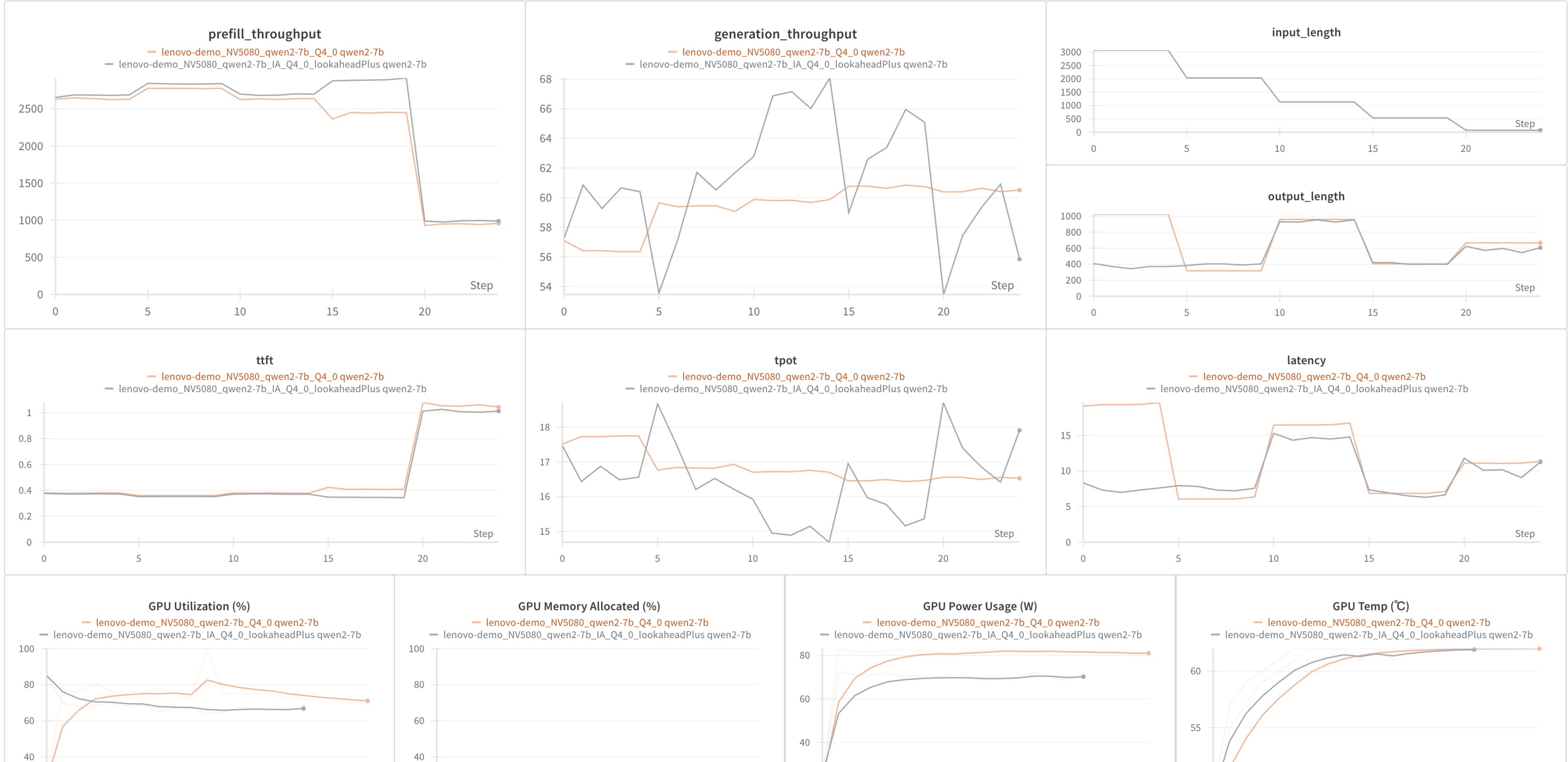
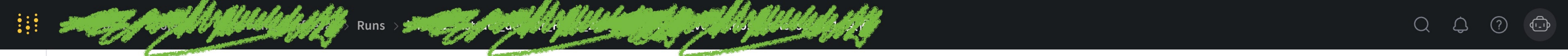


1. 从7~8b的模型来看，开启IA\_Q4\_0和投机后对与GPU的功耗、显存占用大致都为base的85%~90%左右；

## 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

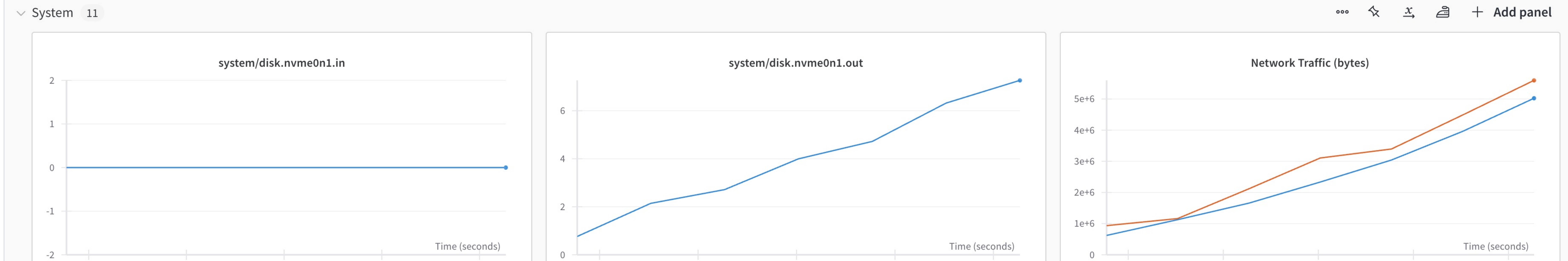
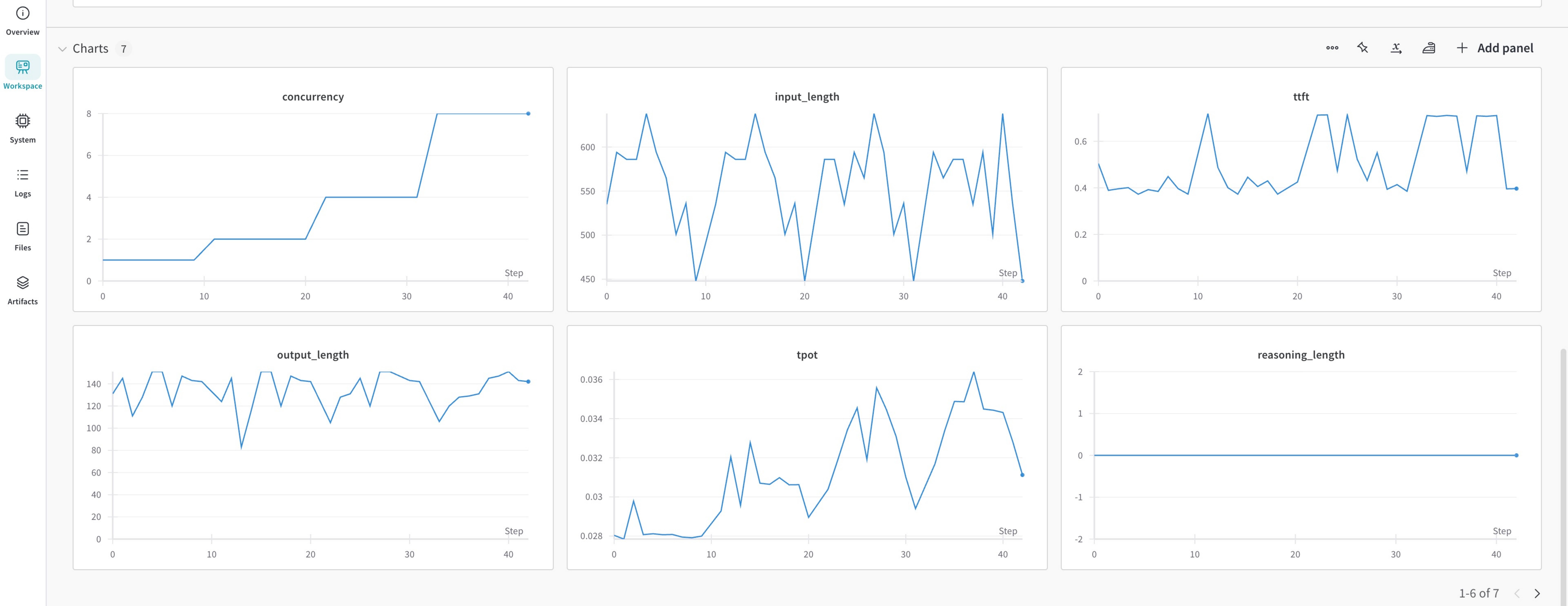
runs.summary["benchmark\_results"]

	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





Overview

Workspace

Runs

Jobs

Automat.

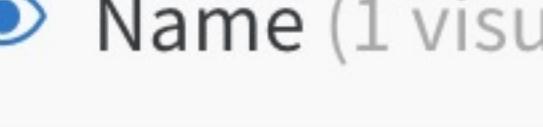
Sweeps

Reports

Artifacts

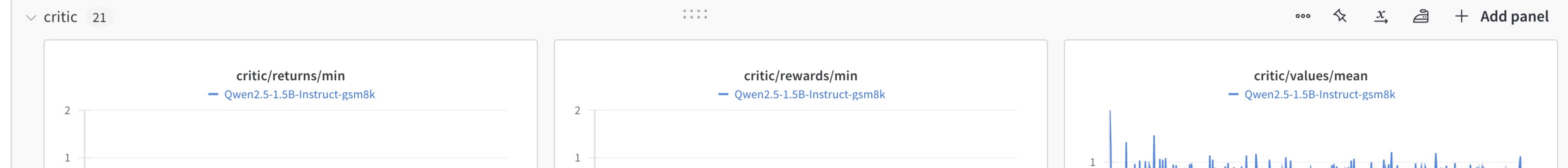
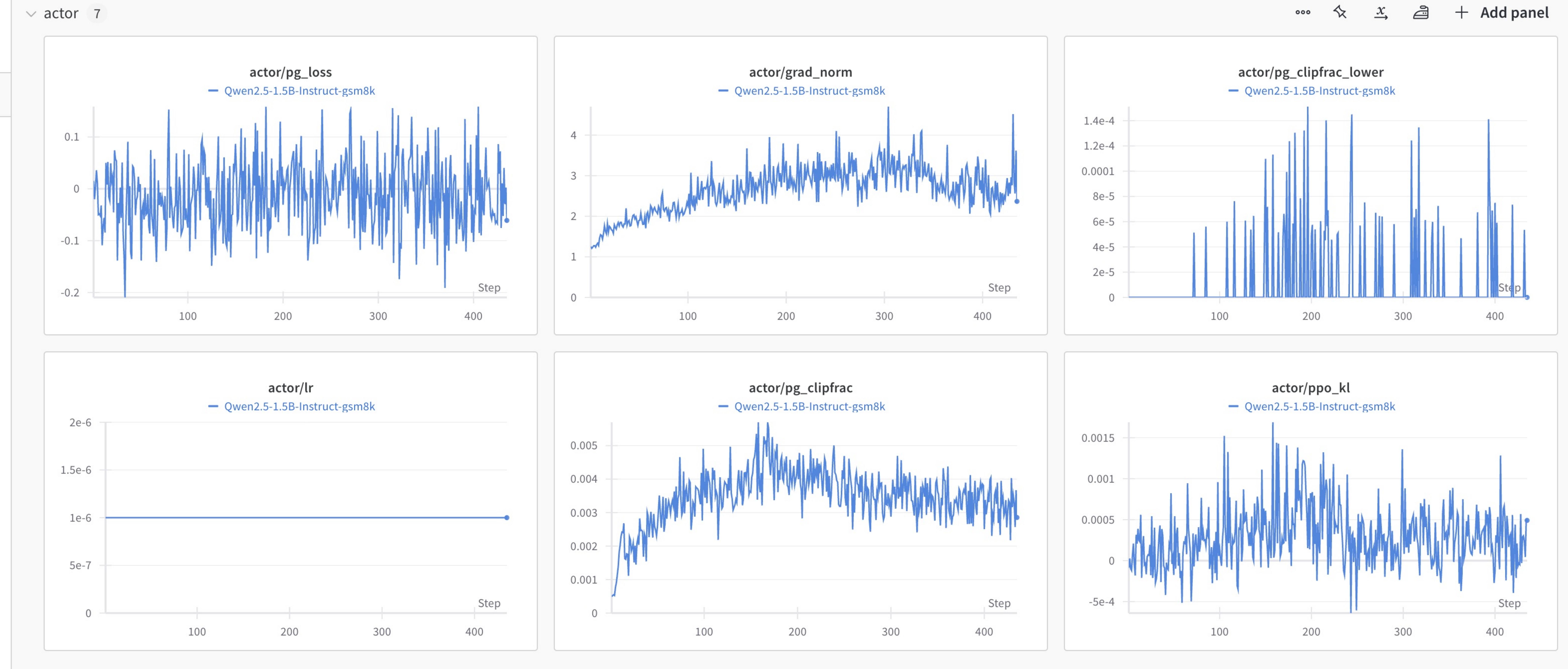
Chenyonghua's workspace Personal workspaceAutosaved just now undo redo more

Runs (1)

Search panels with regex more x print refresh gear Create reportSearch runs \*

Name (1 visualized)

Qwen2.5-1.5B-Instruct-gsm8k

1-1 < >





