



Triton Inference Server Meetup

09/09/2024



Agenda

- Triton vLLM Backend Updates

- Triton Inference Server Updates

- Open Source Improvements and Deliverables

- GenAI Perf

- Triton 3 and Disaggregated Serving



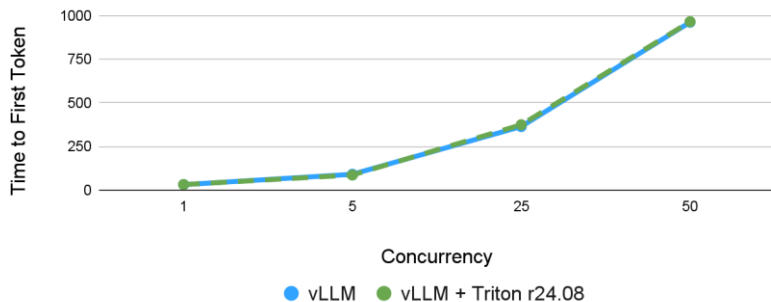
Triton vLLM Backend Updates

vLLM Backend Performance Improvement

Llama2-7B on single A100 40GB with 200 input - 1000 output tokens using vLLM 0.5.3 post

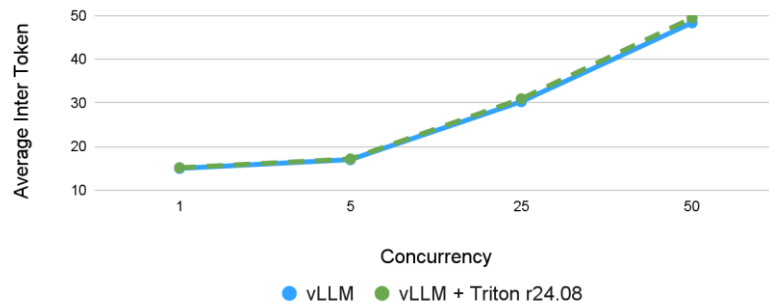
Time to First Token

Lower the Better



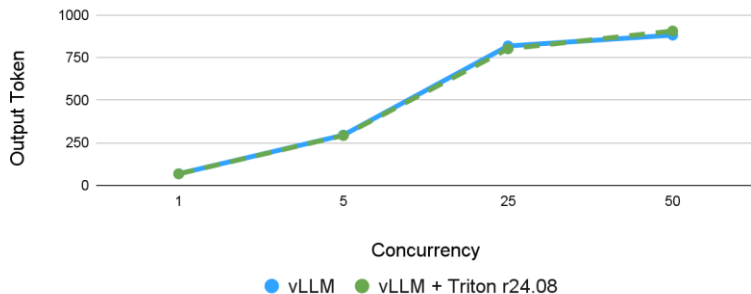
Average Inter Token Latency

Lower the Better



Output Token Throughput

Higher the Better



**Within < 2% to vLLM's performance
for both throughput and latency**

Delegate response sending and cancellation checks to another thread and wait with GIL released, allowing vLLM Engine to have more CPU time.

vLLM Metrics Access Through Triton

In addition to counter and gauge metrics, Triton now supports **histogram metrics**.

The following vLLM metrics will be supported through Triton with 24.08 and 24.09 releases

Supported vLLM metrics in r24.08

- **vllm:prompt_tokens_total** : counter of prefill tokens processed.
- **vllm:generation_tokens_total** : counter of generation tokens processed.
- **vllm:time_to_first_token_seconds** : histogram of time to first token in seconds.
- **vllm:time_per_output_token_seconds**: histogram of time per output token in seconds.

More in upcoming r24.09

- **vllm:e2e_request_latency_seconds** : histogram of end to end request latency in seconds.
- **vllm:request_prompt_tokens** : histogram of prefill tokens processed.
- **vllm:request_generation_tokens** : histogram of generation tokens processed.
- **vllm:request_params_best_of** : histogram of the best_of request parameter.
- **vllm:request_params_n** : histogram of the n request parameter.



Triton Inference Server Updates

Kubernetes (K8s) Multi-Node

Problem

- **Deploy Massive LLMs** (100B - 1T+ parameters) that can't fit on a single node, or make use of older GPUs
- **Automatically Scale and Load Balance** multi-node replicas

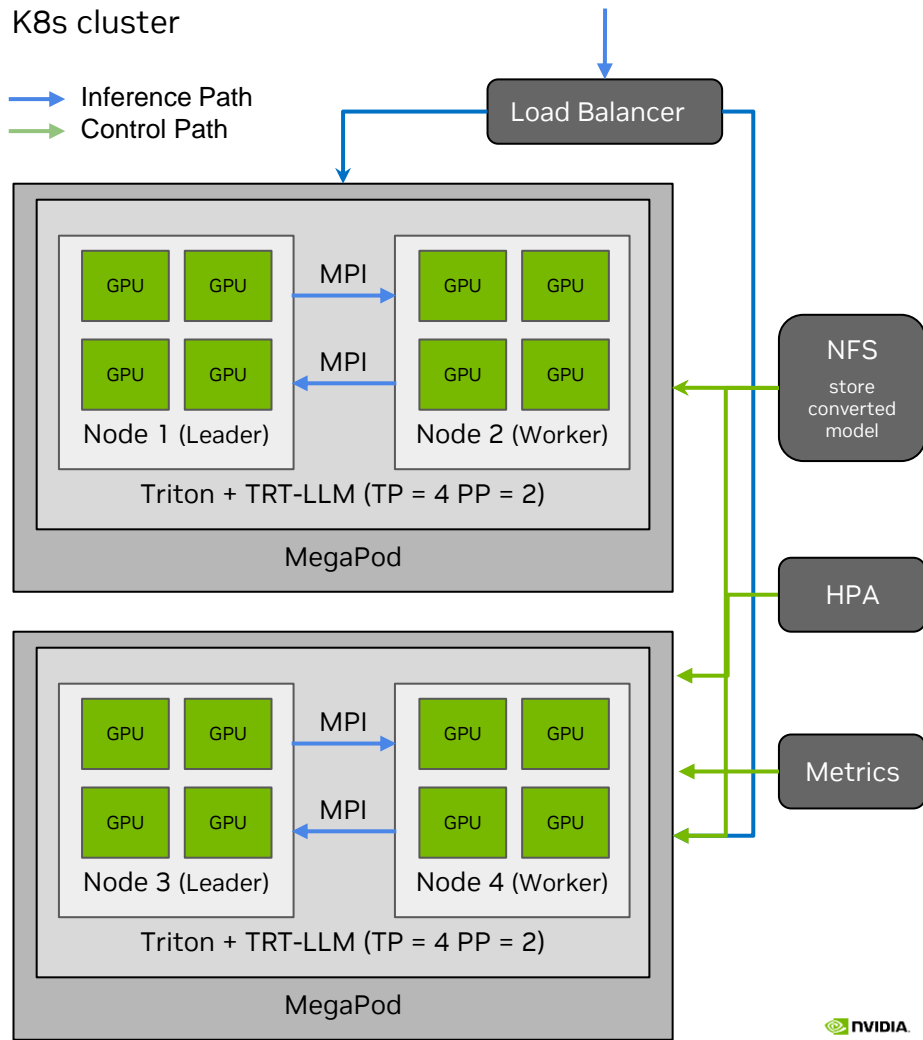
Solution

- **Flexible Kubernetes definitions** that can scale any model supported by TRT-LLM on any GPU.

Key Challenges

- **Multi-Node Communication** (MPI)
- **Deploying Groups of Nodes** (Gang Scheduling)
- **Scaling Groups of Nodes** (LeaderWorkerSet)
- **Leader-Aware Load Balancing**
- **Auto-Scaling Metrics**

K8s cluster

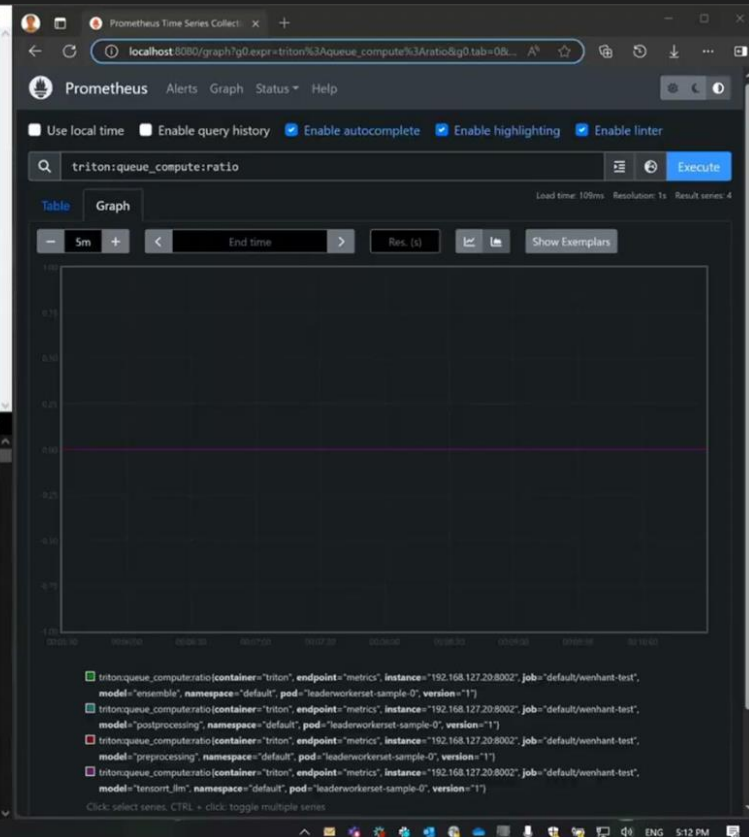


Kubernetes Multi-Node Demo

Llama 3 70B on 2 nodes of 4x A10Gs

```
aws_terminal:~$ kubectl get pods -w
NAME                                READY   STATUS    RESTARTS   AGE
leaderworker-sample-0              1/1     Running   0           95m
leaderworker-sample-0-1            1/1     Running   0           95m
```

```
venhant@NV-HD20283: ~
local_terminal:~$
```



Triton In-Process Python API

- **Python bindings to C APIs** for inference, scheduling, and model management
- **Native Python** support for built-in types like dicts, lists, etc.
- **Non-blocking** inference calls to allow running business logic concurrently
- **Flexible Tensor interoperability** for zero-copy conversions between NumPy, CuPy, PyTorch, DLPack

```
import tritonserver

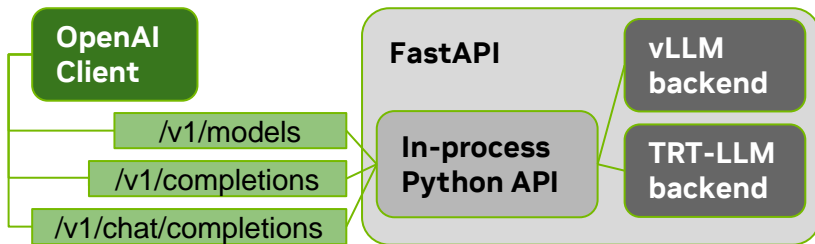
# Configure
options = tritonserver.Options(
    model_repository="/path/to/models",
    log_info=True,
    log_warn=True,
    log_error=True,
)

# Serve
server = tritonserver.Server(options)
server.start(wait_until_ready=True)

# Inference: Non-blocking
inputs = {"text_input": ["Machine learning is"],
          "stream": [True]}
model = server.model("llama-3.1-8b-instruct")
responses = model.infer(inputs=inputs)

# Responses are populated as they are received, so other
# work can be done in parallel while awaiting responses.
for response in responses:
    print(response.outputs["text_output"].to_string_array())
```

OpenAI Compatible API (Beta)



- **Customizable** FastAPI app with support for OpenAI schemas
- **Drop-in replacement** for the supported APIs with OpenAI clients, genai-perf, curl, etc.
- **Triton In-Process Python API** used for managing state and handling inference requests
- **Optimized** Triton backends for vLLM and TensorRT-LLM used for the inference runtime

```
from openai import OpenAI
client = OpenAI(
    base_url="http://localhost:8000/v1", # Self-hosted
    api_key="empty",
)

completion = client.chat.completions.create(
    model="meta-llama/Meta-Llama-3.1-8B-Instruct",
    messages=[
        {"role": "system", "content": "You are a Triton
Inference Server expert."},
        {"role": "user", "content": "Hello vLLM meetup!"}
    ],
    max_tokens=256
)

print(completion.choices[0].message.content)
```

Constrained Decoding and Function Calling

LLM is a component of an **agentic workflow**

Constrained decoding ensures LLMs to enforce a specific formatting, e.g. JSON

- [Constrained Decoding tutorial](#) provides examples for structured generation via prompt engineering and external libraries such as LM Format Enforcer

Function calling enables LLMs to perform complex tasks requiring specific computations or data retrieval

- [Function Calling tutorial](#) provides further insights into controlling model behavior

```
# Specify the custom logits post-processor to use
executor_config.logits_post_processor_map = {
    "<name>": custom_logits_processor
}
self.executor = trtllm.Executor(
    model_path=...,
    model_type=...,
    executor_config=executor_config)
...
# At inference time
request.logits_post_processor_name = "<name>"
```

```
python3 client.py --prompt "How's Nvidia doing?" -o 200 -
-verbose
```

```
{
    "step": "1"
    "description": "Get the current stock price for
NVIDIA",
    "tool": "get_current_stock_price",
    "arguments": {
        "symbol": "NVDA"
    }
}
```

```
=====
Executing function:
get_current_stock_price({'symbol': 'NVDA'})
Function response: 106.38
=====
```

Revamped documentation (24.09)

Triton – TensorRT-LLM Guide

- Run TRT-LLM models with Triton Server
- Advanced Configurations
- Deployment Strategies
- Tutorials for popular LLM models

Triton Scaling Guide

- Multi-Node
- Multi-Instance





AI Agents Guide

- Constrained Decoding
- Function Calling

More in upcoming months

- Versioning
- OpenAI API Guide

Before
Getting Started
Quickstart
User Guide
Deploying your trained model using Triton
Triton Architecture
Model Repository
Repository Agent
Model Configuration
Request Cancellation
Optimization
Ragged Batching
Rate Limiter
Model Analyzer
Model Management
Custom Operations
Decoupled Backends and Models
Triton Response Cache
Metrics
Triton Server Trace
Triton Inference Server Support for Jetson and JetPack
Version 1 to Version 2 Migration
Secure Deployment Considerations

After
Table of Contents
Home
Release Notes 
Compatibility matrix
Getting Started
Quick Deployment Guide (by backend) 
Overview
TRT-LLM
vLLM
Python (HuggingFace)
PyTorch
ONNX
TensorFlow
Openvino
Multimodal model
Stable diffusion
Scaling guide
Multi-Node
Multi-Instance
API Reference
OpenAI API (BETA)
KServe API 
In-Process Triton Server API 



Open Source Improvements and Deliverables

Public CI

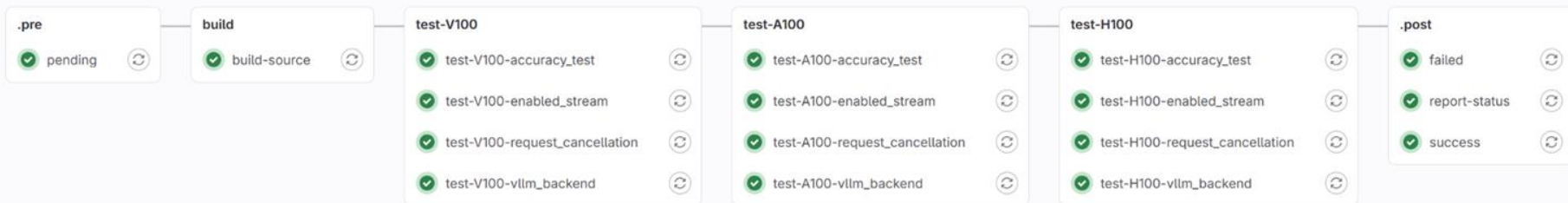
- Enable community members to confidently contribute.
- Verify your changes are not breaking
 - Every new vLLM release to be validated against the latest Triton release
 - Every Pull Request on vllm_backend to be validated against the latest vLLM release
 - Provide coverage across state-of-the-art hardware
- Improve transparency of our build process

vllm_backend / README.md 

License BSD3 Triton 24.08 vLLM 0.5.5 CI Passing V100,A100,H100

6 checks passed		
✓	 mirror_repo	Details
✓	 build (ubuntu-22.04, 3.10)	Details
✓	 pre-commit	Details
✓	 trigger-ci	Details
✓	 ci/gitlab/gitlab-master.nvidia.com Pipeline passed on GitLab	

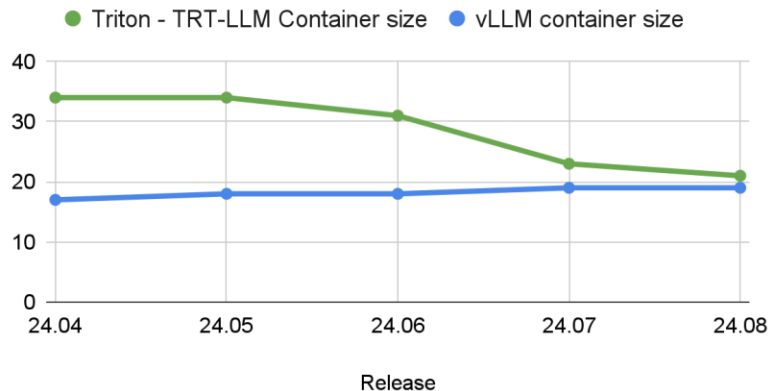
Pipeline Needs Jobs 27 Tests 0



Container Size Reduction and Manylinux Support

Container size reduction

Triton container sizes (GB)



- **Targeting 15GB**
- **40% container size reduction** from 34GB to 21GB, over the course of last three releases
 - Removed unnecessary installations
 - Optimized our docker build

Manylinux

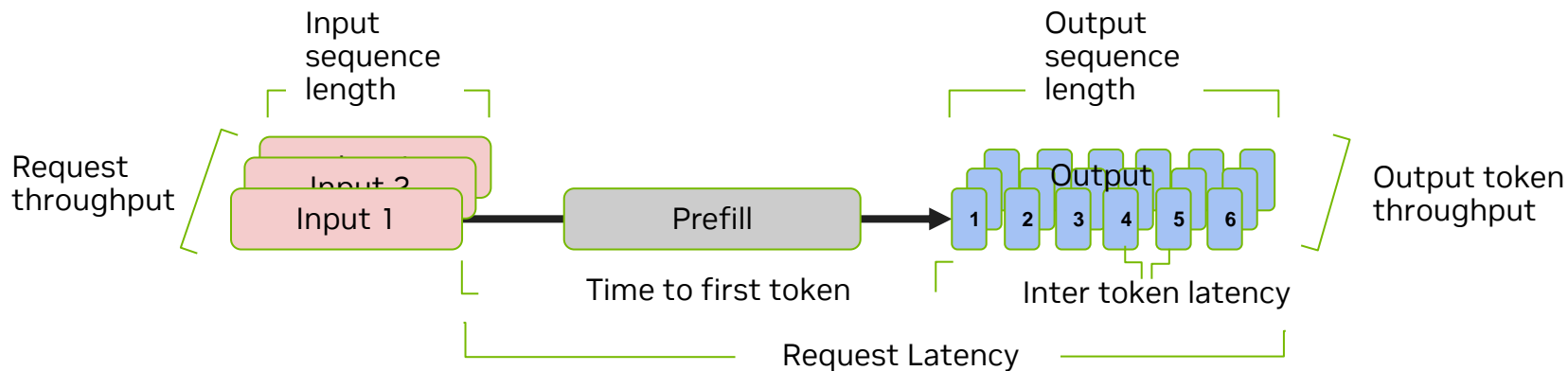
- **vLLM** is coming soon!
- **24.11 GA**: Targeting complete feature parity with our NGC containers
- **Early Access** release available upon request
 - Supported backends
 - PyTorch
 - TensorFlow2
 - ONNX Runtime
 - Supported HW
 - X86 CPU
 - ARM64 CPU
 - GPU
 - Triton server core



GenAI Perf

GenAI Perf

Enables apples to apples comparison for GenAI perf evaluation



Key features

- HW/SW agnostic performance benchmarking via **KServe and OpenAI API**
- Supports LLM, Visual Language Model, embedding, re-ranking, and multi-LoRA
- In-process benchmarking for TRT-LLM

Planned features

- **Customizable front-end** to allow custom benchmarking logic to easily extend GenAI-Perf.
- Access to **raw metric data** for easy post processing and **scripts for generating plots**.

GenAI-Perf Usage and Output

Three Input Formats

- **Synthetic**
- **Bring Your Own Data (File)**
- **External Dataset**

Server-agnostic and model-agnostic

Available via Docker and local build

Example command

```
genai-perf profile \  
-m gpt2 \  
--service-kind openai \  
--endpoint-type completions \  
--synthetic-input-tokens-mean 200 \  
--synthetic-input-tokens-stddev 50 \  
--streaming \  
--output-tokens-mean 100 \  
--output-tokens-stddev 25
```

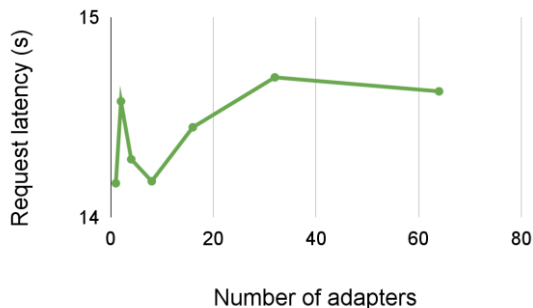
Example output

NVIDIA GenAI-Perf | LLM Metrics

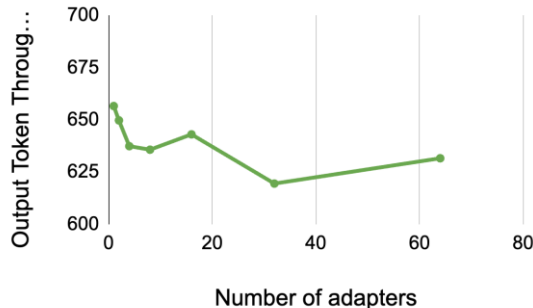
Statistic	avg	min	max	p99	p90	p75
Time to first token (ms)	30.82	18.76	59.07	58.41	52.63	29.02
Inter token latency (ms)	5.29	2.98	15.75	15.71	13.68	4.05
Request latency (ms)	647.73	183.16	2,344.49	2,265.43	1,563.13	582.00
Output sequence length	114.78	49.00	184.00	174.82	146.60	138.50
Input sequence length	206.16	127.00	293.00	287.60	261.20	238.00
Output token throughput (per sec)	176.97	N/A	N/A	N/A	N/A	N/A
Request throughput (per sec)	1.54	N/A	N/A	N/A	N/A	N/A

GenAI-Perf Benchmarks Multi-LoRA Models

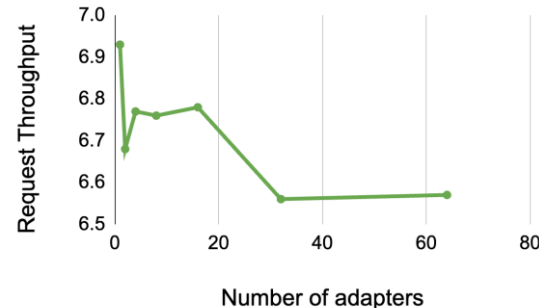
Insight 1 : Stable latency



Insight 2: Stable token throughput



Insight 3 : Lower throughput at > 16 adapters



How many adapters can we deploy without compromising performance?

We can serve up to 64 adapters on one server. For more adapters, we can use GenAI-Perf for further benchmarking.

Benchmarking Setup

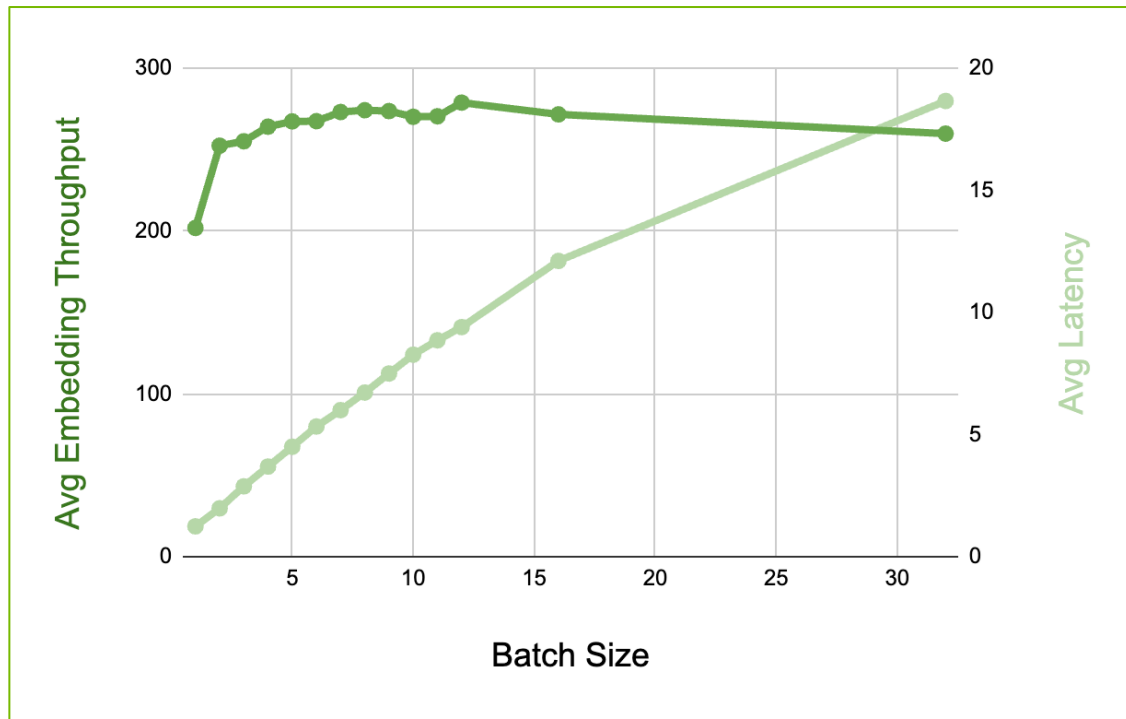
Server: Mistral-7B-v0.1 loaded on a single TGI server running on an NVIDIA RTX 5880

GenAI-Perf: concurrency of 128, random model selection, 500 synthetic prompts

Key arguments to GenAI-Perf:

- `“-m <adapter_name> <adapter_2_name> ... <adapter_n_name>”`
- `“--model-selection-strategy random”`

GenAI-Perf Benchmarks OpenAI API Embeddings Models



At what batch size do we maximize our embedding throughput?

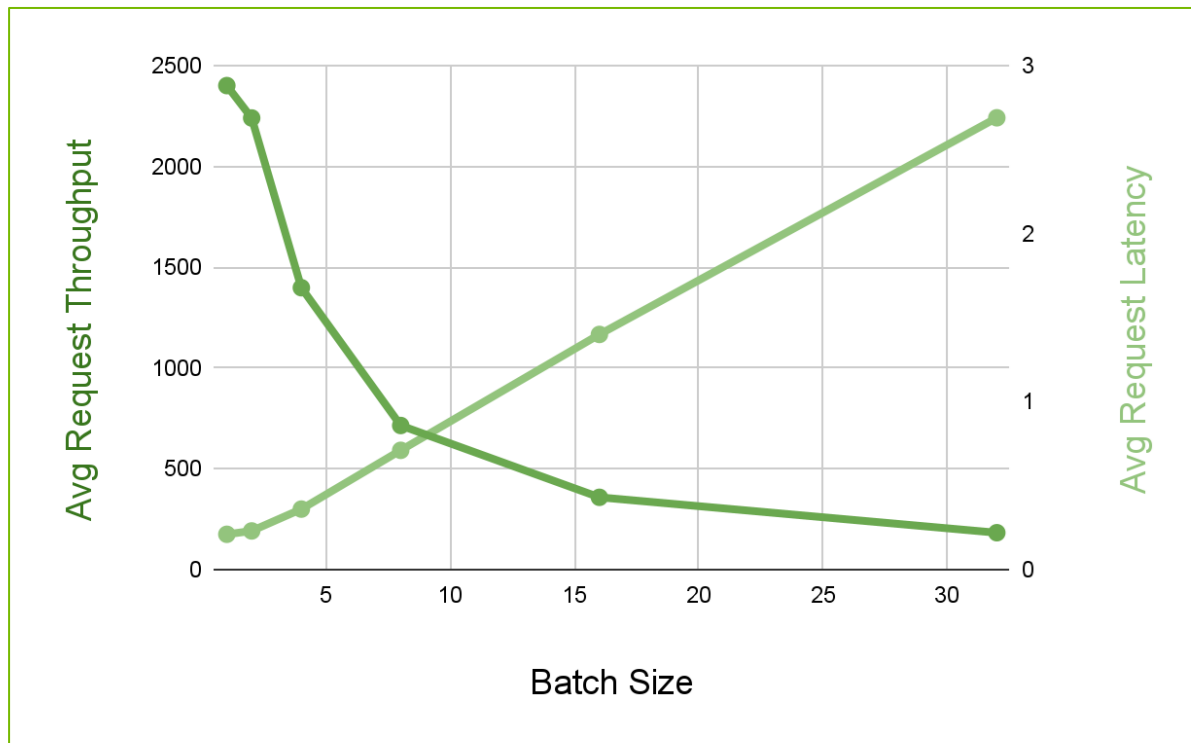
Embedding throughput is maximized at a batch size of 12. Latency increases continually.

Benchmarking Setup

Server: intfloat/e5-mistral-7b-instruct model loaded on a single vLLM OpenAI server running on an NVIDIA RTX 5880

GenAI-Perf: concurrency of 256

GenAI-Perf Benchmarks Hugging Face Re-Ranker API Models



What is the impact of increasing the number of passages per re-ranking request?

The number of passages per request significantly impacts latency and throughput.

Benchmarking Setup

Server: BAAI/bge-reranker-large model loaded on a single Hugging Face TEI server running on an NVIDIA RTX 5880

GenAI-Perf: concurrency of 512

Vision Language Model

Rising popularity of Vision Language Models (VLM) over the past year

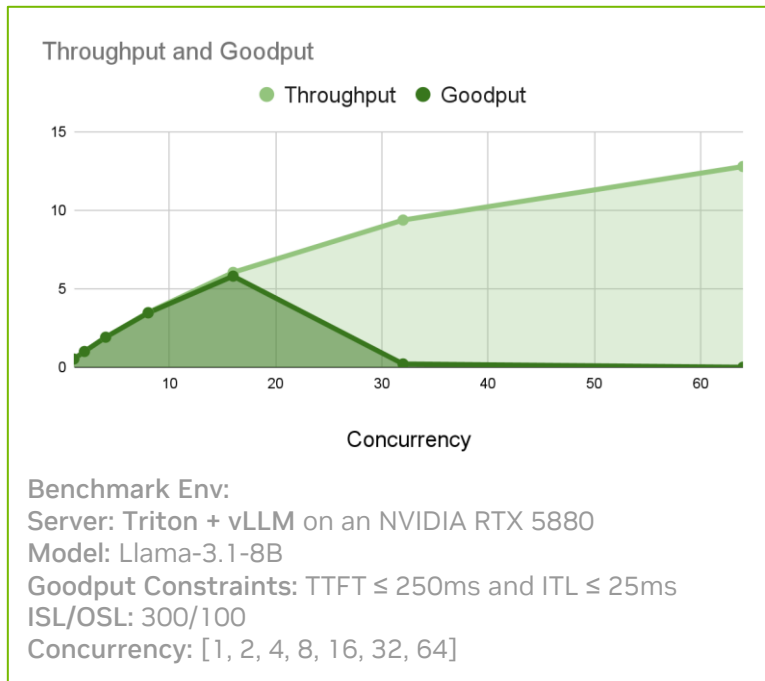
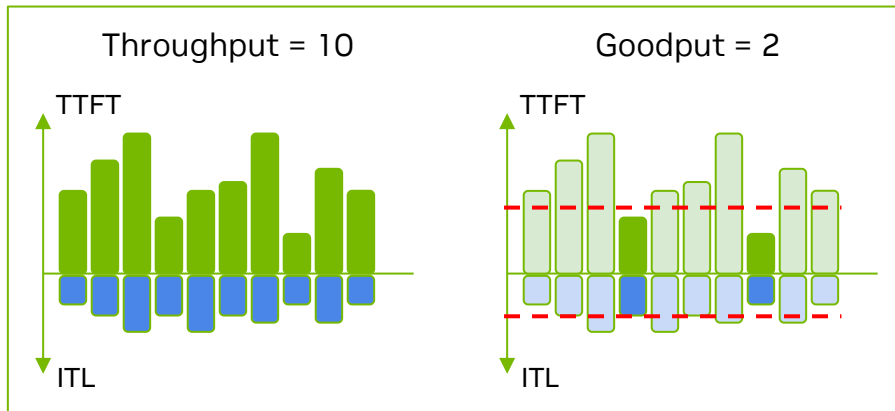
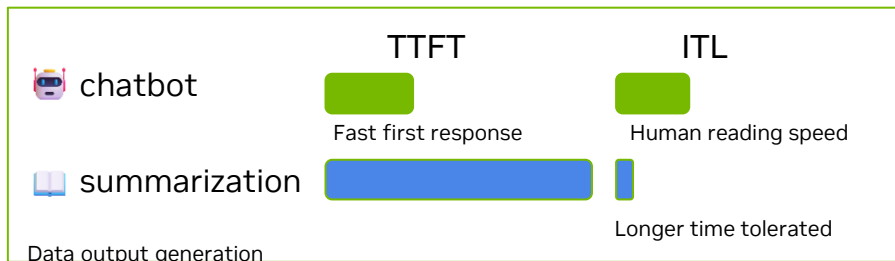
- **OpenAI multi-modal chat API** for VLM inference
 - contains text and image (base64)
- Supports 2 input data sources:
 - **Synthetic data generation**
 - Image resolution (height / width)
 - Image format
 - Prompt length (or, input sequence length)
 - **Bring your own data (BYOD)**

```
payload = {  
  "model": "llava-hf/llava-v1.6-mistral-7b-hf",  
  "messages": [  
    {  
      "role": "user",  
      "content": [  
        {  
          "type": "text",  
          "text": "What's in this image?",  
        },  
        {  
          "type": "image_url",  
          "image_url": {  
            "url":  
f"data:image/jpeg;base64,{base64_image}"  
          }  
        }  
      ]  
    },  
  ],  
  "max_tokens": 100  
}
```

Goodput

Benchmarking GenAI performance from user's perspective

Goodput is defined as the number of completed requests per second that meet certain user-defined service level objectives (SLO) such as time to first token (TTFT) and inter token latency (ITL)



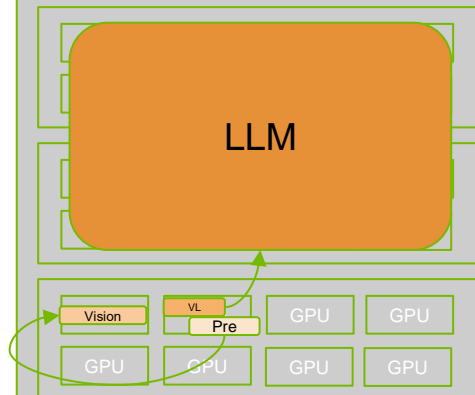


Triton 3.0

Disaggregated Serving

GenAI Models and Workloads Require Data Center Scaling

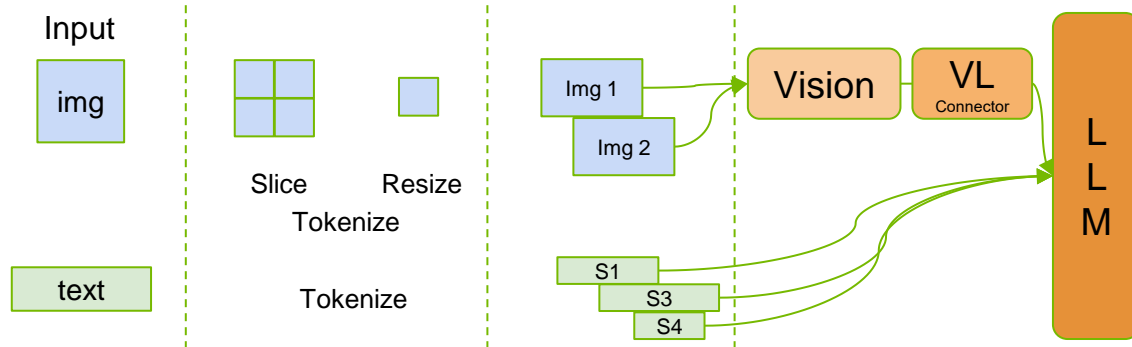
Data center scale



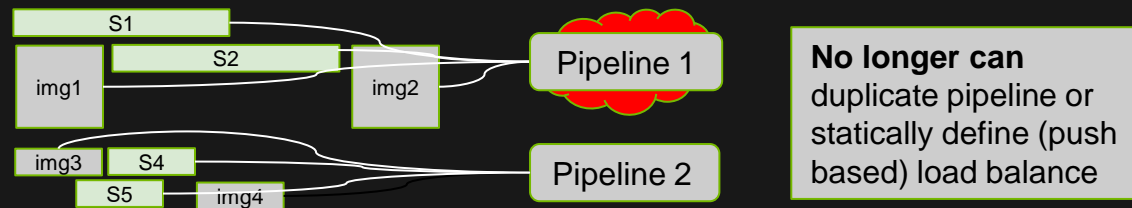
Hetero multistage pipeline

- Multi-node & out of process
- Independent scaling
- Pull (capacity) based load balancing

Simple visual language model (loosely Llava 1.5 HD) example



Bottlenecks & severe under utilization in pipeline (Adobe Firefly, LinkedIn)

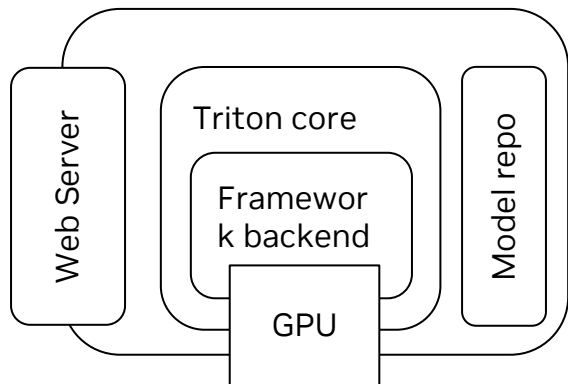


No longer can
duplicate pipeline or
statically define (push
based) load balance

Triton Inference Server

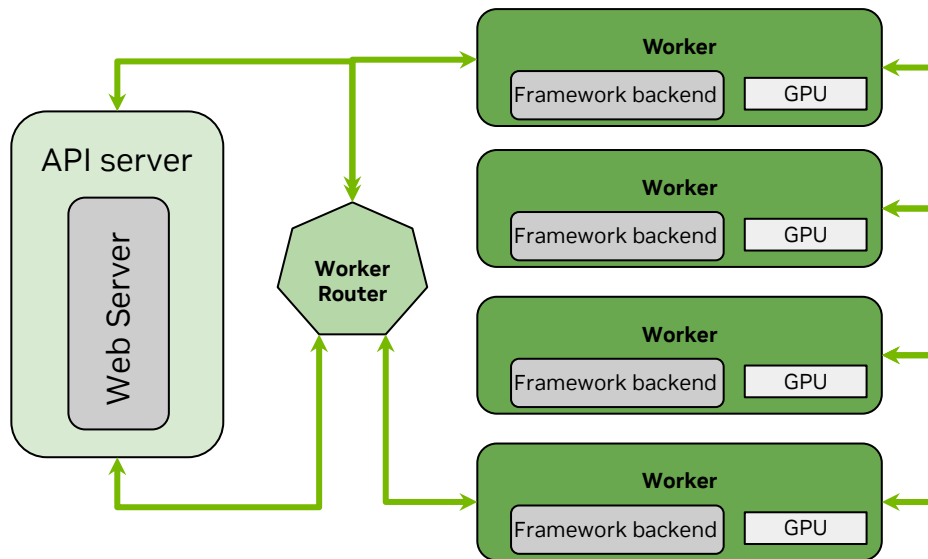
Optimized for Data Center Scale Distributed Workflows

Triton 2.0



- **Optimized for Single Node Performance**
- **Optimized for Web Service Dev Cycle: Develop, Optimize, Deploy**

Triton 3.0

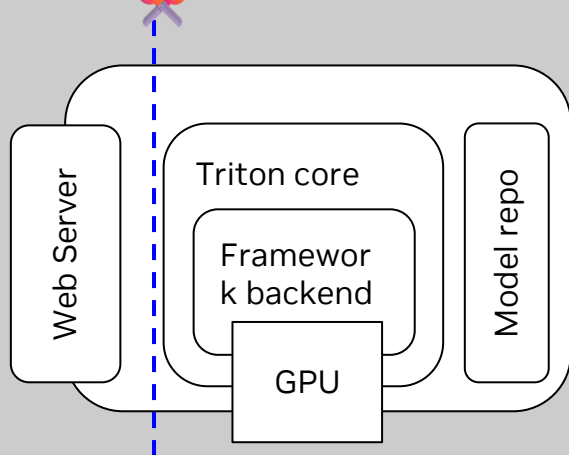


- **Optimized for Gen AI**
- **Optimized for Data Center Scalability**
- **Optimized for Distributed Workflows**
- **Optimized for Quick Iteration**

Triton 3.0 Component Architecture

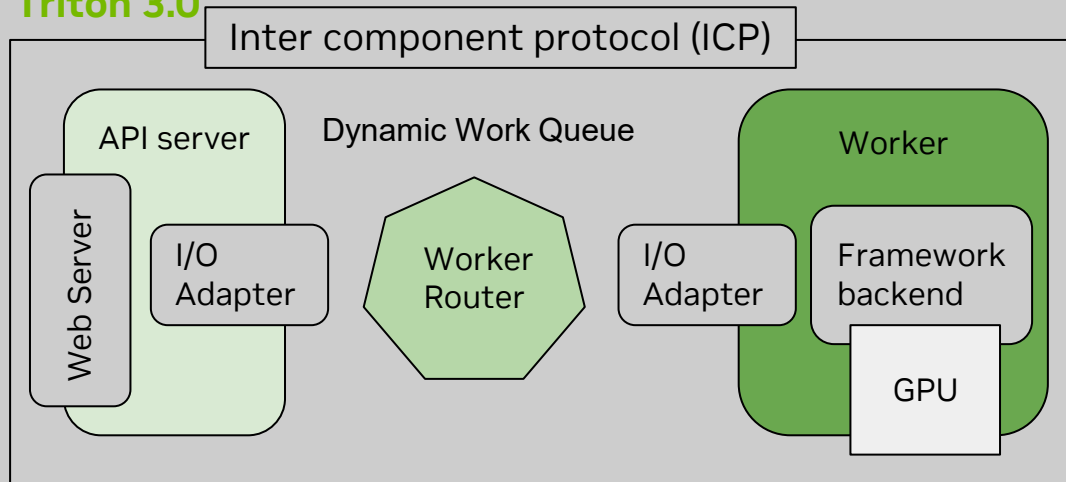
Split Triton Into Reusable Components

Triton 2.0



- Just enabled K8 multi-node autoscaling
- But needs new foundation to improve flexible scalability

Triton 3.0



Interchangeable
API Server (NIM,
RayServe, etc)

Coordinate
data to workers
w.r.t capacity

Create/destroy on demand
Scale up for throughput
Single model or ensemble.

Workers Pull Work Based on Capacity

The diagram illustrates a multi-stage pipeline with capacity-based load balancing. It consists of three stages: Stage 1, Stage 2, and Stage 3. Stage 1 has two workers, each with a GPU. Stage 2 has three workers, each with a GPU. Stage 3 has one worker with a GPU. Workers are connected to a 'Dynamic Work Queue' and a 'Worker Router' in each stage. The diagram shows how work is pulled based on capacity, with arrows indicating the flow of data and work between stages and workers.

Manifold

- Redirects downstream requests with upstream payloads

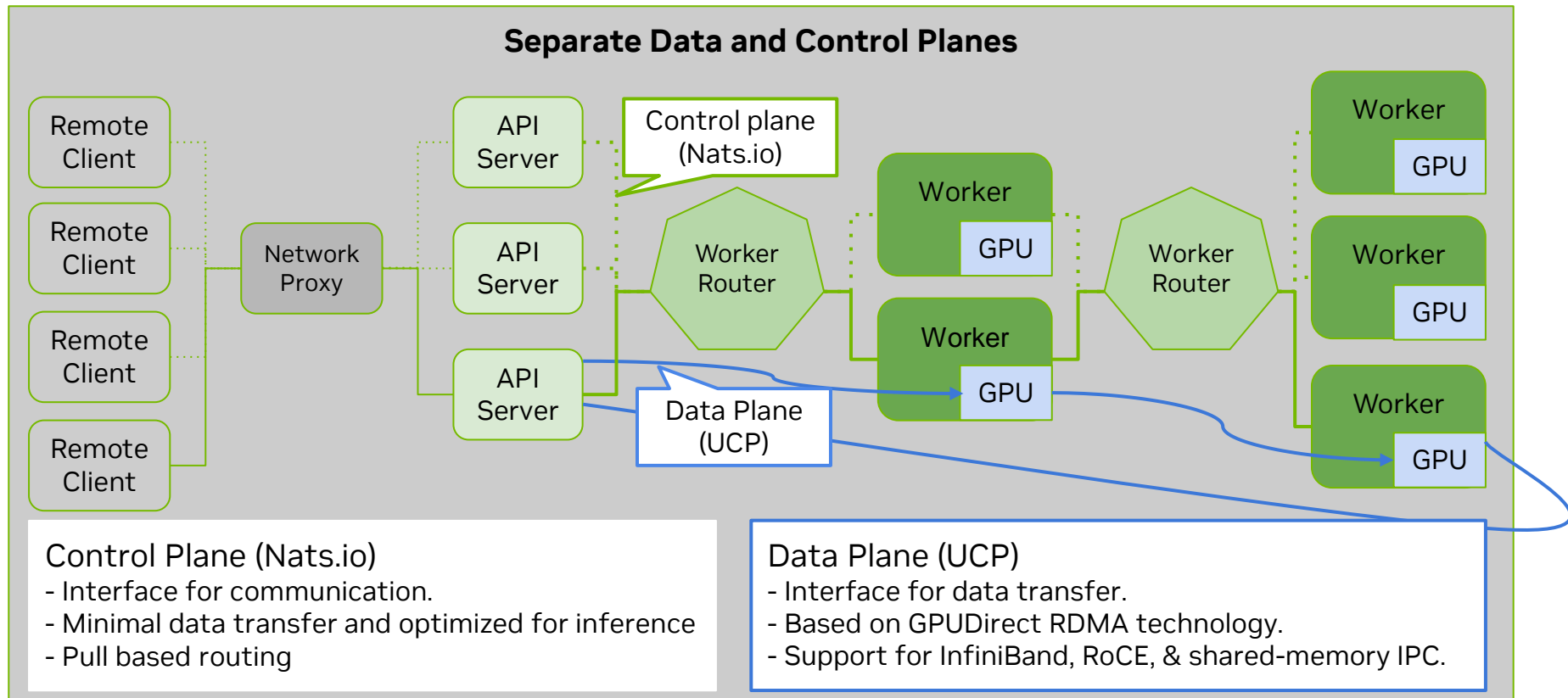
Worker

- Microservice assigned separate compute resources.

Multiple stage pipeline

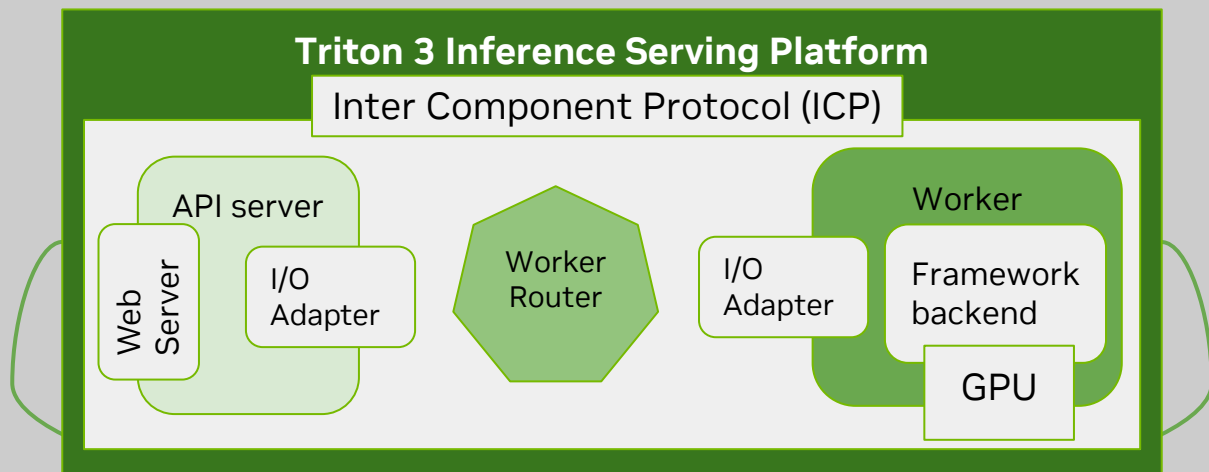
- Stages can be scaled on demand individually
- **Capacity based load balancing** across workers & stages.

Triton 3.0 Inter Component Protocol



Triton 3.0 Component Reuse

Components can be adopted in **modular** fashion and will have first class **Python support**.



Domain Specific APIs:

In addition to inference APIs, include LlamaStack, SGLang, LangGraph, etc

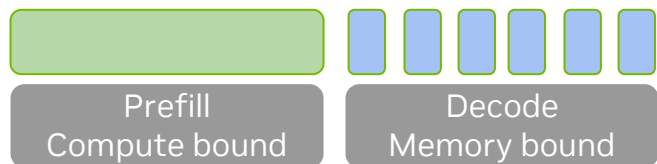
KV cache Transfer & Model Management

- Disaggregated serving
- Domain specific KV cache routing

Perf Benchmarking & Tuning

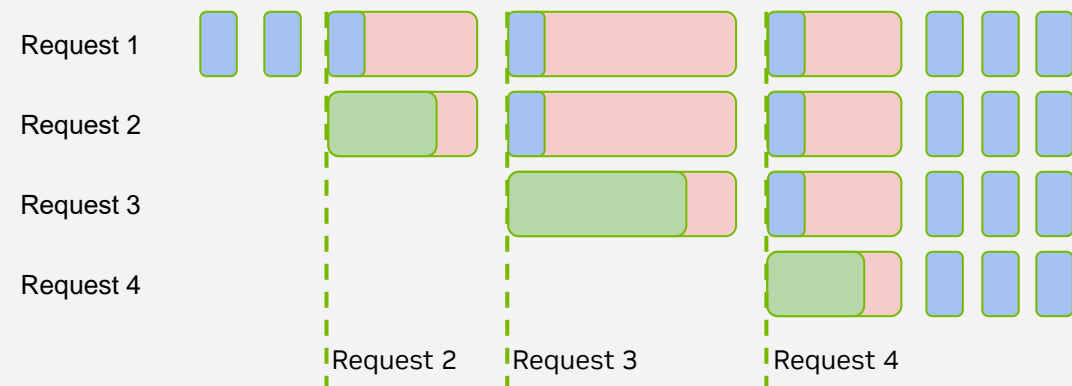
- GenAI Perf & Model navigator
- Worker to worker traffic

Use Case: Disaggregated Serving



Colocating prefill & decoding causes interference

4 requests processed in single GPU with continuous batching



Colocating prefill & decoding couples resource allocation and parallelism

- Different optimal strategies for prefill & decode
- **Prefill:** tensor parallelism for low latency
- **Decode:** data or pipeline parallelism for high throughput



Disaggregate prefill and decode

Leverage high BW NVLink to minimize the KV cache transfer overhead

Use Case: Disaggregated Serving

Dynamic Disaggregated Serving

Dynamically change prefill and decode based on ISL/OSL or TTFT/ITL requirements

Worker 1
Prefill



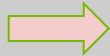
Worker 1
Prefill

Worker 2
Prefill



Worker 2
Prefill

Worker 3
Decode



Worker 3
Decode

Worker 4
Decode



Worker 4
Prefill

Disaggregated Serving on Mixed SKUs

- Simulated results on H100 (prefill) and H20 (decode) show cost benefits

Disaggregated Mixture of Experts (1.8B)

- Can make bigger models more performant

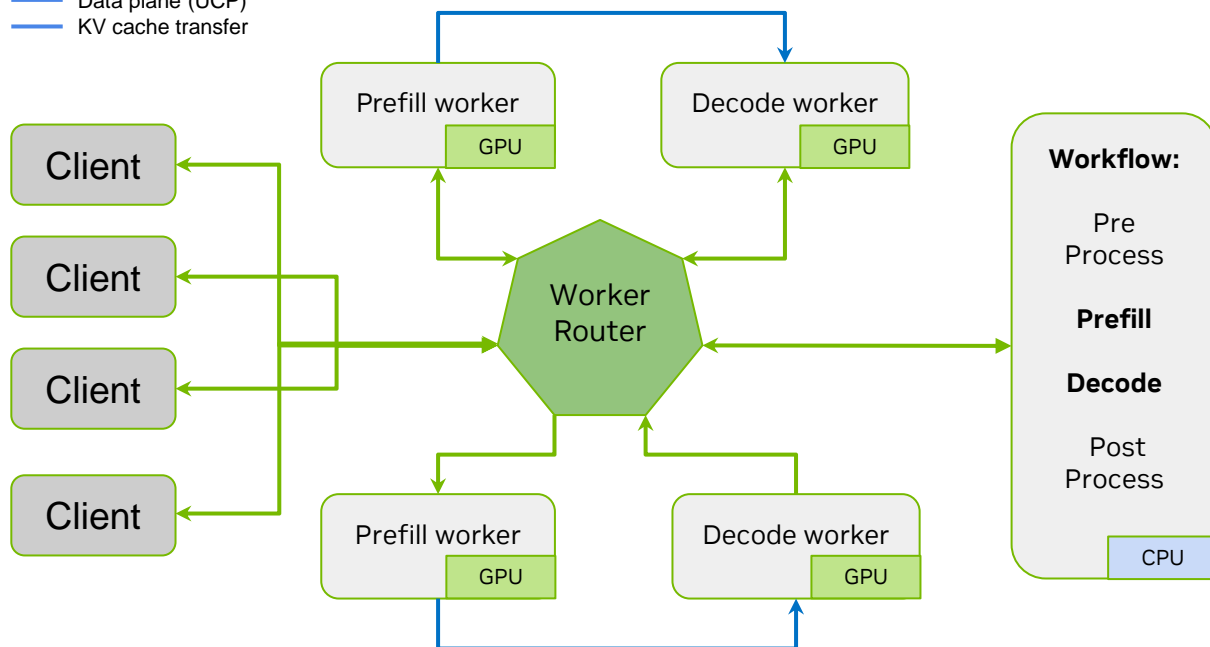
NVLink/NVSwitch => High BW KV cache transfer across workers & nodes

- Enables easier programmability for placing disaggregation.

Better utilization of one gigantic GPU system (Blackwell)

Disaggregated Serving Prototype

↔ Control plane
— Data plane (UCP)
— KV cache transfer



TRT-LLM

- Llama 3 8b Instruct
- TP 1
- 4x A100 GPUs
- Pull Based Routing

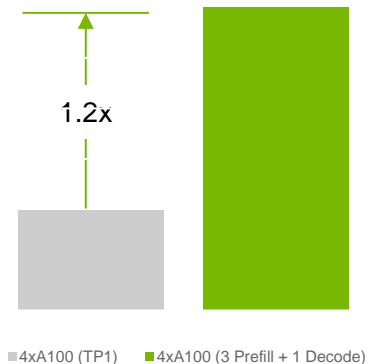
vLLM

- Llama 3.1 70B Instruct
- TP 2
- 8x H200 GPUs
- Pull Based Routing

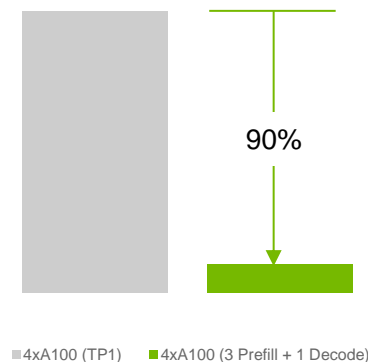
Disaggregated Serving Prototype - Triton 3 with TRT-LLM backend

Llama3 8B, ISL/OSL: 4096/512, TRT-LLM 0.12, Concurrency:10

Output Token Throughput
(Higher is Better)

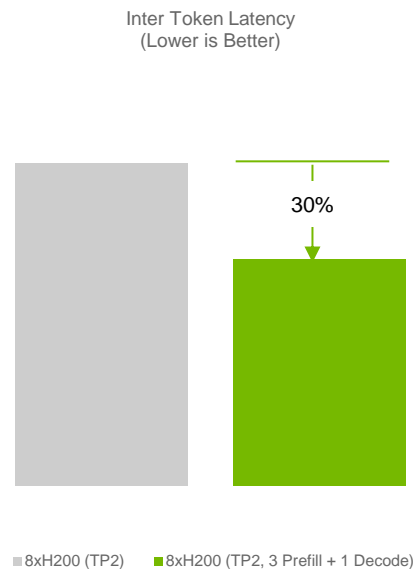
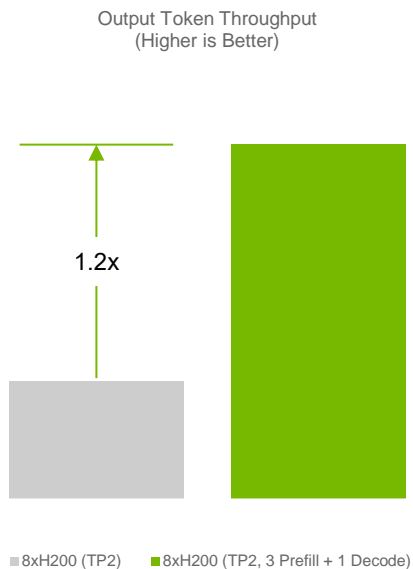


Inter Token Latency
(Lower is Better)



Disaggregated Serving Prototype - Triton 3 with vLLM Backend

Llama3.1-70B, vLLM 0.5.3, 8x H200, ISL/OSL: 2K/128, Concurrency:10



The background features a series of overlapping, wavy, light green bands that create a sense of depth and movement, sloping upwards from left to right. A solid, vibrant green vertical bar is positioned along the left edge of the frame.

Thank You

Getting started with NVIDIA AI Inference

Bring inference to production with performance, ease, and cost savings



[NVIDIA Triton Product Page](#)

[NVIDIA Launch Pad](#)

[Download Triton on NGC](#)

[Explore More Resources for Development](#)

[For more information on NVIDIA Triton Inference Server](#)

[Open-source GitHub repository:](#)

[Latest release information](#)

[Quick start guide:](#)

NVIDIA Deep Learning Institute. Self-paced online course. 4 hours.
[Deploying a Model for Inference at Production Scale](#)