

Portable High-Performance LLM Serving: *A Triton Backend for vLLM*

—
Jan van Lunteren
Burkhard Ringlein
IBM Research
AI Platform Group, Zurich

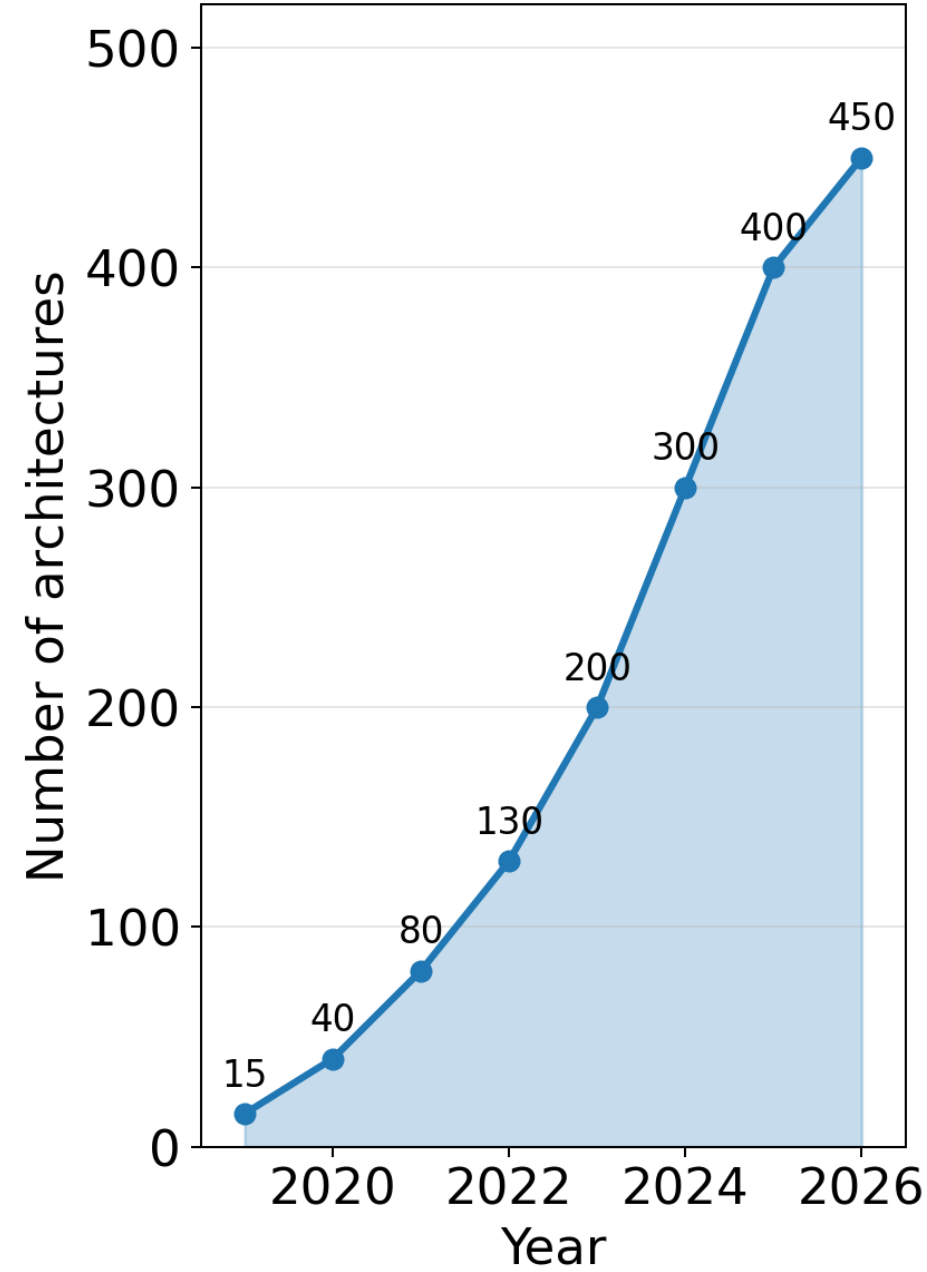
2026-04-08, PyTorch Conference, Paris



Open Source AI is growing...

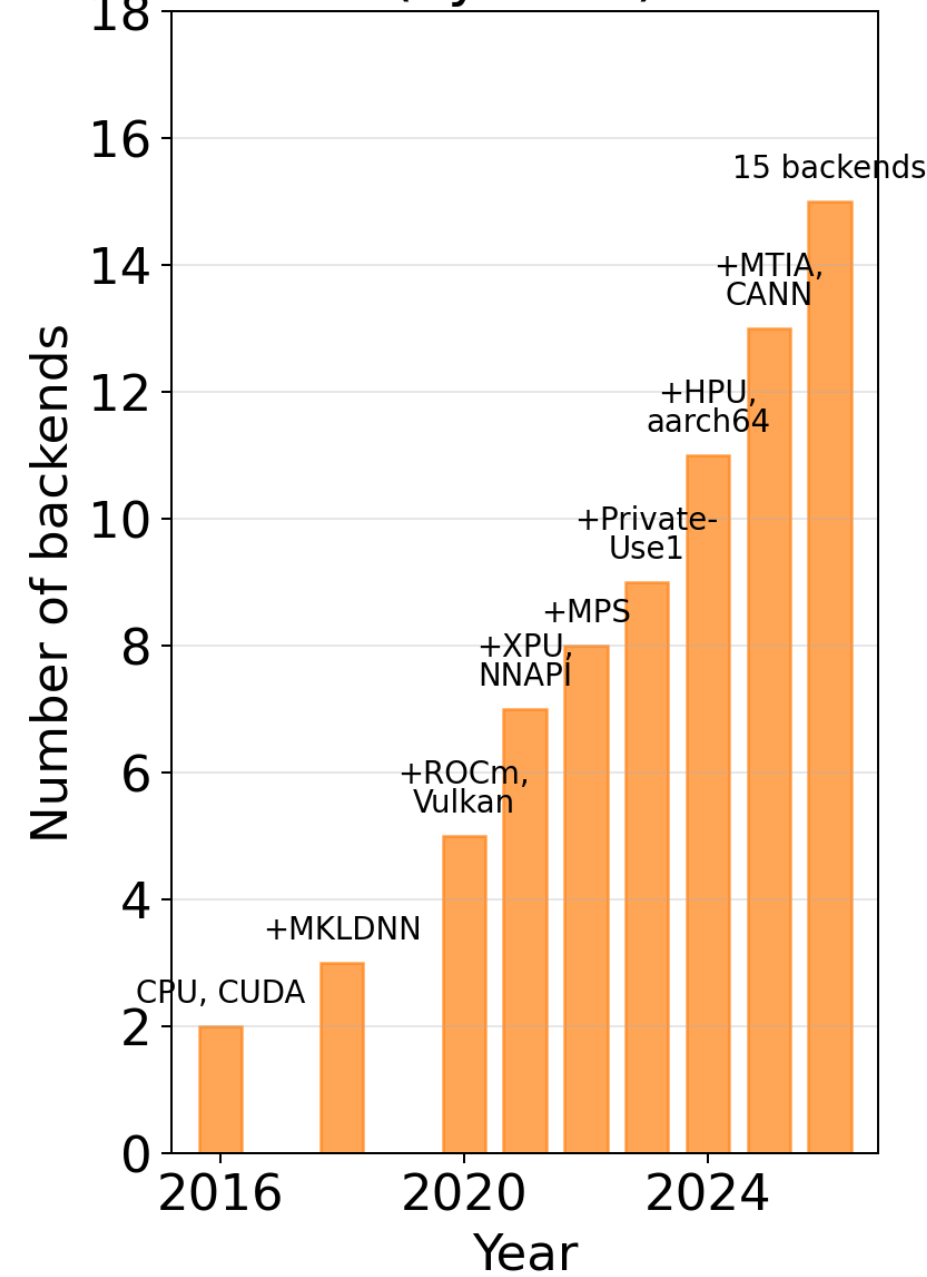
Specialized Kernel Optimizations (Open-Source LLM Inference Ecosystem)

Model Architectures
(HuggingFace Transformers)

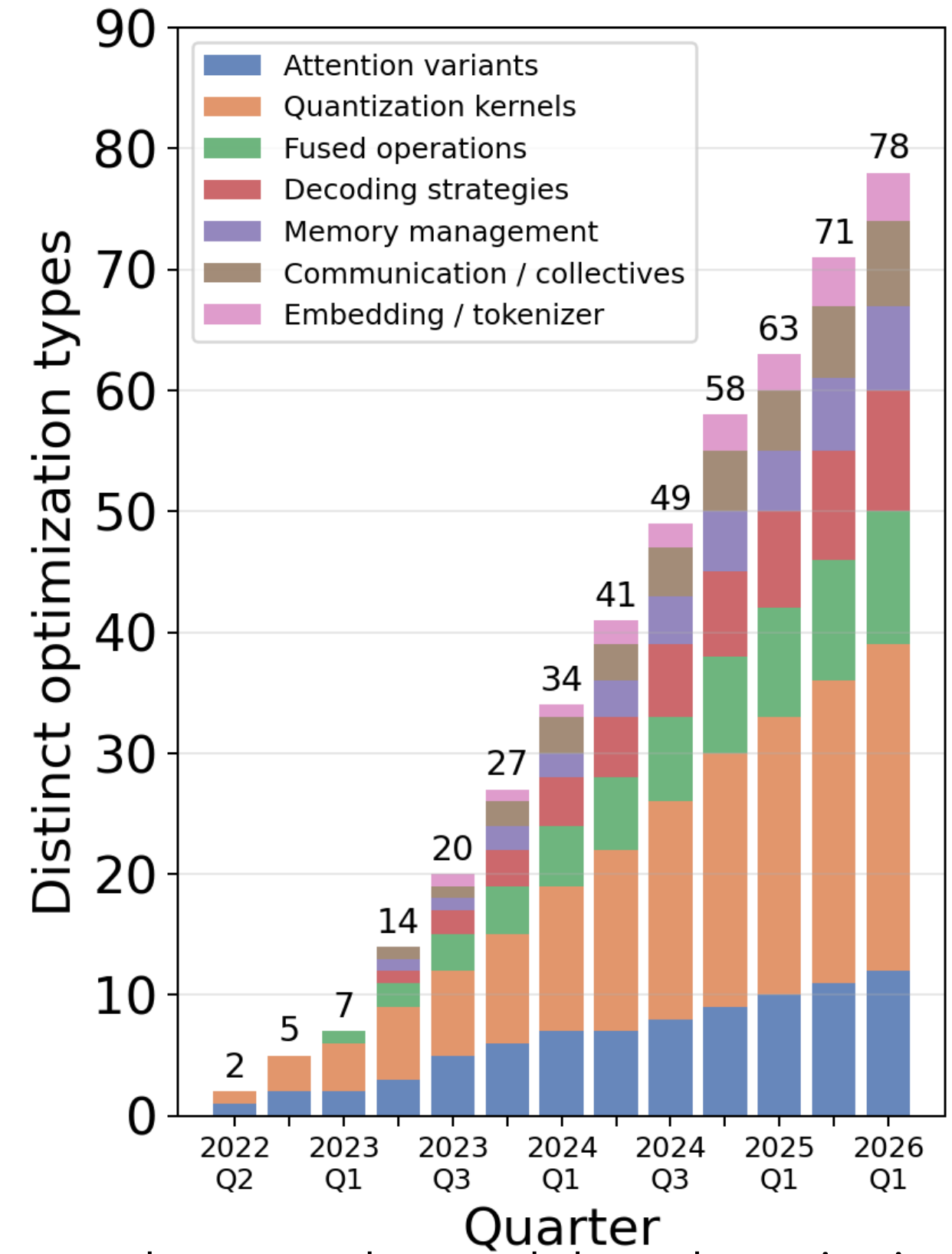


More models...

Hardware Backend Targets
(PyTorch)



...on more hardware...



...with much more kernel-level optimizations?

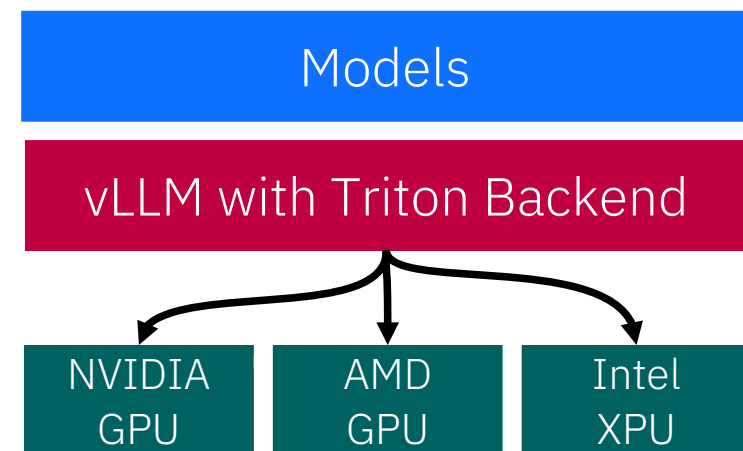
How can **LLM** keep up?

→ performance portability!?!

Agenda: How to achieve performance portability for LLM kernels

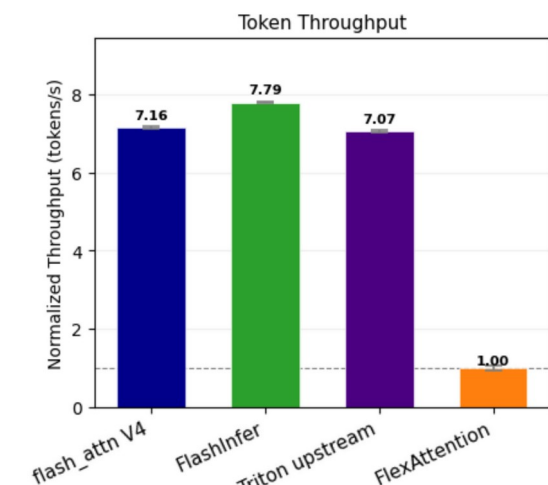
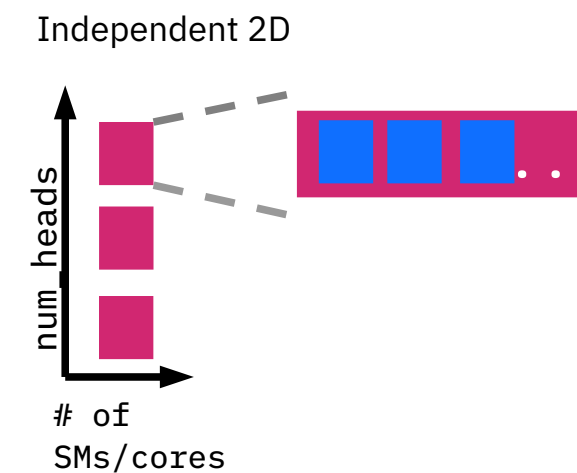
(Or, how to vLLM can keep up!)

1. Some background:
Triton & performance
portability



2. Triton-only Attention
Backend in vLLM

3. Deep Dive: How to
write portable paged
attention



4. Results

Conclusion?

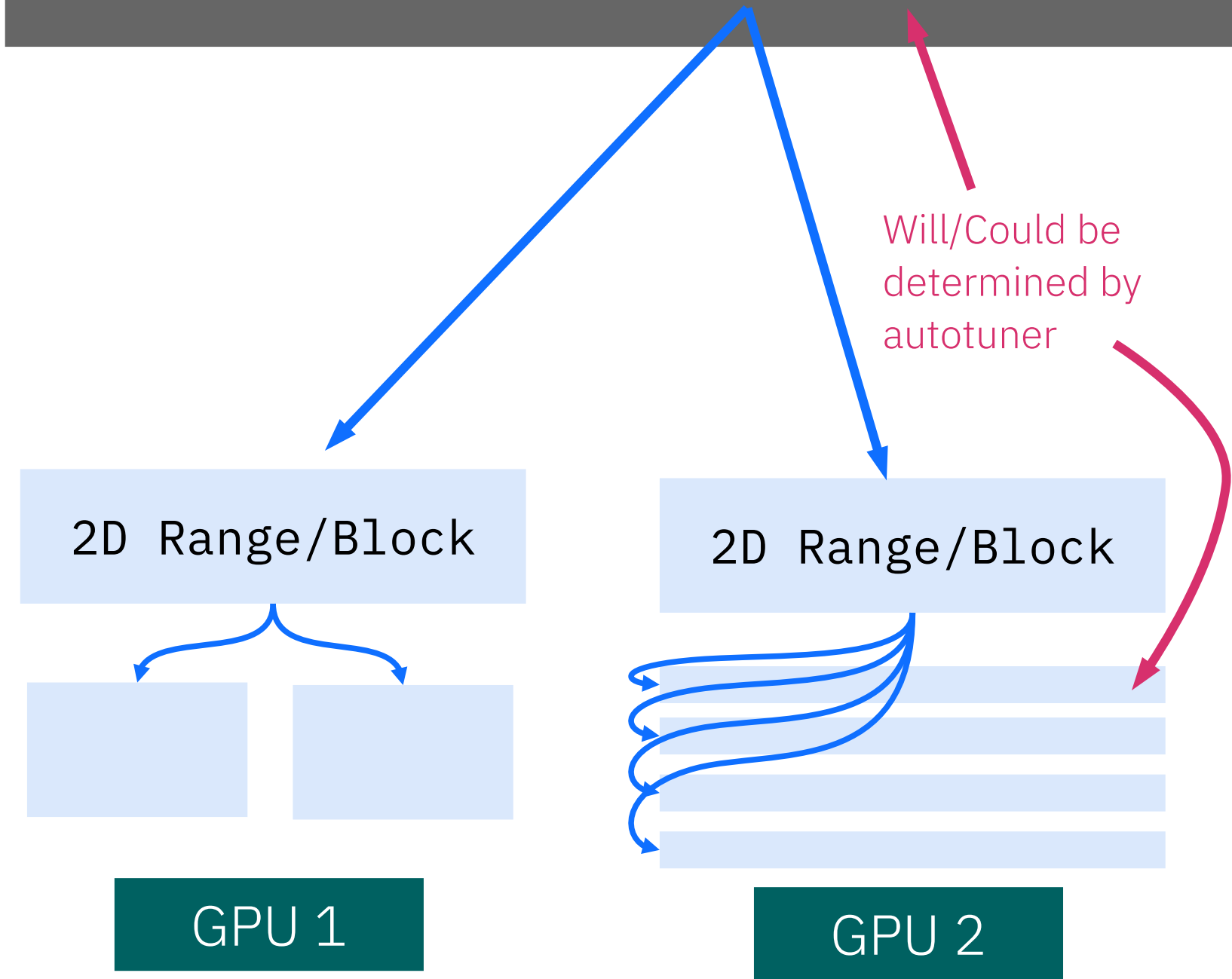
→ **Goal:** (1) Sketch the **way to performance portability** using the 100% open-source solutions
(2) ...getting you excited about the Triton Attention Backend in vLLM! 😊



Performance Portability: What is it and how does Triton help?

- “Performance portability refers to the ability of computer programs and applications to [operate effectively across different platforms.](#)” (wikipedia.org)
- In a Nutshell: Tritons [tiling programming model](#) & abstraction level:
 - Allows to express complicated optimizations (*we’ll come to that...*)
 - While staying hardware agnostic
- Hyperparameters (called “configurations”) allow simple but effective adaptations for different hardware (maybe also with **autotuning**)

```
# triton  
my_tile = index + tl.arange(BLOCK_SIZE_N)[: , None]  
               + tl.arange(BLOCK_SIZE_M)[None, :]
```



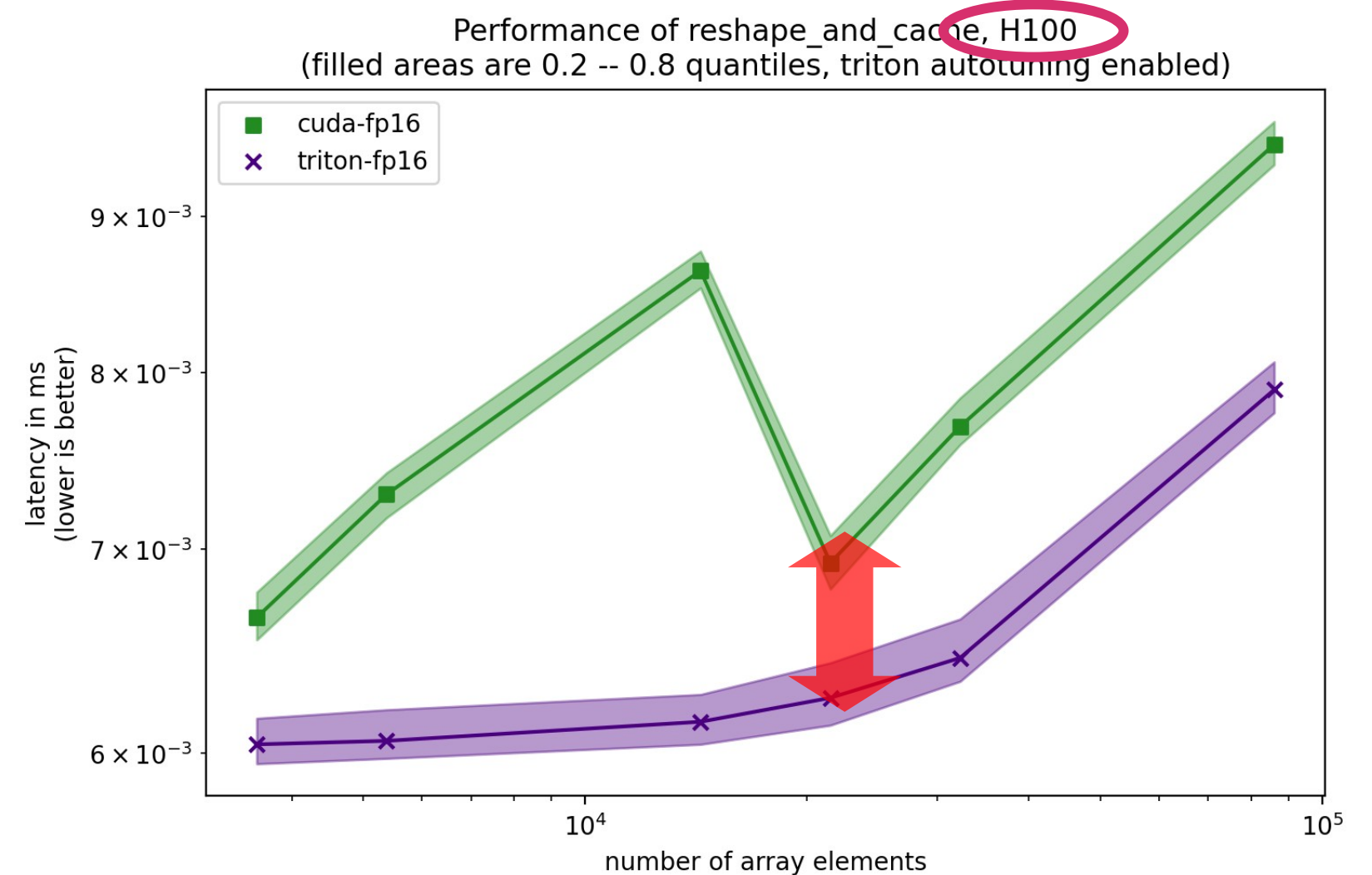
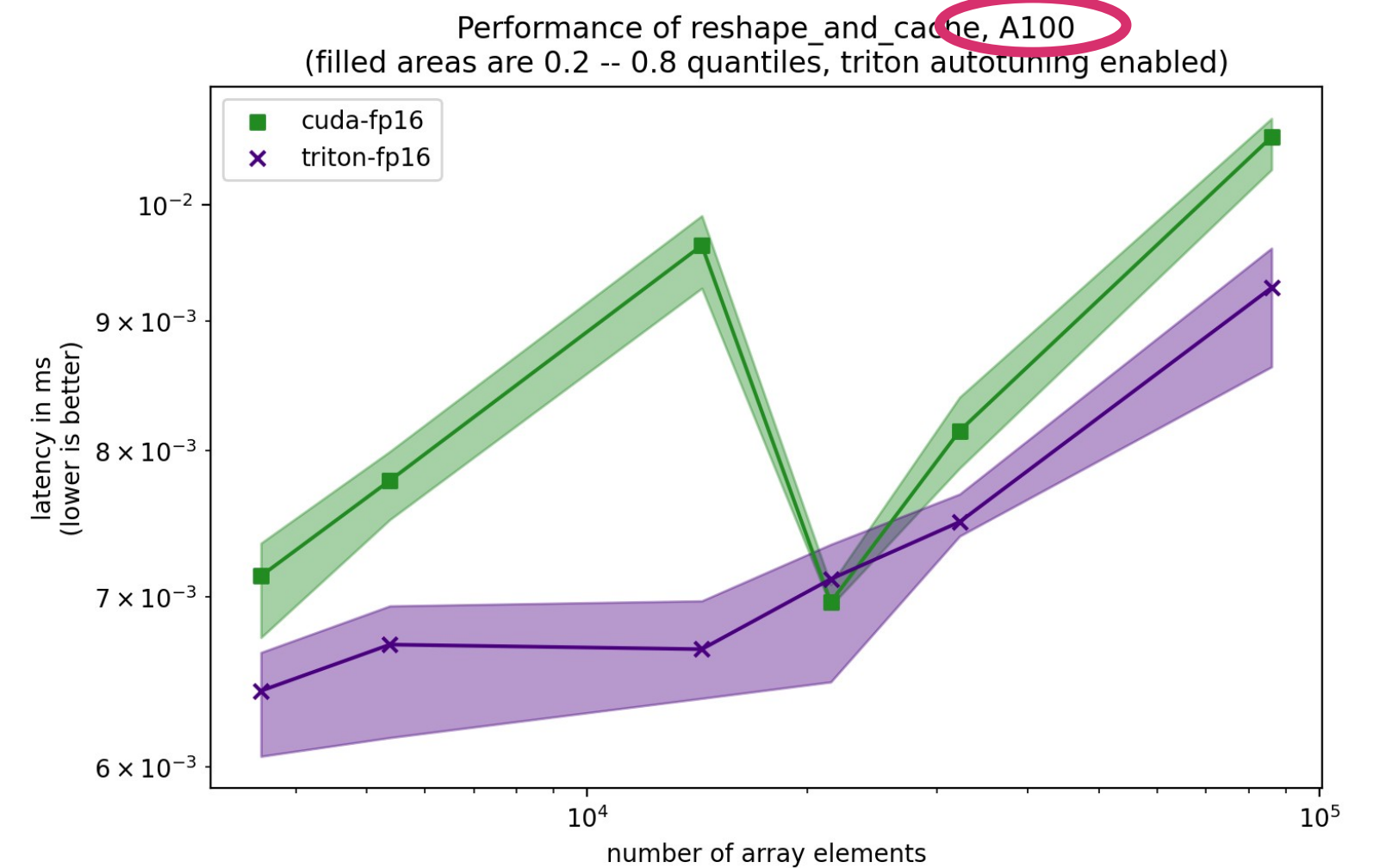
Limited performance portability? → one example

- Reshape and cache: simple vLLM kernel to populate the paged cache after linear layers
- Even within GPUs of the same vendor, CUDA kernels cannot achieve peak performance, without adaption

– Here: Triton outperformed CUDA by 16% (due to inferior blocking)

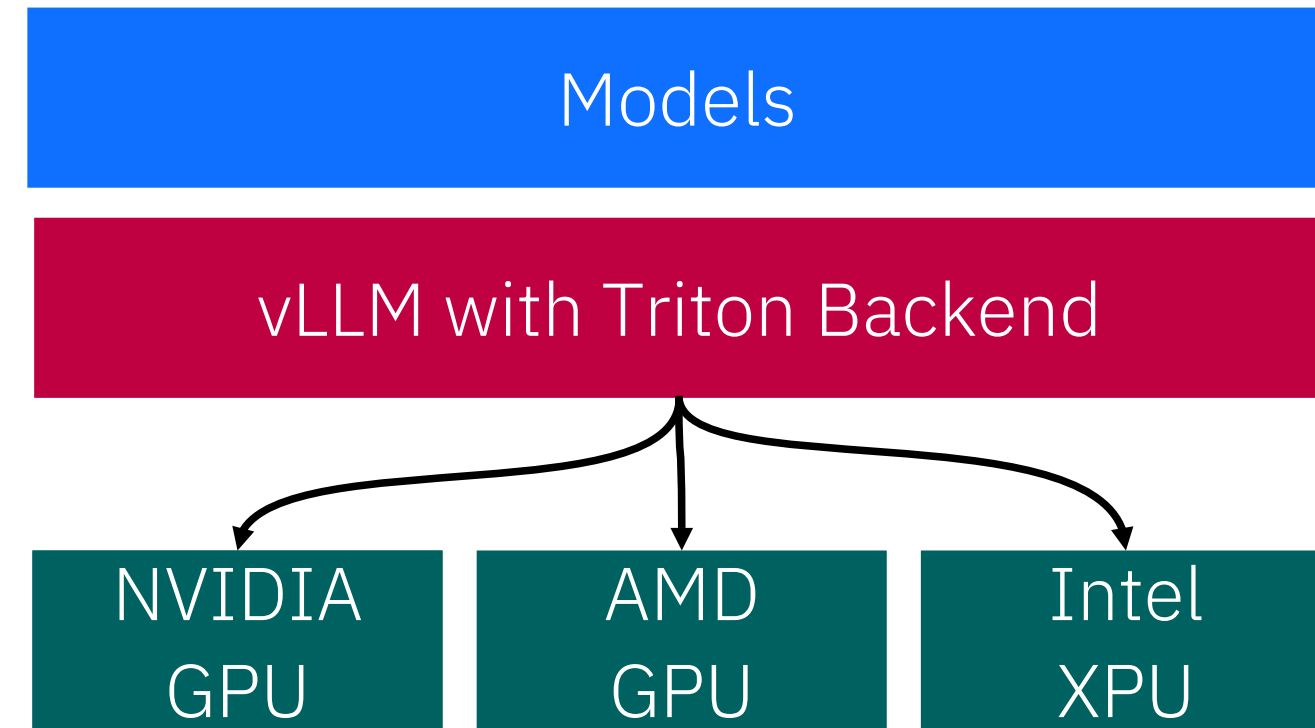
→ Because writing high-performing code requires knowledge of the hardware (e.g. tensor core size, cache size, access patterns)

- ▶ Hardware changes → performance of same low-level code degrades
- ▶ Besides Just-in-Time compilation and autotuning can take more runtime information into account



Introducing: `vllm-triton-backend`

- We (IBM Research & RedHat) built a production-ready vLLM attention backend written entirely in Triton
 - Contains multiple kernels for prefill & decoding (paged attention)
 - Packaged together with heuristics to adapt for different scenarios
- (a Backend in vLLM encapsulates different attention-implementations, usually wrappers for external libraries)
- → First (and only) `triton-only backend in vLLM`
- `Runs on NVIDIA, AMD & Intel`
 - “its always there”
- `Delivering state-of-the-art performance`
 - Is the default for AMD-based deployments, or GTP-OSS on pre-Hopper GPUs, etc.
 - (PRs of latest changes still pending)
- In the meantime: *a lot of help from the community!*
(Thanks a lot!)



Installation:

```
# integrated in vllm
pip install vllm
      (no further dependencies)
```

Run:

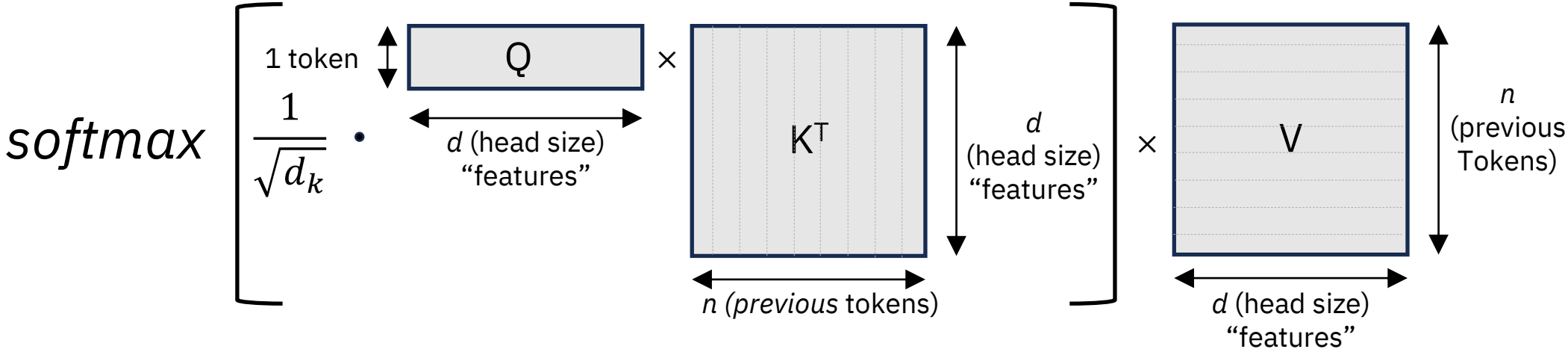
```
vllm serve ibm-granite/granite-4.0-small \
  --attention-backend=TRITON_ATTN
# the configuration only necessary where it
  is not the default
```

Attention

- Original attention formula:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

- For a single *query token*, *sequence*, and *query head* and causal masking:

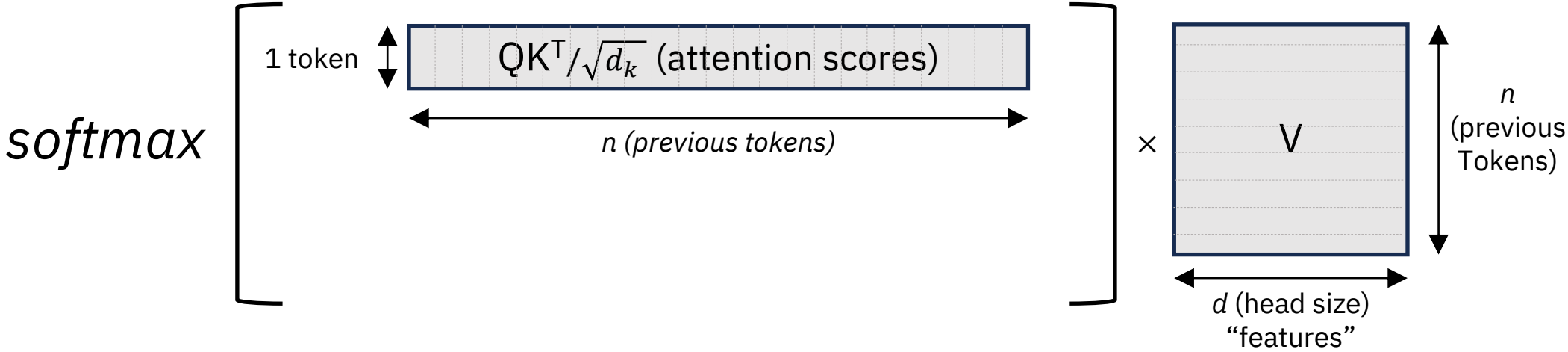


Attention

- Original attention formula:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

- For a single *query token*, *sequence*, and *query head* and causal masking:



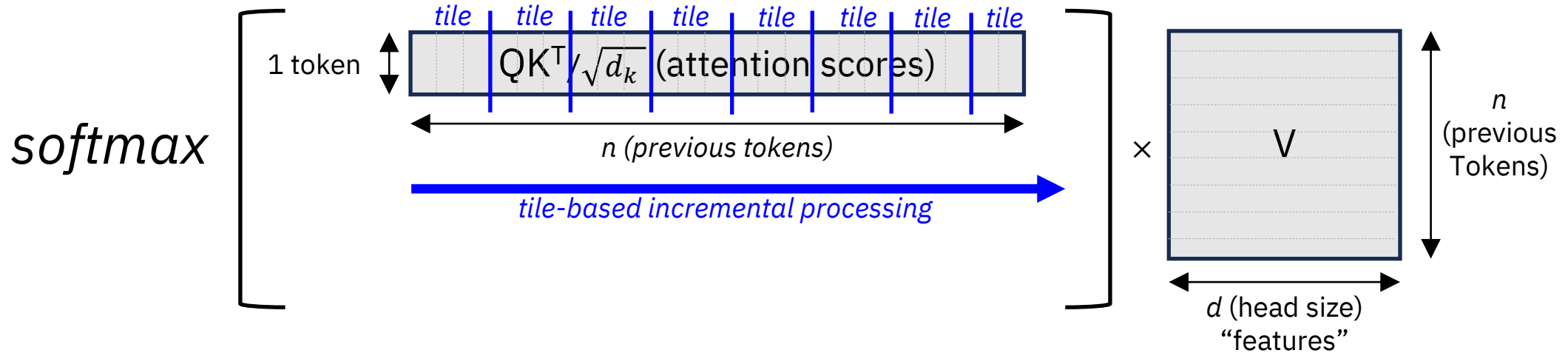
- Numerically stable softmax:

$$softmax(s_i) = \frac{e^{s_i - \max_j s_j}}{\sum_{k=1}^n e^{s_k - \max_j s_j}}$$

- Softmax can be calculated efficiently using a *tiled* approach (a.k.a. *online softmax*)

Attention

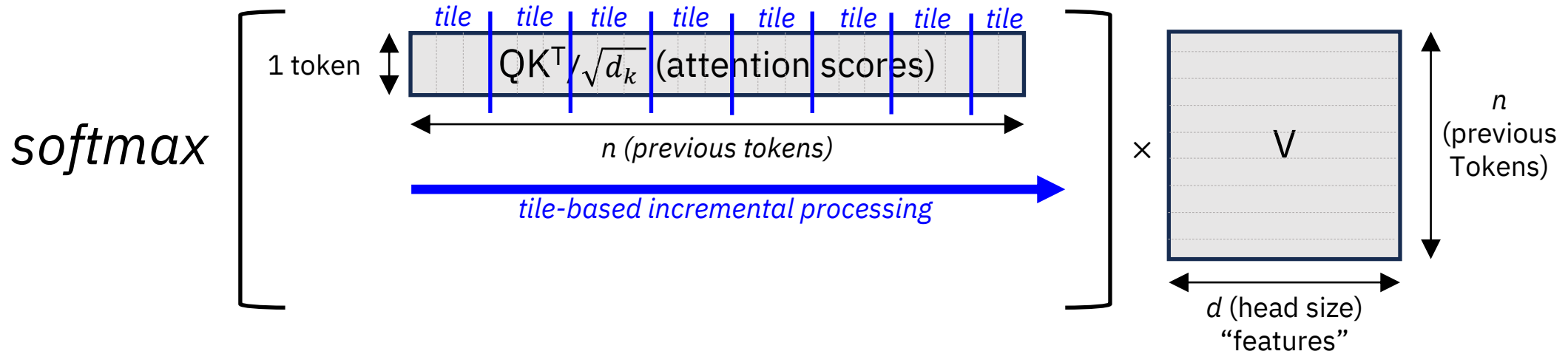
- Original attention formula:
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$
- For a single *query token*, *sequence*, and *query head* and causal masking:



- Numerically stable softmax:
$$\text{softmax}(s_i) = \frac{e^{s_i - \max_j s_j}}{\sum_{k=1}^n e^{s_k - \max_j s_j}}$$
- Softmax can be calculated efficiently using a *tiled* approach (a.k.a. *online softmax*)
 - split row into tiles, process incrementally
 - track running maximum and sum, rescale if maximum changes, normalize at the end

Attention

- Original attention formula:
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$
- For a single *query token*, *sequence*, and *query head* and causal masking:

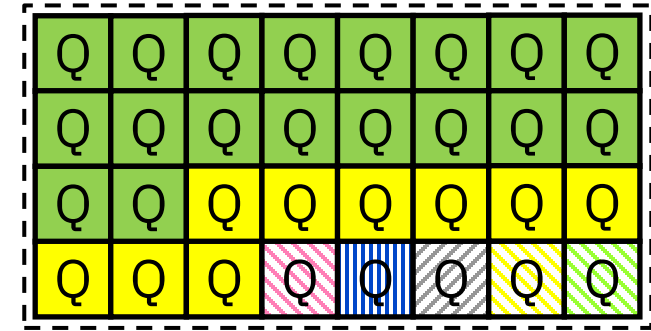


- Numerically stable softmax:
$$\text{softmax}(s_i) = \frac{e^{s_i - \max_j s_j}}{\sum_{k=1}^n e^{s_k - \max_j s_j}}$$
- Softmax can be calculated efficiently using a *tiled* approach (a.k.a. *online softmax*)
 - fused* with tiled matrix multiplications QK^T and V
 - core optimization in *FlashAttention*

Triton Attention Kernel

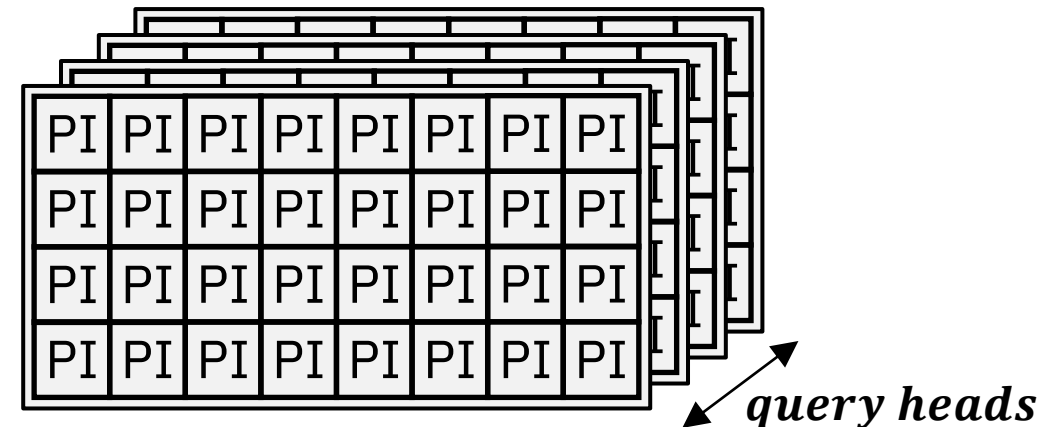
vLLM

- Example: scheduled *continuous batch* with chunked prefill
 - prefills (solid color)
 - decodes (hatched color)
- belong to different sequences with *separate KV contexts*



Naïve Approach

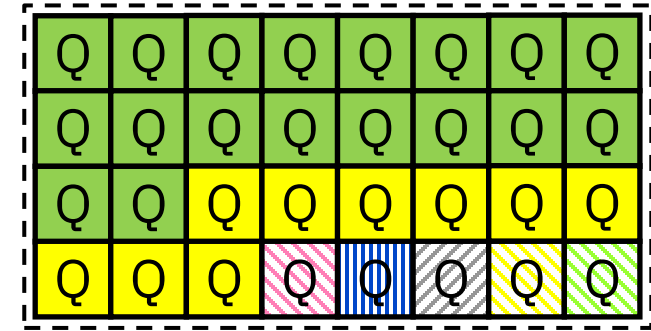
- Launch Triton attention kernel with grid: *query tokens* × *query heads*
 - Each *program instance* (PI)
 - handles a *single* query token and query head
 - attends over the KV context of its sequence
- *redundant* KV reads across PIs
- for query tokens from the same prompt
 - for query heads mapped to the same KV head (in case of Grouped Query Attention, *GQA*)



Triton Attention Kernel

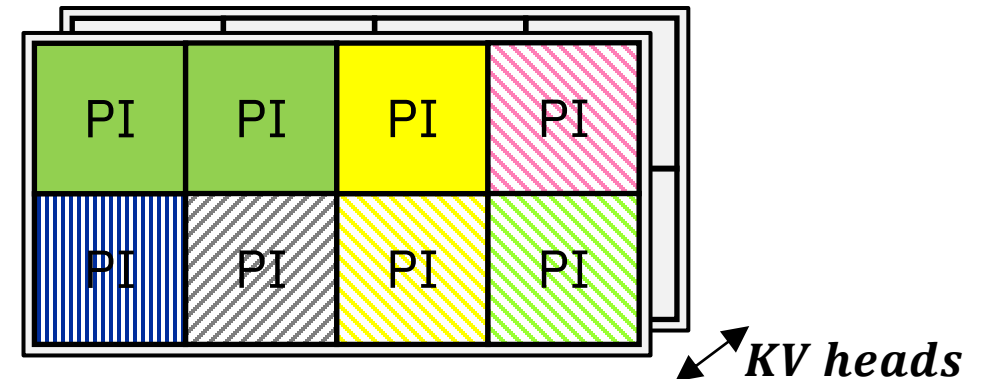
vLLM

- Example: scheduled *continuous batch* with chunked prefill
 - prefills (solid color)
 - decodes (hatched color)
- belong to different sequences with *separate KV contexts*



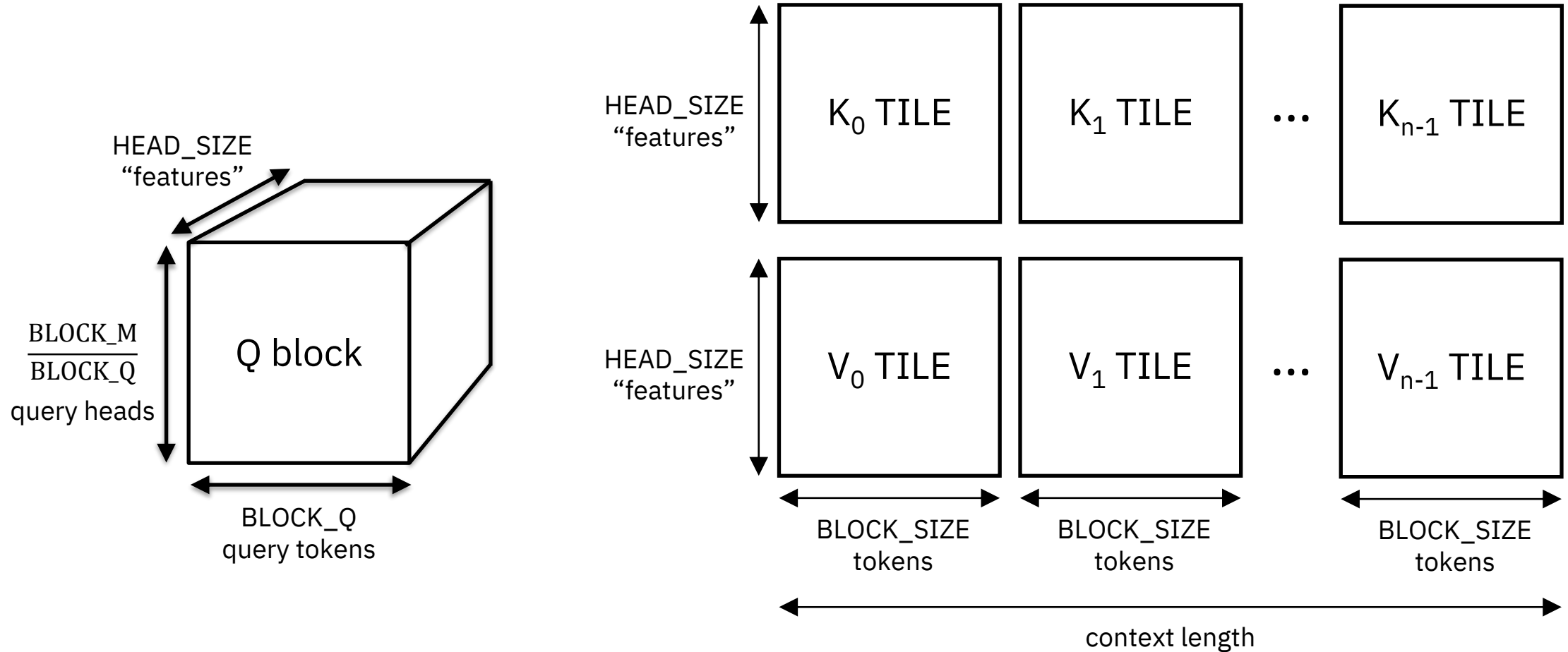
Q Block Optimization

- Launch Triton attention kernel with grid: ***Q blocks*** × ***KV heads***
 - Each *program instance* (PI) handles a *Q block*
 - multiple (prefill) or one (decode) query token(s)
 - multiple query heads mapped to the same KV head (*GQA*)
 - attends over the KV context of its sequence
- better *KV reuse*, lower memory bandwidth



Triton Attention Kernel

- Each *program instance* (PI) handles one Q block of size $BLOCK_M$, covering
 - $BLOCK_Q$ query tokens (prefill) or 1 query token (decode)
 - $\frac{BLOCK_M}{BLOCK_Q}$ query heads mapped to the same KV head (GQA)



Triton Attention Kernel

Challenge: Few Q blocks Lead to Small Launch Grid $Q \text{ Blocks} \times KV \text{ heads}$

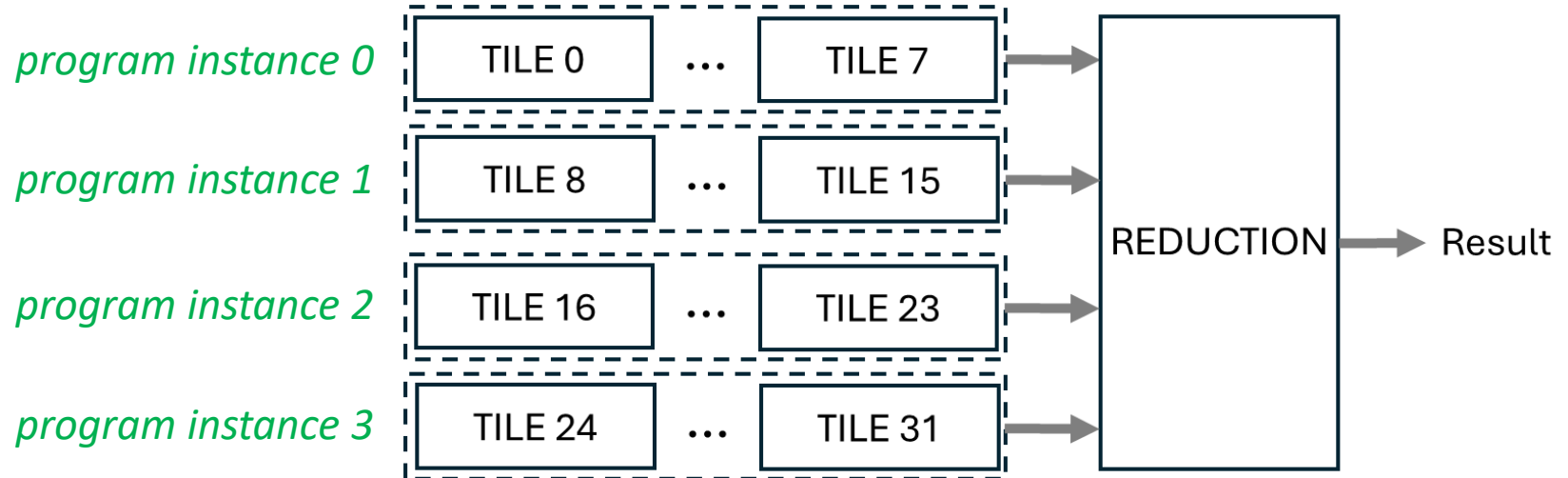
- Example: single-sequence decode involves one Q block per KV head
- Small *decode* batches \rightarrow few program instances \rightarrow *GPU underutilization*

Solution

- *Parallelize* tiled softmax across multiple program instances
 - effective for small batches and long sequences
 - key optimization in *FlashAttention* and *FlashDecoding*

- Multiple program instances
 - *sequential* tile processing within each instance
 - *parallel* tile processing across instances

\rightarrow requires *intermediate buffers* and a *reduction kernel*



Triton Attention Kernel

Kernels

- **2D:** $Q \text{ blocks} \times KV \text{ heads}$ exploits Q block parallelism
 - **3D + Reduce:** $Q \text{ blocks} \times KV \text{ heads} \times \text{tiles}$ exploits tile parallelism when Q parallelism is low
- *heuristic-based* runtime selection

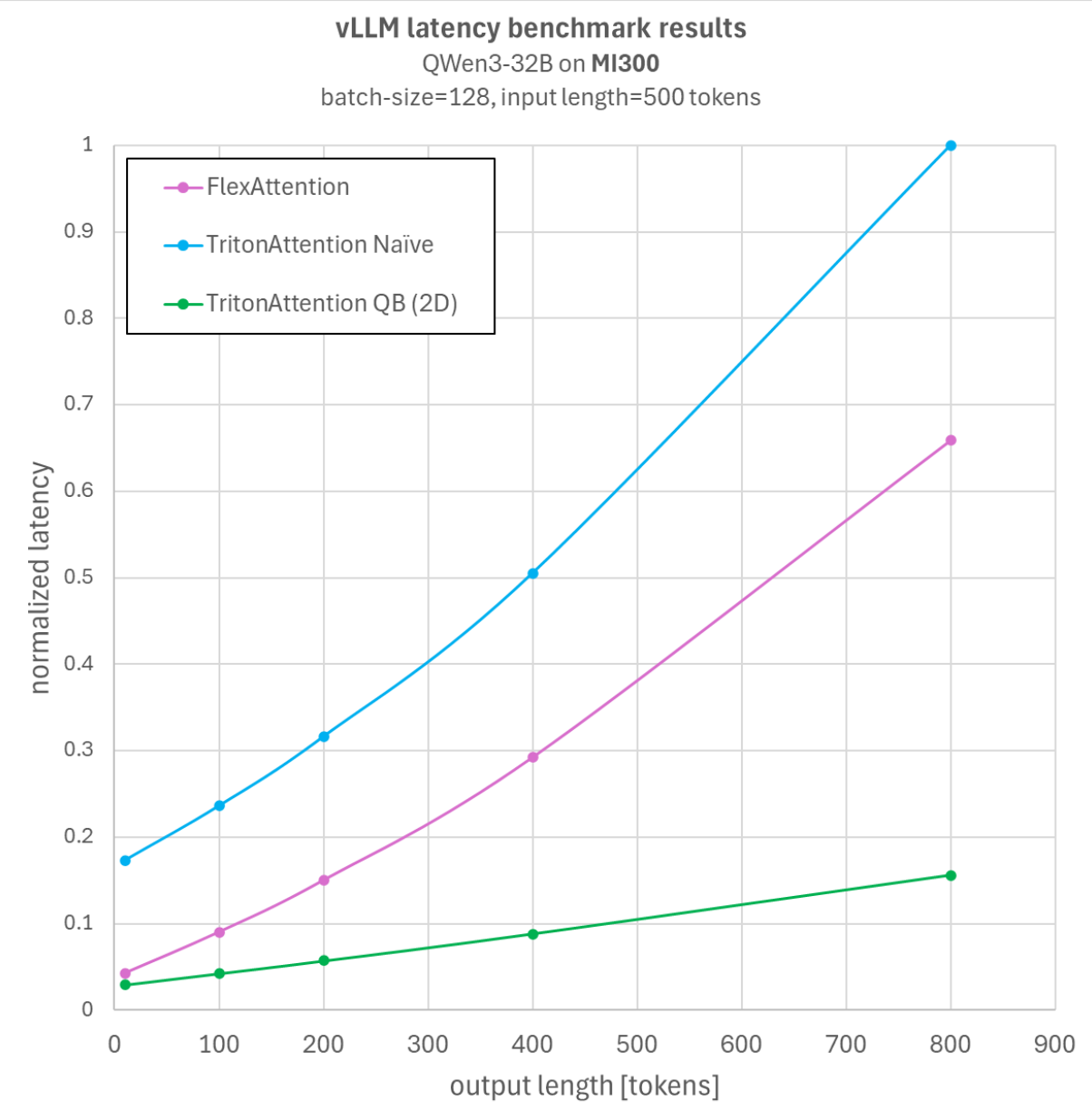
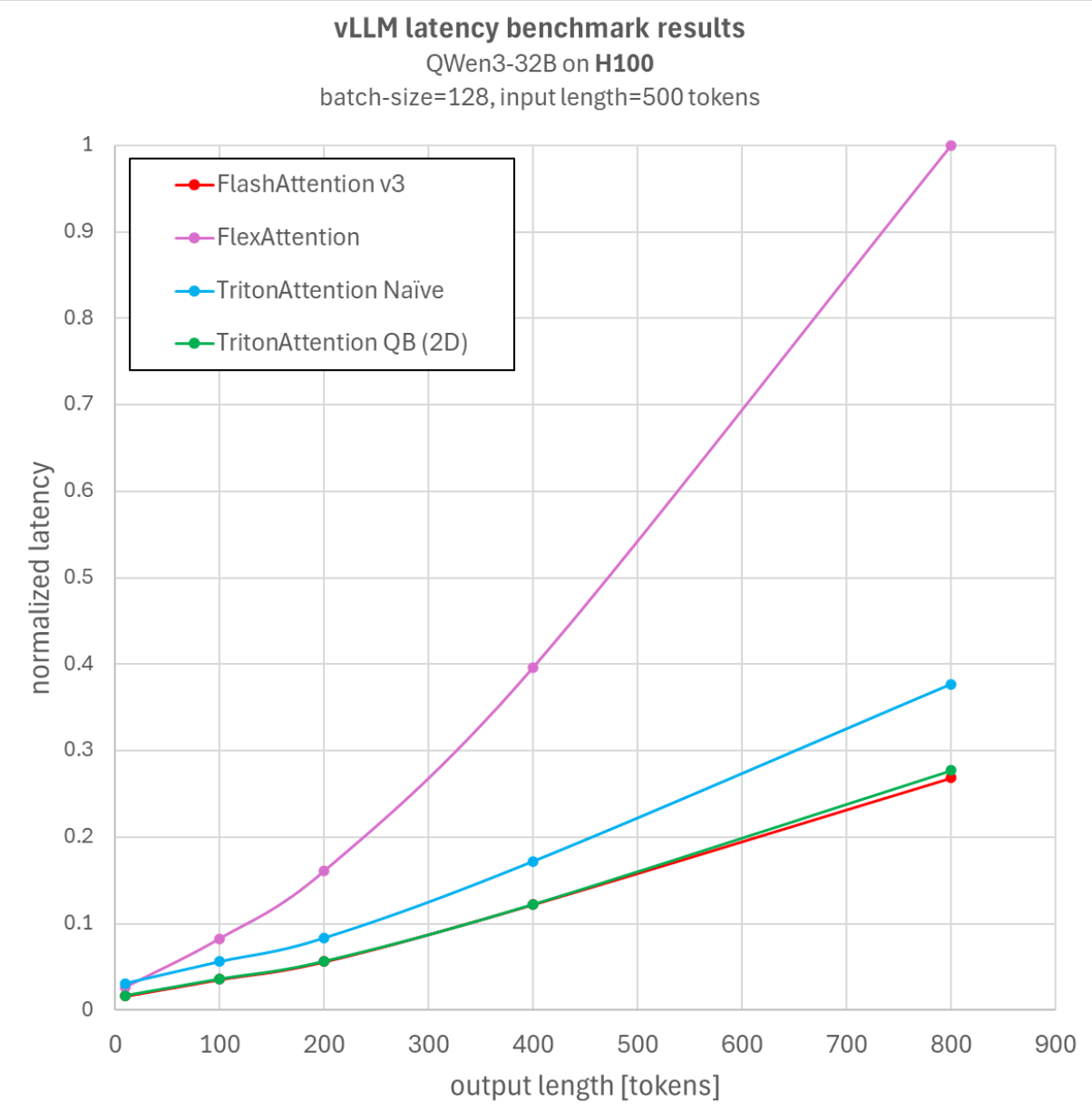
Tile Size

- Decoupled from **BLOCK_SIZE** (vLLM KV-cache paging granularity)
 - compute tiles independent of memory layout
 - enables tuning and support for hybrid/SSM models

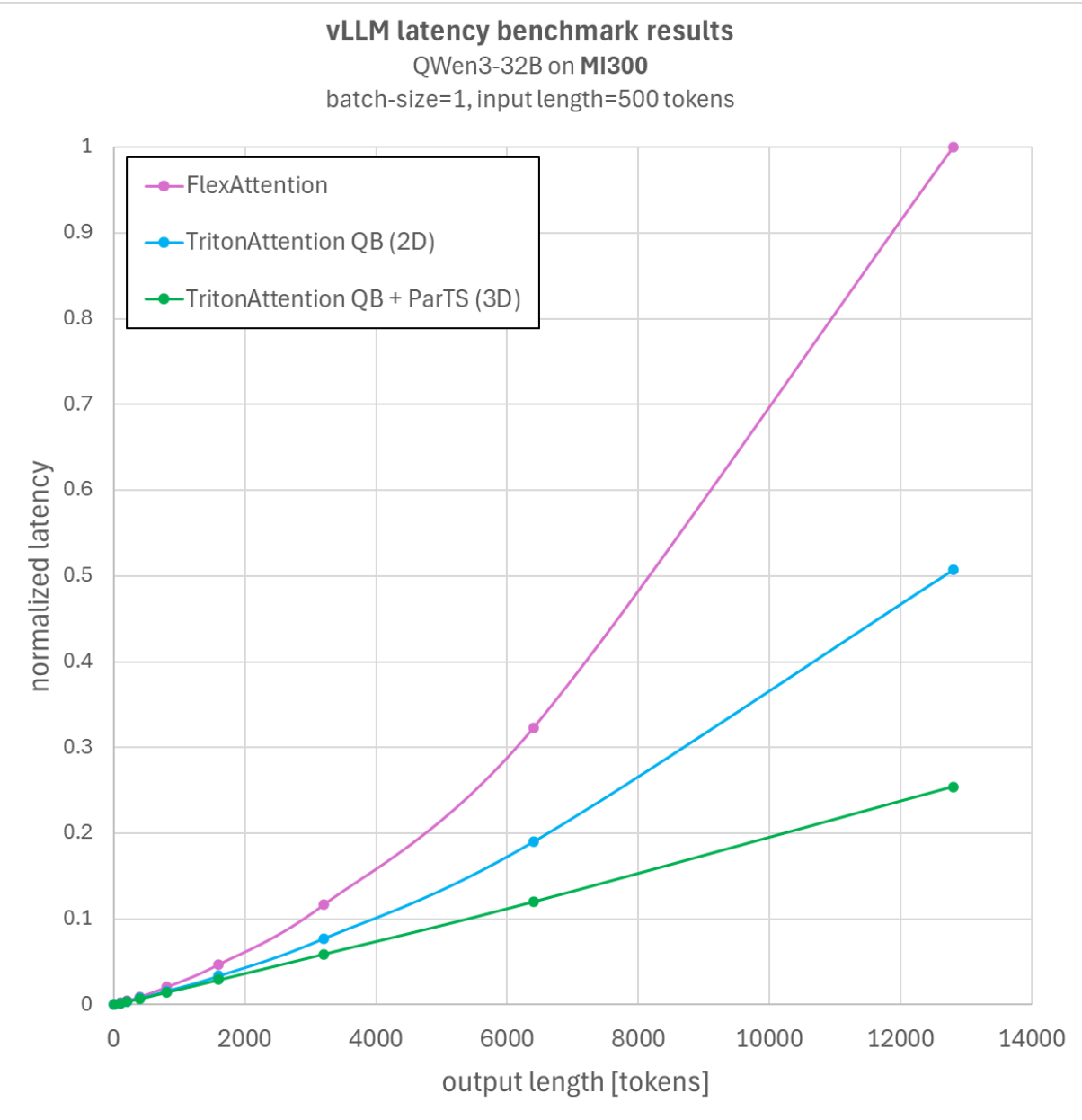
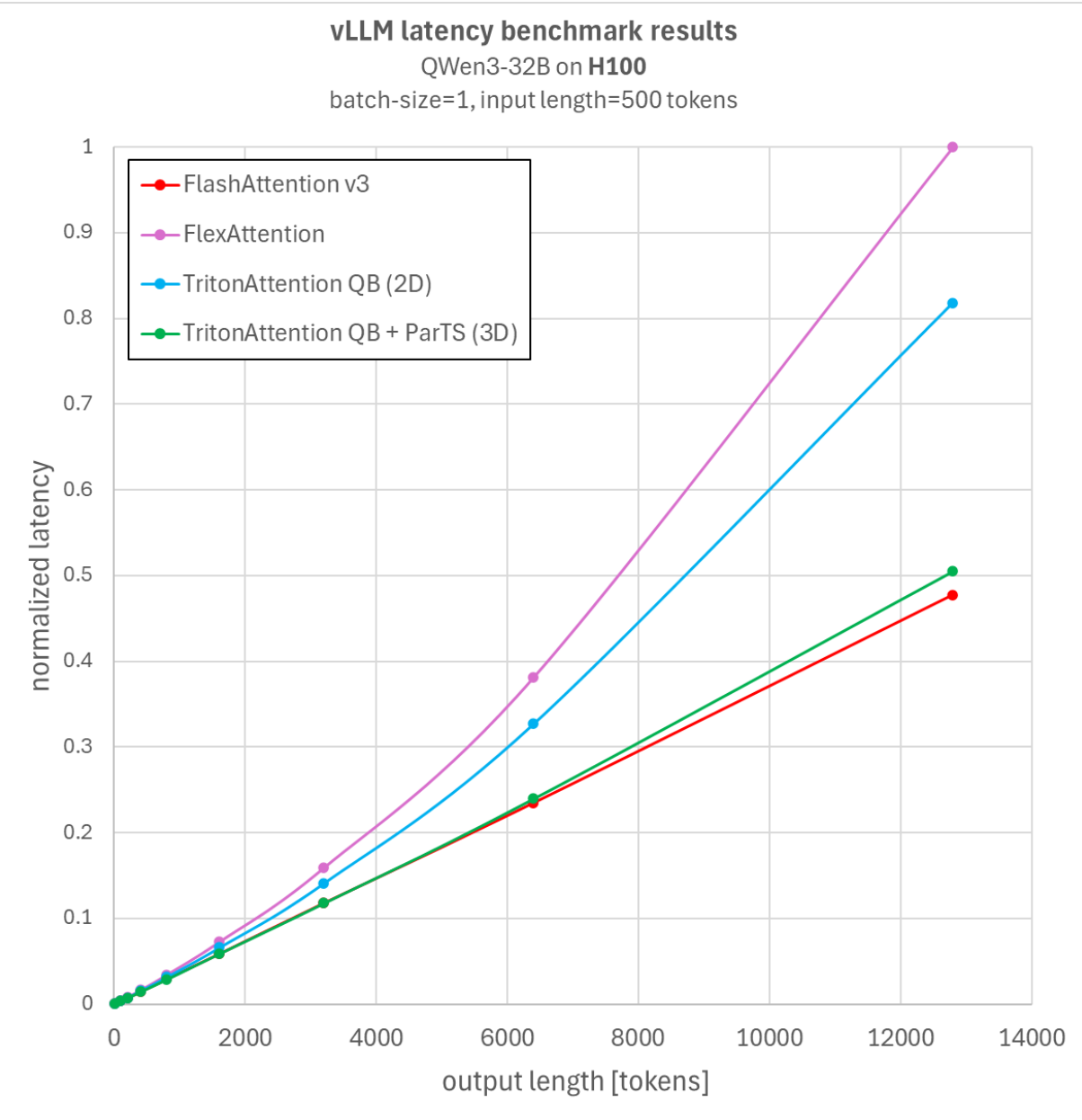
Persistent 3D Kernel Implementation (*PR in progress*)

- Static launch grid: $SMs \text{ per } KV \text{ head} \times KV \text{ heads}$
 - enables *persistent* execution (e.g., H100: 132 SMs, MI300: 304 CUs)
- *Per-tile* work distribution (not per sequence)
 - tiles evenly distributed across SMs
 - CPU-generated scheduling tensors map each SM to (Q block, tile) ranges
- *Compact* intermediate buffers
 - only sequences split across SMs require additional buffer slots and reduction
 - $total \ slots \leq max \ batch \ size \ 3D + SMs \ per \ KV \ head$

Performance Evaluation 2D Kernel

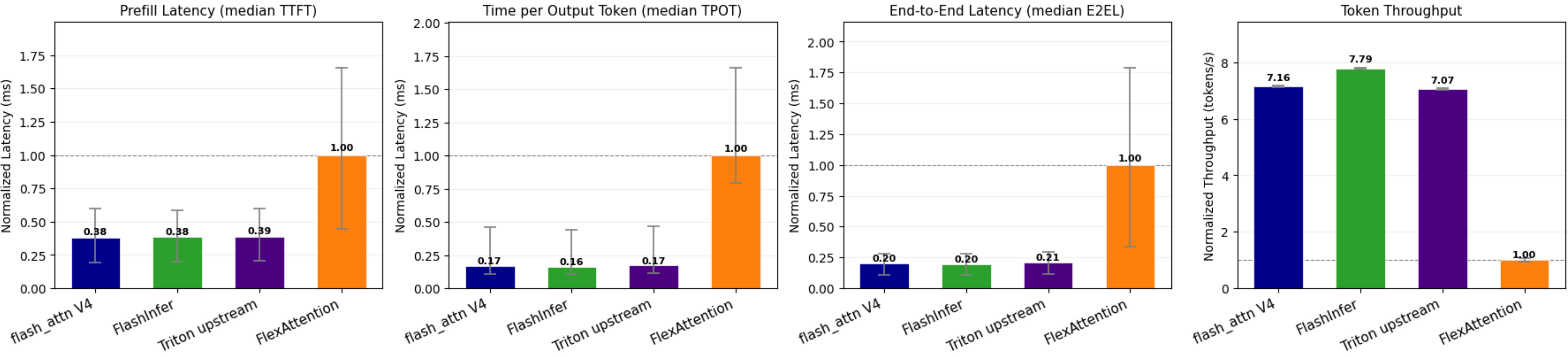


Performance Evaluation 3D Kernel



Measuring end-to-end with ShareGPT: B200

vllm bench serve with ShareGPT dataset — Qwen3-32B on B200
(normalized to FlexAttention)
(end2end benchmarking of vLLM, no prefix caching)

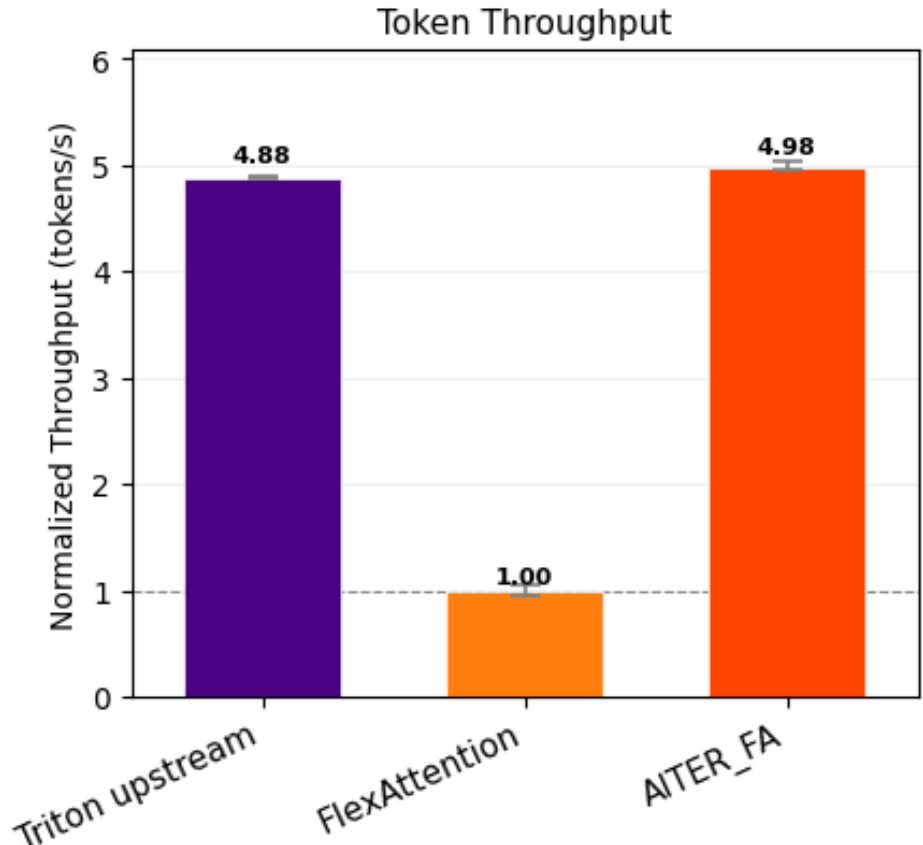
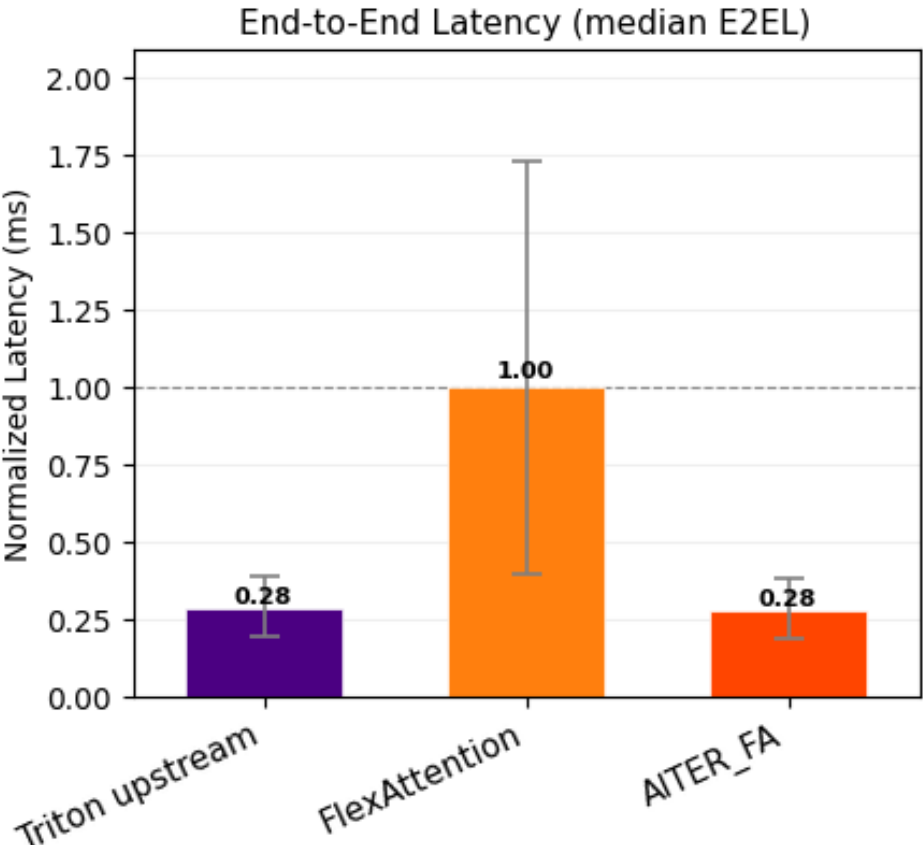
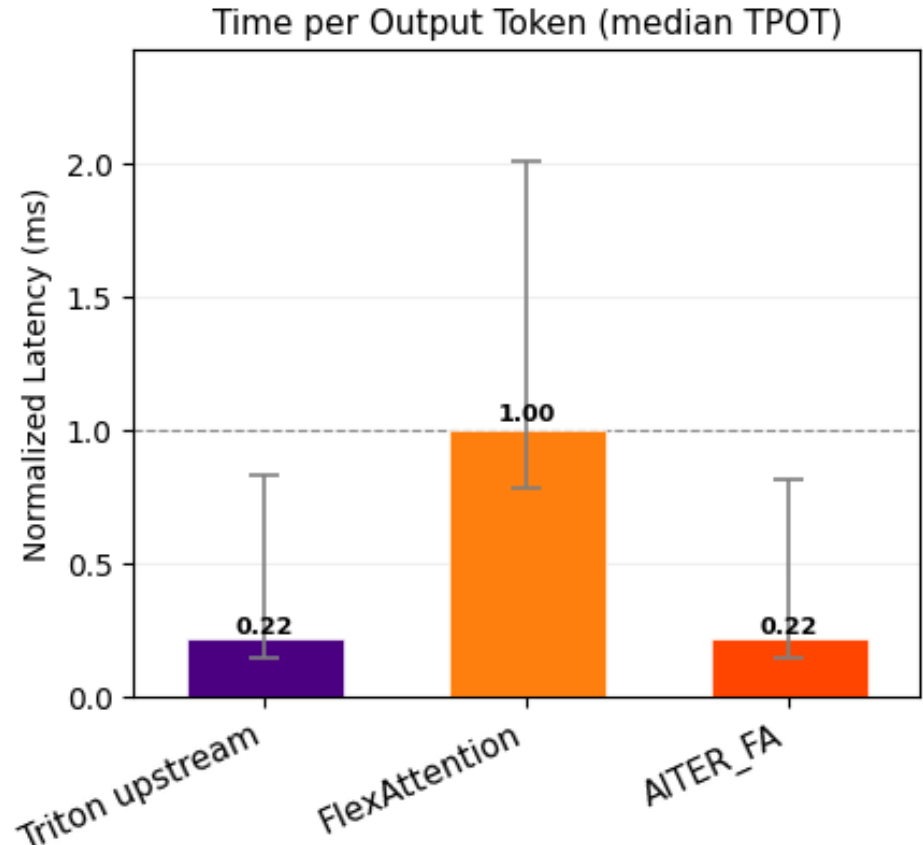
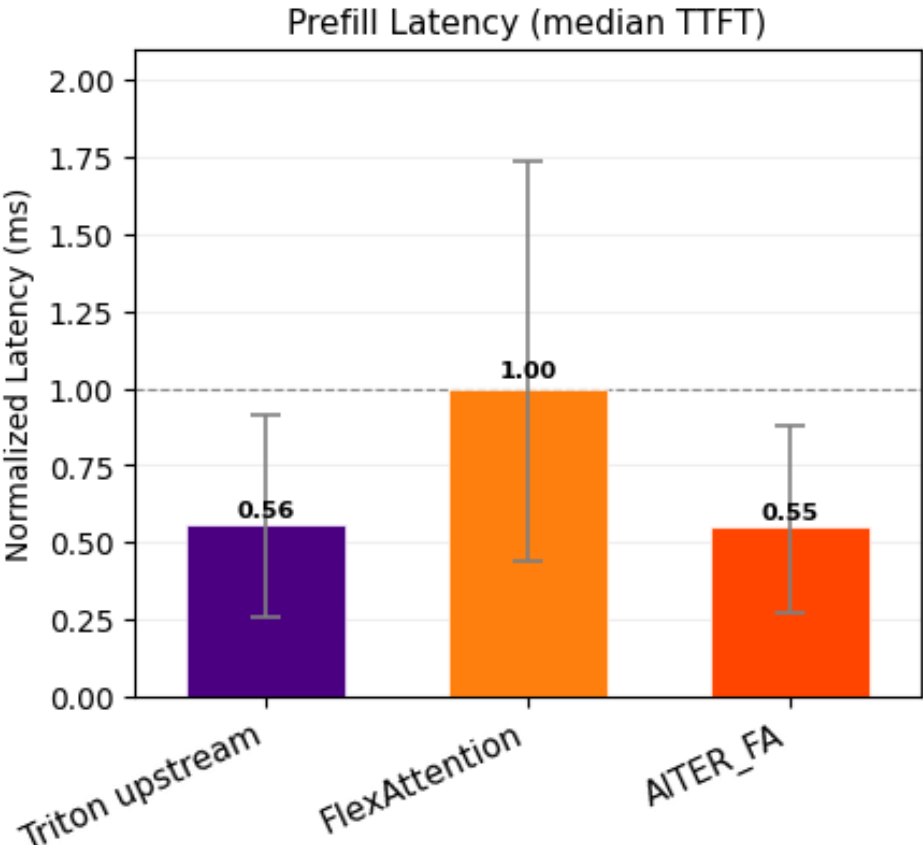


- errorbars are p20/p80 for latencies, min/max for throughput
 - all values are average of 3 measurements

warmup: 2 iterations

Measuring end-to-end with ShareGPT: MI300X

**vllm bench serve with ShareGPT dataset — Qwen3-32B on MI300X
(normalized to FlexAttention)
(end2end benchmarking of vLLM, no prefix caching)**



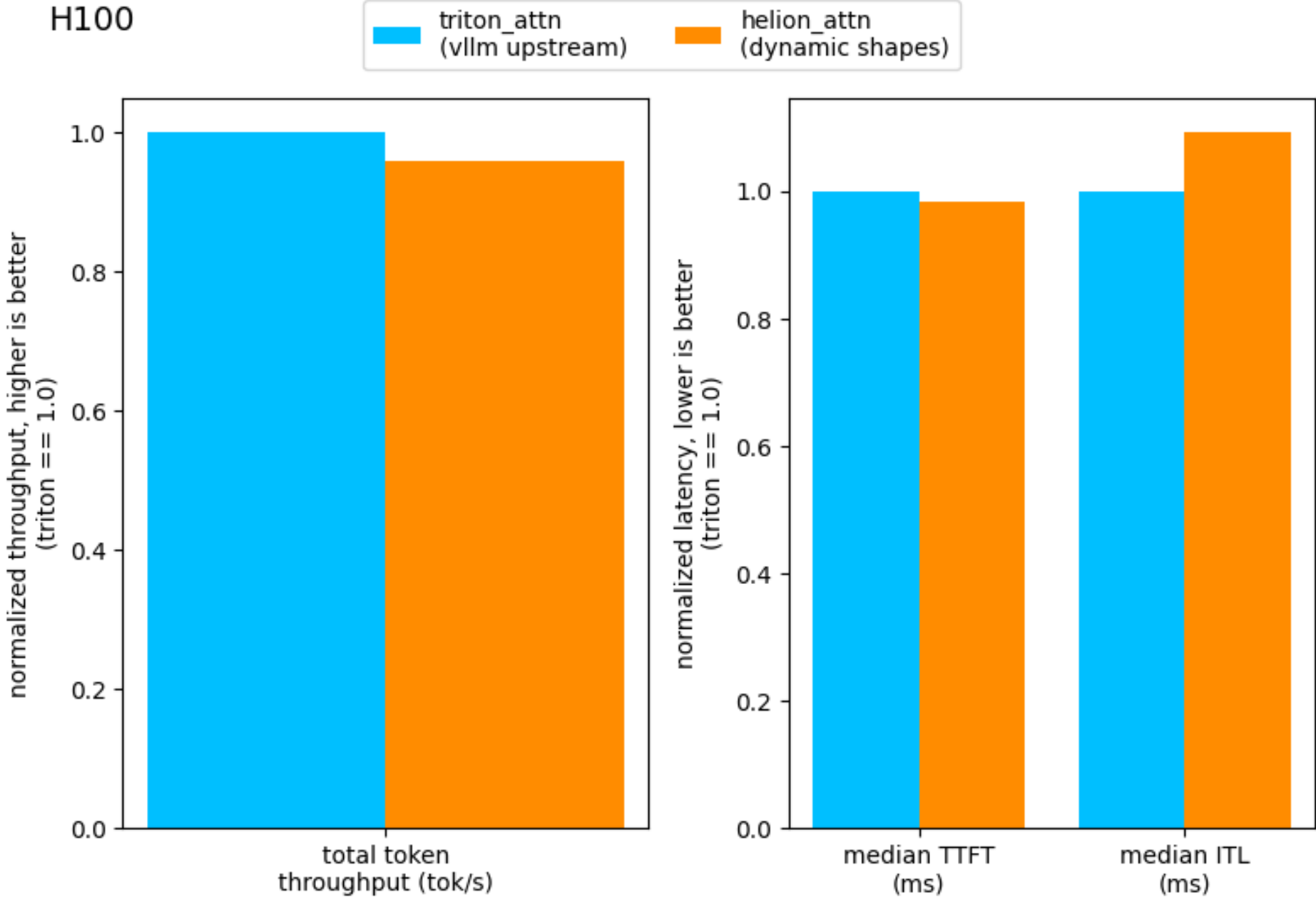
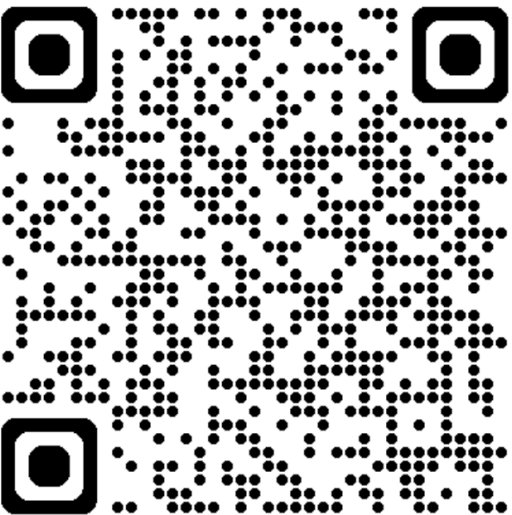
- errorbars are p20/p80 for latencies, min/max for throughput
- all values are average of 3 measurements

warmup: 2 iterations

Outlook: Helion, the new kid on the block?



- Helion as new DSL from pytorch
- “Tiled pytorch” or “higher-level Triton”
- We implemented a (simple) paged attention kernel in Helion
 - → <https://pytorch.org/blog/portable-paged-attention-in-helion/>
 - Already getting close....
- Code as vLLM draft [PR#27293](#)



Conclusion?

- We developed a **triton-only feature-complete attention backend for vLLM**
 - that achieves state-of-the-art performance on **AMD & NVIDIA** (and also runs on Intel)
 - with a single-source kernel
- **Community-maintained kernels in an open-source high-level language** help to “catch the train” in LLM kernel development by **enabling performance portability**

```
pip install vllm
# no further
# dependencies
```

- **Big thanks to the team:** Thomas Parnell, Chih-Chieh Yang, Sara Kokkila Schumacher, Sage Moore, Lukas Wilkinson, Luka Govedič

Want to know more?

- **Paper:** → With all details and further lessons learned

→ **Paper:** “*The Anatomy of a Triton Attention Kernel*”



- Pytorch blog:
 - pytorch.org/blog/portable-paged-attention-in-helion
 - pytorch.org/blog/enabling-vllm-v1-on-amd-gpus-with-triton
- VLLM blog: vllm.ai/blog/vllm-triton-backend-deep-dive
- **... Contribute!** (Play around...and we/vLLM welcomes PRs)

...looking forward to all your questions!

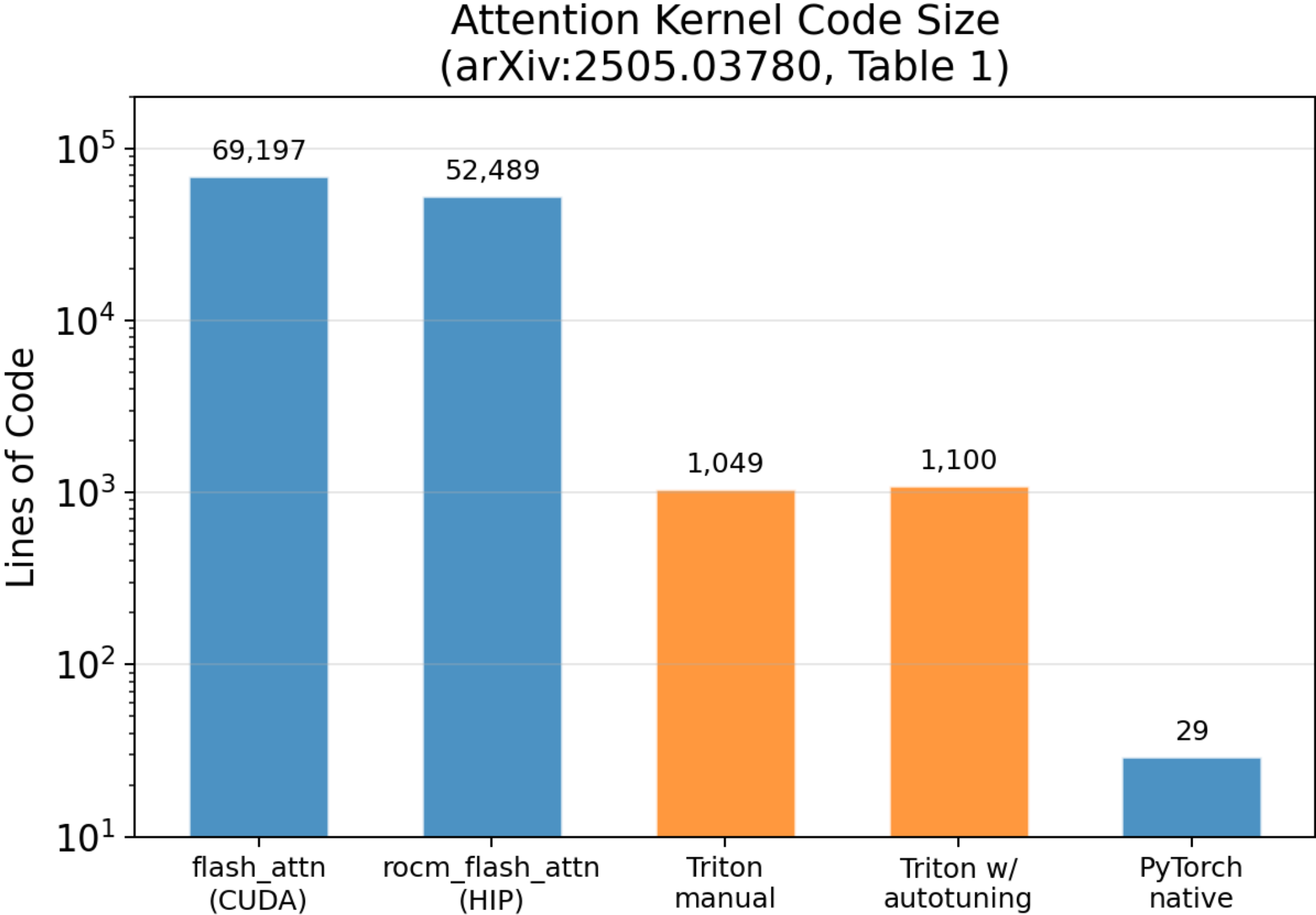
Dr. Burkhard Ringlein **Dr. Jan van Lunteren**

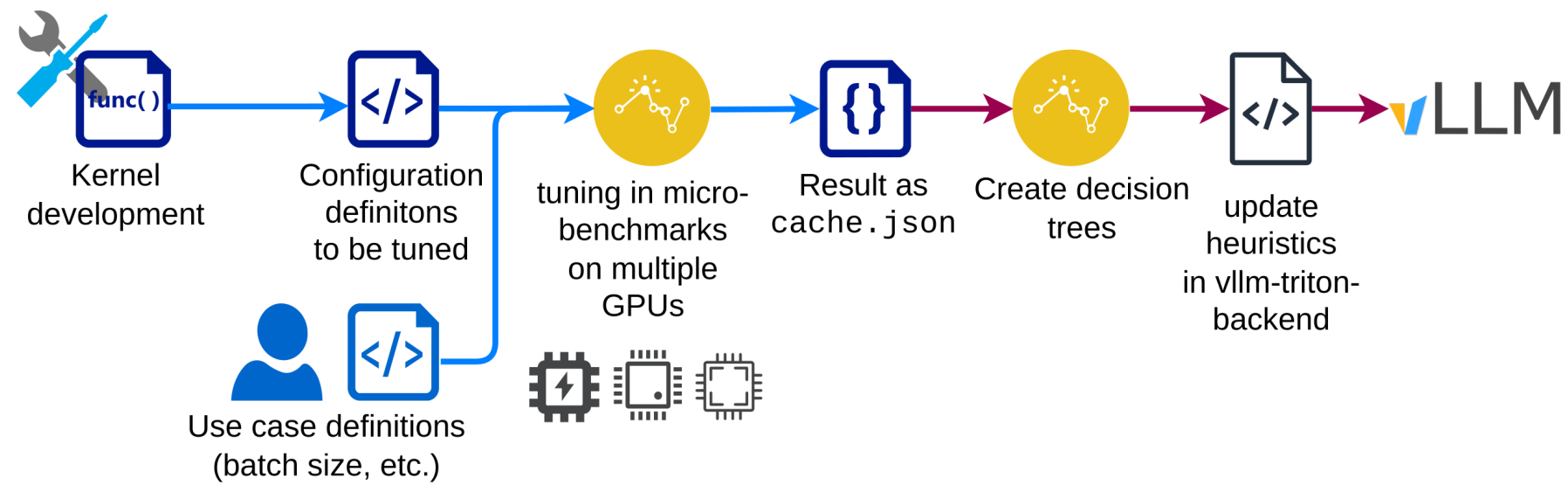
 ngl@zurich.ibm.com jvl@zurich.ibm.com

 [burkhard-ringlein](https://www.linkedin.com/in/burkhard-ringlein) [janvanlunteren](https://www.linkedin.com/in/janvanlunteren)

Appendix

Portability in LoC...

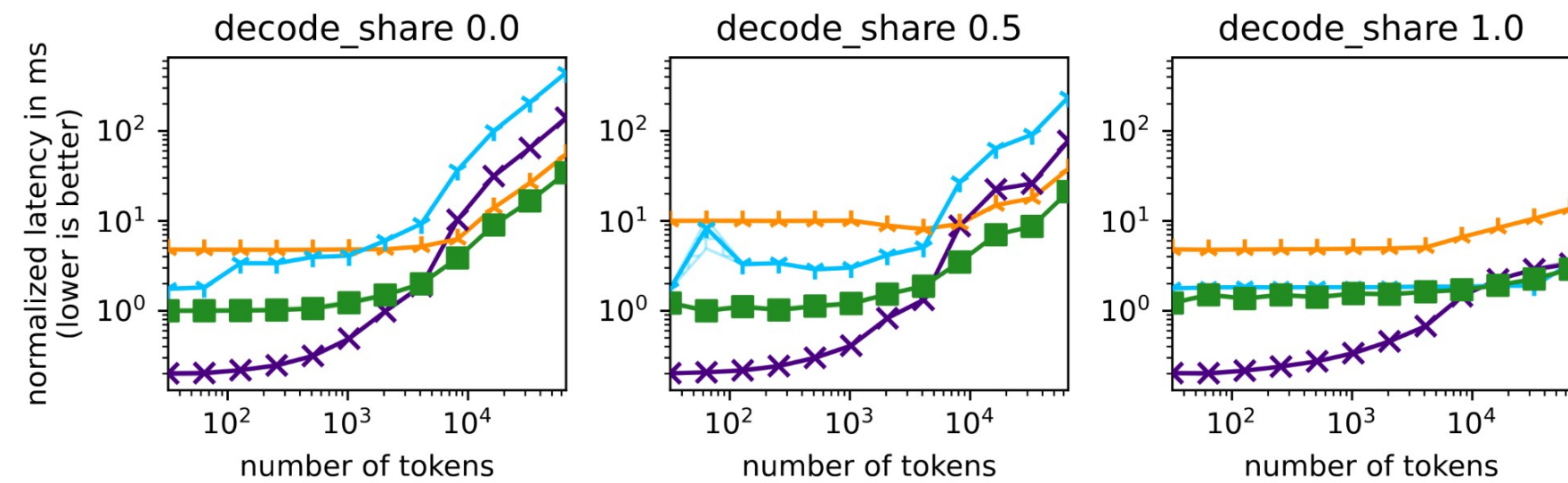
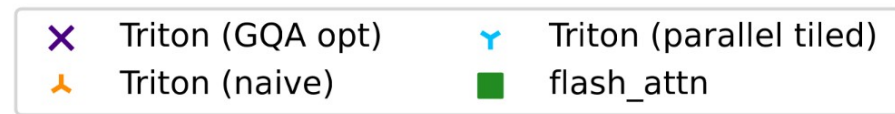




```

if torch.version.hip:
    TILE_SIZE_PREFILL = 64
    NUM_STAGES_PREFILL = 1
    NUM_WARPS_DECODE_3D = 4
else: # cuda platform
    TILE_SIZE_PREFILL = 16
    NUM_STAGES_PREFILL = 4
    NUM_WARPS_DECODE_3D = 2

```



Key-Insights: Autotuning & Heuristics

- Tuning block and tile sizes for different models & GPUs can increase performance
- However:
 - Tuning still has a lot of overhead (even if done ahead-of-time)
 - Looking up best config before every kernel call adds latency
 - Ahead-of-time tuning can barely cover all possible use cases
- → Train **heuristics based on extensive micro-benchmarks**
 - Simplifies kernel launch
 - Generalizes well to unseen workloads
 - (simple if-else trees for now)
 - Compatible with CUDA graphs
 - **Use extensive Micro-benchmarks**
- Future work: Implement ML-based decision Trees

IBM