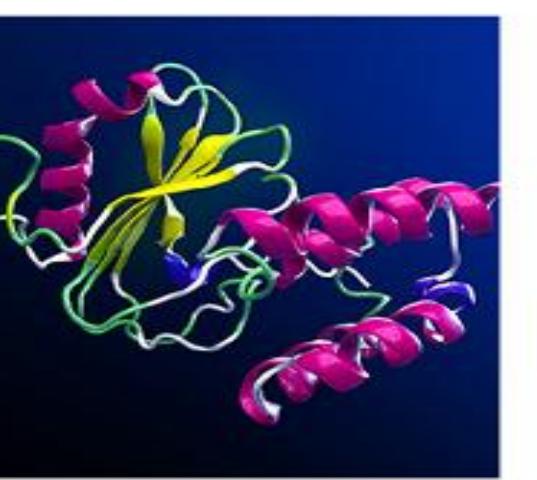
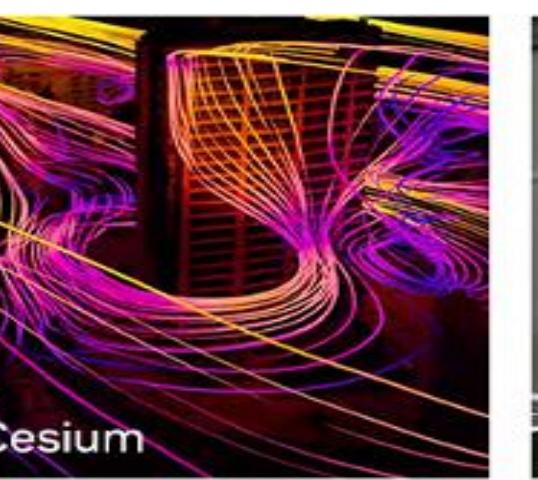
Thank You !

GTC 2025- Watch
on Demand



**What's Next in AI
Starts Here**

March 17-21, 2025



Demo Listing

ARG

```
BASE_IMAGE=nvcr.io/nvidia/tritonserver:2
5.03-trtllm-python-py3
FROM ${BASE_IMAGE} AS base
```

```
FROM base AS dev
```

Install additional packages

```
RUN pip install transformers==4.51.3 && \
    pip install flash-attn --no-build-
isolation && \
    pip install
/opt/tritonserver/python/triton*.whl
```

Clone and set up OpenAI frontend

```
RUN git clone https://github.com/triton-
inference-server/server.git && \
    cd server/python/openai/ && \
    pip install -r requirements.txt
```

Set working directory

```
WORKDIR /workspace
```

```
FROM dev AS final
```

-

```
from transformers import AutoModelForCausalLM, AutoTokenizer
import torch
```

```
# Load model and tokenizer from local path
# model_path = "/models/models--meta-llama--Llama-4-Scout-17B-16E-
# Instruct/current_snapshot"
model_path = "/models/Llama-3.1-8B-Instruct"
print(f"Loading model from local path: {model_path}")
```

```
try:
    # Use trust_remote_code=True for custom model code
    model = AutoModelForCausalLM.from_pretrained(
        model_path,
        device_map="auto",
        torch_dtype=torch.float16, # Use float16 as specified in the
        build_command
        local_files_only=True
    )
    tokenizer = AutoTokenizer.from_pretrained(
        model_path,
        local_files_only=True
    )
    print("Model and tokenizer loaded successfully!")
    print(f"Model type: {type(model).__name__}")
    print(f"Tokenizer type: {type(tokenizer).__name__}")

    # Test inference
    print("\nRunning test inference...")
    test_prompt = "What is the capital of France?"
    print(f"Test prompt: {test_prompt}")
```

```
# Tokenize input
inputs = tokenizer(test_prompt,
return_tensors="pt").to(model.device)
```

```
# Generate response
with torch.no_grad():
    outputs = model.generate(
        **inputs,
        max_new_tokens=50,
        temperature=0.7,
        top_p=0.9,
        do_sample=True
    )
```

```
# Decode and print response
response = tokenizer.decode(outputs[0], skip_special_tokens=True)
print(f"Model response: {response}")
```

```
except Exception as e:
    print(f"Error loading model: {str(e)}")
```

Demo Listing 2



```
pip install --upgrade
huggingface-hub
huggingface-cli download meta-
llama/Llama-3.1-8B-Instruct --
local-dir
/mnt/shared/work/Llama-3.1-8B-
Instruct
docker compose up --build
triton-trt-llm-build
docker compose up trt-llm
docker compose up trt-llm-build
docker compose up triton-deploy
docker compose up triton-sdk-
genai-perf
```

```
export HOME=$(pwd)
python3 -
m venv venv && . venv/bin/act
ivate

# install dynamo
pip install ai-dynamo[all]

# run dynamo
dynamo run out=vllm deepseek-
ai/DeepSeek-R1-Distill-Llama-
8B
```

