



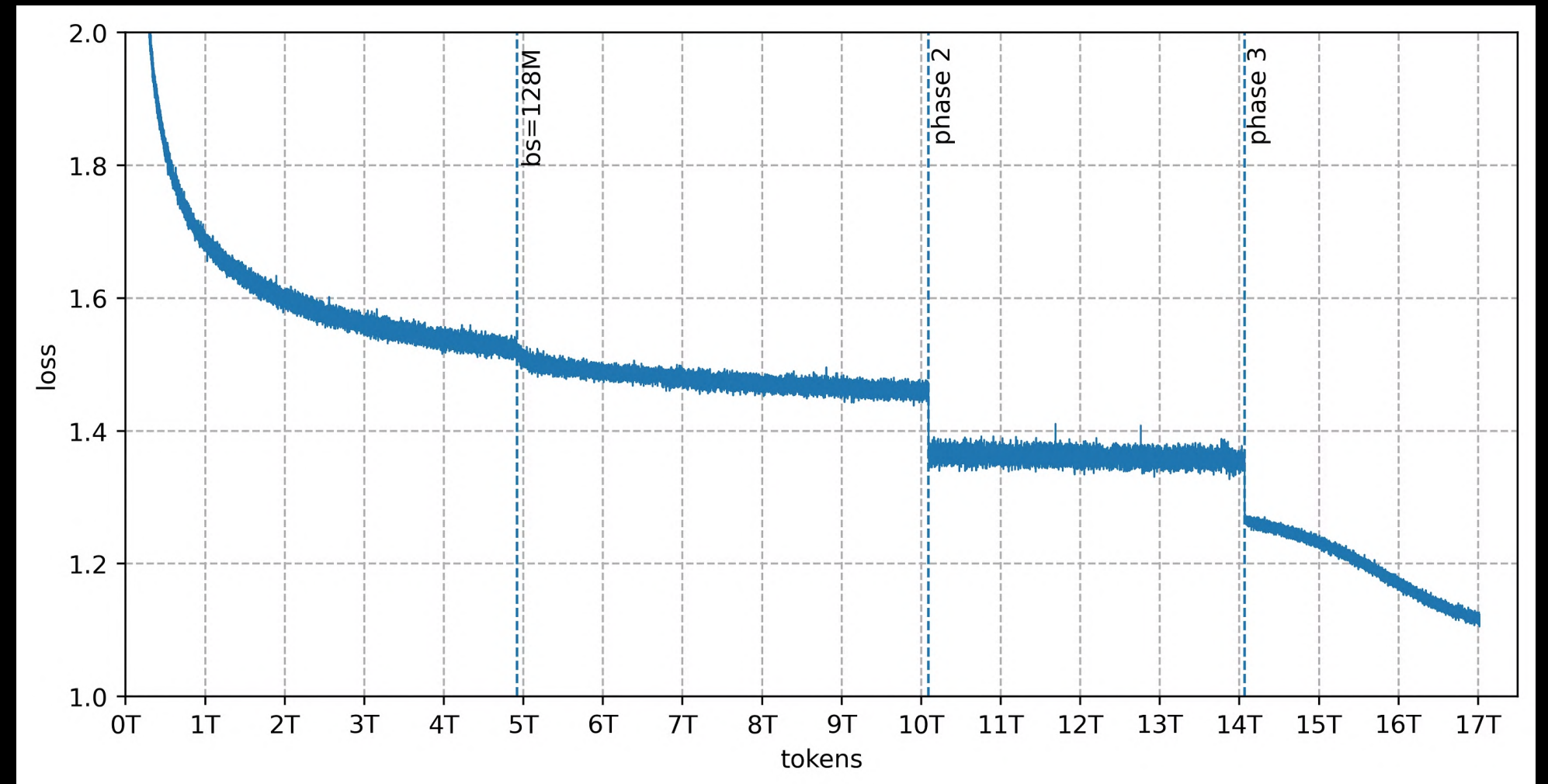
Trinity Large & torchtitan

Matej Sirovatka

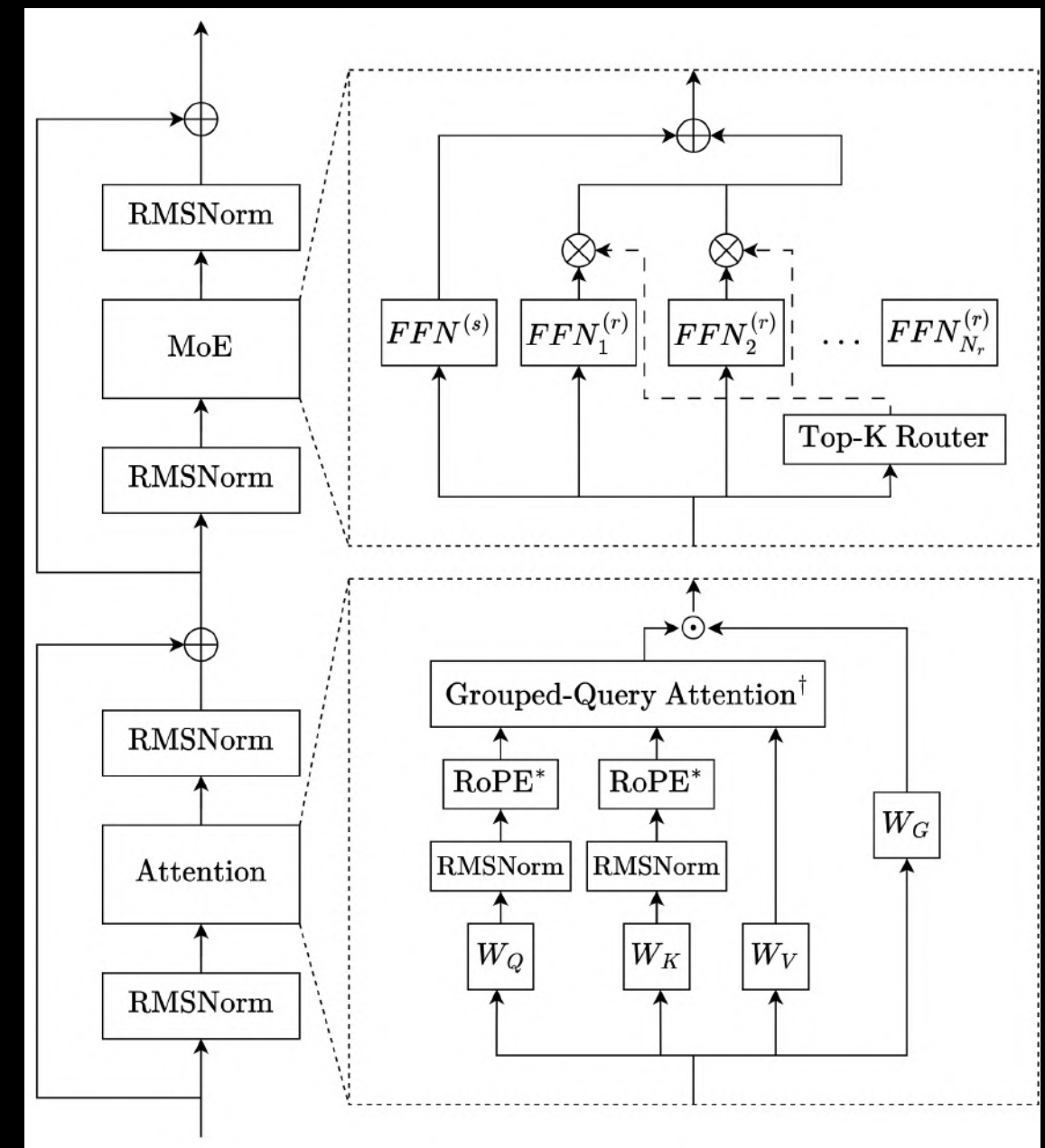
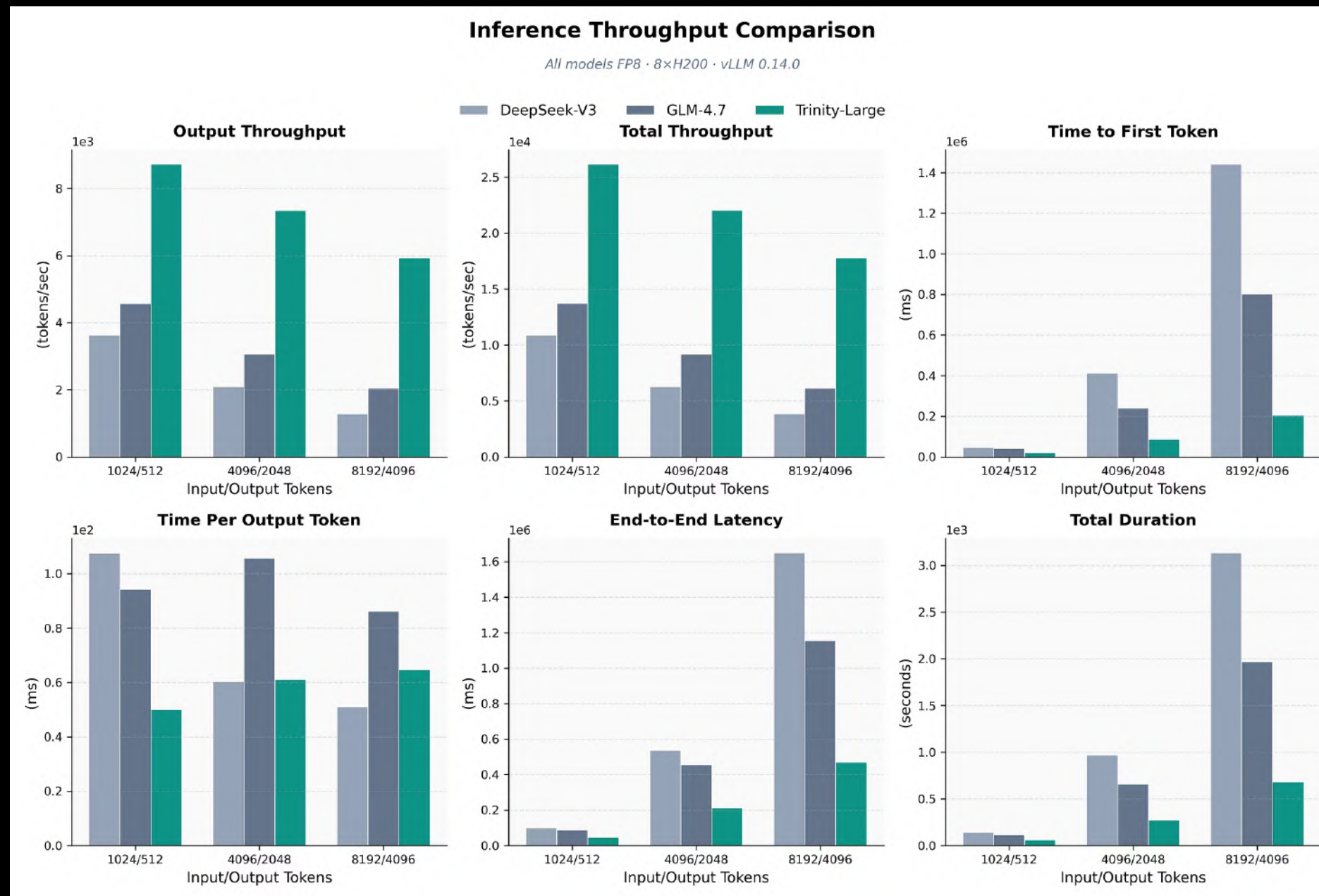


TRINITY LARGE

- 400B/13B
- 4/256 experts
- 1 month
- 2048 B300s
- 20T tokens
- torchtitan only

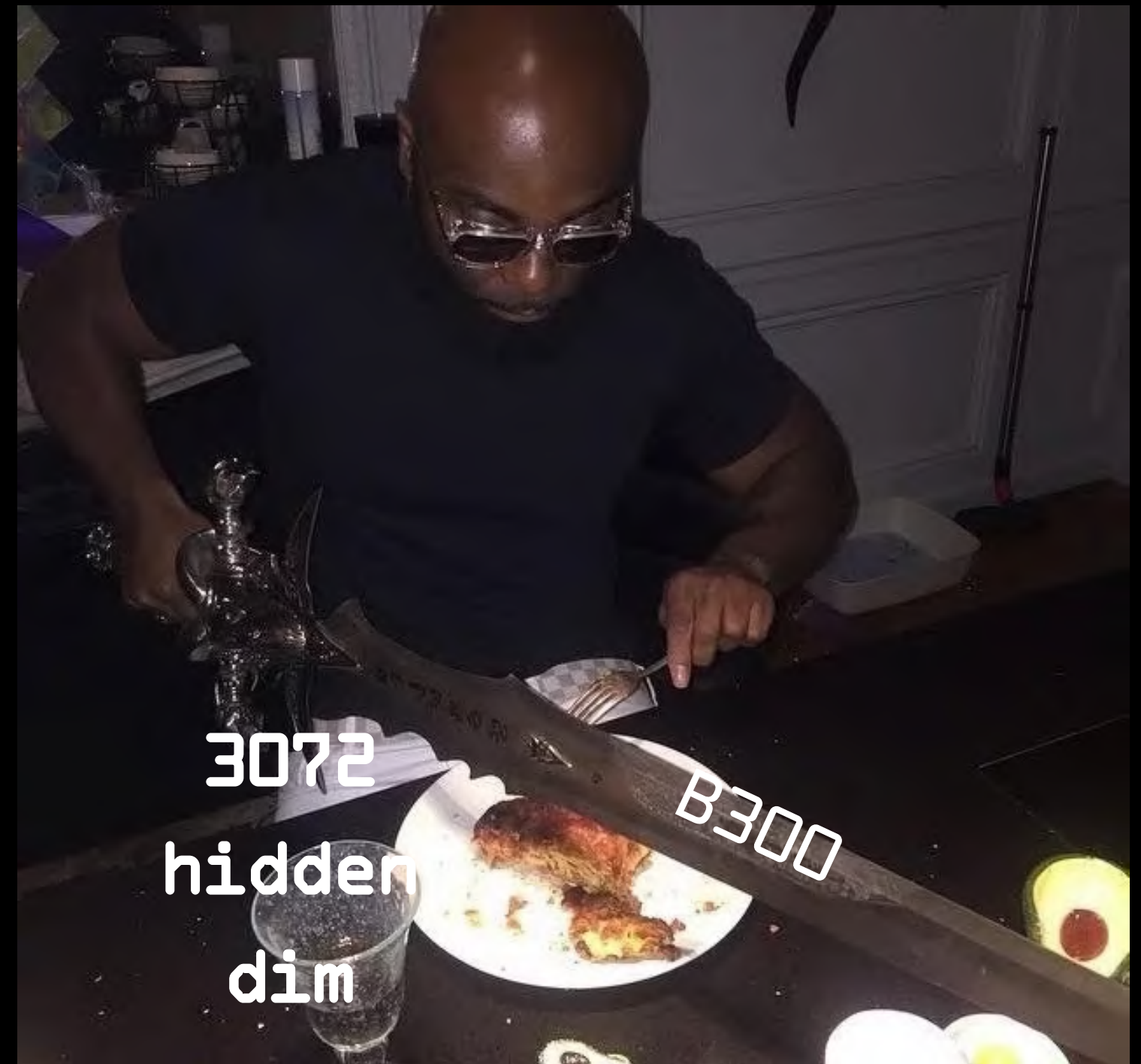


INFERENCE FRIENDLY



WHAT PROBLEMS CAN WE POSSIBLY HAVE

- problem 1 - 2k B300s
- problem 2 - 3072 hidden dim
- problem 3 - 4/256 experts



HOW TO PARALLELIZE



Federico Cassano ✓

@ellev3n11

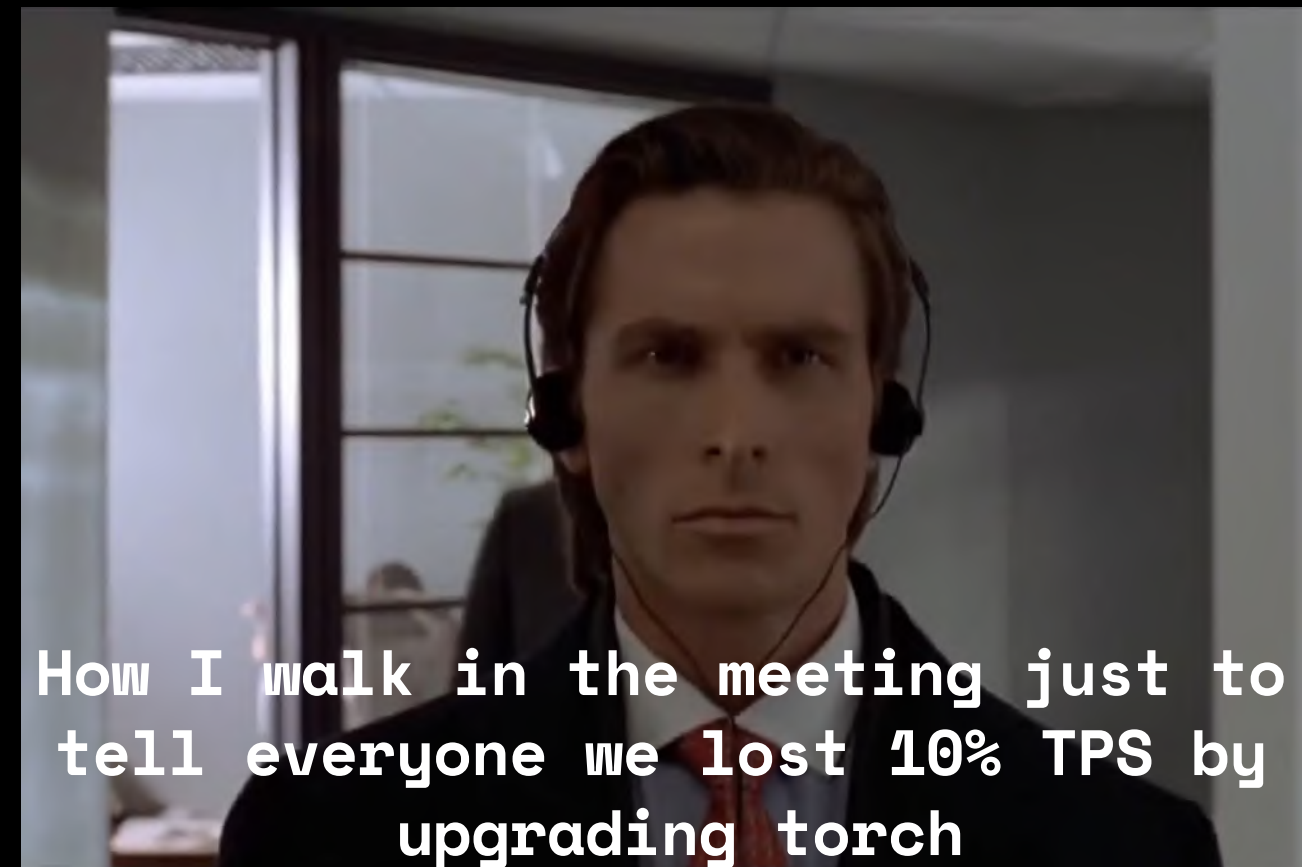
tensor parallelism is a scam



“Was it wise to pipeline? As we now know, pipelining is not wise.”

SETTING UP ENV IS ACTUALLY HARD

- you upgrade torch and lose 10% MFU
- how to actually enable FA?



How I walk in the meeting just to
tell everyone we lost 10% TPS by
upgrading torch

FAULT TOLERANCE

- being first to use a cluster is not so nice in general
- torchFT - nice in theory
- heartbeat monitoring & cold-starts



Me when a betterstack notification pops up

NOW WE CAN JUST TRAIN THE MODEL, RIGHT?

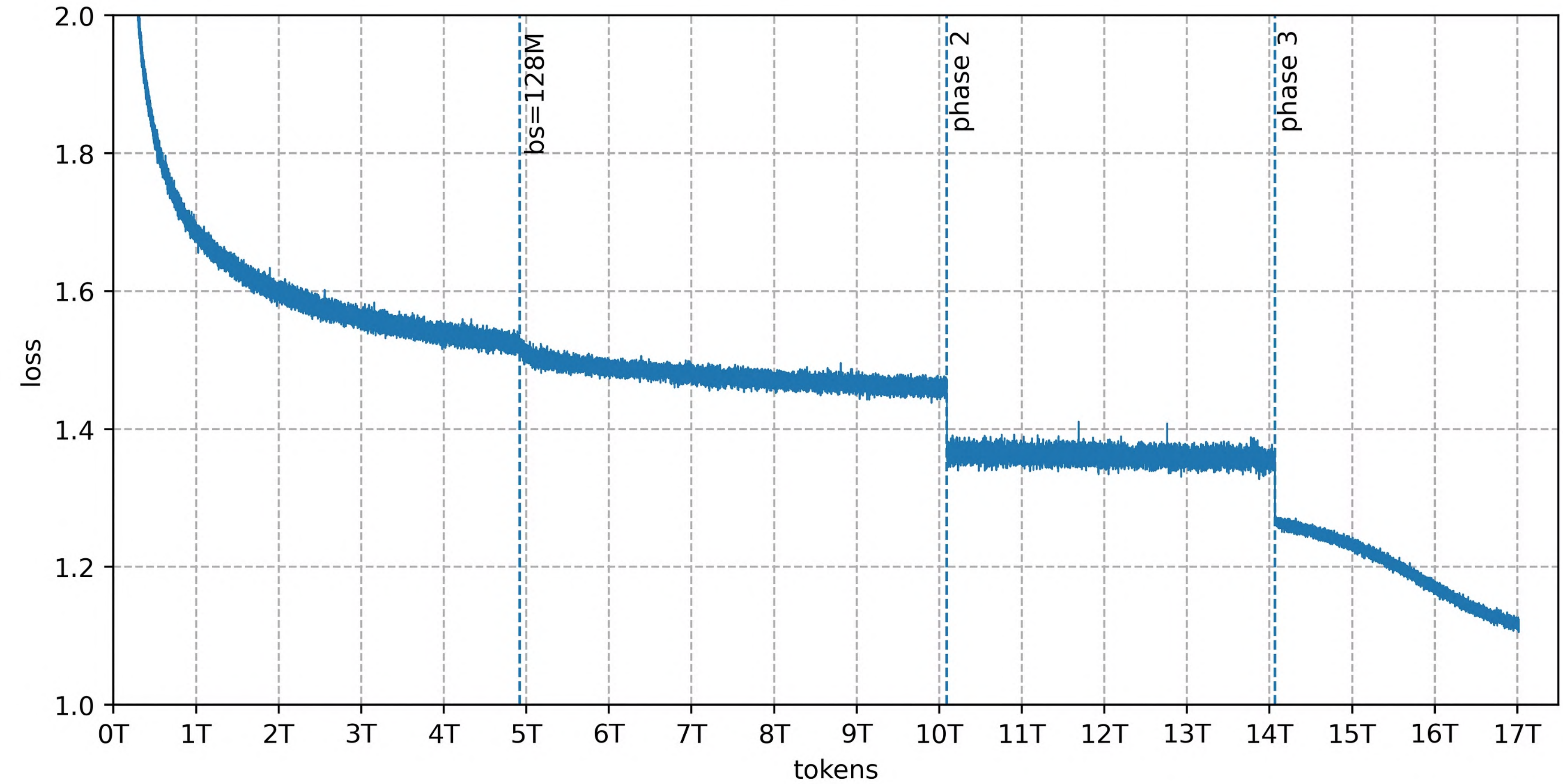
- random expert collapse
- random grad norm spikes

Solution?

- one golden goose run

CONCLUSION?

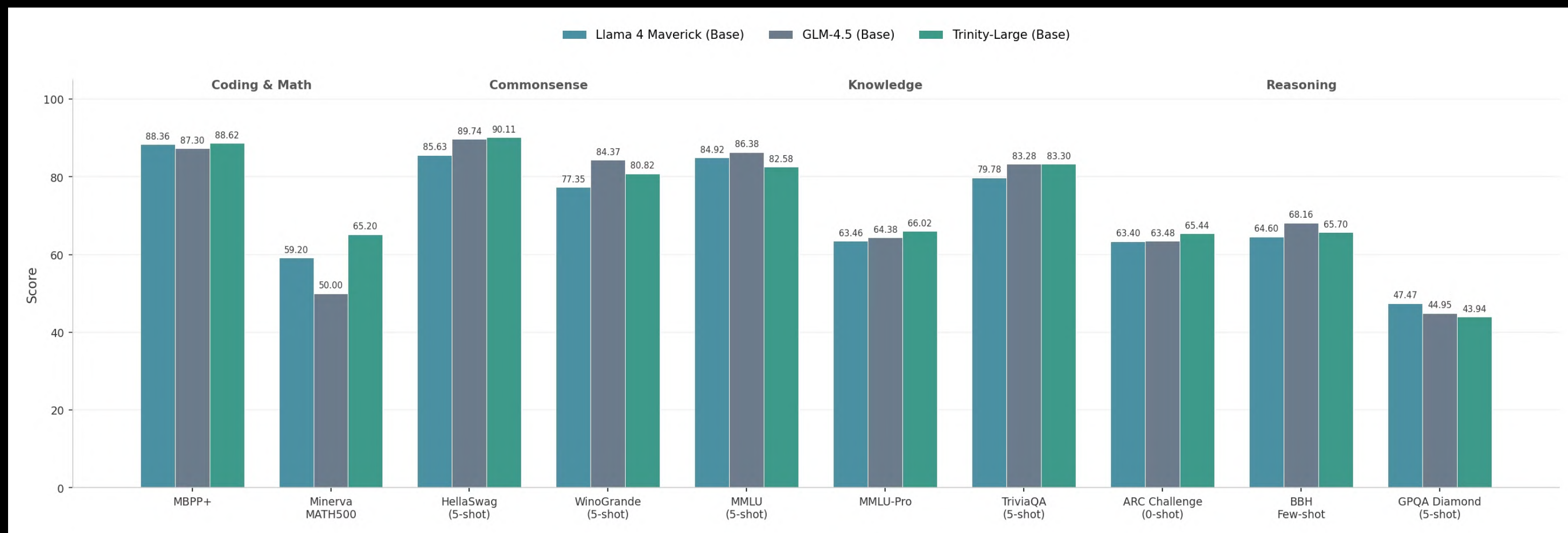
- we managed to train a very sparse model, on a not so stable cluster, with limited time and resources, with a codebase we didn't write
- and it was actually very pleasant



3T tokens on open-router in 50 days

top 4 model on openclaw usage

successful SOTA post-train done by Arcee





GPQA-D

Tau2-Airline

Tau2-Telecom

PinchBench

AIME25



BFCLv4

MMLU-Pro

IFBench

SWE-bench Verified