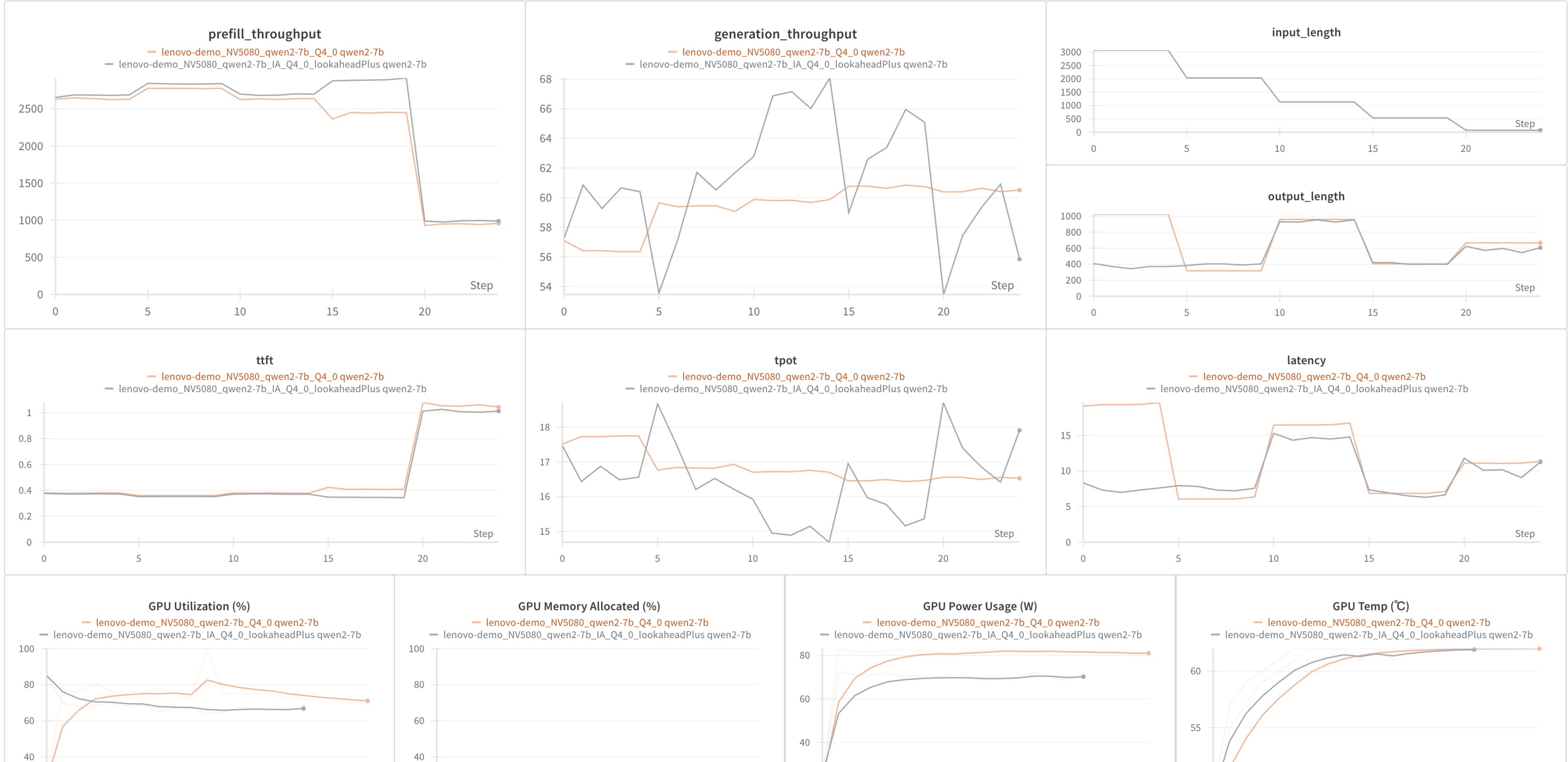
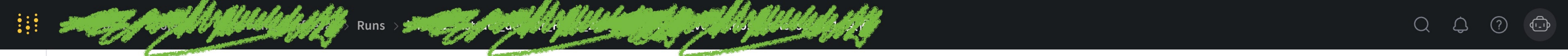


1. 从7~8b的模型来看，开启IA_Q4_0和投机后对与GPU的功耗、显存占用大致都为base的85%~90%左右；

指标看板

下方表格勾选指定tab筛选模型 (默认勾选ds-2-7b)





☰ Chenyonghua's run workspace Personal workspace

Autosaved just now ⌂ ⌂ ⌂

Overview

Workspace

System

Logs

Files

Artifacts

kimi-k2-instruct_default_1

Search panels with regex

Create report

Tables 2

Add panel

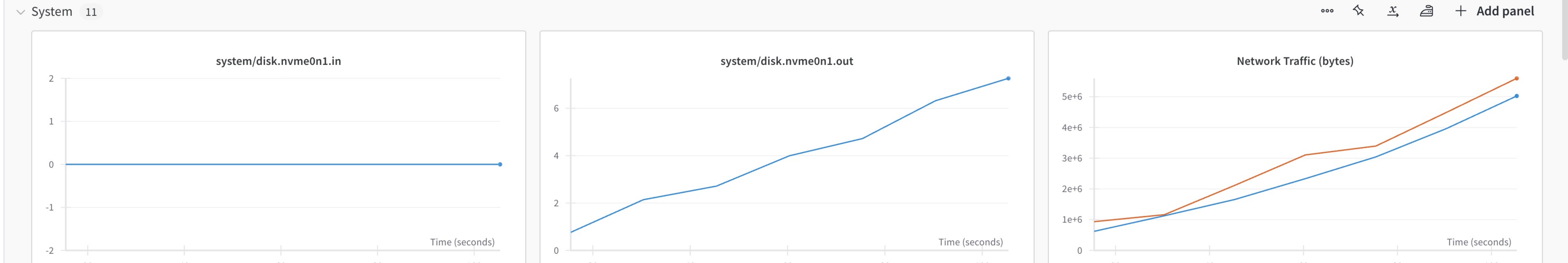
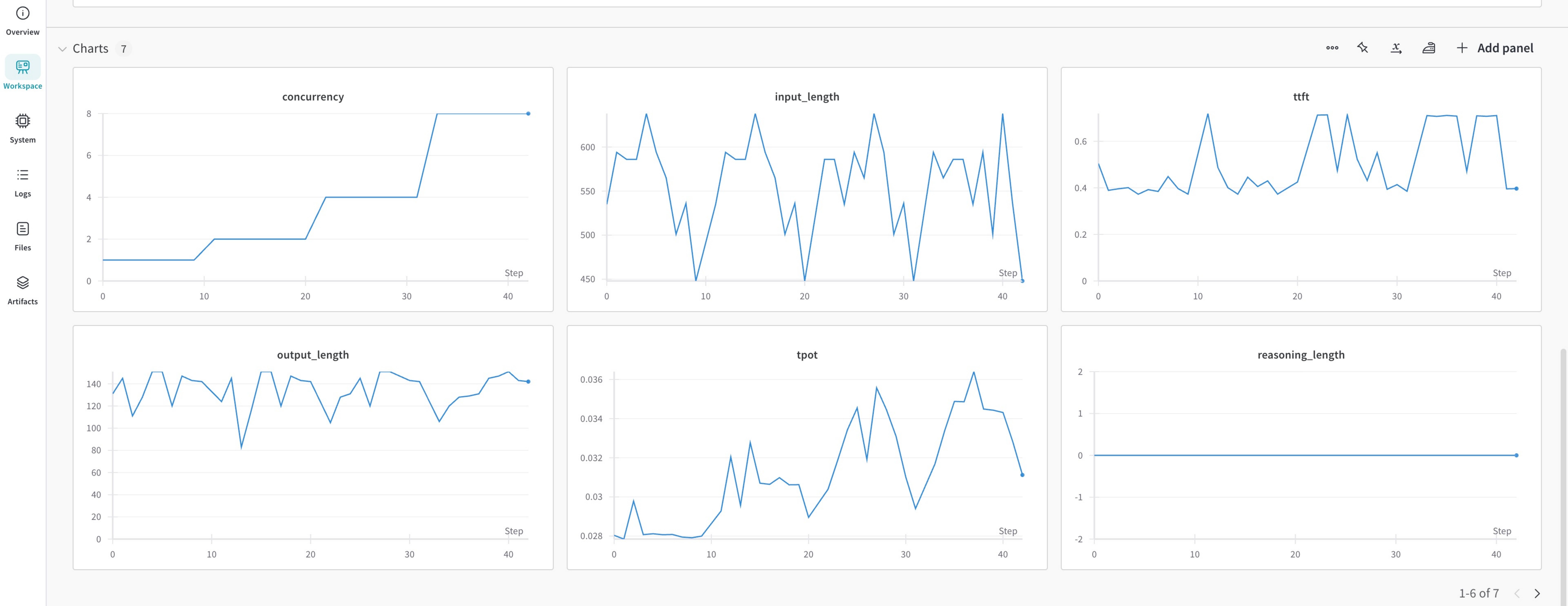
runs.summary["benchmark_results"]

	芯片规格	模型	软件框架	集群规格	数据集	请求总量	请求速率	请求成功率(%)	审核命中	并发数	输入总token量	输出总token量	思维链总token量	总耗时(s)	总请求吞吐(req/s)	总输入吞吐(token/s)	总输出吞吐(token/s)	TTFT_avg
1	kimi-k2-instruct	-	-	-	random	10	串行	100.0	0	1	5583	1369	0	42.32	0.24	131.92	32.35	405.71
2	kimi-k2-instruct	-	-	-	random	10	串行	100.0	0	2	5583	1322	0	22.82	0.44	244.66	57.93	445.84
3	kimi-k2-instruct	-	-	-	random	10	串行	100.0	0	4	5583	1363	0	14.18	0.71	393.79	96.14	530.99
4	kimi-k2-instruct	-	-	-	random	10	串行	100.0	0	8	5583	1342	0	9.47	1.06	589.56	141.71	622.53

Export as CSV Columns... Reset table

Lengths Distribution

Scenario	Percentile	Prompt Input	Actual Output	Reasoning
random Req:10 Rate:-1.0 Con:8	0%	448	106	0
random Req:10 Rate:-1.0 Con:8	25%	535.25	128.25	0
random Req:10 Rate:-1.0 Con:8	50%	575.5	136.5	0
random Req:10 Rate:-1.0 Con:8	75%	592	144.5	0
random Req:10 Rate:-1.0 Con:8	100%	638	151	0
random Req:10 Rate:-1.0 Con:8	Avg	558.3	134.2	0



☰ Chenyonghua's workspace Personal workspace

autosaved just now ⏪ ⏩ ⋮

Runs (1)

 Search runs

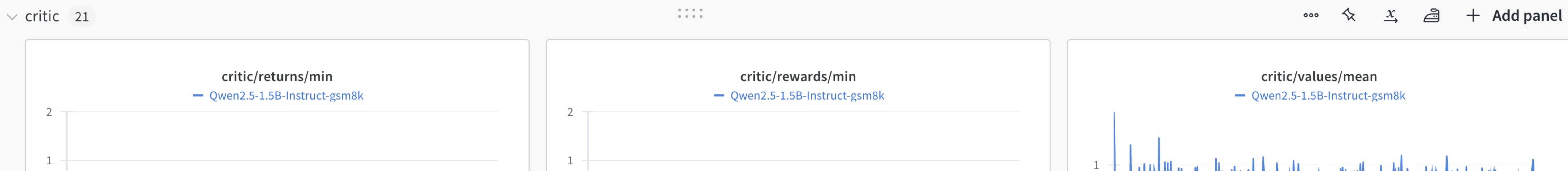


 Name (1 visualized)

 Qwen2.5-1.5...struct-gsm8k



-6 of 7 < >





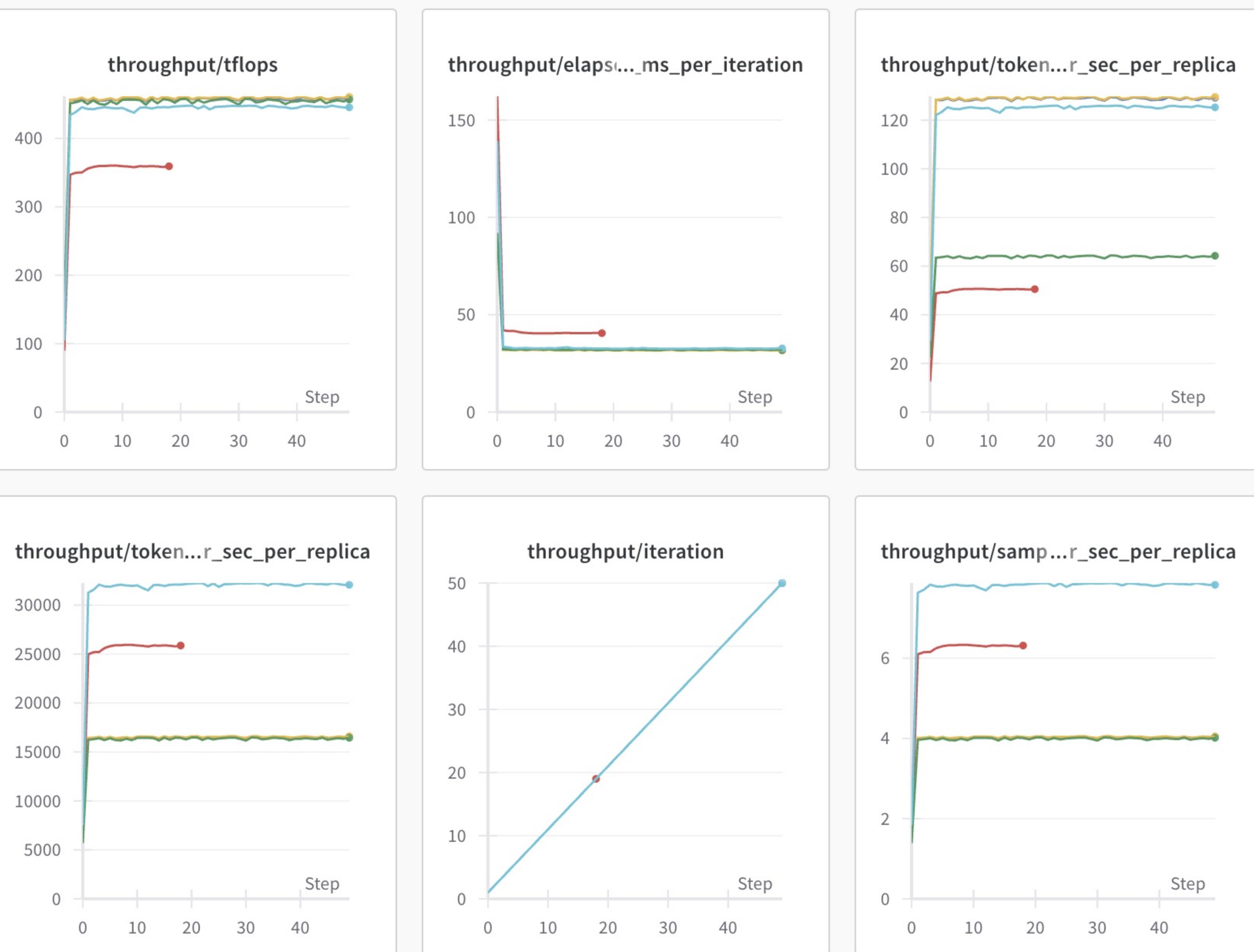




Name (5 visualized)

- pretrain-llama-70B-WS512-tp4-pp8-gbs4096-mbs1-seqlen4096
- pretrain-llama-70B-WS128-tp4-pp4-gbs1024-mbs1-seqlen4096
- pretrain-llama-70B-WS128-tp4-pp4-gbs1024-mbs1-seqlen4096
- pretrain-llama-70B-WS128-tp4-pp4-gbs1024-mbs1-seqlen4096
- pretrain-llama-70B-WS256-tp2-pp8-gbs1024-mbs1-seqlen4096
- pretrain-llama-70B-WS256-tp2-pp8-gbs2048-mbs1-seqlen4096
- pretrain-llama-70B-WS256-tp2-pp4-gbs128-mbs1-seqlen4096
- pretrain-llama-70B-WS256-tp4-pp4-gbs4096-mbs1-seqlen4096
- pretrain-llama-70B-WS256-tp4-pp4-gbs2048-mbs1-seqlen4096
- pretrain-llama-70B-WS256-tp4-pp8-gbs4096-mbs1-seqlen4096
- pretrain-llama-70B-WS256-tp4-pp4-gbs1024-mbs1-seqlen4096
- pretrain-llama-70B-WS256-tp4-pp8-gbs1024-mbs1-seqlen4096

▼ throughput 11



Search runs

*

≡

☰

↑

Name (9 visualized)

MR-V100_vllm_Qwen2.5-32B

MR-V100_vllm_Qwen2.5-14B

MR-V100_vllm_Llama-2-13b-hf

MR-V100_vllm_Qwen1.5-14B

MR-V100_vllm_Qwen1.5-14B

MR-V100_vllm_Llama-2-70b-hf

MR-V100_vllm_Llama-2-70b-hf

MR-V100_vllm_Llama-2-70b-hf

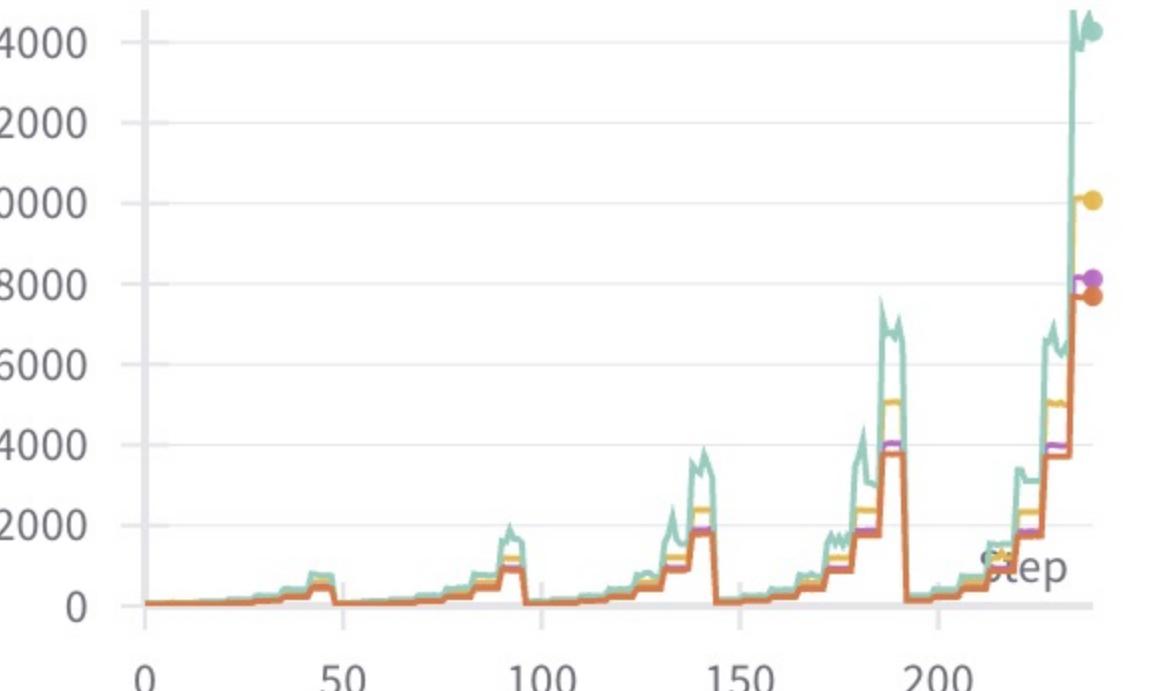
MR-V100_vllm_Llama-2-70b-hf

performance 15

... ✎ ✖ + Add panel

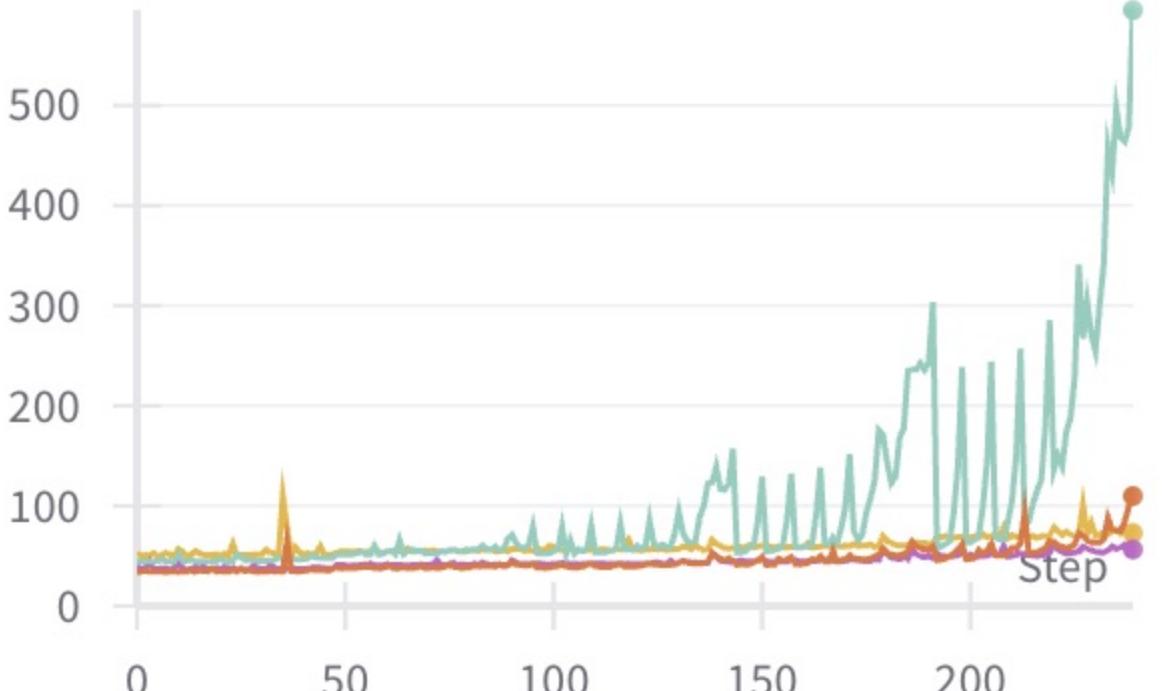
performance/avg_ttft_ms

MR-V100_vllm_Qwen2.5-32B
MR-V100_vllm_Qwen2.5-14B
MR-V100_vllm_Llama-2-13b-hf
MR-V100_vllm_Qwen1.5-14B



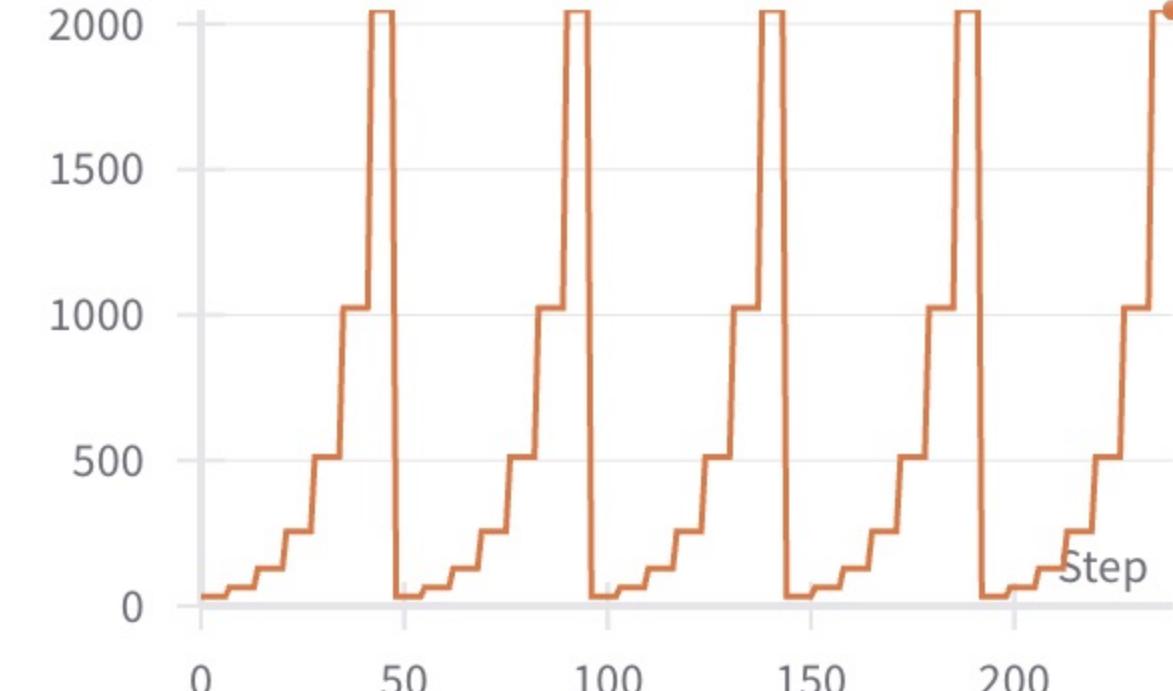
performance/max_tpot_ms

MR-V100_vllm_Qwen2.5-32B
MR-V100_vllm_Qwen2.5-14B
MR-V100_vllm_Llama-2-13b-hf
MR-V100_vllm_Qwen1.5-14B



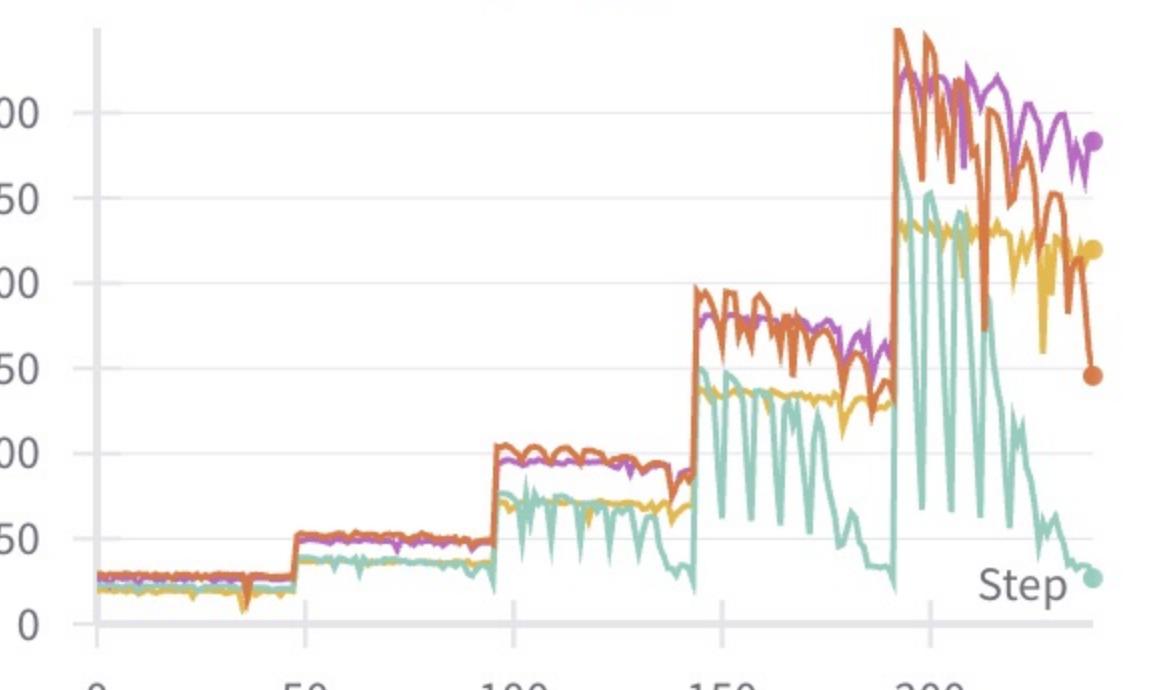
performance/input_length

MR-V100_vllm_Qwen2.5-32B
MR-V100_vllm_Qwen2.5-14B
MR-V100_vllm_Llama-2-13b-hf
MR-V100_vllm_Qwen1.5-14B



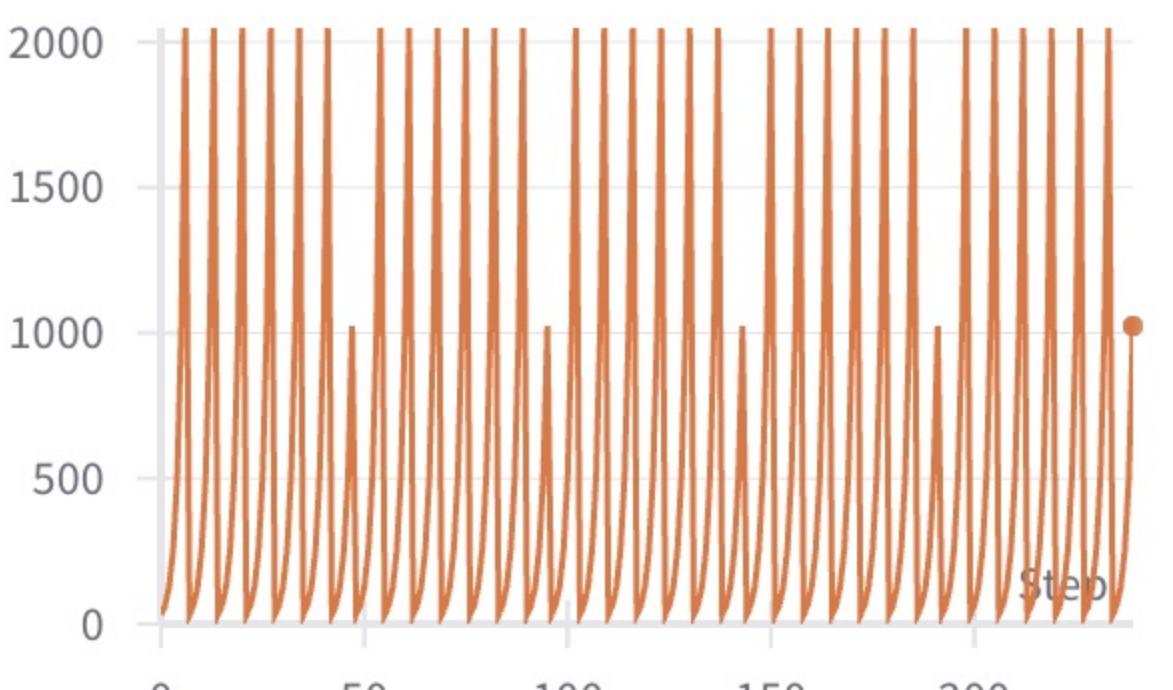
performance/min_incremental_throughput

MR-V100_vllm_Qwen2.5-32B
MR-V100_vllm_Qwen2.5-14B
MR-V100_vllm_Llama-2-13b-hf
MR-V100_vllm_Qwen1.5-14B



performance/output_length

MR-V100_vllm_Qwen2.5-32B
MR-V100_vllm_Qwen2.5-14B
MR-V100_vllm_Llama-2-13b-hf
MR-V100_vllm_Qwen1.5-14B



performance/avg_incremental_throughput

MR-V100_vllm_Qwen2.5-32B
MR-V100_vllm_Qwen2.5-14B
MR-V100_vllm_Llama-2-13b-hf
MR-V100_vllm_Qwen1.5-14B

