

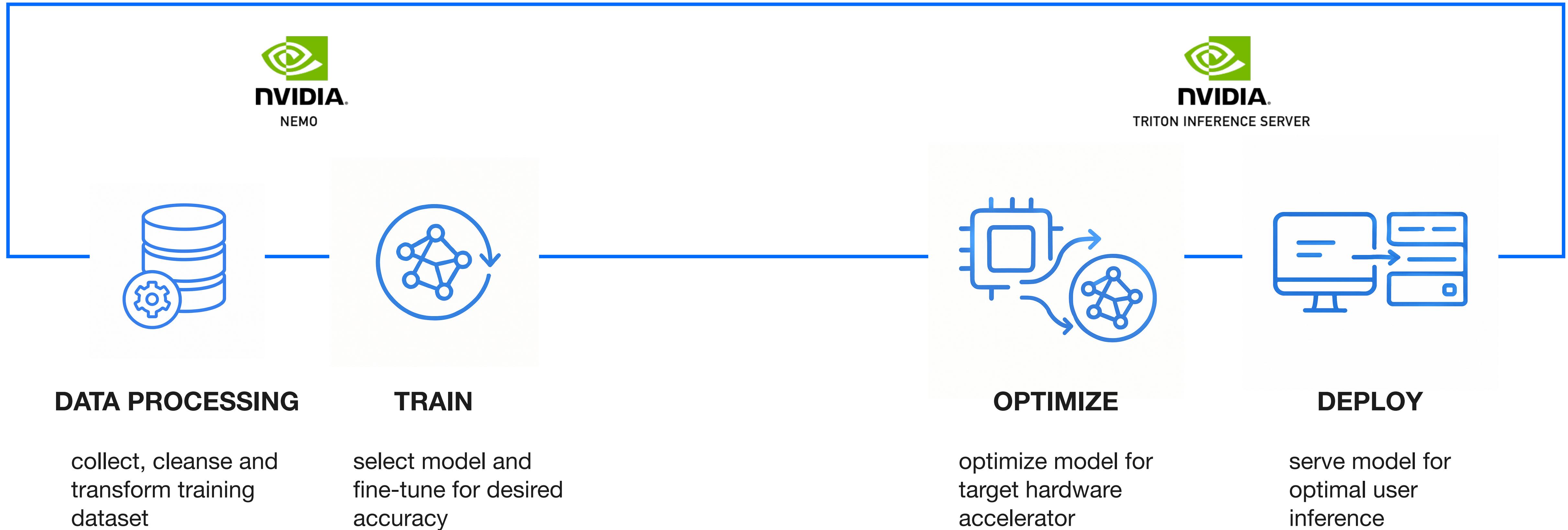
Fine-Tuning and Deploying LLMs with NeMo Framework and Triton Inference Server

Emiliano Molinaro, Ph.D.

SDU eScience Center

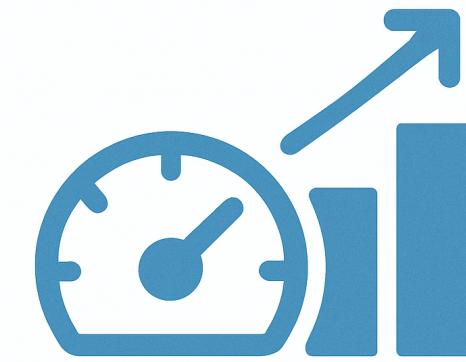
From Training Pipelines to Real-Time Inference

AI Training



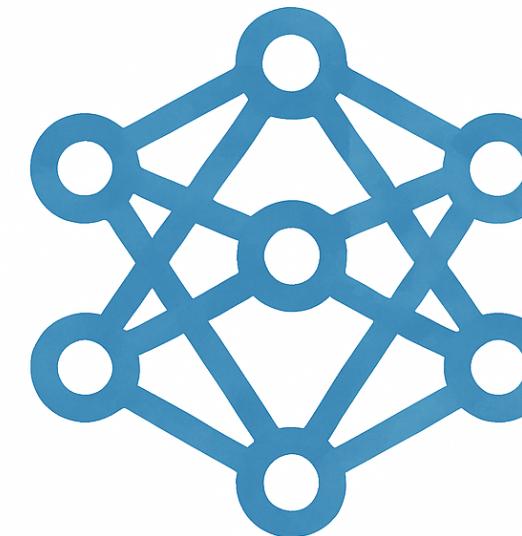
Overview of NeMo Framework

Performance & Scalability



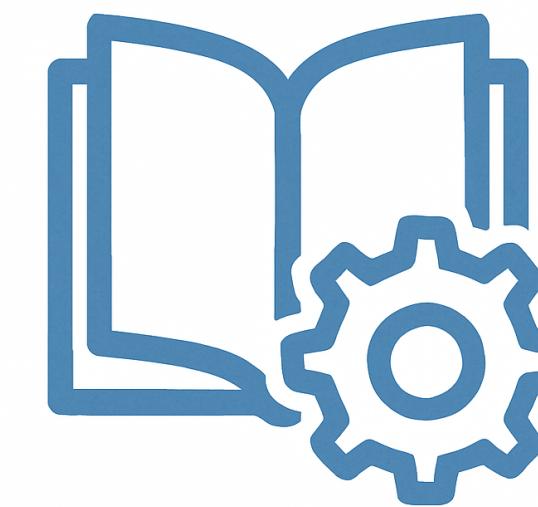
- Data Parallelism
- Tensor Parallelism
- Pipeline Parallelism
- Context Parallelism
- Expert Parallelism
- GPU Acceleration

Model Coverage



- Supports Several Model Families:
LLMs, SSMs, MoEs, SD, VLMs, ...

Training Algorithms



- PEFT: LoRA, P-tuning, IA3, QLoRA, Adapters
- Model Alignment: RLHF PPO, DPO, KTO, IPO, RLAIF, SteerLM, Rejection Sampling

Usability



- Python-based API
- Config-Driven
- Modular
- Production-Ready

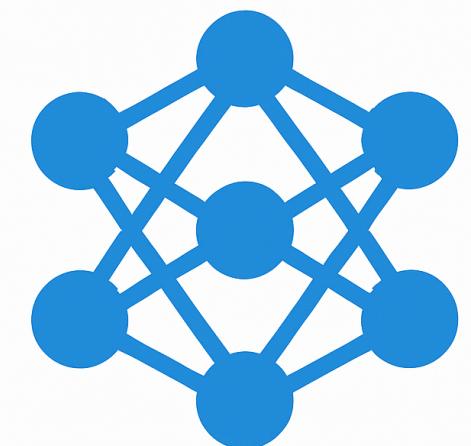
Parameter-Efficient Fine-Tuning

Efficient Fine-Tuning Methods

Technique	How it works	Parameters trained
Prompt Tuning	◆ Add prompts to input	<input checked="" type="checkbox"/> Prompt embeddings only
Prefix Tuning	◆ Insert prompts at every layer	<input checked="" type="checkbox"/> Layer-wise prompts only
P-Tuning	◆ Generate prompts with encoder (LSTM)	<input checked="" type="checkbox"/> Small encoder only
LoRA	◆ Low-rank adaptation in attention layers	<input checked="" type="checkbox"/> Tiny subset of layer parameters

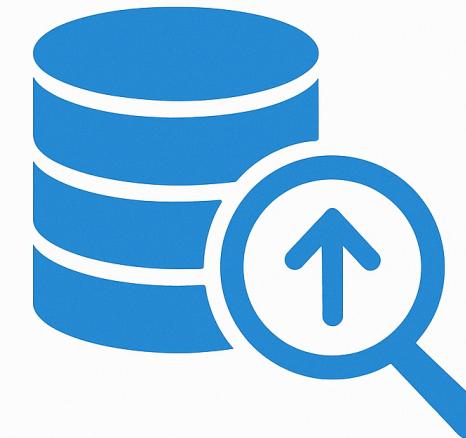
Overview of Triton Inference Server

Any Framework



Multiple Backends:
TensorFlow, PyTorch,
TensorRT, XGBoost,
ONNX, Python,
TensorRT-LLM, vLLM

Any Query Type



Optimised for Real Time,
Batch, Audio/Video Streaming,
Ensemble Inferencing

High Performance

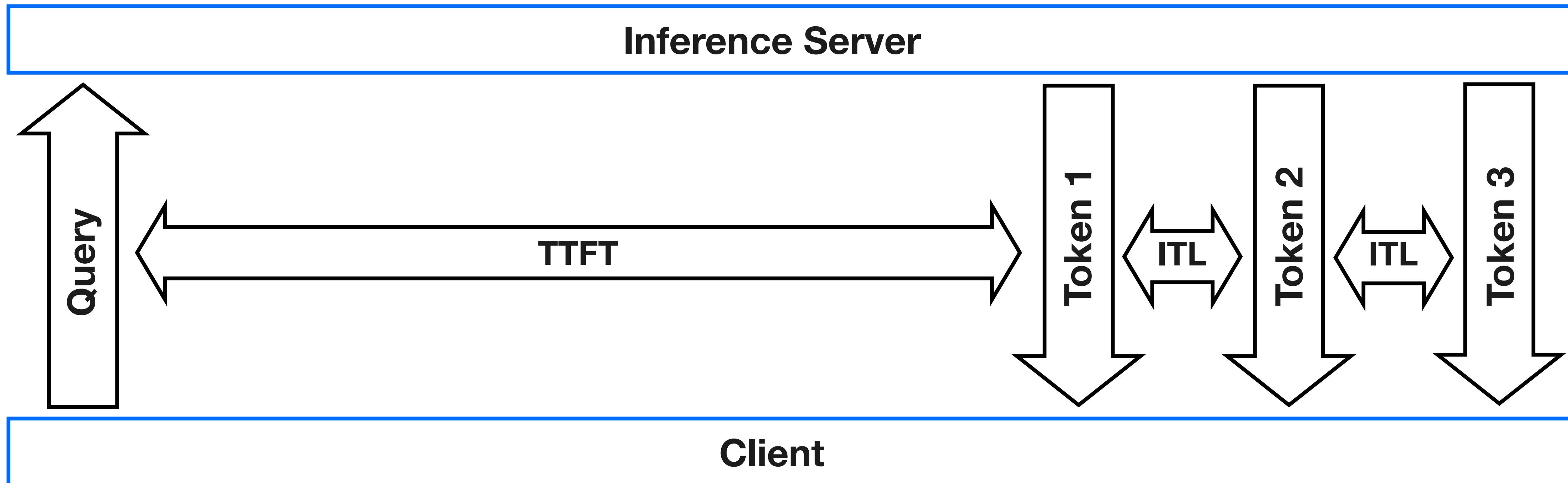


Model Analyser for
Optimal Configuration

Optimised for High GPU/CPU
Utilisation, High Throughput,
Low Latency

Latency Metrics

- **TTFT (Time To First Token):**
The duration from the request initiation to the generation of the first token in the response.
- **End-to-End Latency (E2E):**
The total time taken by the model to generate the entire response from start to finish.
- **Inter-Token Latency (ITL) / Time Per Output Token (TPOT):**
The average time between the generation of consecutive tokens in streaming scenarios, reflecting client-side token wait times.



Inflight Batching

Maximize GPU Utilisation during LLM Serving

Benefits of Inflight Batching:

- **Enhanced Throughput:**
 - Processes more requests per GPU operation.
 - Reduces per-request overhead, leading to higher overall throughput.
- **Optimized Resource Utilization:**
 - Maximizes GPU usage by maintaining high occupancy rates.
 - Cost-effective scaling by leveraging existing hardware efficiently.
- **Minimal Latency Impact:**
 - Configurable batching parameters ensure that latency remains within acceptable bounds.
 - Balances between batch size and response time to meet application needs.

[Figure from NVIDIA](#)

Batch Elements	Iteration									...
	1	2	3	4	5	6	7	8	9	
R ₁	Context	Gen	EoS	NoOp	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	Context	Gen	EoS	NoOp	END	R ₅				...
R ₂			END	R ₆						...
R ₃					END	R ₇			END	R ₉
R ₄									END	R ₈

Inflight Batching

Context Gen EoS NoOp

Inflight batching and the additional kernel-level optimizations improve GPU usage and doubles throughput which helps reduce energy costs and minimize TCO.

Hands-On Session

Implement end-to-end workflow on **UCloud - Interactive HPC**:

- Deploying an LLM in production for optimized inference
- Creating a synthetic medical Q&A dataset
- Setting up and fine-tuning a LLM with PEFT (LoRA)

Notebooks available in this GitHub repository:

<https://github.com/emolinaro/ucloud-workshop-28-05-2025>

