

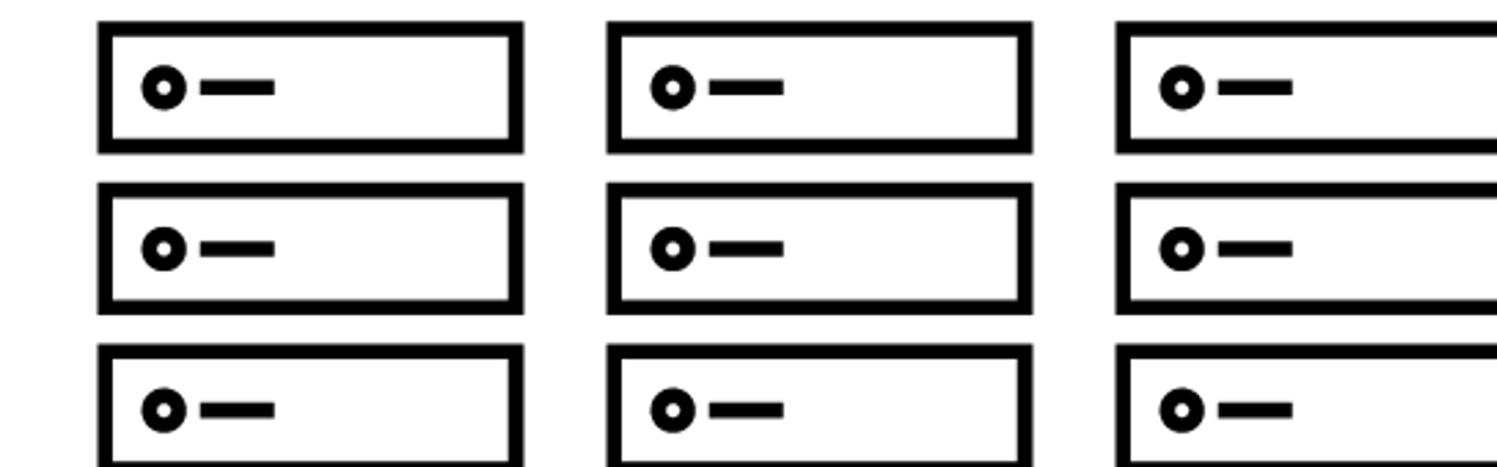
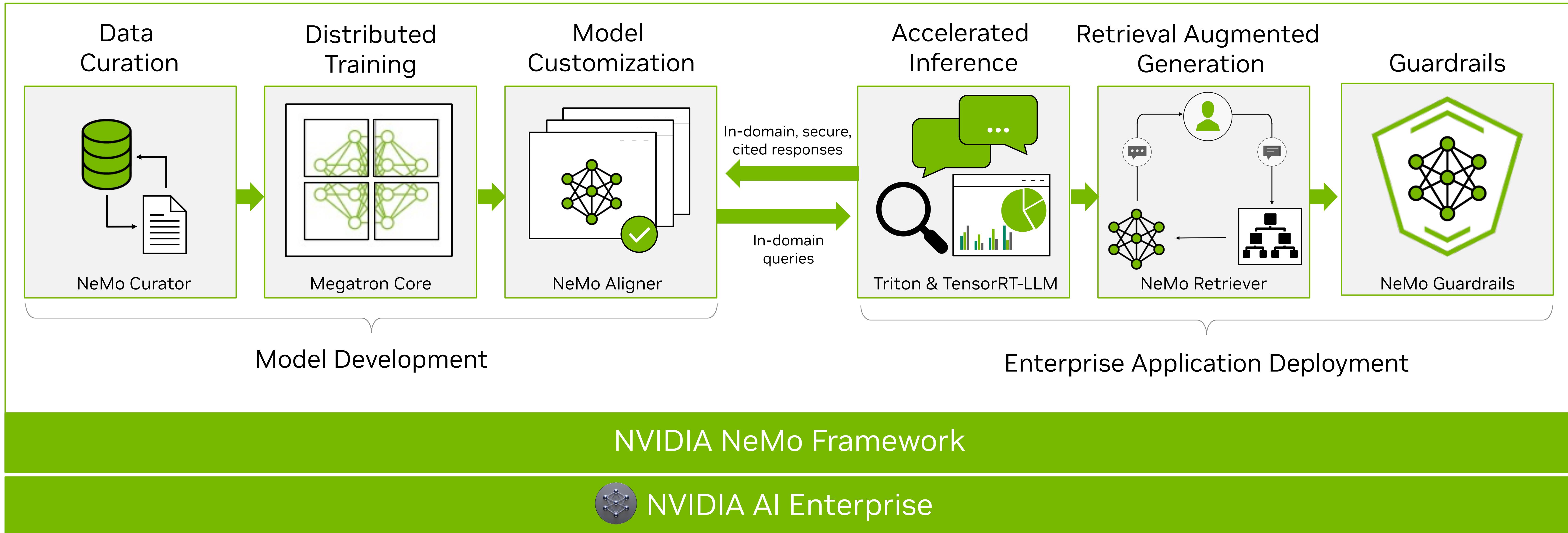


NVIDIA NeMo Framework

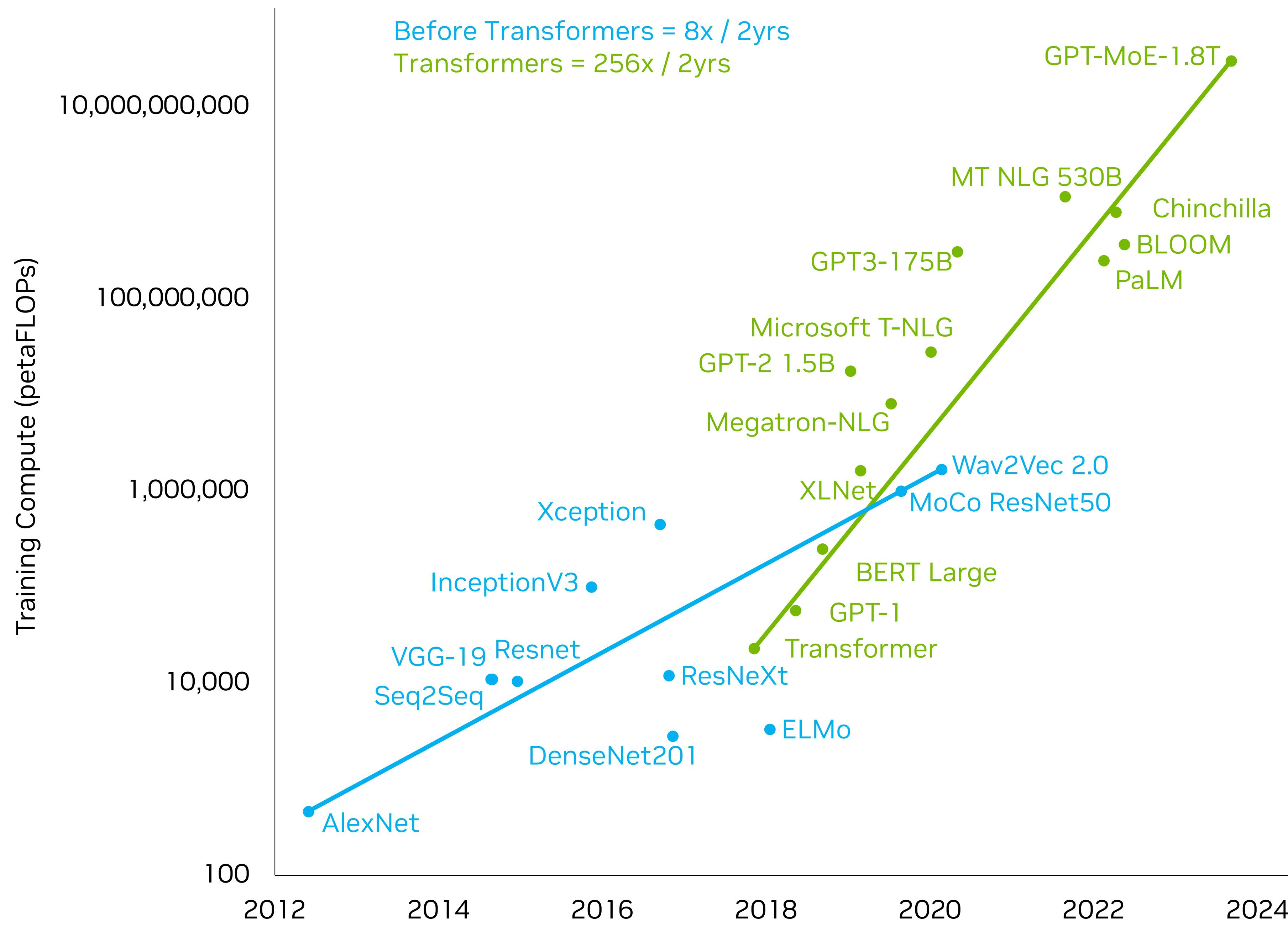
Accelerate Enterprise AI Development

Building Generative AI Applications for the Enterprise

Build, customize and deploy generative AI models with NVIDIA NeMo

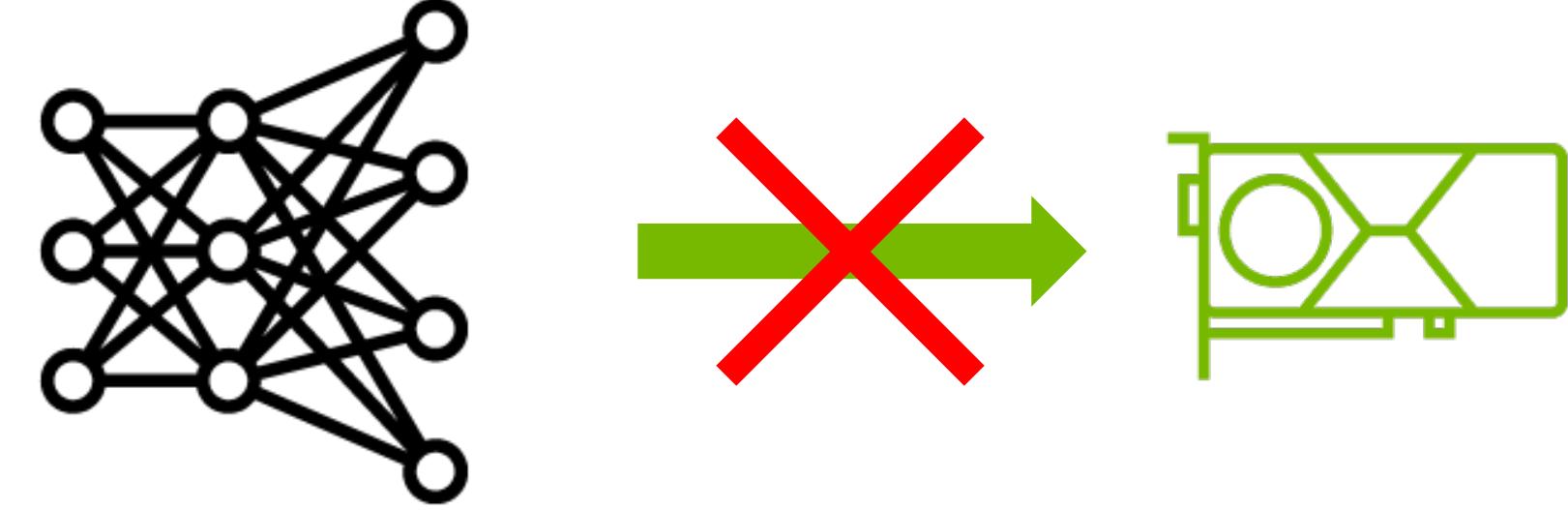


Explosive Growth in AI Computational Requirements



	Time (GPU hours)	Power Consumption (W)
Llama 3 8B	1.3M	700
Llama 3 70B	6.4M	700
Total	7.7M	

LLM

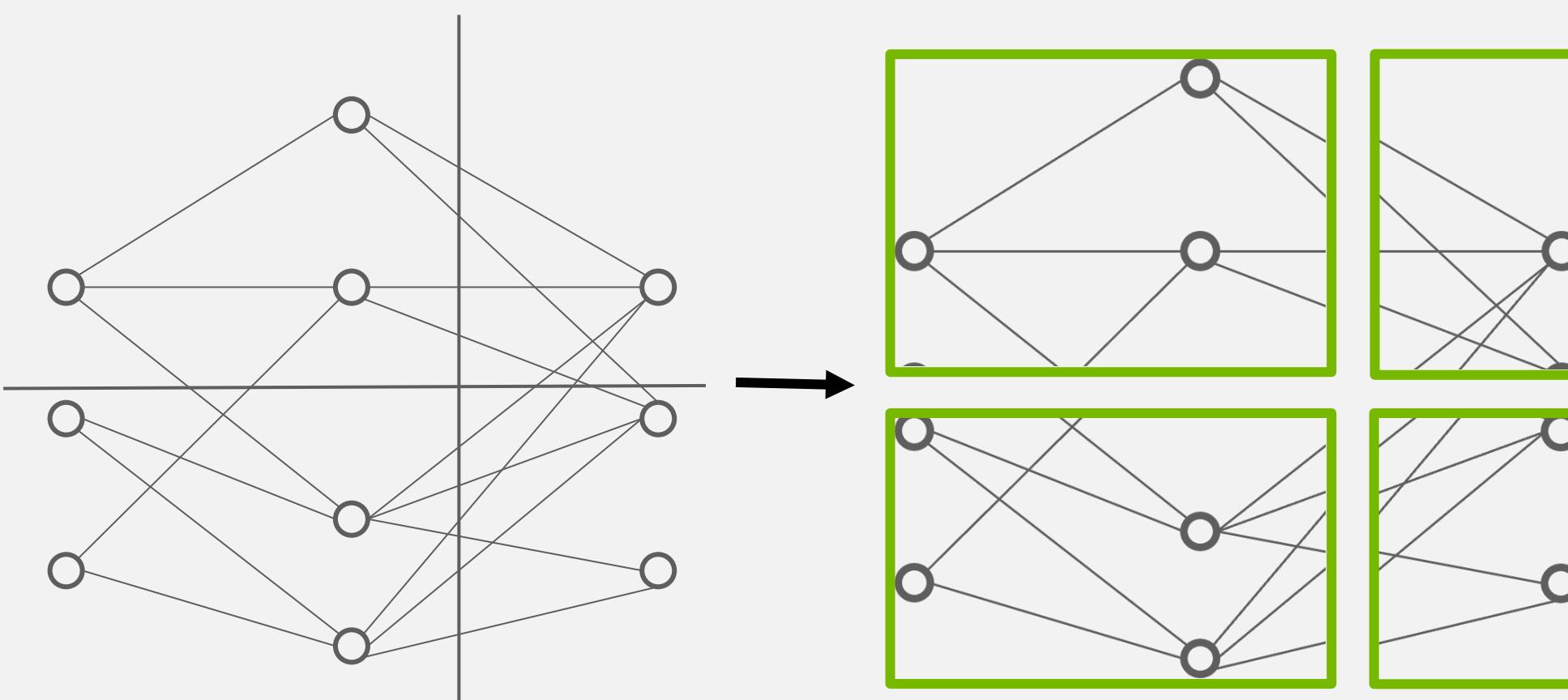


Can't fit into single GPU

Building Generative AI Foundation Models

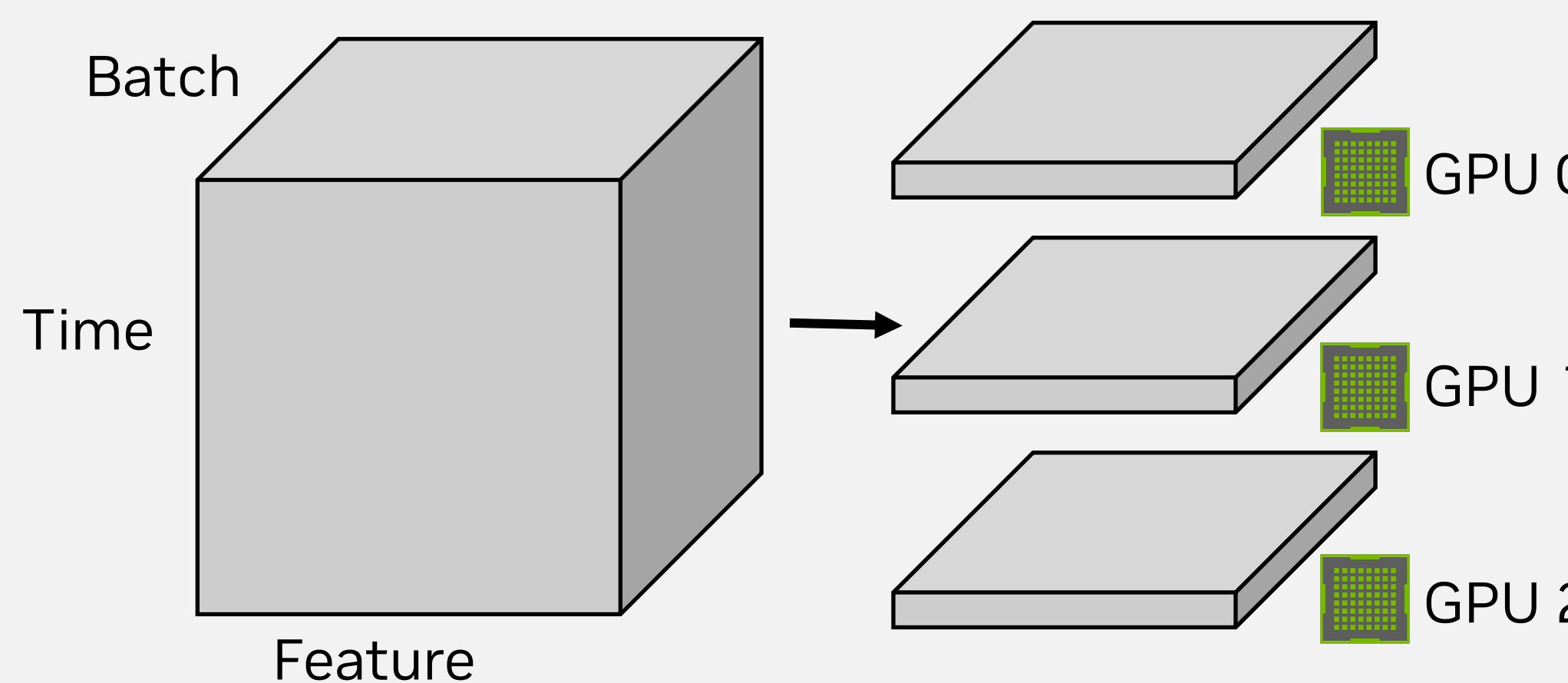
Efficiently and quickly training models using NVIDIA NeMo

Tensor & Pipeline Parallelism



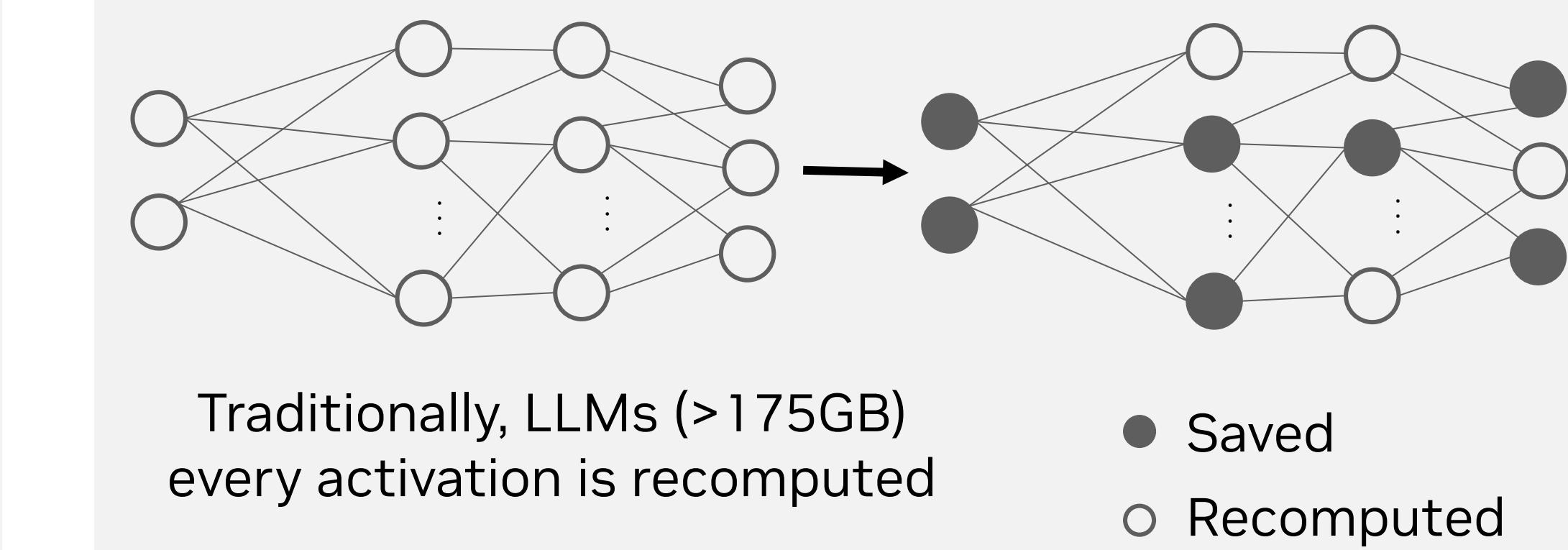
Reduced memory footprint and allows for large-scale training of LLMs across accelerated infrastructure

Sequence Parallelism



Working with tensor processing to increase the batch size that can be support for training

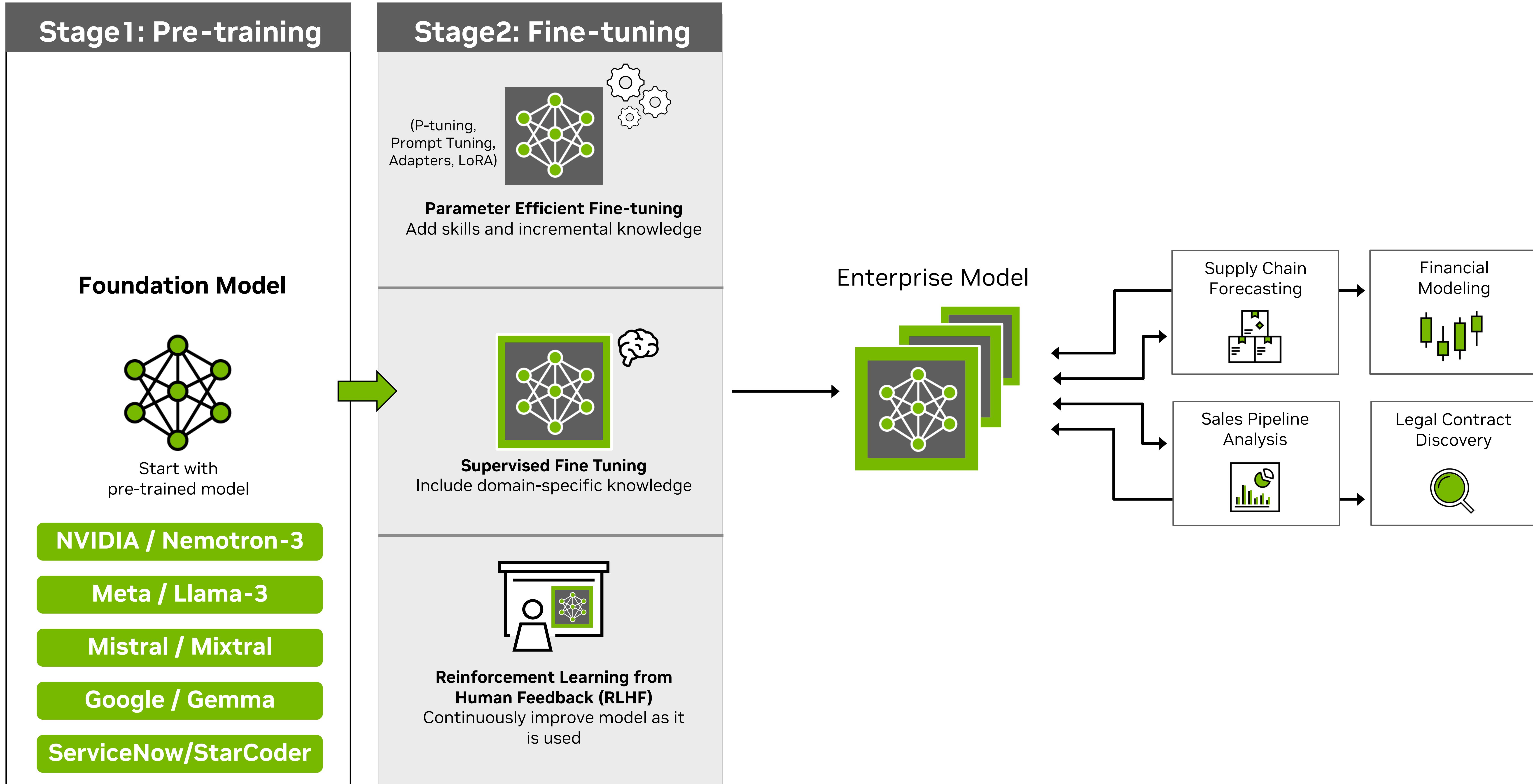
Selective Activation Recomputation



Smart activation checkpointing provides greatest trade-off between memory and recomputation

Model Customization for Enterprise Ready LLMs

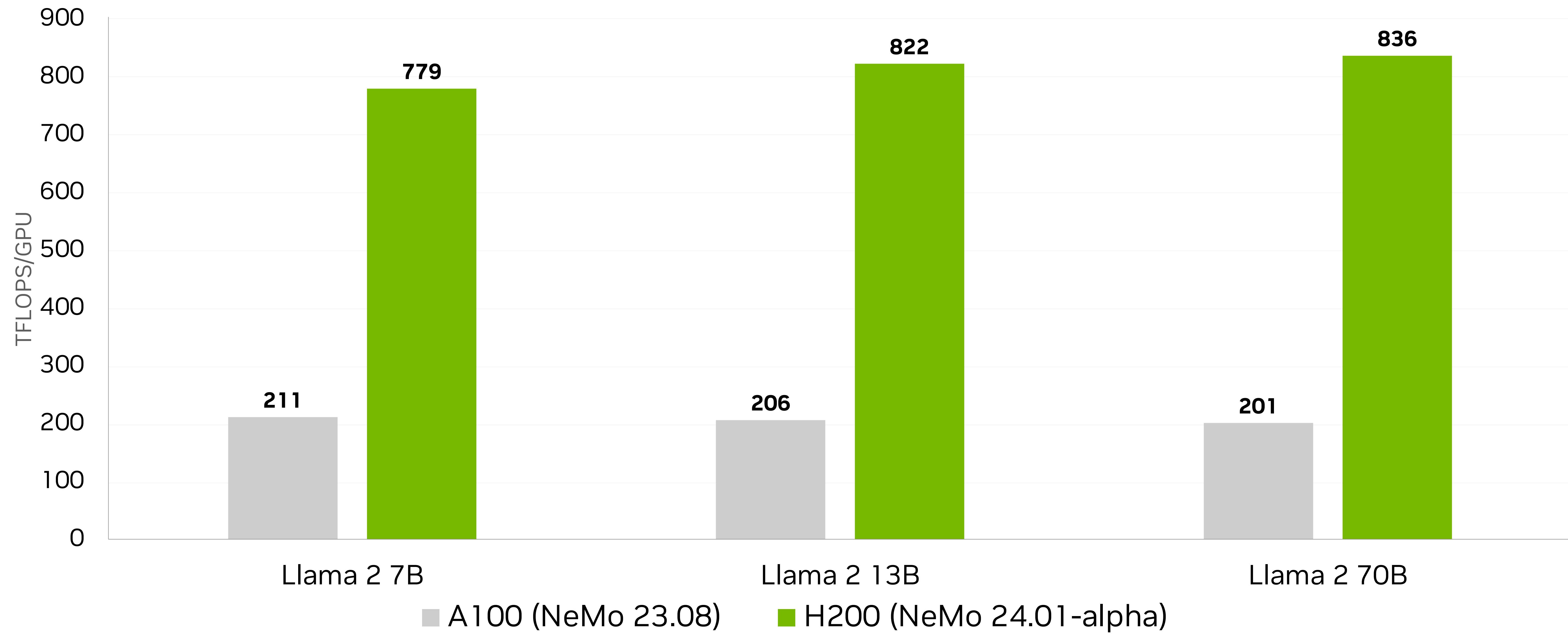
Customization techniques to overcome the challenges of using foundation models



NVIDIA H200 with NeMo Delivers Stunning Performance

H200 and NeMo Achieve Up to 4.2X Faster Llama 2 Pre-Training and SFT

Higher is Better



Measured performance per GPU. Global Batch Size = 128.

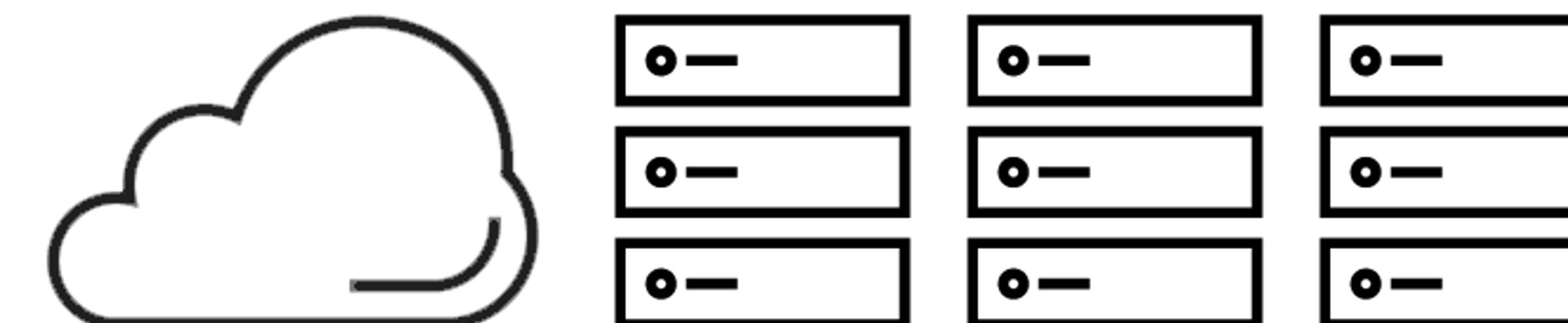
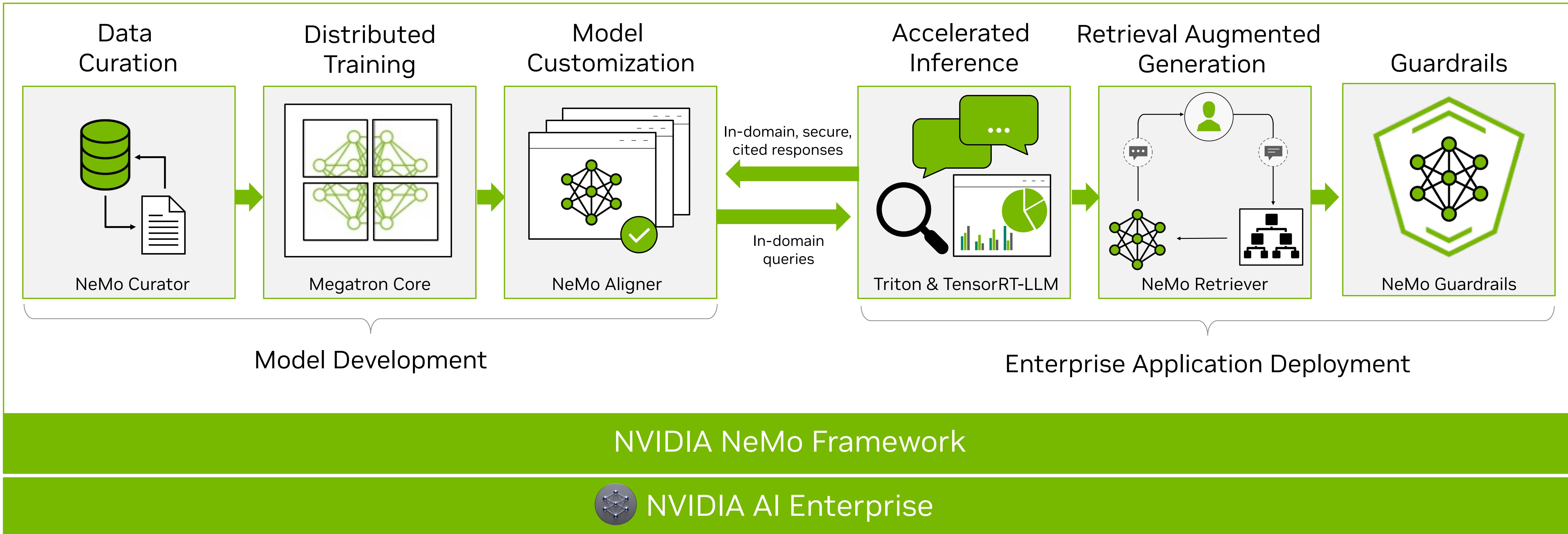
Llama 2 7B: Sequence Length 4096 | A100 8x GPU, NeMo 23.08 | H200 8x GPU, NeMo 24.01-alpha

Llama 2 13B: Sequence Length 4096 | A100 8x GPU, NeMo 23.08 | H200 8x GPU, NeMo 24.01-alpha

Llama 2 70B: Sequence Length 4096 | A100 32x GPU, NeMo 23.08 | H200 8x GPU, NeMo 24.01-alpha

Building Generative AI Applications for the Enterprise

Build, customize and deploy generative AI models with NVIDIA NeMo



Triton & TensorRT-LLM

NVIDIA AI Inference Platform

- **Triton Inference Server**

- Triton Inference Server deploy AI model from multiple deep learning and machine learning frameworks. It supports inference across cloud, data center, edge and embedded devices on NVIDIA GPUs, x86 and ARM CPU, or AWS Inferentia. Triton Inference Server delivers optimized performance for many query types, including real time, batched, ensembles and audio/video streaming.

- **TensorRT-LLM**

- For deployments with high amounts of customization or control required
 - New or customized model architectures
 - Fine tuned control over optimizations, quantization, & performance
- Model developers & optimizers desiring deep control of optimization
- Requires higher layers for serving (Triton) and functionality (Guardrails, RAG, etc.) for final deployments

Triton Inference Server

High Performance Inference Serving for AI model production deployments

TensorRT-LLM backend

Speed-of-Light LLM optimization

TensorRT

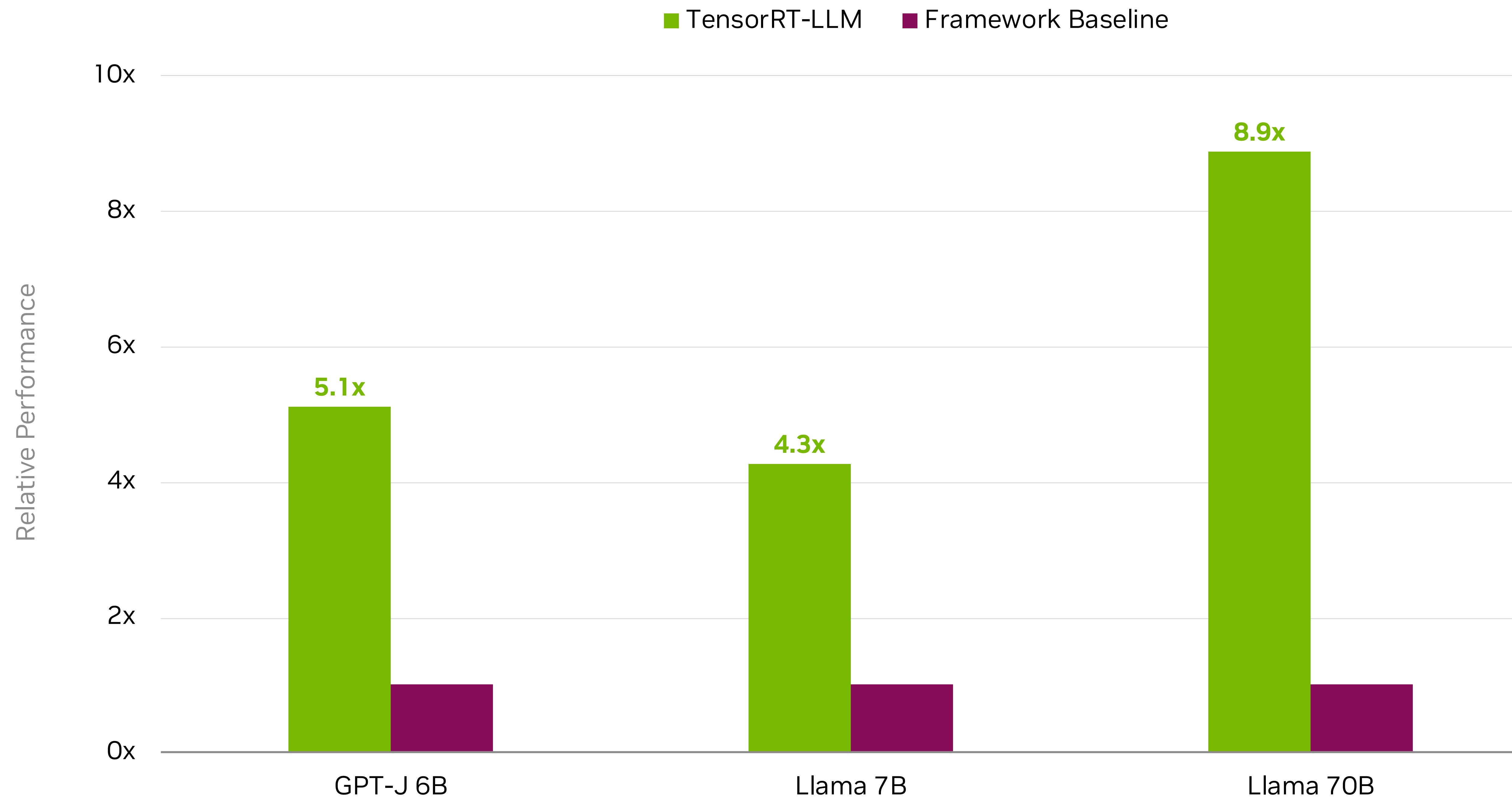
NV Internal kernel libs

NVIDIA LLM Inference Stack

Recommended that most developers start with NeMo Framework

TensorRT-LLM Performance Improvement

Up to 9x better throughput than implementations in DL frameworks

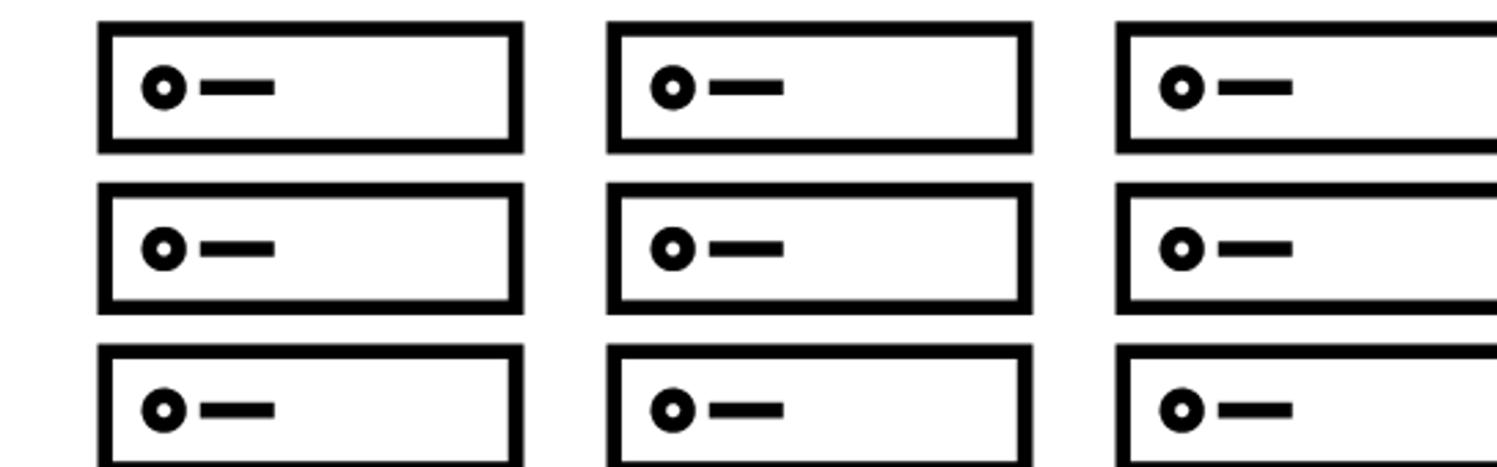
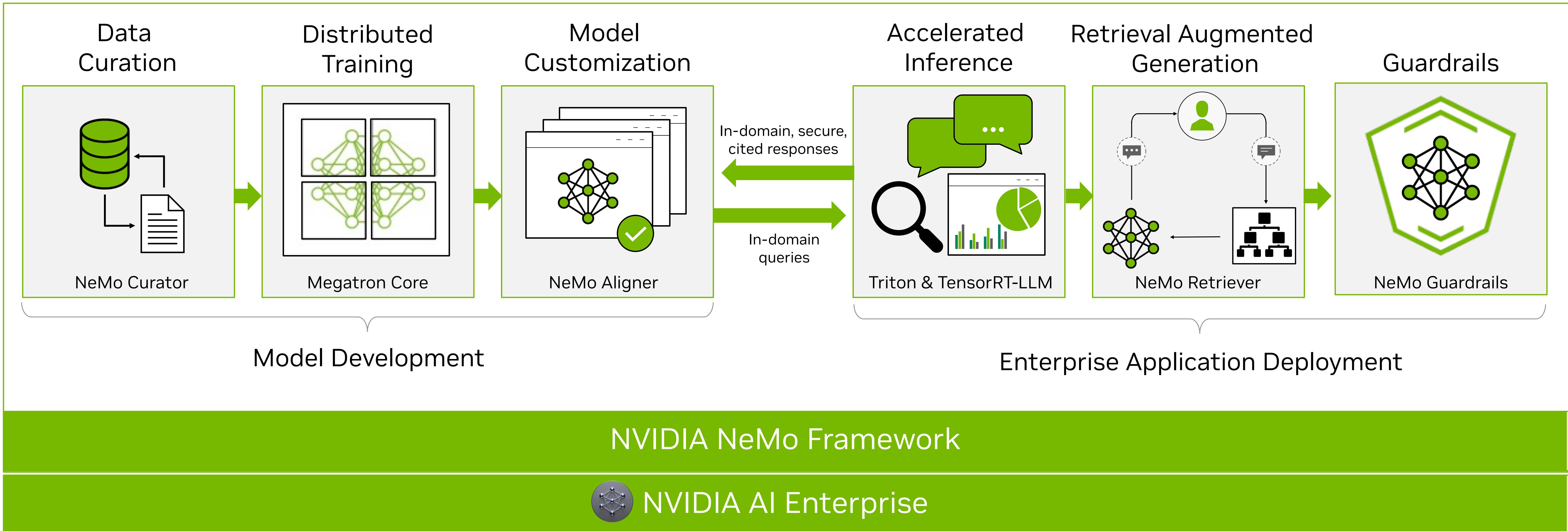


TensorRT-LLM v0.5.0 internal build. HF Accelerate. Tokens/s/GPU relative improvement
DGX H100. TensorRT-LLM FP8, HF Accelerate FP16.
Max batchsize up to 64. Input & output sequence length 128:128
TensorRT-LLM Llama 70B TP2, HF Accelerate PP4



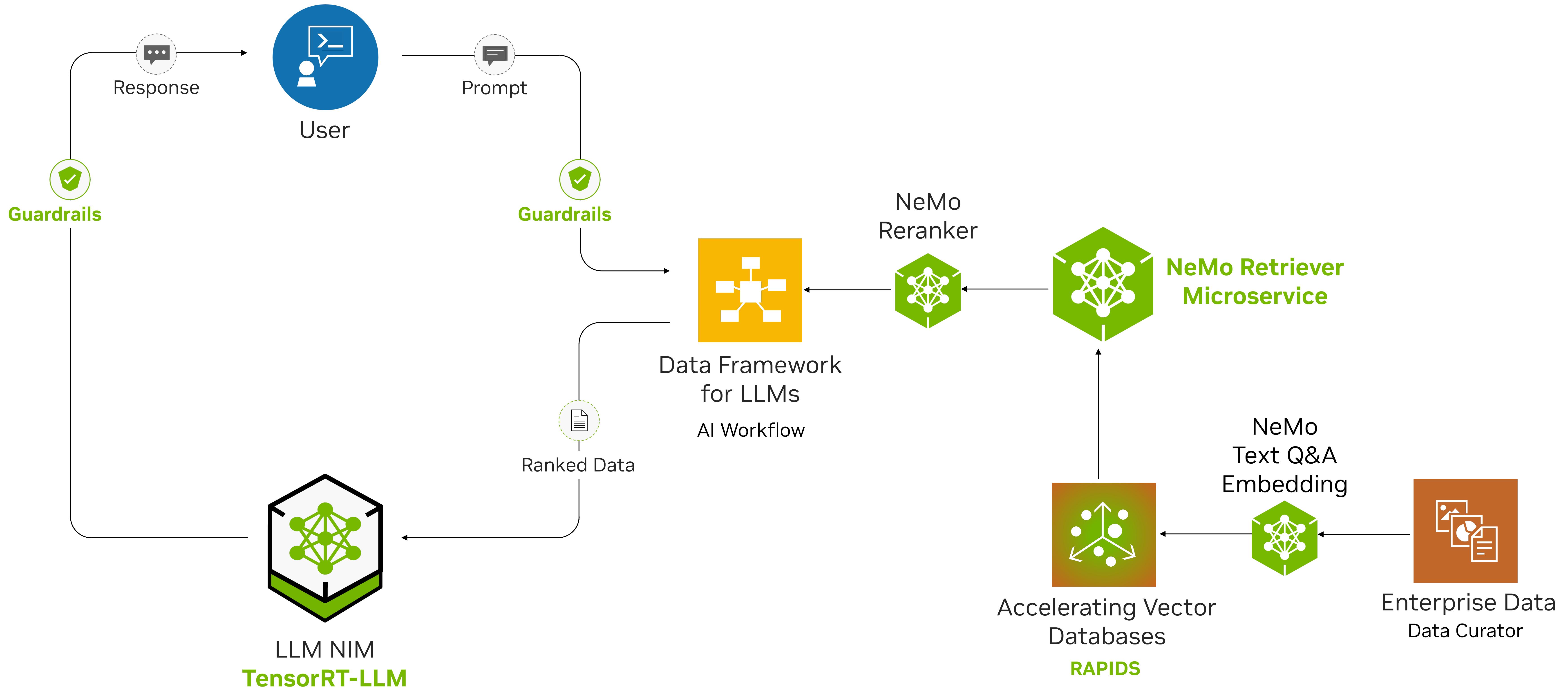
Building Generative AI Applications for the Enterprise

Build, customize and deploy generative AI models with NVIDIA NeMo



NVIDIA Provides Optimized Retrieval Augmented Generation

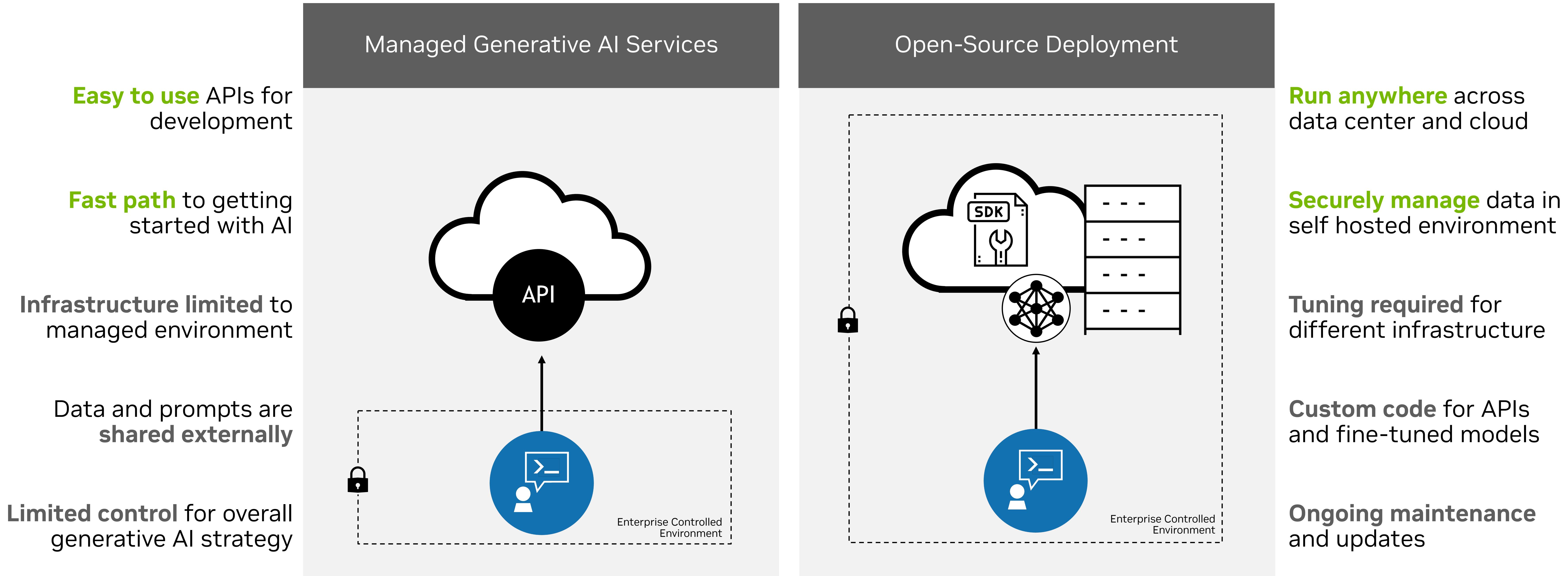
Commercially viable, optimized embedding, reranking, and personalization to deliver highest accuracy and performance



NeMo Microservice

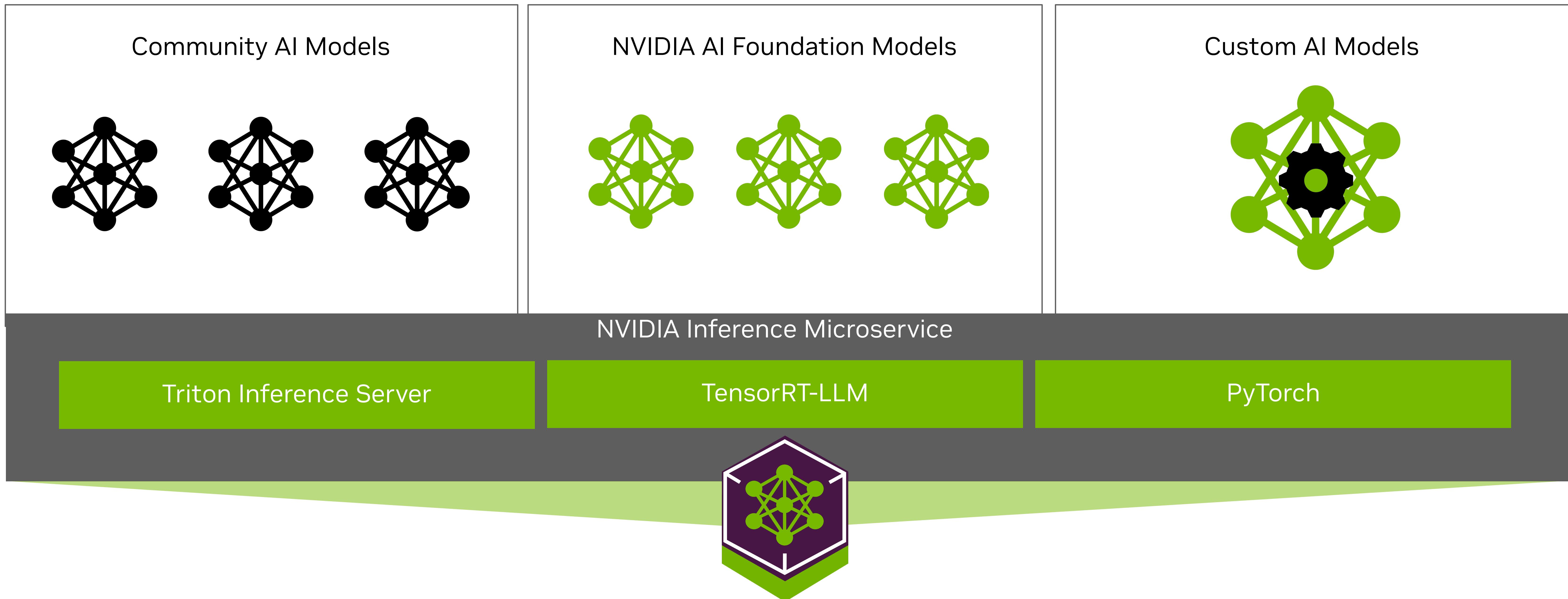
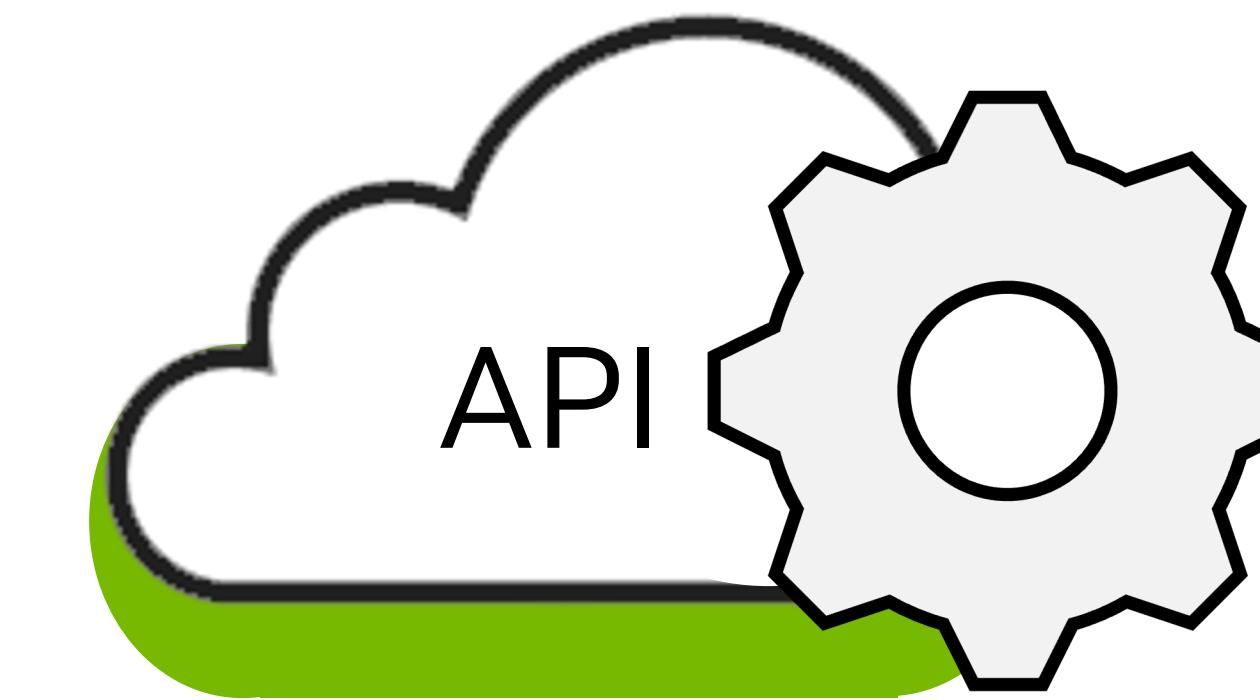
Enterprises Face Challenges Experimenting with Generative AI

Organizations must choose between ease of use and control



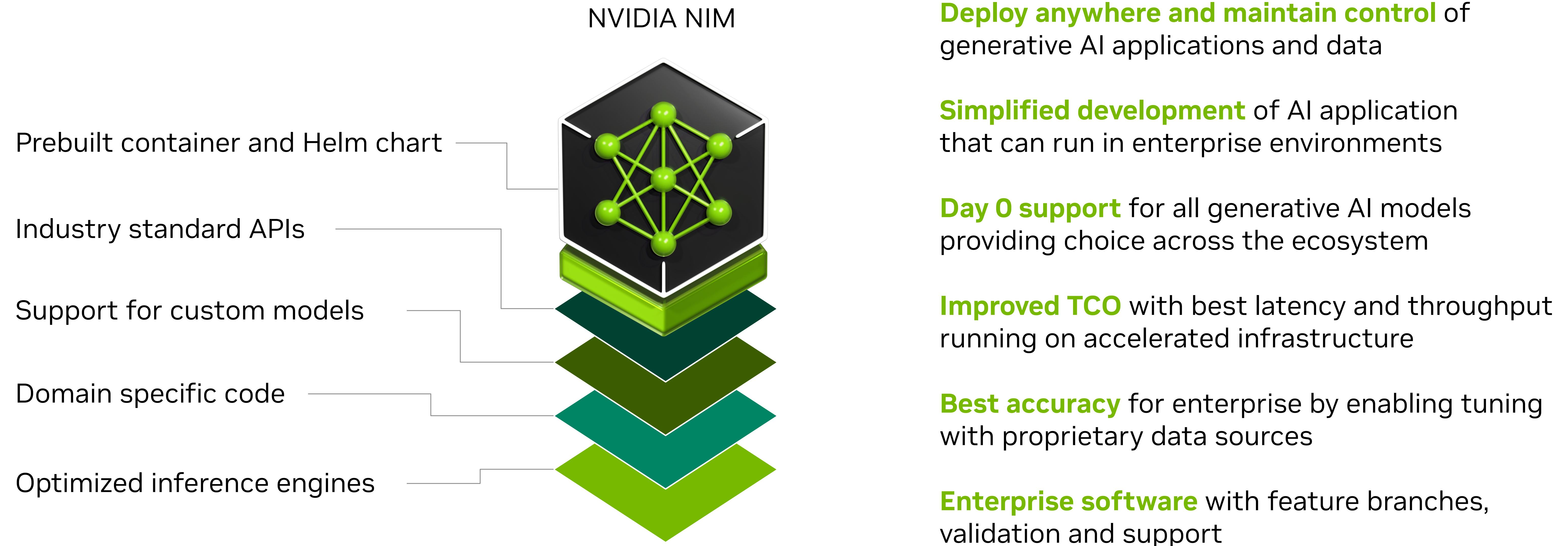
NVIDIA Inference Microservice

Easy-to-use microservice to run generative AI in production at scale anywhere



NVIDIA NIM Optimized Inference Microservices

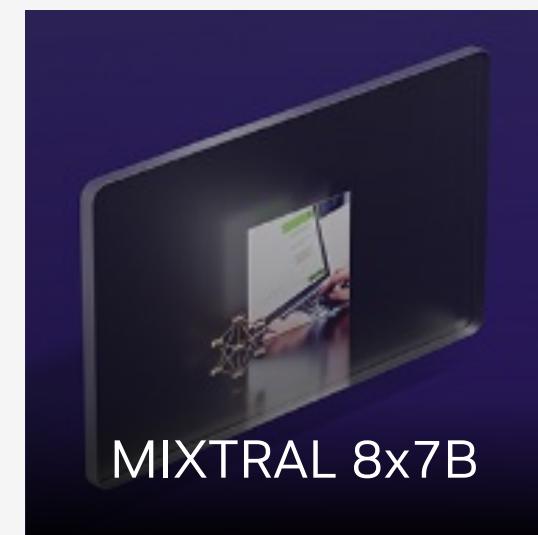
Accelerated runtime for generative AI



Inference Microservices for Generative AI

NVIDIA NIM is the fastest way to deploy AI models on accelerated infrastructure across cloud, data center, and PC

NVIDIA API Catalog



MIXTRAL 8x7B



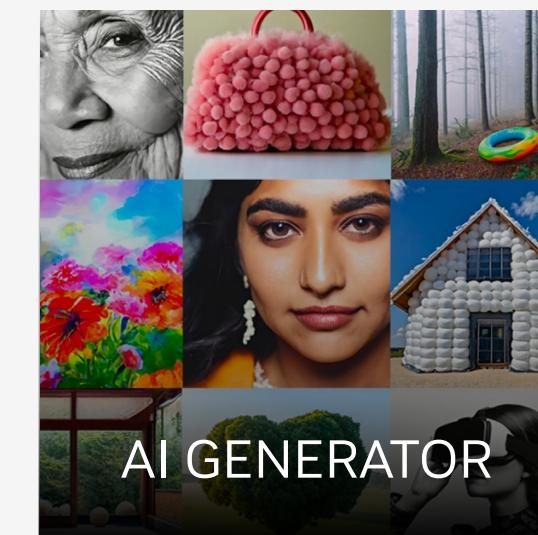
GEMMA 7B



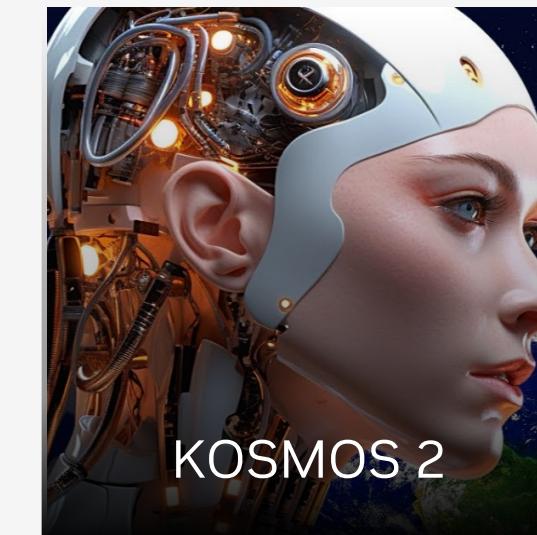
FUYU



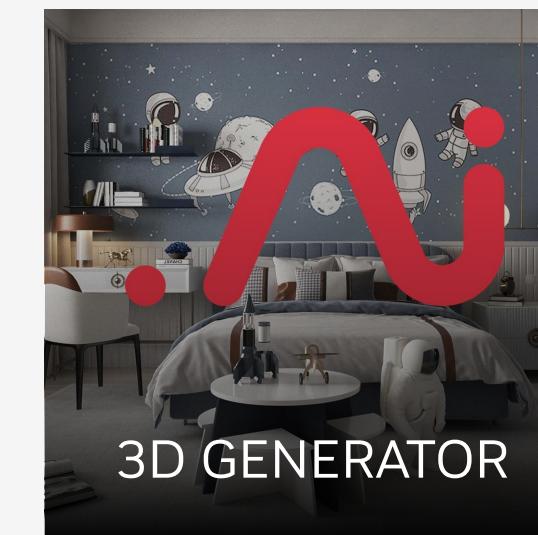
NEMO RETRIEVER



AI GENERATOR



KOSMOS 2



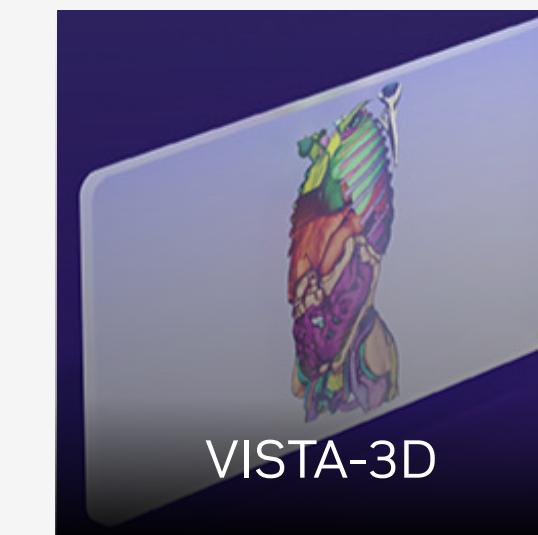
3D GENERATOR



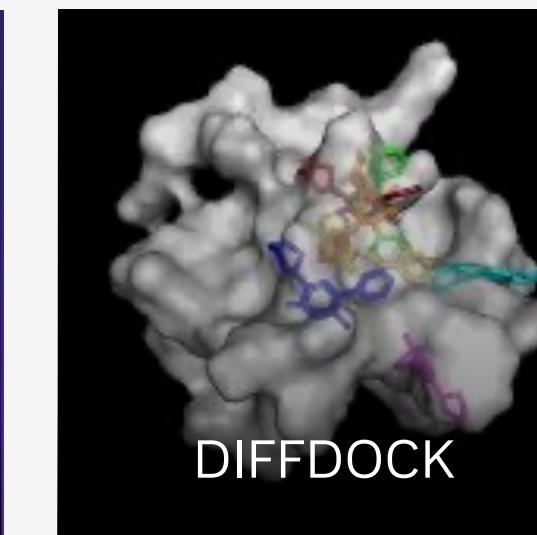
AUDIO2FACE



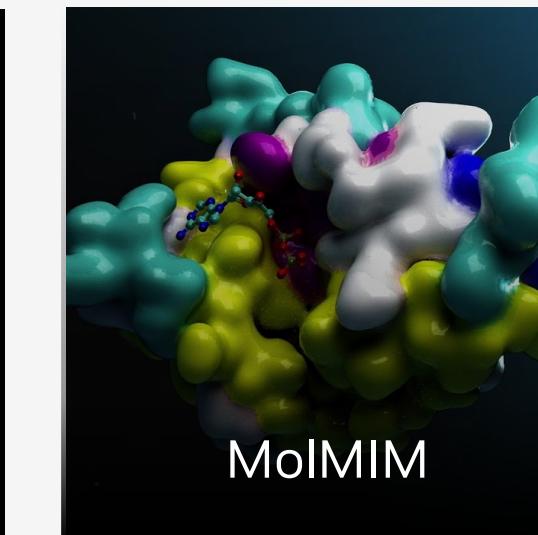
ESM FOLD



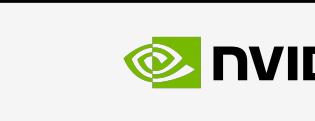
VISTA-3D



DIFFDOCK



MolMIM



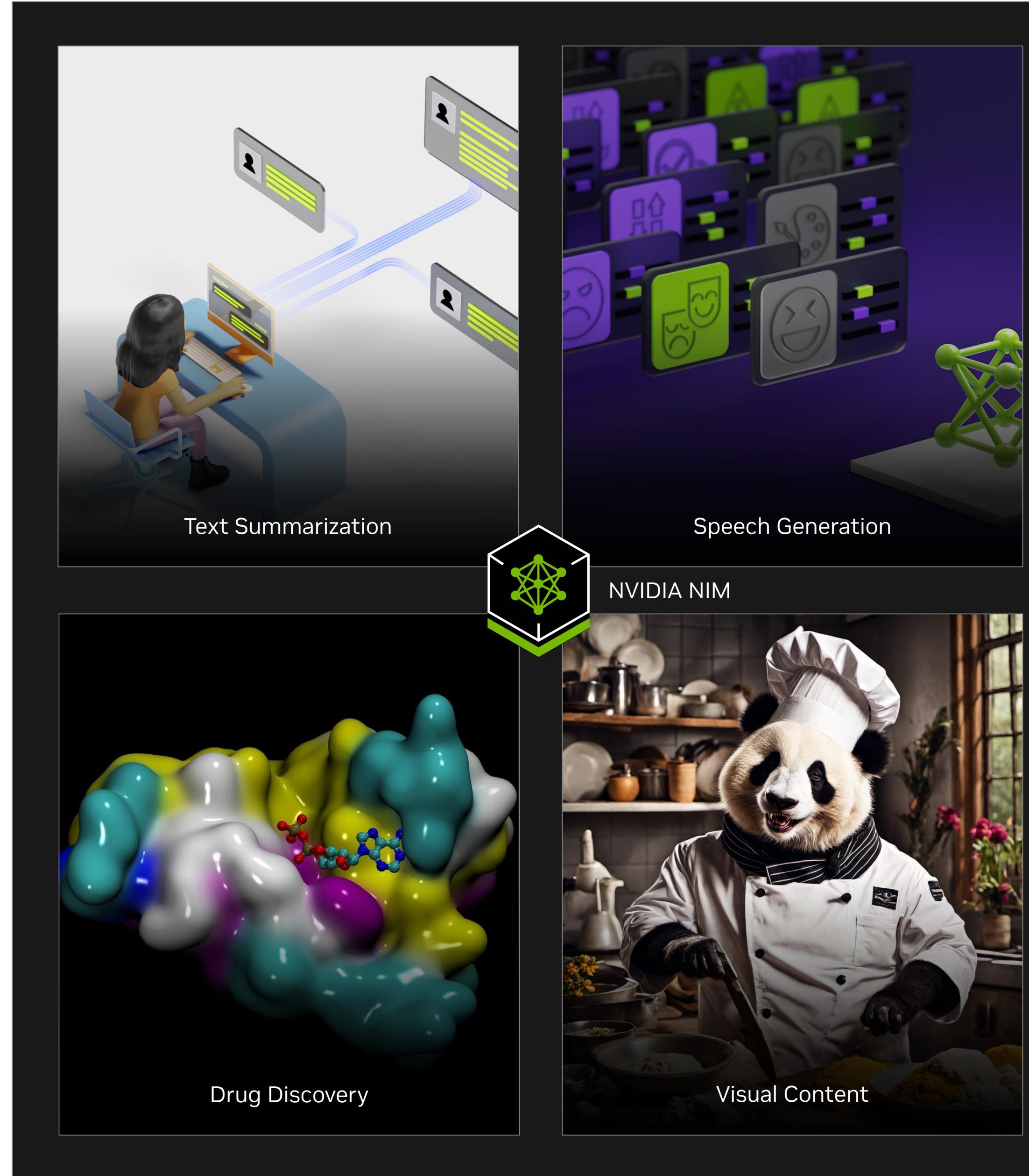
NVIDIA NIM is the Fastest Path to AI Inference

Reduces engineering resources required to deploy optimized, accelerated models

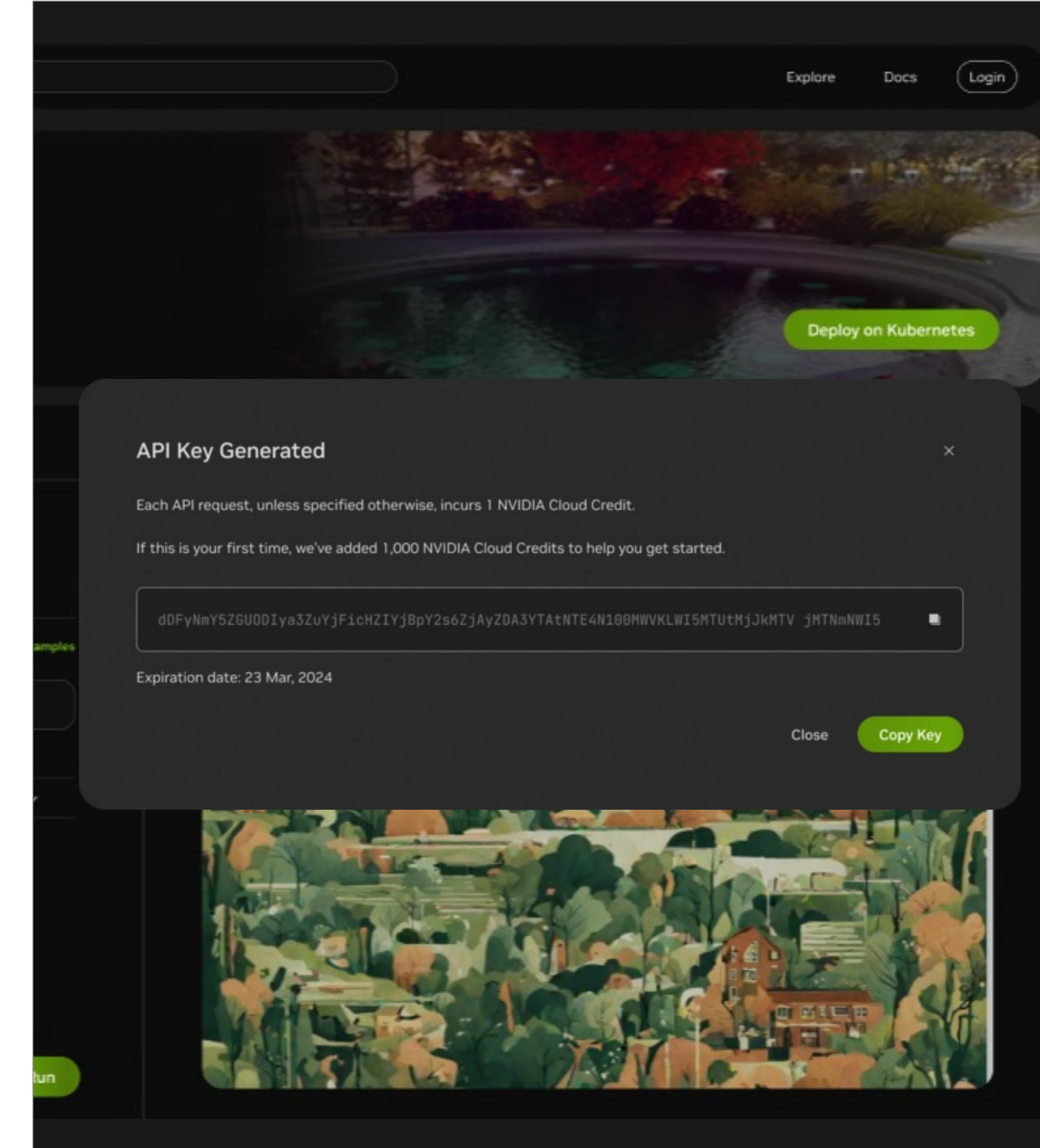
	NVIDIA NIM	Triton + TRT-LLM Opensource
Deployment Time	5 minutes	~1 week
API Standardization	Industry standard protocol OpenAI for LLMs, Google Translate Speech	User creates a shim layer (reducing performance) or modify Triton to generate custom endpoints
Pre-Built Engine	Pre-built TRT-LLM engines for NV and community models  MISTRAL AI_  Llama 2  starcoder  Nemotron	User converts checkpoint to TRT-LLM format and creates and runs sweeps through different parameters to find the optimal config
Triton Ensemble/ BLS Backend	Pre-built with TRT-LLM to handle pre/post processing (tokenization)	User manually sets up + configures
Triton Deployment	Automated	User manually sets up + configures
Customization	Supported – P-tuning and LORA, more planned	User needs to create custom logic
Container Validation	Pre-validated with QA testing	No pre-validation
Support	NVIDIA AI Enterprise - Security and CVE scanning/patching and tech support	No enterprise support

Experience and Run Enterprise Generative AI Models Anywhere

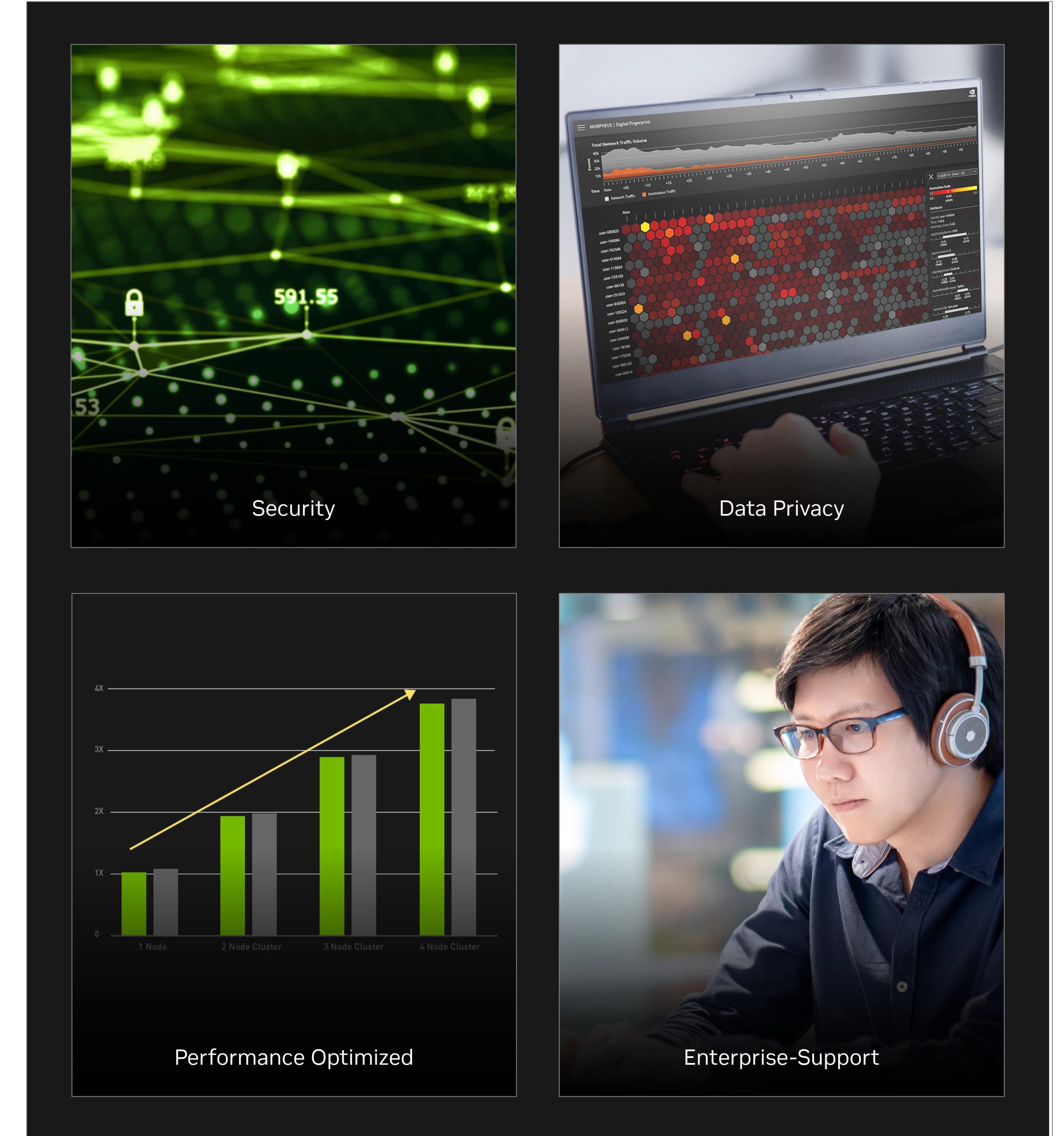
Seamlessly integrate AI in business applications with NVIDIA AI APIs



Experience Models



Prototype with APIs



Deploy with NIMs

Building Generative AI Applications for the Enterprise

Build, customize, and deploy generative AI models with NVIDIA NeMo.

