# Sizing LLM Inference Systems

Notebook 1: Understanding Batching Strategies
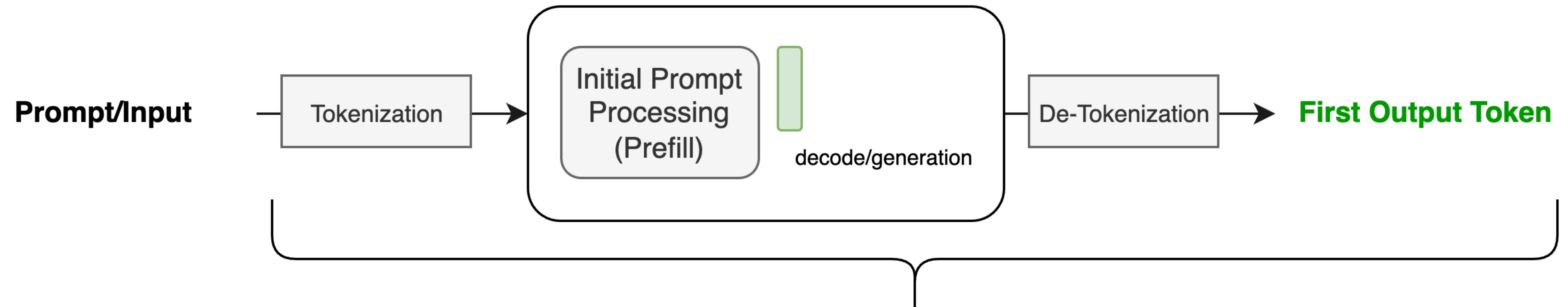
# Metrics

# LLM Inference Measurement – Time to First Token

The time it takes to process the input and generate the first token



Time to First Token (TTFT)

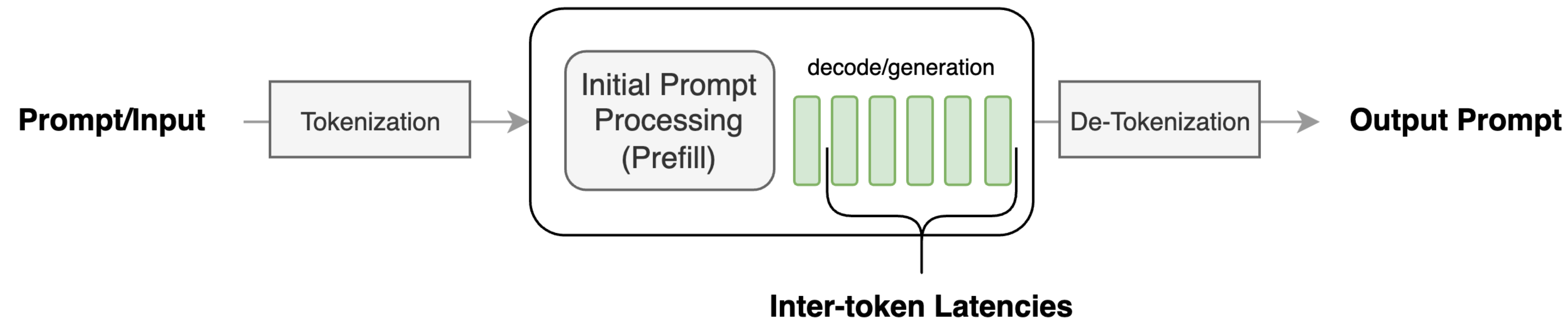Signifies time it takes to process the prompt and generate the first token
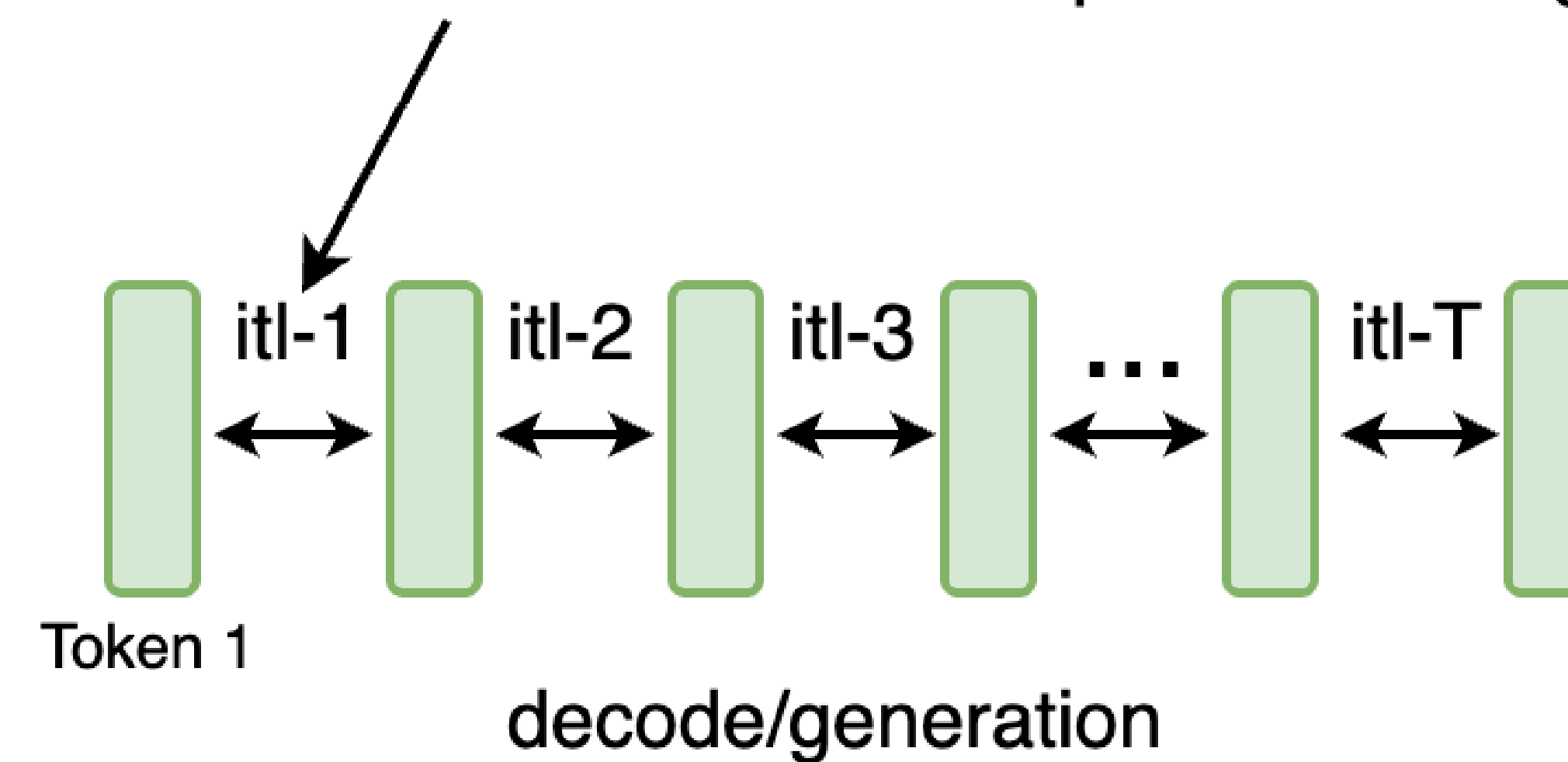
Synonyms:
- First Token Latency (FTL)

# LLM Inference Measurement – Inter-token Latency

Measuring generation performance when the system is under load



Inter-token Latencies

**Inter-token Latencies** - Should be near constant over all token generations.
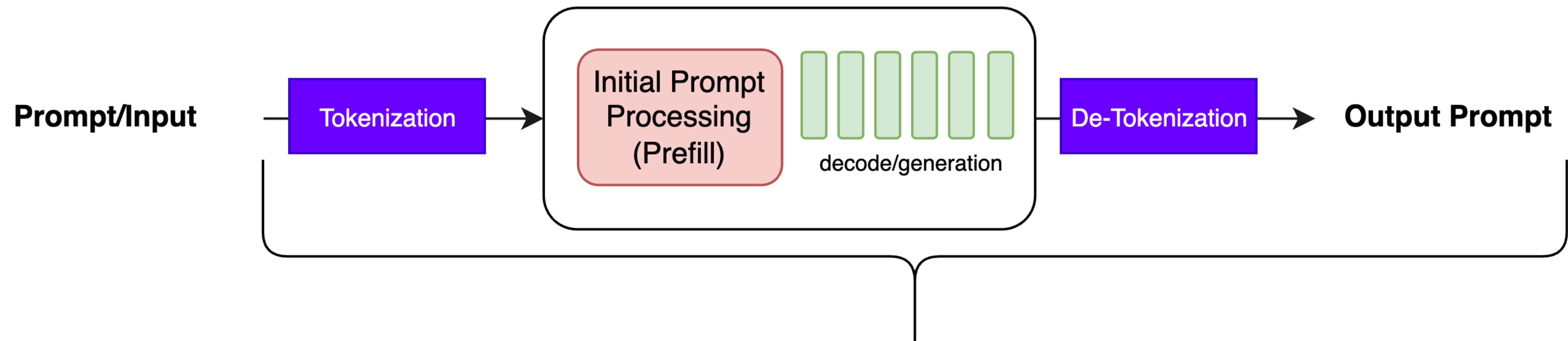
Near constant time implies efficient generation

Synonyms:
- Time Per Output Token (TPOT)

# LLM Inference Measurement – Total Time to Generation

The time it takes to process a prompt and generate all the tokens



**Total Time to Generation**

Measures to total time to process input and generate all tokens
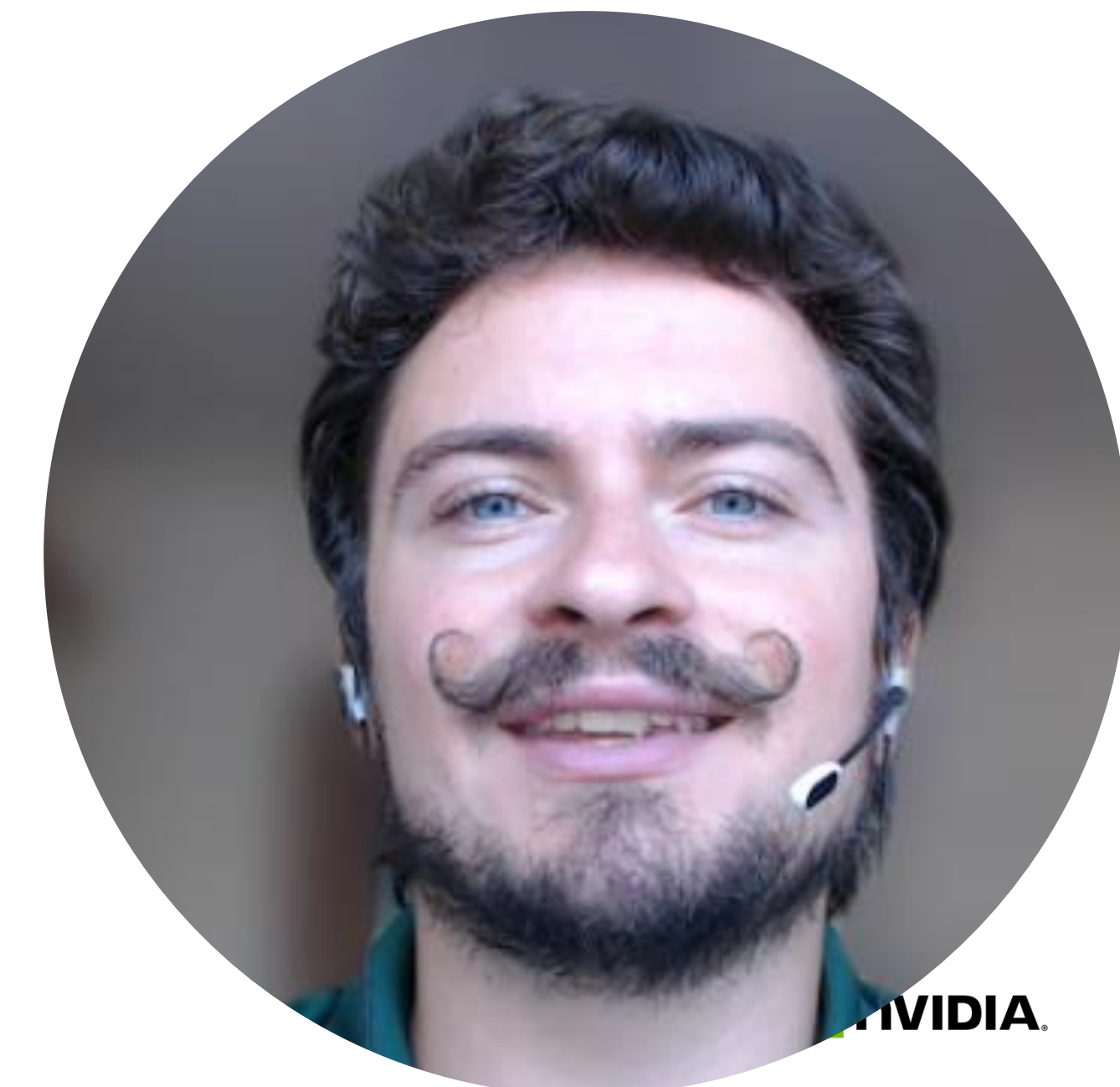
Synonyms:
- End-to-End Latency (E2E latency)
- Time to last Token

# Throughput Metrics

- Always specify
  - Model
  - Precision
  - Input Length
  - Output Length
  - Concurrency
  - TP

- The most unambiguous metric
  - to measure is **requests/second/instance**
  - to use in sizing **requests/second/GPU**
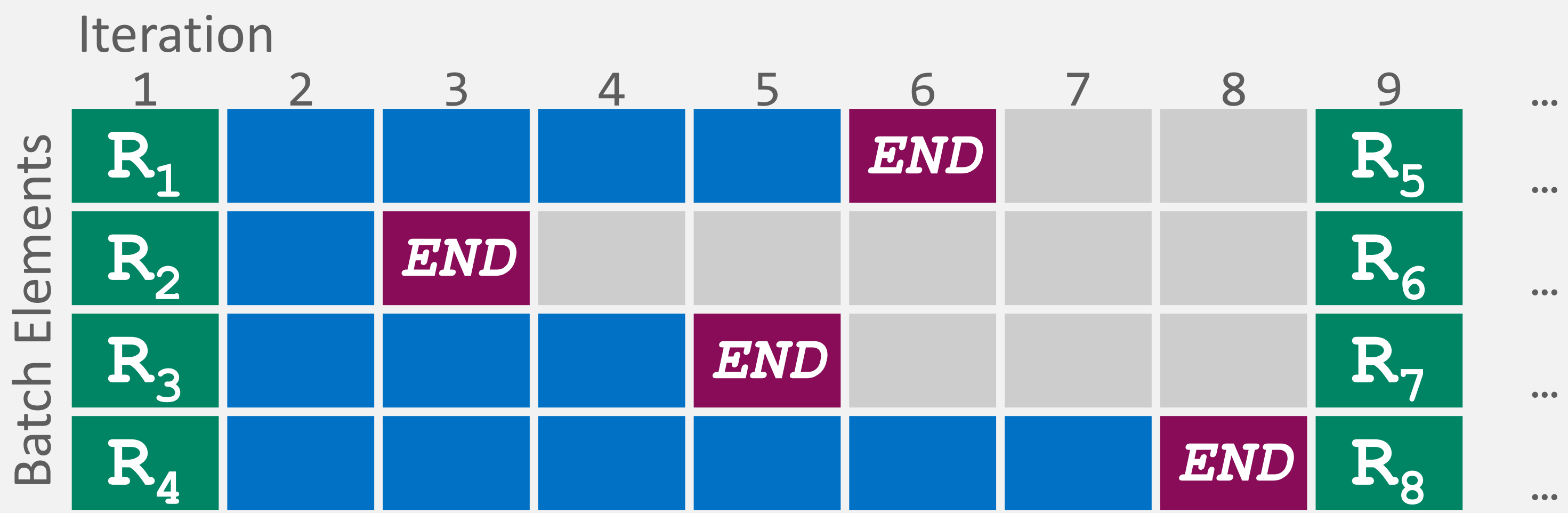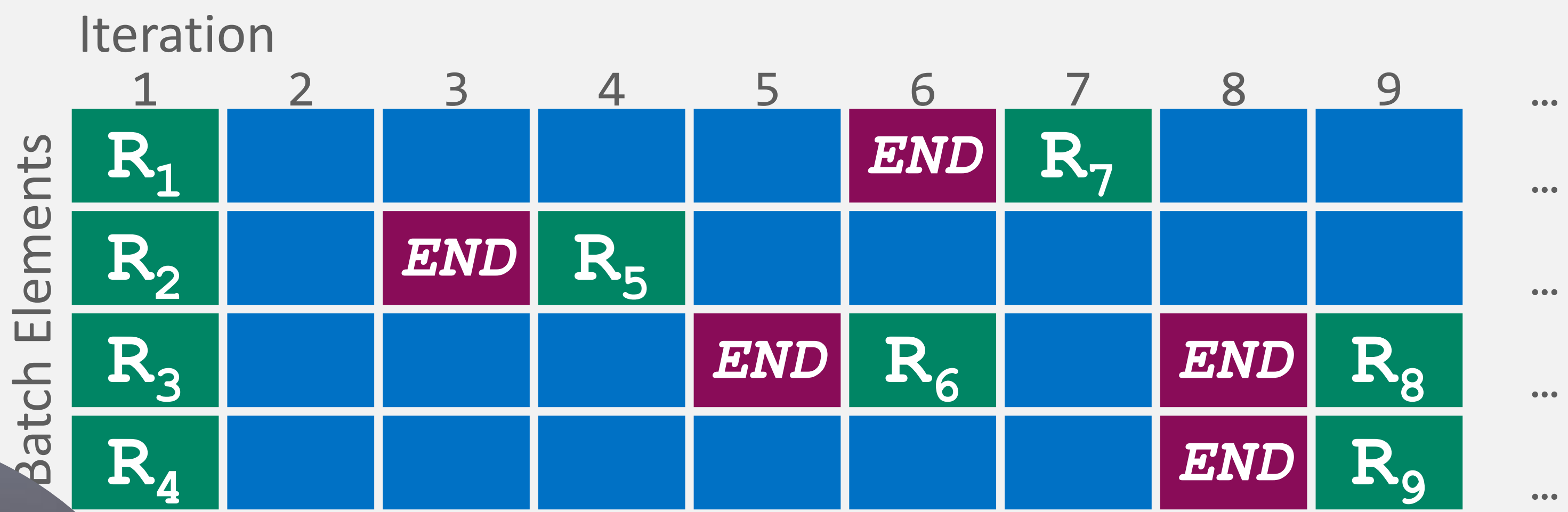
# Inference Optimizations

# Inflight Batching

## Maximing 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 utilizaiton

- Integrated directly into the TensorRT-LLM Triton backend
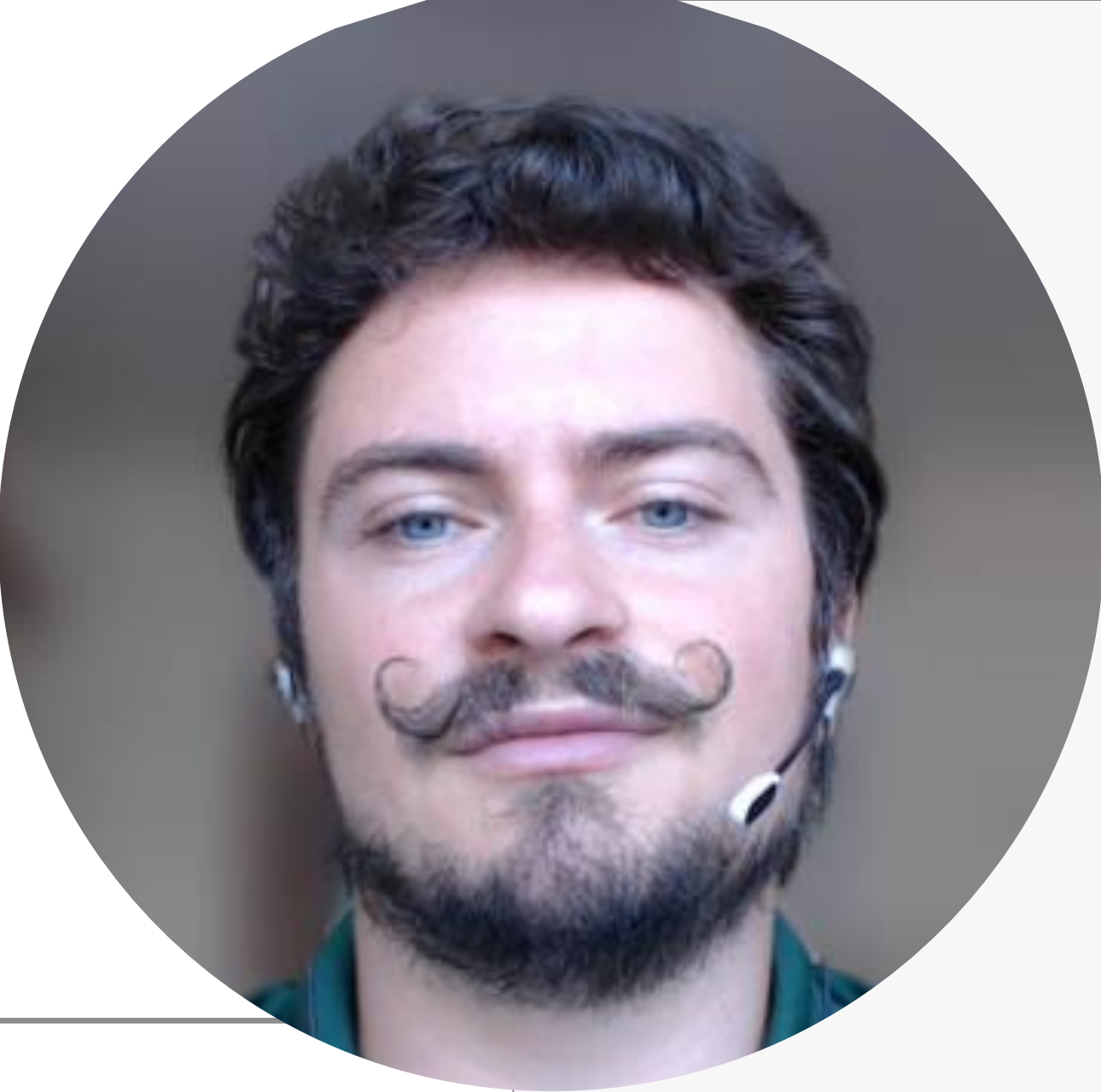
- Accessible though the TensorRT-LLM Batch Manager



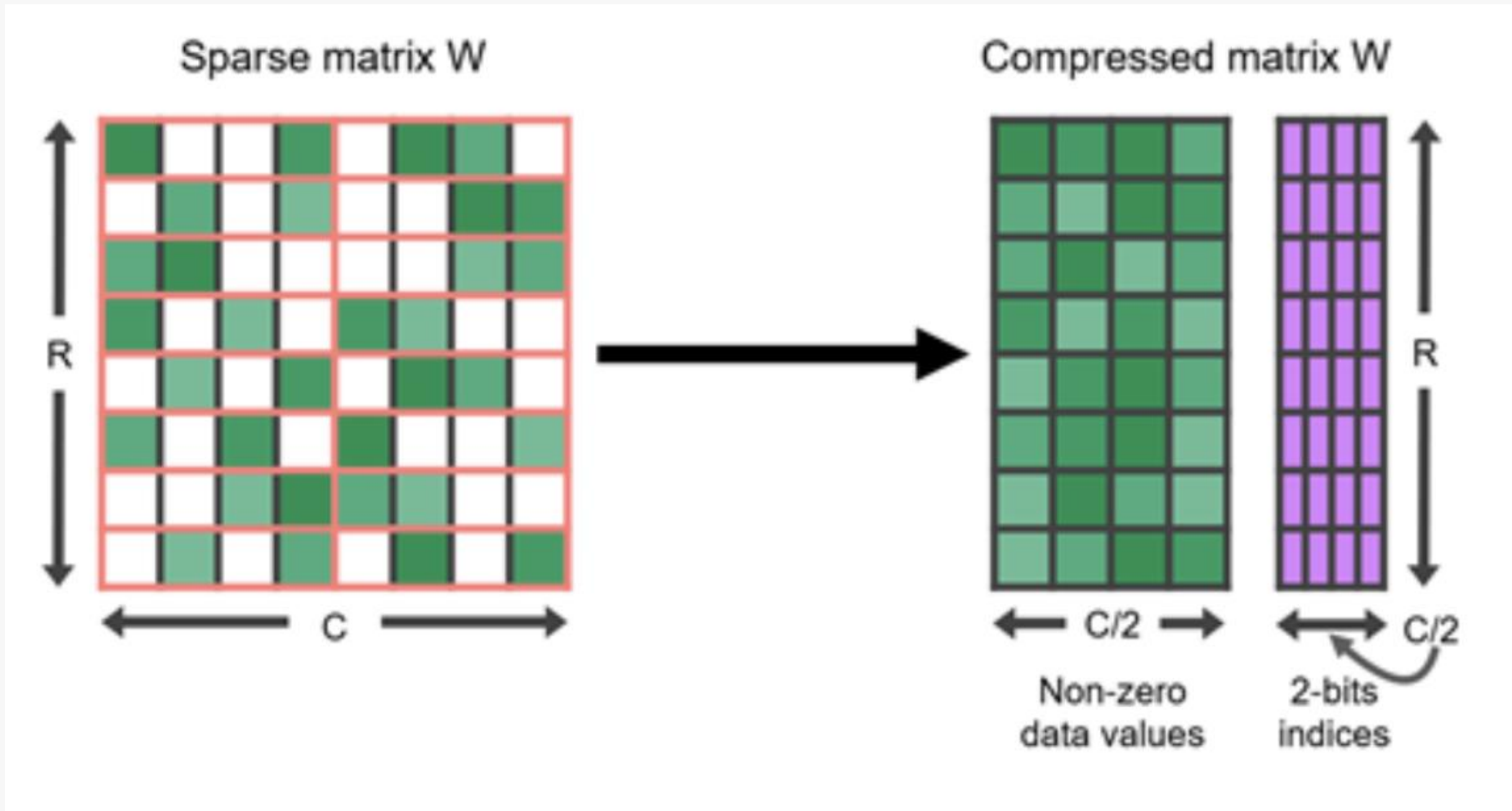Static Batching



**Inflight Batching**

Context | Gen | EoS | NoOp

# Model Optimizer - Sparsity



A 2:4 structured sparse matrix W, and its compressed representation

Model Optimizer offers Sparsity (with fine-tuning or post-training sparsity) to reduce the memory footprint and accelerate inference.

It supports NVIDIA 2:4 Sparsity sparsity pattern and various sparsification methods, such as (NVIDIA ASP) and (SparseGPT).

- **MLPerf-Inference v4.0 Results** (blog)

    - Uses 2:4 sparsity on a Llama2-70B.

    - Post-training sparsification (PTS) with SparseGPT yields no accuracy drop.

    - PTS reduces model size by 37%, facilitating deployment on a single H100 SXM with TP=1, PP=1, and achieving a 1.3x speedup.

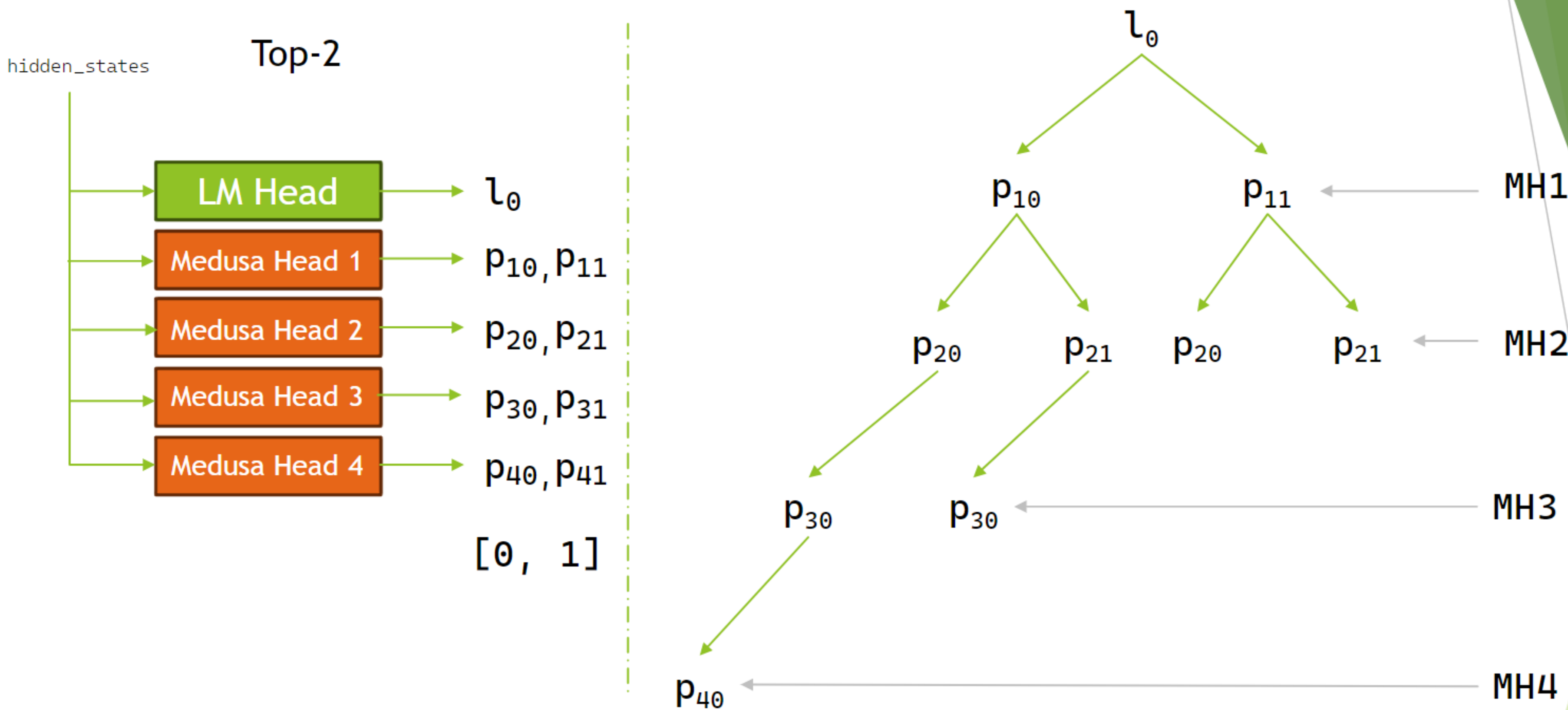| Model | Batch_size | Inference Speedup (Compared to the FP8 dense model with the same batch size) |
|---|---|---|
| Sparsified Llama 2-70B | 32 | 1.62x |
| | 64 | 1.52x |
| | 128 | 1.35x |
| | 896 (MLPerf) | 1.30x |

Performance improvement with sparsity

# Speculative Decoding

- Generate of more than one token per forward pass iteration
- Helpful if only if the GPU is underutilized due to small batch sizes.
- Perform more computations, keeping memory access almost the same
- Supports: separate draft model, Medusa, ReDrafter

hidden_states

Top-2

| LM Head | $\to$ | $l_0$ |
| Medusa Head 1 | $\to$ | $p_{10}, p_{11}$ |
| Medusa Head 2 | $\to$ | $p_{20}, p_{21}$ |
| Medusa Head 3 | $\to$ | $p_{30}, p_{31}$ |
| Medusa Head 4 | $\to$ | $p_{40}, p_{41}$ |

[0, 1]

$l_0$

$p_{10}$     $p_{11}$     MH1

$p_{20}$     $p_{21}$     $p_{20}$     $p_{21}$     MH2

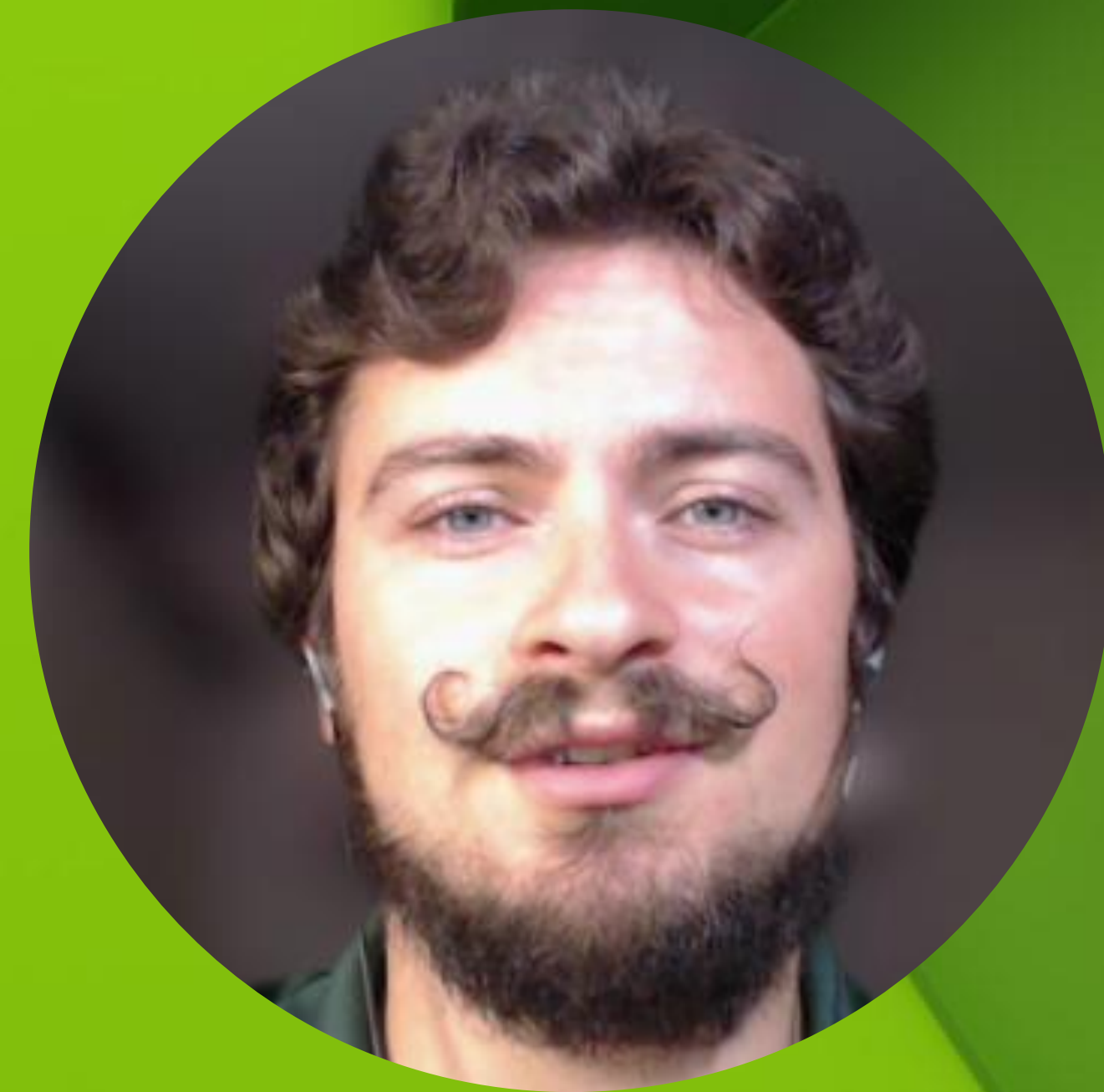$p_{30}$     $p_{30}$     MH3

$p_{40}$     MH4

# of paths: 4 (16 if all MHs have Top-2)
Paths: $l_0 p_{10} p_{20} p_{30} p_{40}$, $l_0 p_{10} p_{21} p_{30}$, $l_0 p_{11} p_{20}$, $l_0 p_{11} p_{21}$

# of candidates: 10 ($l_0$, $l_0 p_{10}$, $l_0 p_{10} p_{20}$, $l_0 p_{10} p_{20} p_{30}$, $l_0 p_{10} p_{20} p_{30} p_{40}$, $l_0 p_{10} p_{21}$, $l_0 p_{10} p_{21} p_{30}$, $l_0 p_{11}$, $l_0 p_{11} p_{20}$, $l_0 p_{11} p_{21}$)
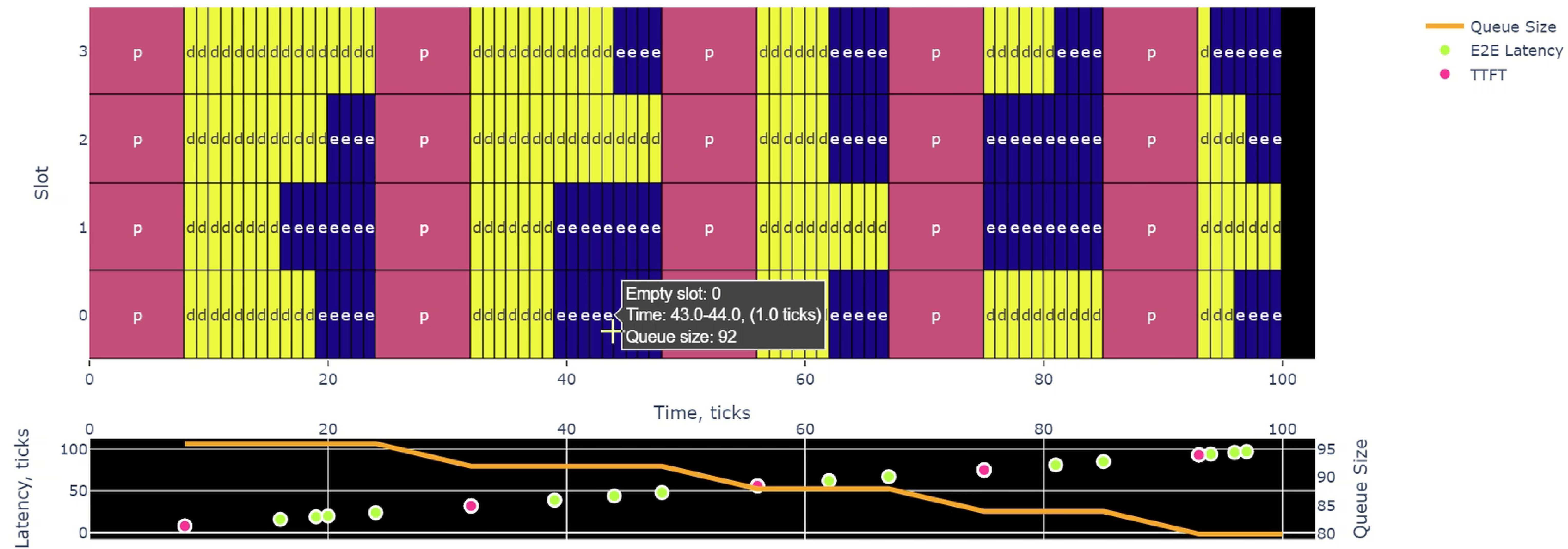
NVIDIA.

# The Simulator

# The Simulator



Prefill and Decoding Simulation

# Objectives of this notebook

1. Understand and measure time to first token (TTFT), end-to-end latency (E2E Latency), and inter-token latency (ITL).

2. Analyze throughput metrics and simulate their dependencies on various factors.

3. Explore the impact of batching and inflight batching on GPU utilization and performance.

4. Investigate the effects concurrency settings on latency and throughput.