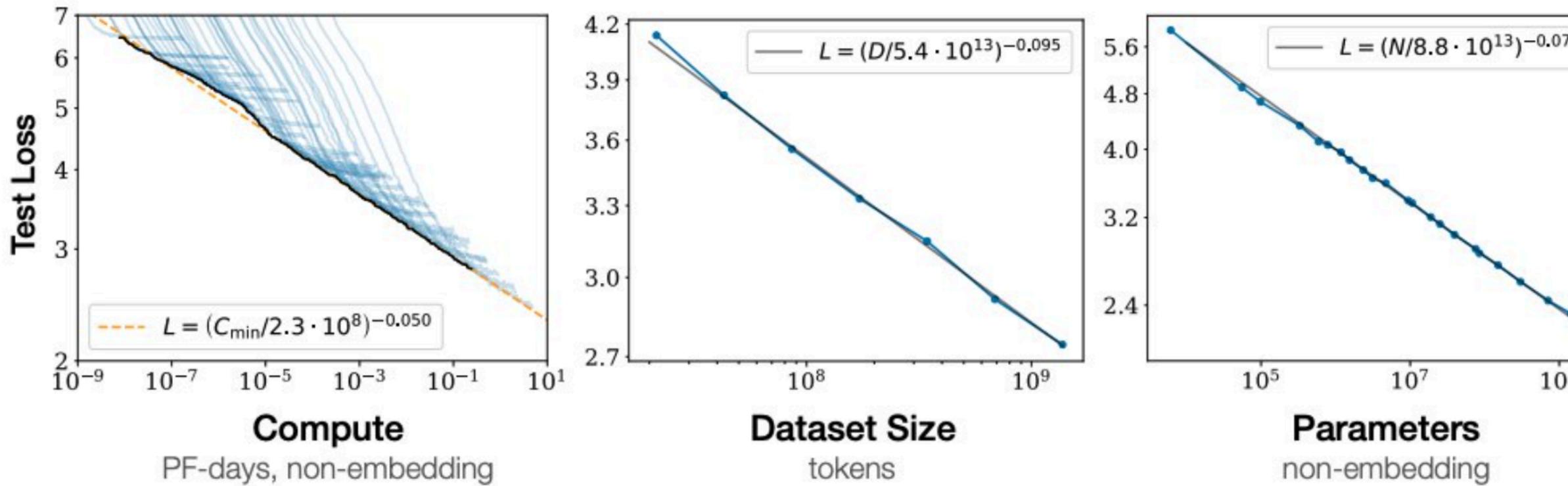


Decentralised Learning

AKA How to train really large models on consumer
devices over the Internet

MIKA
SENGHAAS

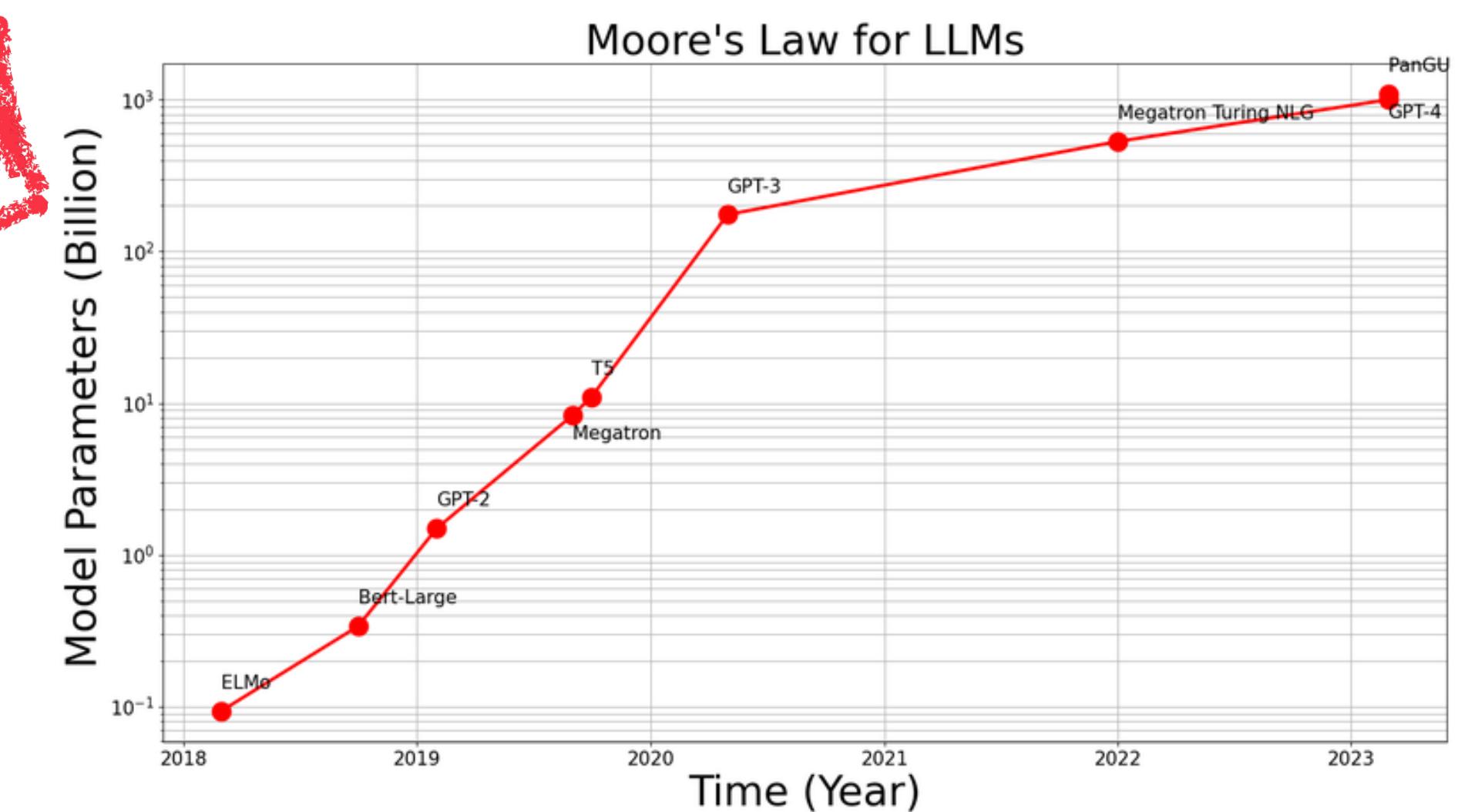
Scale



Kaplan et.al (2020),
Hoffmann et.al (2022)

LLMs improve as function of scale
predictably

We train larger models
↳ 100B+ (1000B+?)

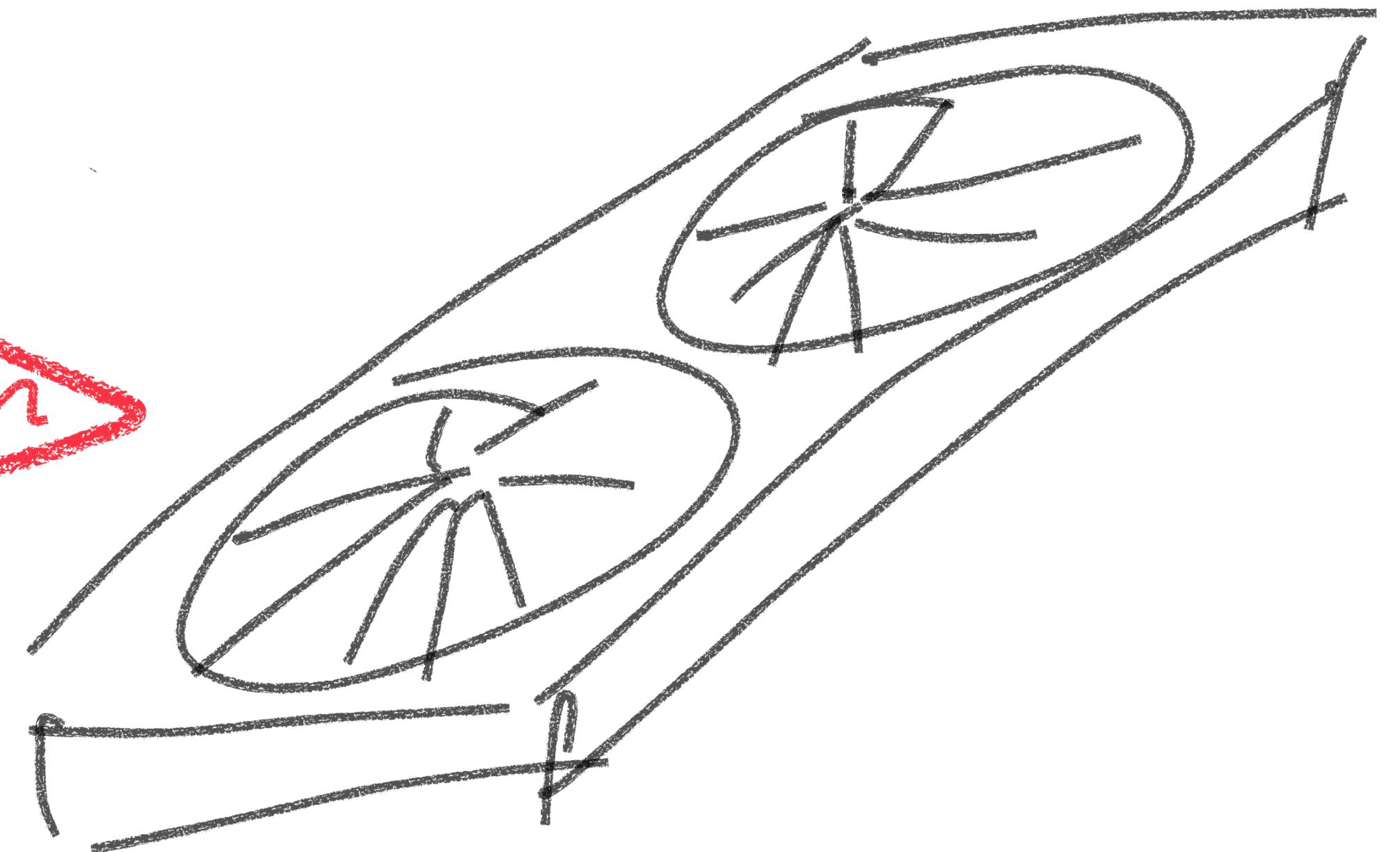


Scale

Llama - 3.1
405B

But really
anything > 1GB

Doesn't fit
~~1000~~

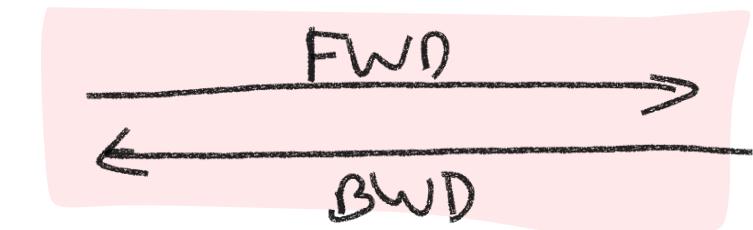


A100 / H100

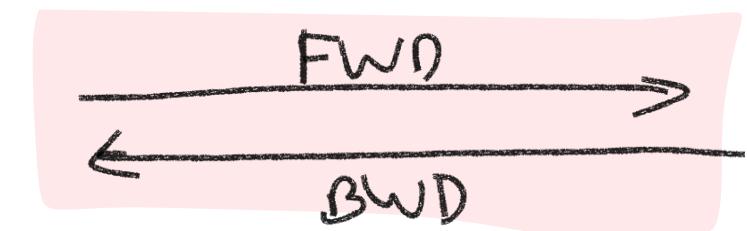
10B at BF/FPI6 \approx 20GB · 3 \approx 60GBt
+ Activations

Parallelism Paradigms

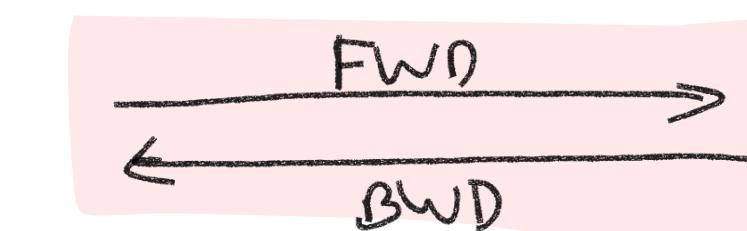
D_1



D_2



D_3



P_1



P_2

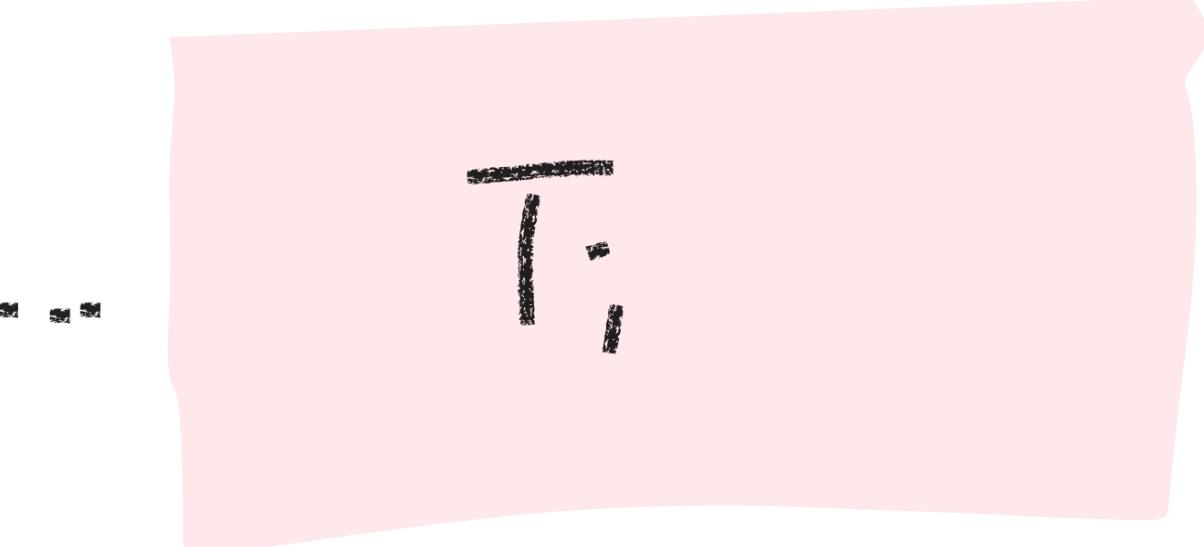


P_3



...

T_i



....

$$A \times B_1 B_2 =$$

$$\begin{matrix} AB_1 \\ \oplus \\ AB_2 \end{matrix}$$

Data Parallel

DP

Pipeline Parallel

PP

Tensor Parallel

TP

Parallelism Paradigms

Data Parallel
DP

Pipeline Parallel
PP

Tensor Parallel
TP

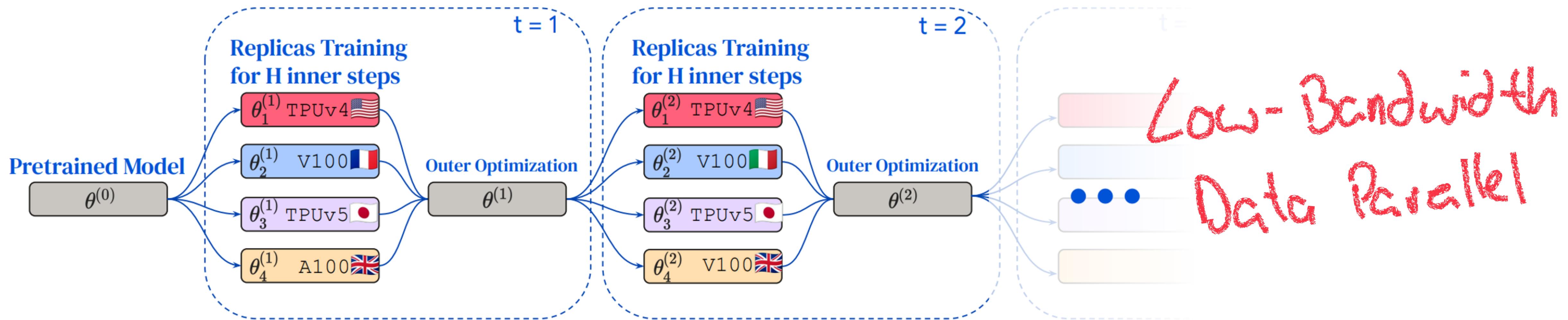
WORK IN **DATA CENTRE**
BUT NOT
DECENTRALISED

Decentralised Setting

- ① Heterogeneous Hardware
- ② Heterogeneous Network
- ③ Reliability
- ④ Scale



DILoCo (Dovilas et al., 2023)

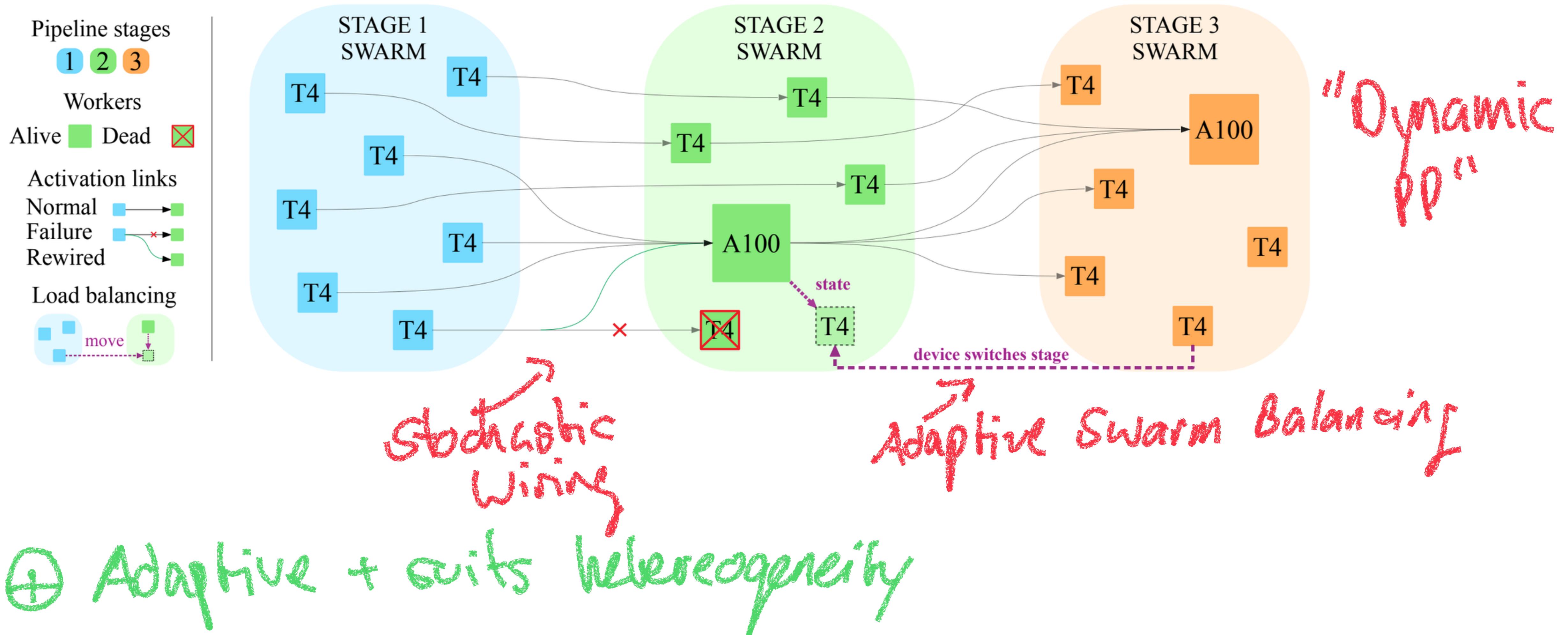


- ⊕ Works for islands of data centres
- ⊕ Open source

- ⊖ Model has to fit in min. GPU mem
- ⊖ Synchronous global update (\neq heterogeneity)

SWARM

(Ryabinin et.al, 2023)



Research Directions

- ① Benchmarking
- ② Battle-test SWARM
- ③ Async DiLoCo
- ④ Multi-parallel Trainings

More Ambitions
Less % Success