# NVIDIA Inference Stack Demystified
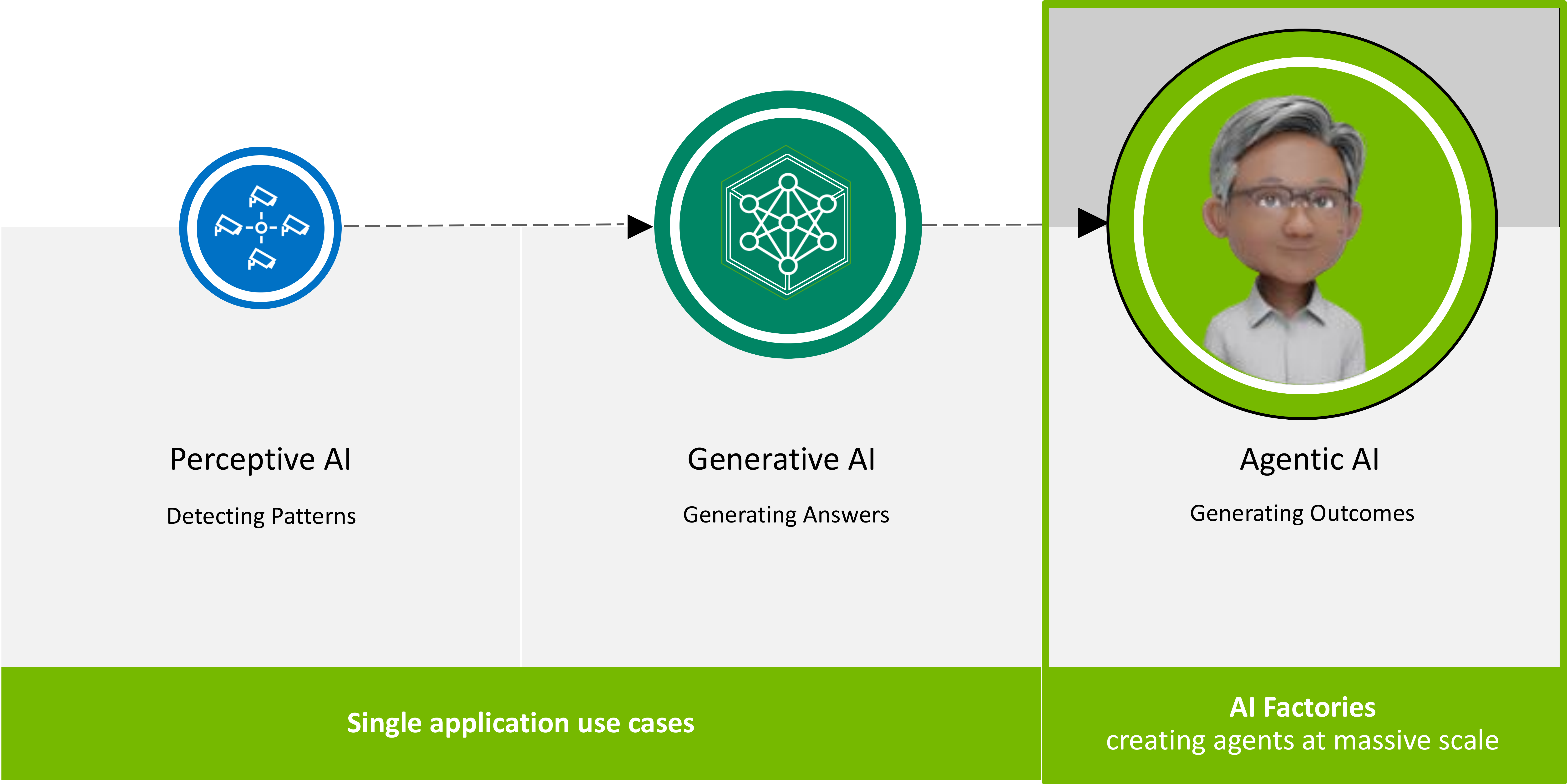
Amit Kumar | Manager | Industries | Solutions Architecture & Engineering | Generative AI | NVIDIA
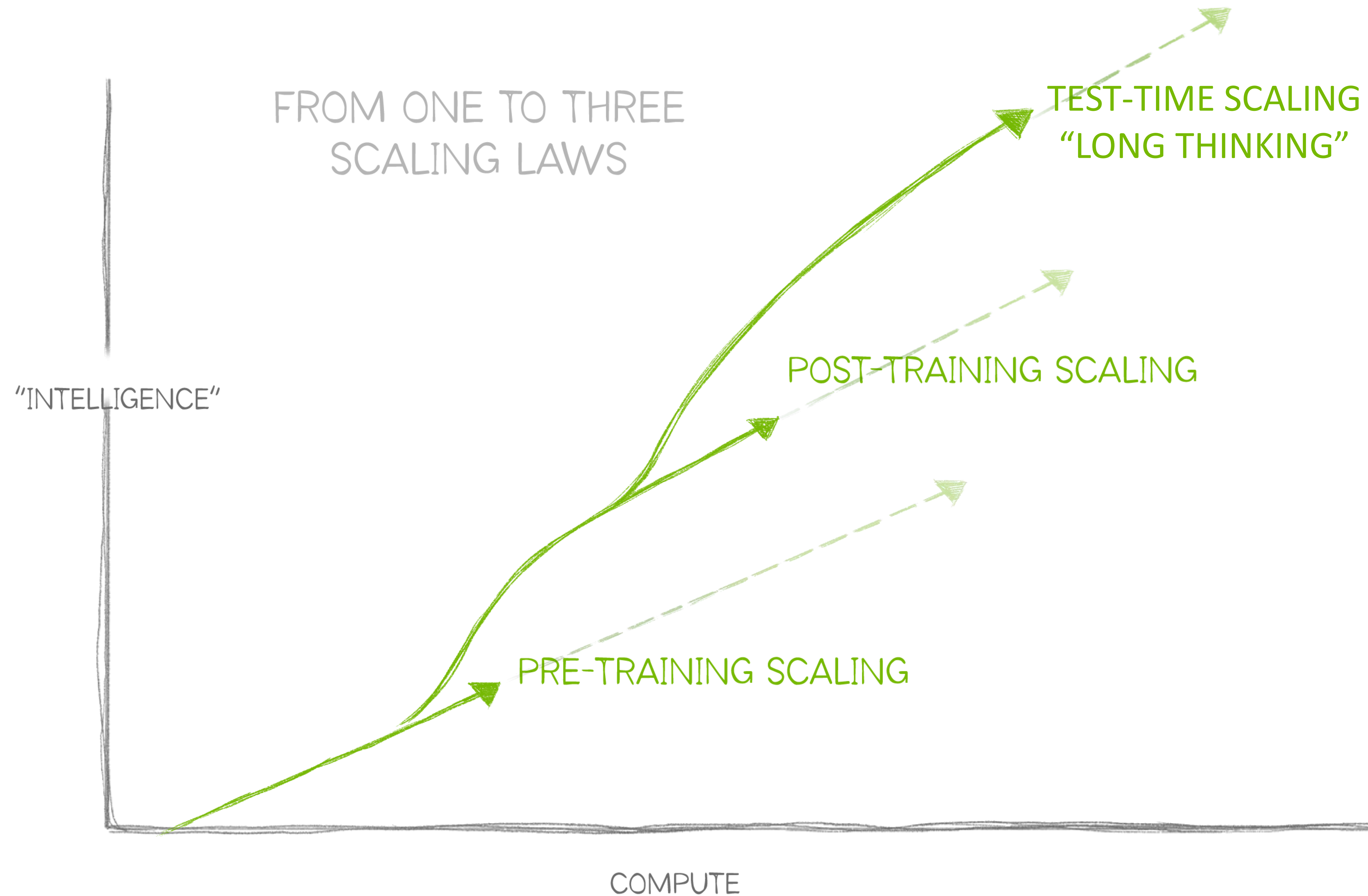
https://www.linkedin.com/in/amit-kumar-b4b43ab/

# AI Scaling Laws Drive Exponential Demand for Compute

New "long thinking" supercharges inference scaling

# Inference Compute Requirements Scaling Exponentially

Fueled by reasoning models and AI agents
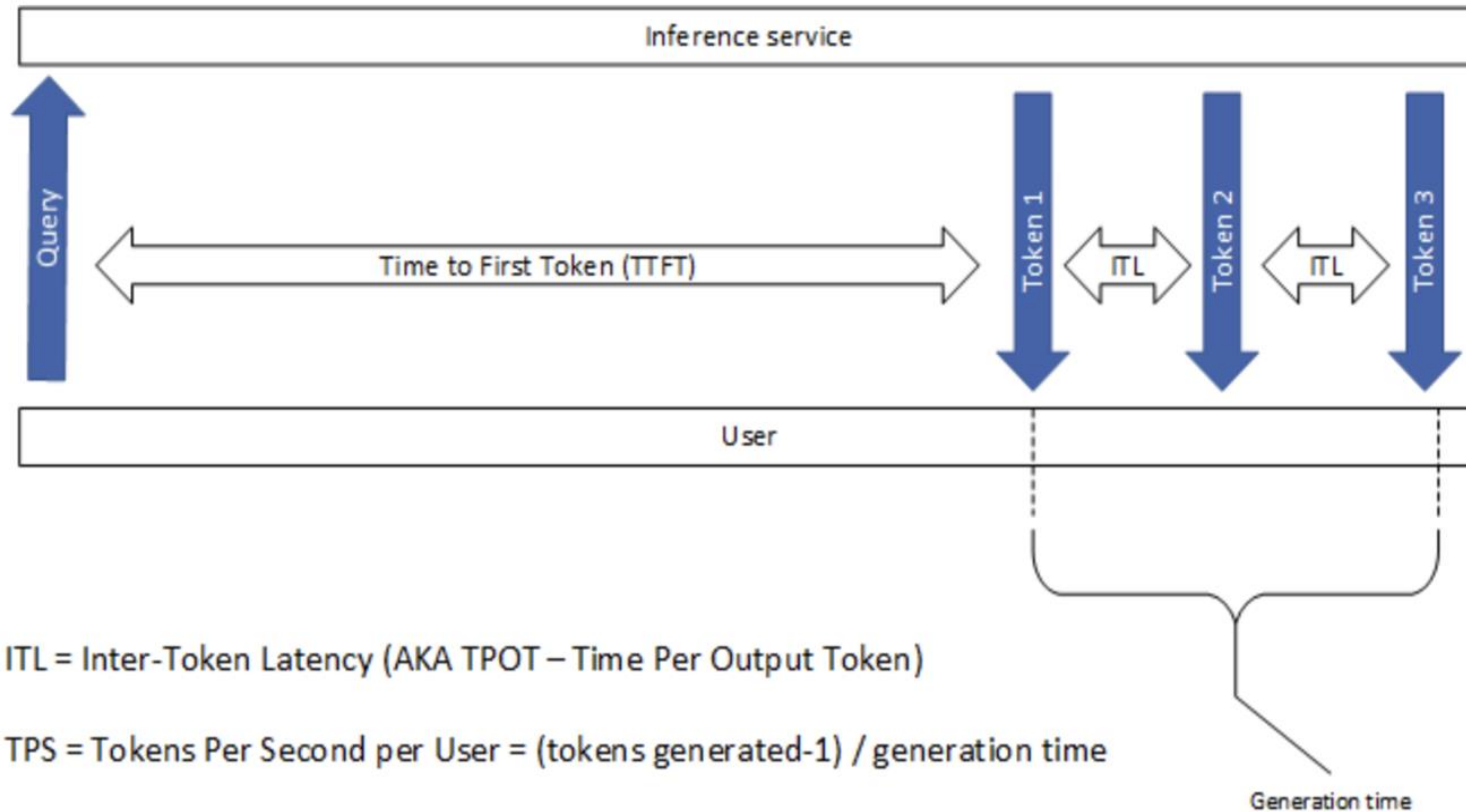
Hundreds of billions of parameters

100x more thinking tokens

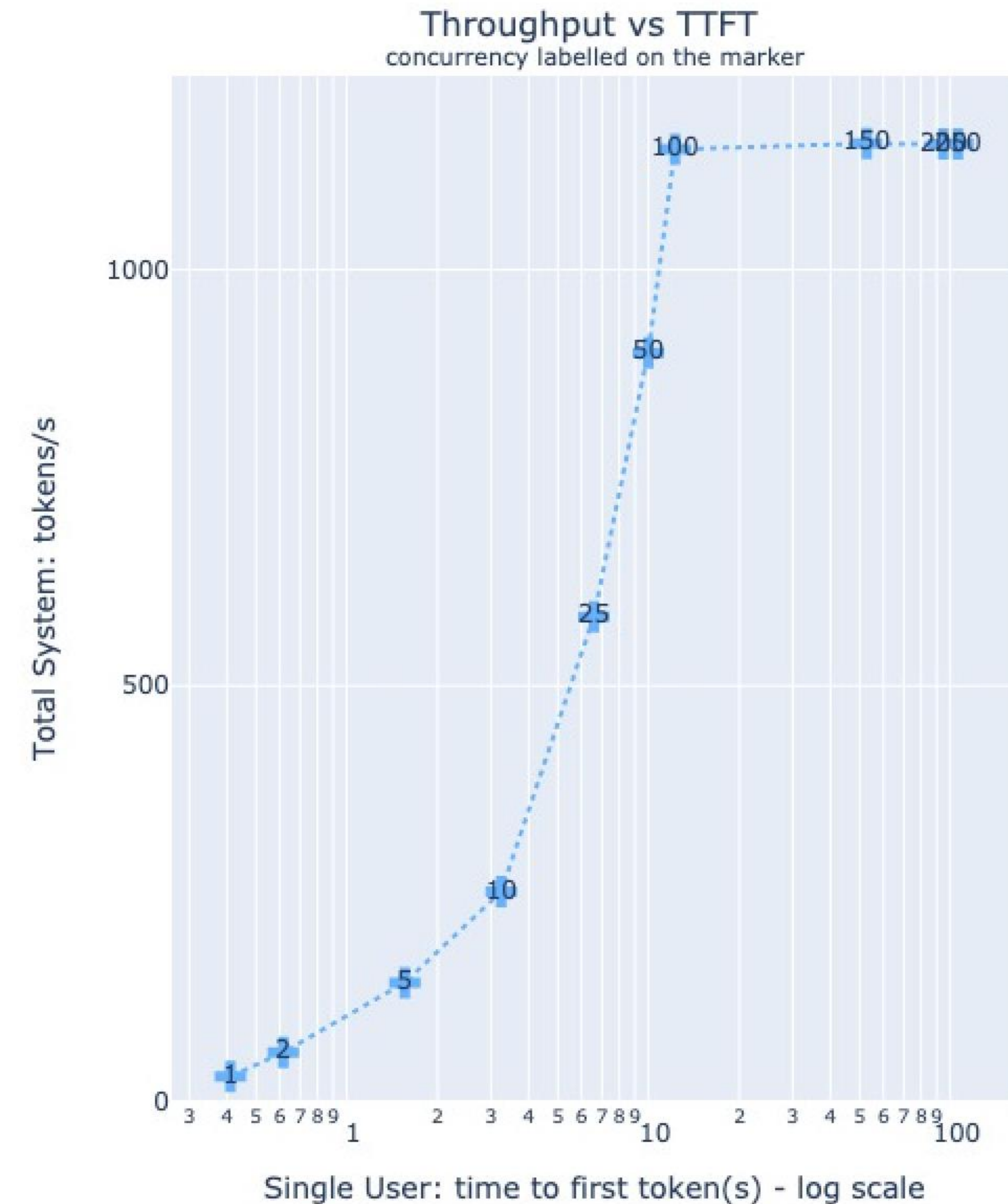Millions of input tokens

NVIDIA.

# Metrics Definition



ITL = Inter-Token Latency (AKA TPOT – Time Per Output Token)

TPS = Tokens Per Second per User = (tokens generated-1) / generation time

Request latency (e2e) = TTFT + ITL*(tokens generated-1)

# A word about Latency and Throughput Trade-off across concurrecies



X-axis: Latency (Time to first token)

y-axis: Throughput (total system tokens/s)

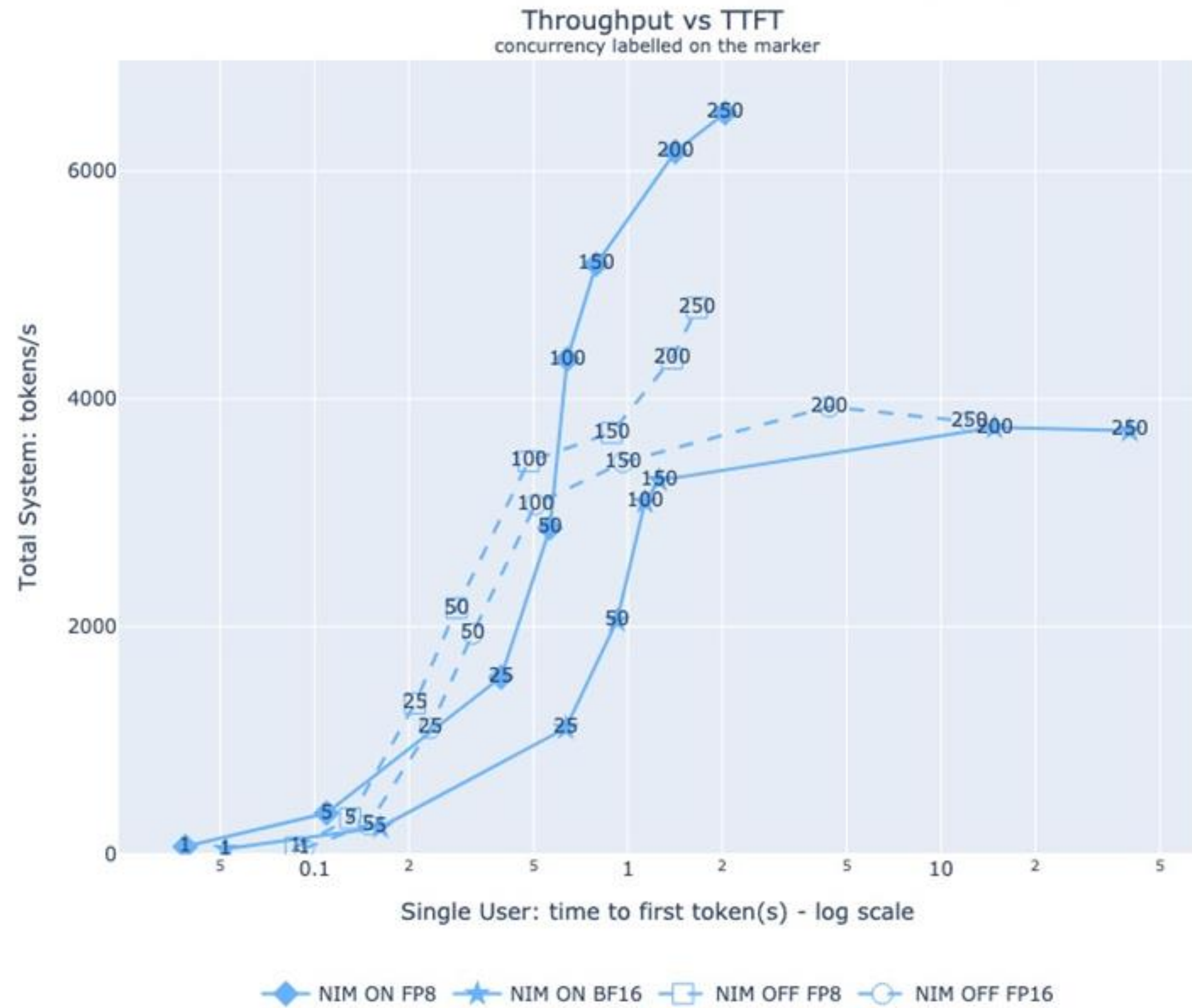Values on the marker: concurrency values

**Uses for the chart**

- Find Maximum Throughput Under a Latency Budget

- Find the Throughput and Latency under a known concurrency value

- Find the throughput saturation threshold beyond which only TTFT increases when concurrencies increase.

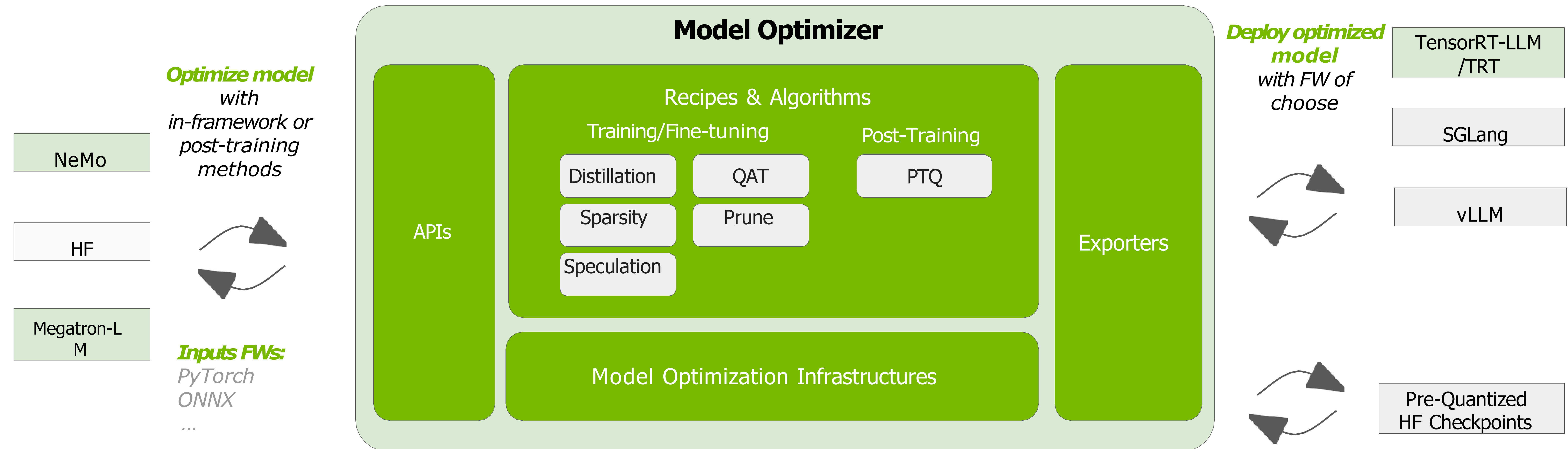- Similar chart for Throughput vs ITL

# 500/2000



Meta-Llama3.1-70B-Instruct Inference 500/2000 on H100 TP4

# NVIDIA TensorRT Model Optimizer
## Product Overview

**Model Optimizer**

**Optimize model** *with in-framework or post-training methods*

NeMo

HF

Megatron-LM

**Inputs FWs:** *PyTorch ONNX ...*

APIs

**Recipes & Algorithms**

Training/Fine-tuning

Post-Training

Distillation

QAT

PTQ

Sparsity

Prune

Speculation

Model Optimization Infrastructures

Exporters

**Deploy optimized model** *with FW of choose*

TensorRT-LLM /TRT

SGLang

vLLM

Pre-Quantized HF Checkpoints

Links:
TensorRT-Model-Optimizer
HF Checkpoints

Nvidia AI SW Stack

# 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
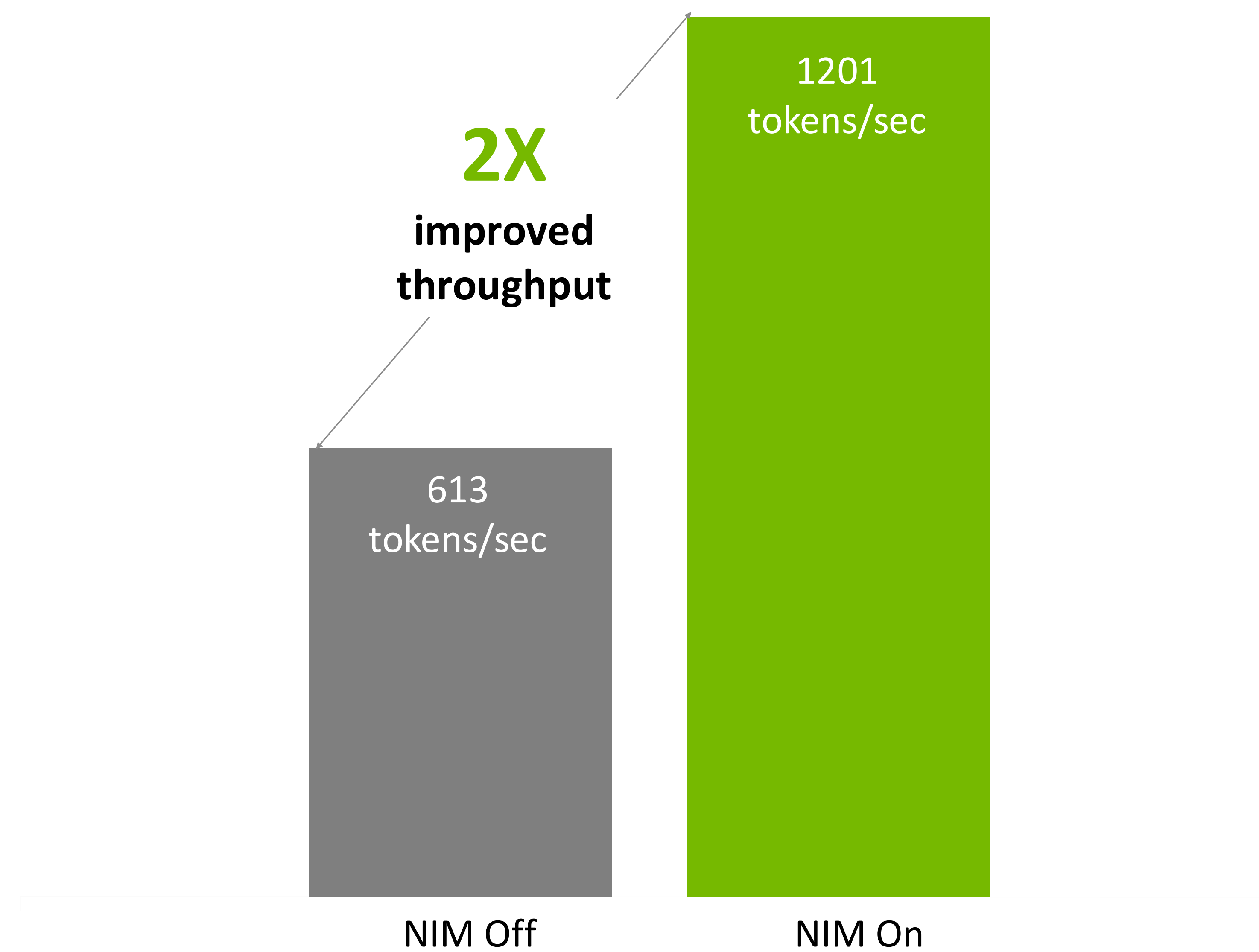
Microsoft Azure · aws · Google Cloud · ORACLE · DGX & DGX Cloud · CISCO · DELL Technologies · Hewlett Packard Enterprise · Lenovo · SUPERMICRO

NVIDIA

# Optimized Efficiency Out of the Box

Improved performance on the same infrastructure with every release

**2X**

**improved
throughput**

1201
tokens/sec

613
tokens/sec

NIM Off                    NIM On

**ⓩ NVIDIA.**

# NIM Accelerates the AI Factory

Maximize token generation, or minimize total cost of ownership

**2X**

**improved throughput**

$100M

NVIDIA AI Enterprise adds only 7% cost to the AI factory, while drastically increasing throughput

NVIDIA AI Enterprise

$107M

1.92M tokens/sec

980K tokens/sec

Infrastructure & Energy

Infrastructure & Energy

**Also includes:**
Support SLA
Security patching
Lifecycle management

NIM Off

NIM On

Infrastructure includes hardware systems with GPU, CPU, Networking, Storage, as well as infrastructure operations. Energy includes 4-year operating costs.
Tokens/second based on Llama 3.1-8B-Instruct running on 1x H100 SXM compared to leading open-source inference solutions

NVIDIA.

# NIM Deployment

Easily manage cost-efficient autoscaling across NVIDIA accelerated cloud infrastructure

NIM on baremetal, VMI, AKS, Azure ML and Azure AI Studio

NIM Deploy - Azure

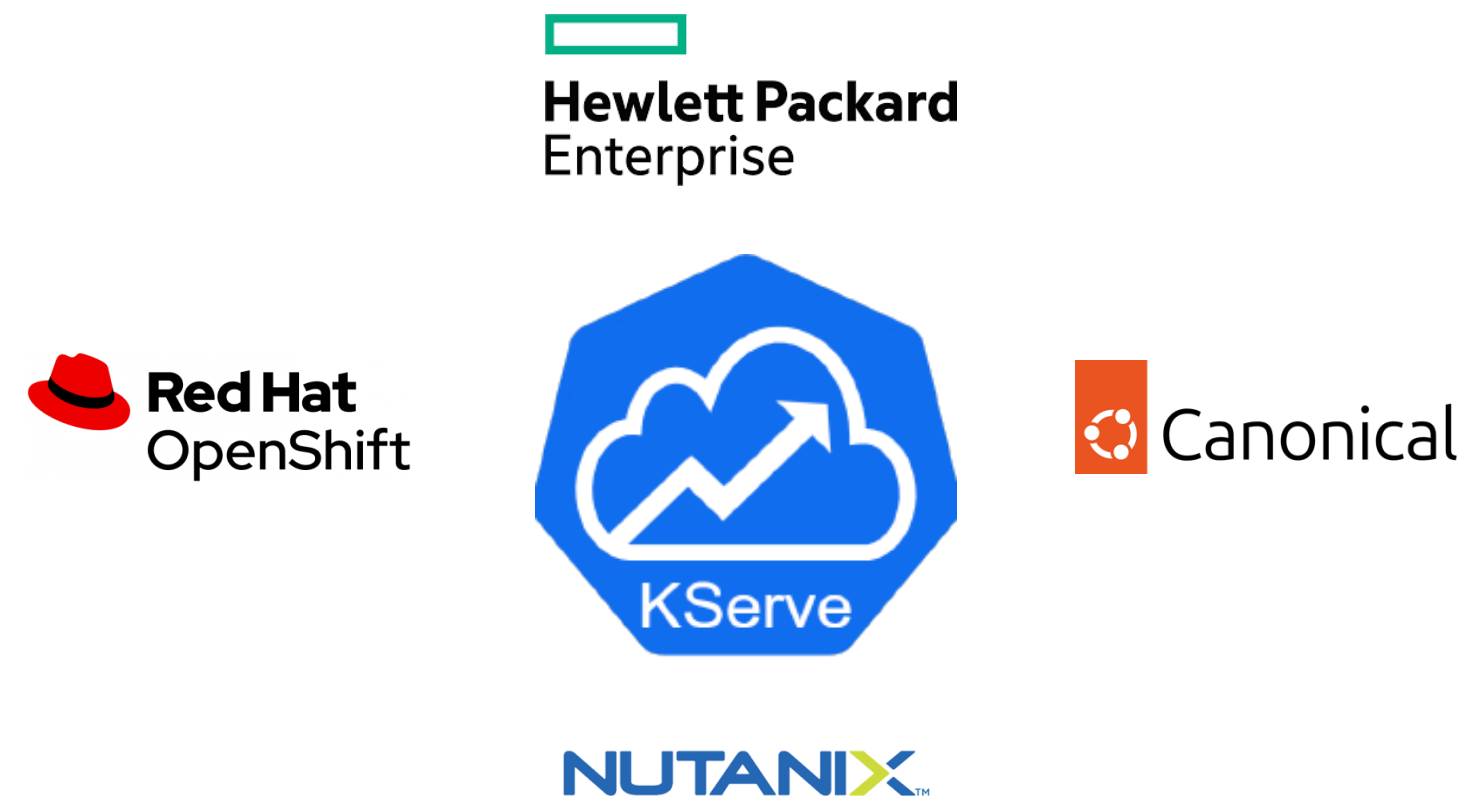NIM on baremetal, AMI, EKS and Amazon SageMaker

NIM Deploy - AWS

NIM on baremetal, VMI, GKE and Google Cloud Vertex AI

NIM Deploy – Google Cloud

NIM on baremetal, CMI, OKE, and OCI Data Science Service

How-to Blog

NIM on KServe Inference Platform on Kubernetes

NIM Deploy - KServe

Simplified NIM Deployment on Kubernetes

NIM Deploy - Helm

NIM Deployment Lifecycle Management on Kubernetes

NIM Operator

# NVIDIA Agent Intelligence (AI-Q) Toolkit

An open-source library for building enterprise-ready agentic systems
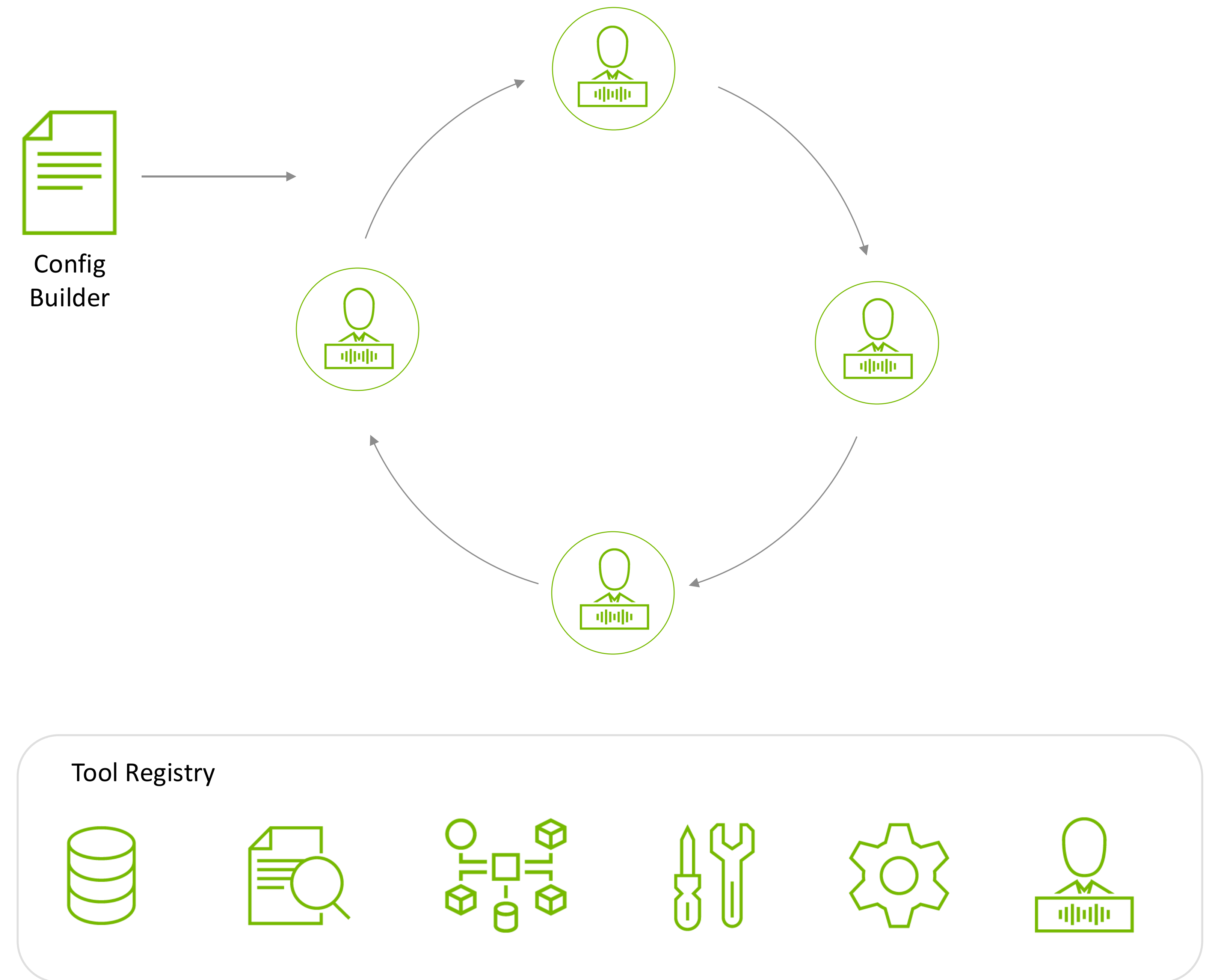
## Agent Interconnect

- Universal descriptors for agents, tools, and workflows across frameworks
- Reusable Agent/Tool registry
- Workflow Configuration/Builder

## Profiling & Optimizations

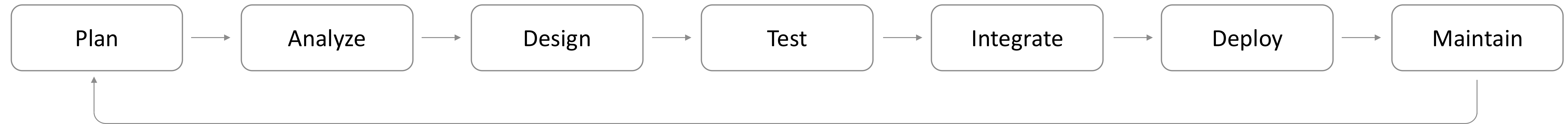- Fine-grained AI workflow telemetry collected can be used to implement agentic system accelerations.

## Evaluation & Observability

- Evaluate system level accuracy
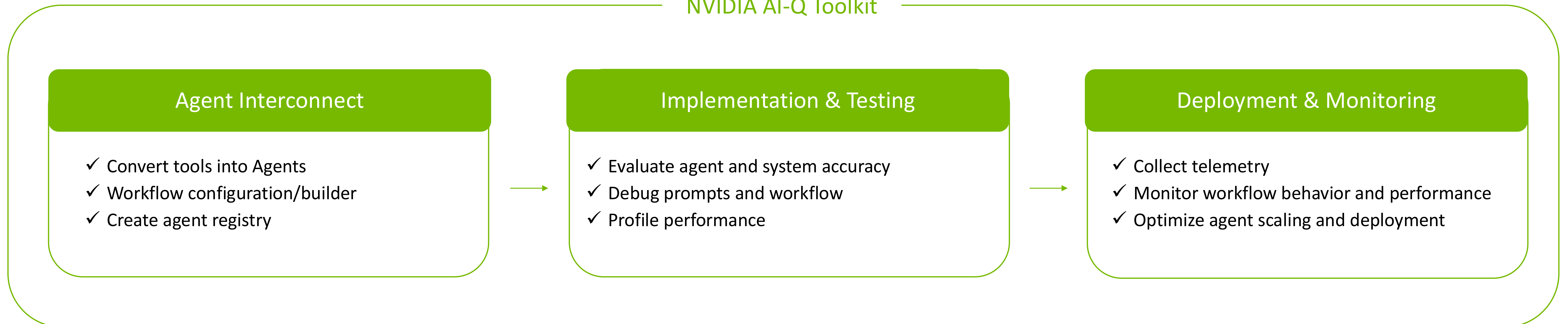- Understand and debug inputs and outputs for each component in the AI workflow

Config Builder

Tool Registry

# Enabling Software Development Lifecycle for AI Agents

## NVIDIA AI-Q Toolkit

Plan → Analyze → Design → Test → Integrate → Deploy → Maintain

### NVIDIA AI-Q Toolkit

**Agent Interconnect**
- ✓ Convert tools into Agents
- ✓ Workflow configuration/builder
- ✓ Create agent registry

**Implementation & Testing**
- ✓ Evaluate agent and system accuracy
- ✓ Debug prompts and workflow
- ✓ Profile performance

**Deployment & Monitoring**
- ✓ Collect telemetry
- ✓ Monitor workflow behavior and performance
- ✓ Optimize agent scaling and deployment

# NVIDIA AI-Q Toolkit

## Accelerate AI Agents and Streamline Agentic Workflow Optimization

### SAVE TIME

**Simplify** the development of agentic systems

➢ Flexibly choose, and connect, agent frameworks best suited for each task

➢ Easily reuse existing and new RAG pipelines, different Agentic workflows, and tools across your Enterprise

➢ Quickly elevate existing Gen AI workflows to Agentic AI workflows

### REDUCE COSTS

**Accelerate** agent responses—do more with what you have

➢ System level optimizations provide accelerated Agentic AI performance

➢ NVIDIA AgentIQ collects telemetry that provides opportunities for optimization, driving efficiency for an agentic workflow.
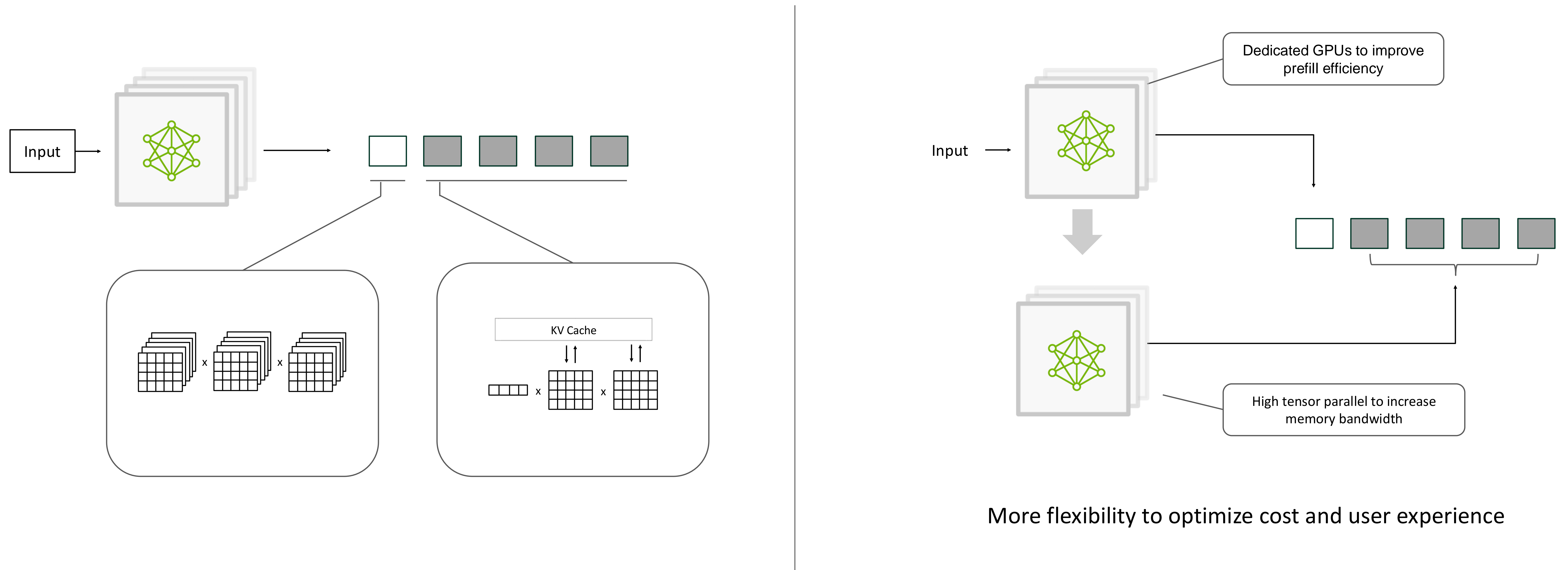
### IMPROVE BUSINESS OUTCOMES

**Increase** agentic system accuracy

➢ Evaluate agentic system response accuracy

➢ Understand and debug inputs and outputs for each component in the system

➢ Traceability and auditing of agent communications

# New Inference Optimization Techniques to Boost Inference

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



Input

KV Cache

Input

Dedicated GPUs to improve prefill efficiency

High tensor parallel to increase memory bandwidth

More flexibility to optimize cost and user experience

NVIDIA.

# Announcing NVIDIA Dynamo

## AI Inference Software for Reasoning Inference at Scale
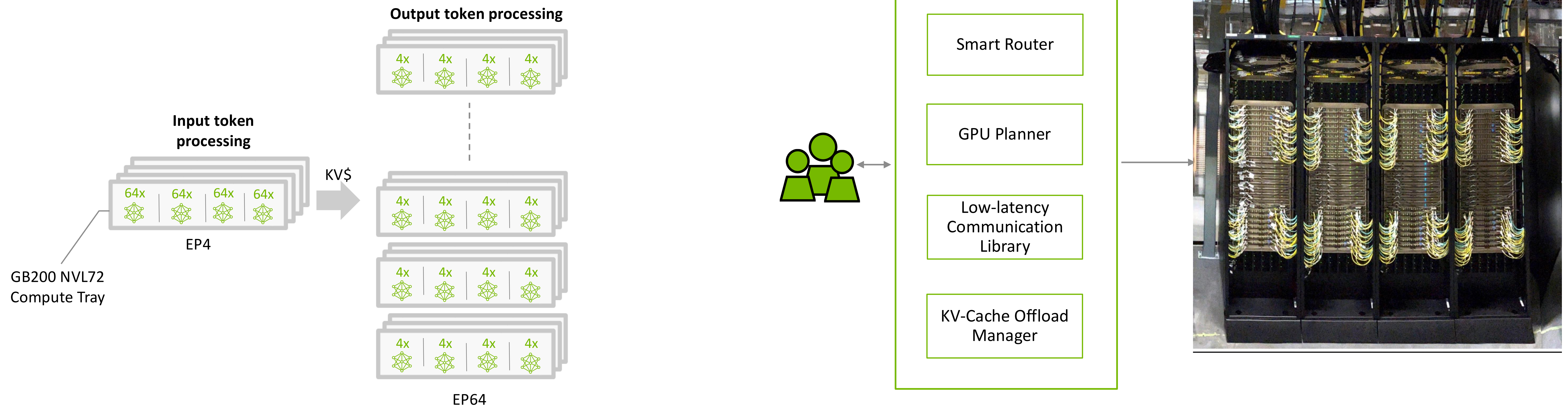
**30X**

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
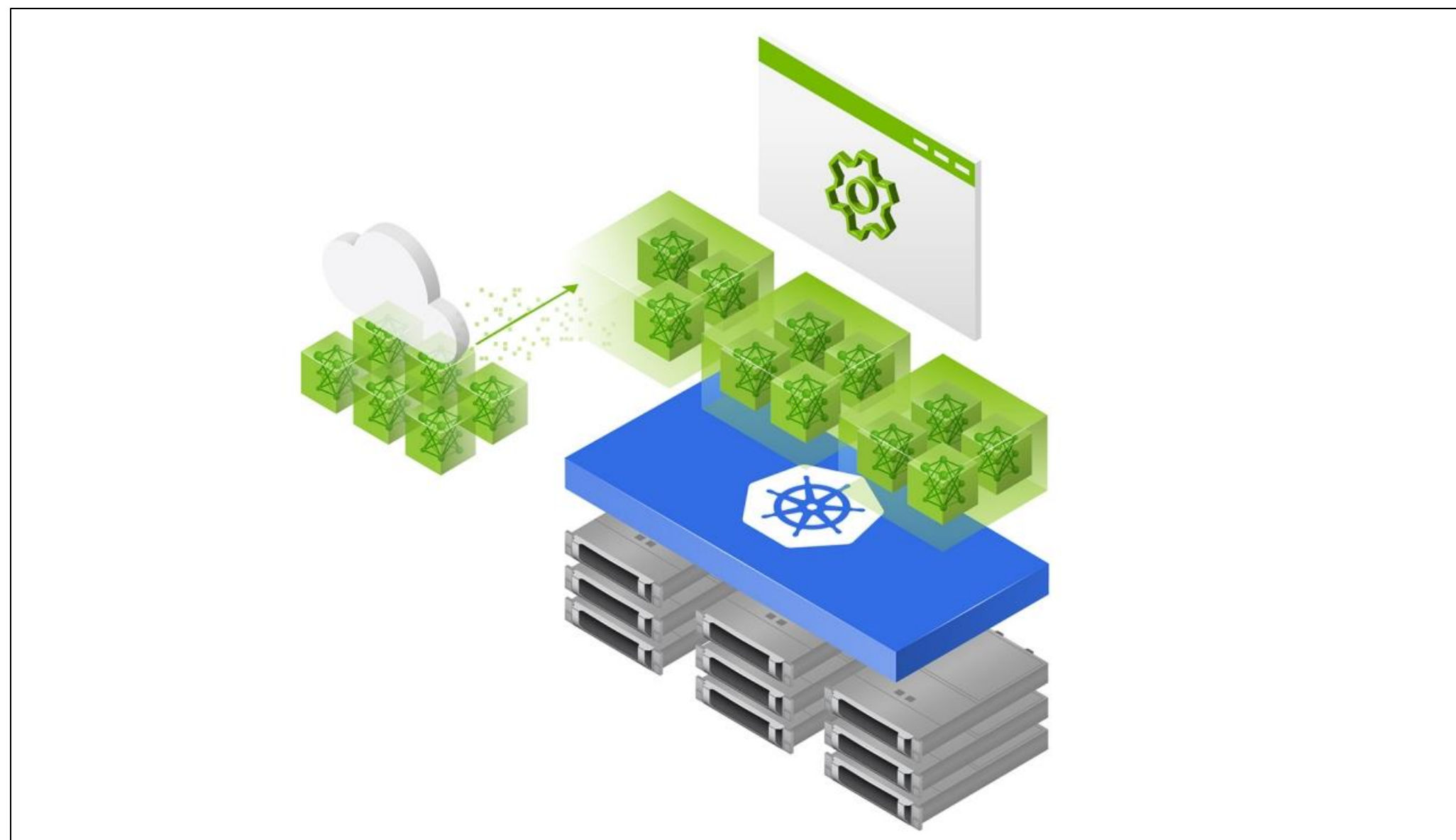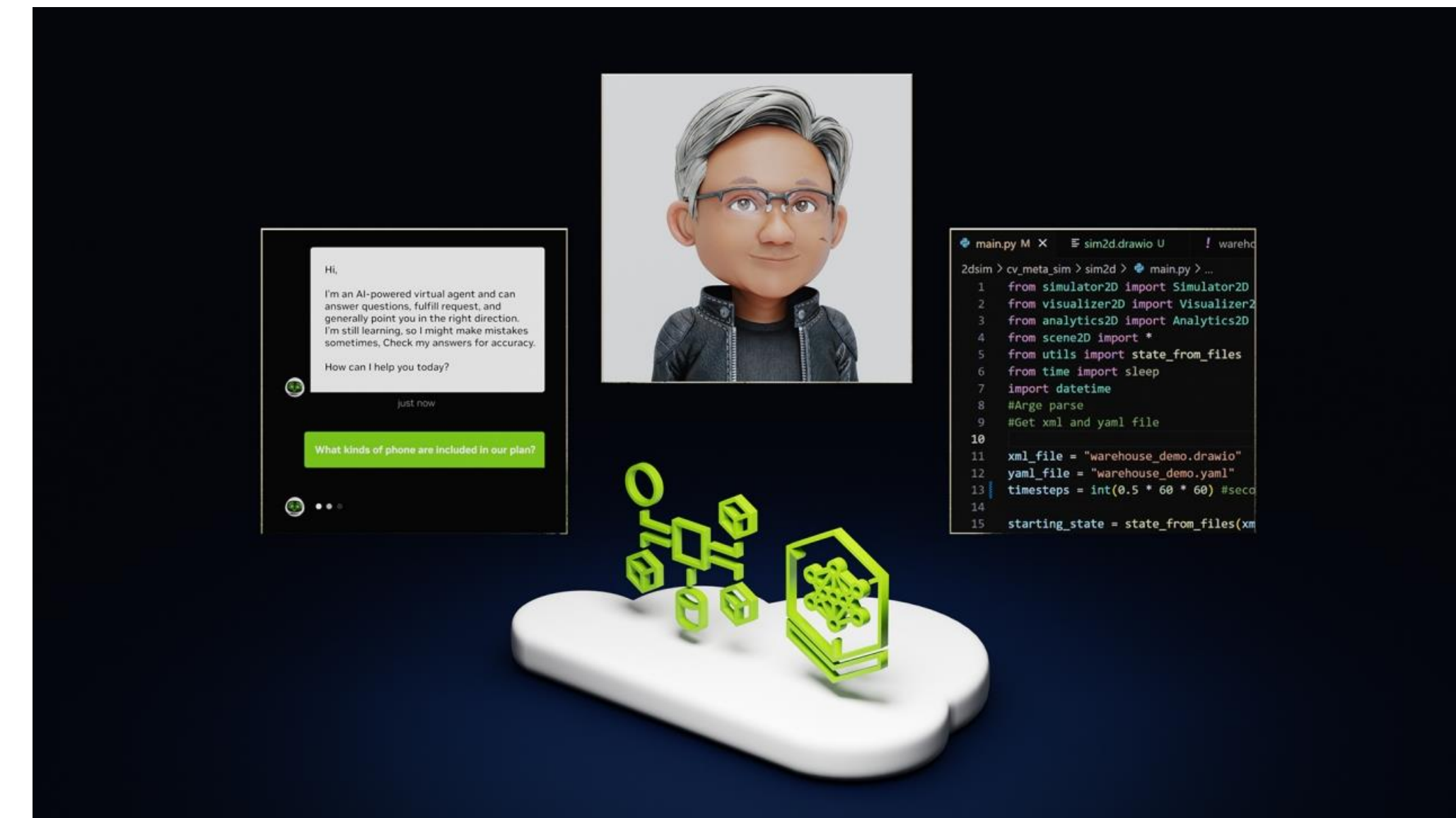
### DeepSeek R1

**Output token processing**

| 4x | 4x | 4x | 4x |

**Input token processing**

| 64x | 64x | 64x | 64x |

KV$

EP4

| 4x | 4x | 4x | 4x |

| 4x | 4x | 4x | 4x |

| 4x | 4x | 4x | 4x |

EP64

GB200 NVL72
Compute Tray

## NVIDIA Dynamo

Smart Router

GPU Planner

Low-latency
Communication
Library

KV-Cache Offload
Manager

# NVIDIA Dynamo Use Cases

Unlock the full protentional of reasoning models and AI agents
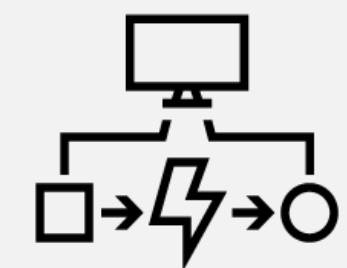


**Serving Reasoning Models**



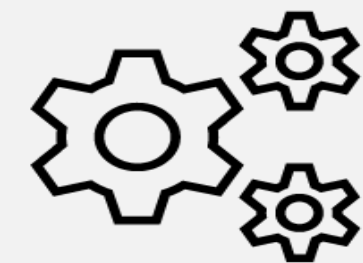**Code Generation**

# NVIDIA Dynamo Breakthrough Features

A modular generative AI inference server designed for distributed and disaggregated serving

**NVIDIA Dynamo**

**Distributed Inference Serving**

Seamlessly scale LLMs from a single GPU to thousands of GPUs

**GPU Planning & Scheduling**

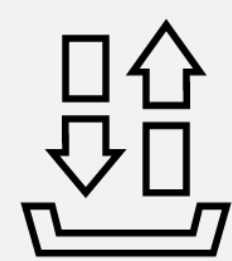Meet changing demand patterns w/o over or under provisioning of resources

**Smart Request Router**

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

**Low-latency Inference Data Transfer Library**

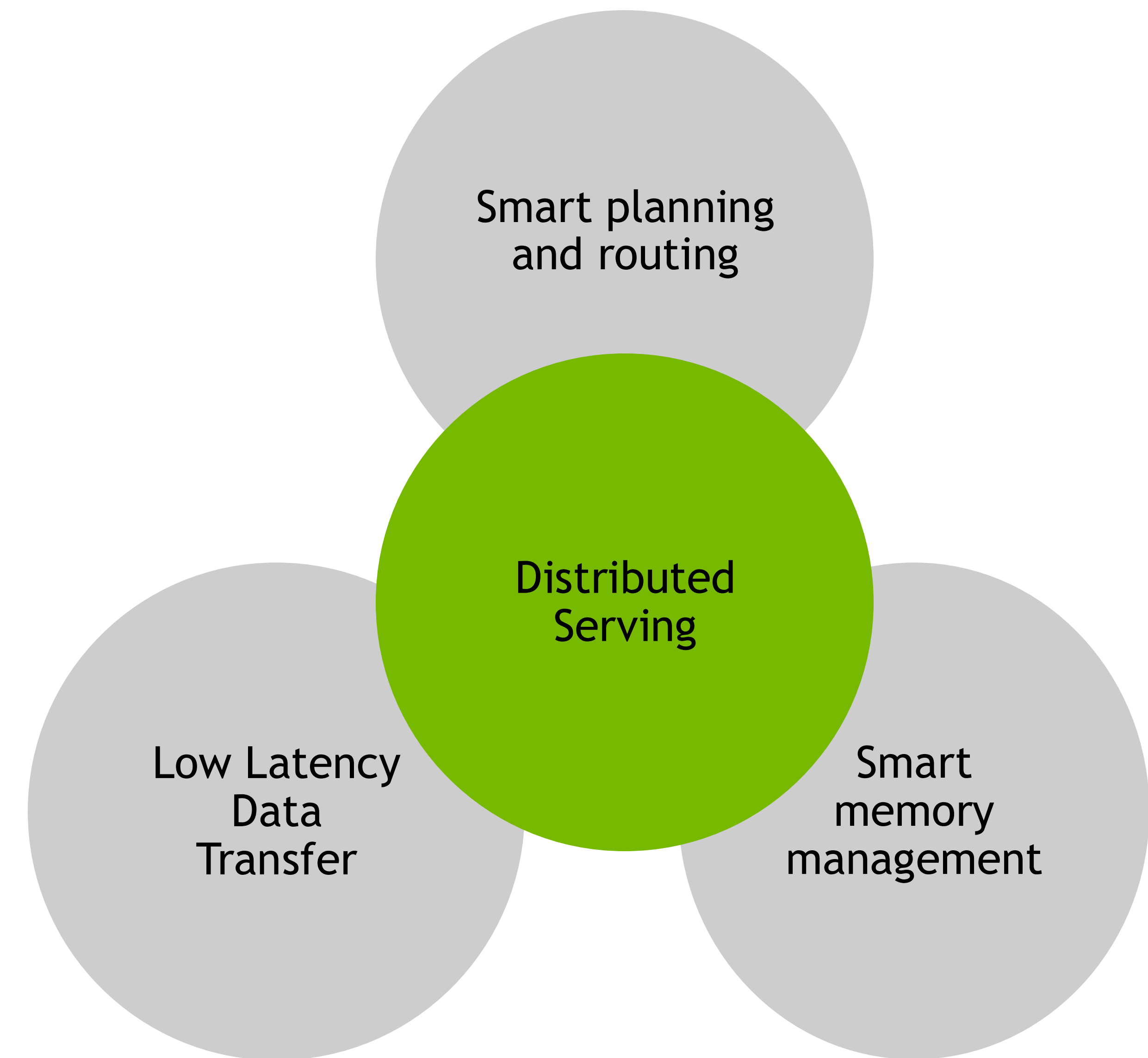Accelerate GPU-to-GPU communication to enhance user experience

**KV Cache Manager**

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

# NIM with NVIDIA Dynamo

Easy button for high performance inference at data center scale

- **Turbocharge agentic AI outcomes** with 30X AI reasoning throughput on NVIDIA accelerated datacenter GPUs

- **Unlock game-changing use cases** with distributed scale-out of complex requests across 1000+ GPUs

- **Maximize token revenue generation** for AI factories with distributed serving without bottlenecks

- **Push button deploy** in 3 simple steps in 5 minutes or less without coding

Smart planning and routing

Distributed Serving

Low Latency Data Transfer

Smart memory management

NVIDIA.

# Accelerating AI Ecosystem

Fully open source and supports all major AI frameworks

# Learn with NVIDIA (DLI)

# Register free for NVIDIA GTC