

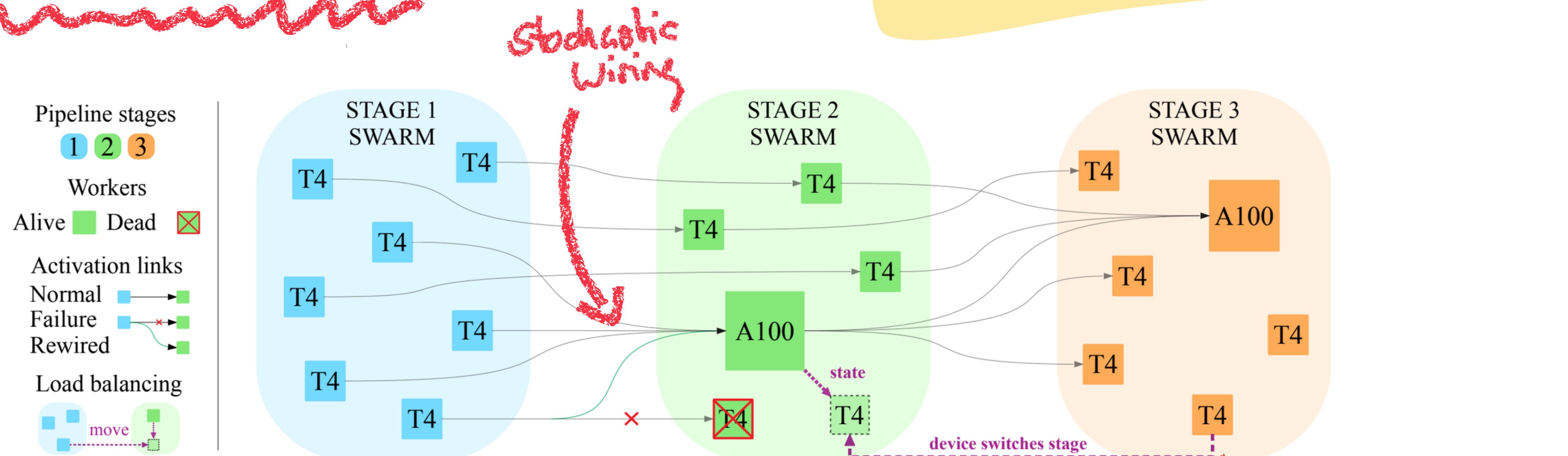
BATTLE-TESTING SWARM

A RESEARCH OUTLINE

MIKA
SENGHAAS

EPFL

SWARM



DYNAMIC + REDUNDANT PP

Decentralised Setting

① Heterogeneous Hardware

② Heterogeneous Network

③ Reliability

④ Scale

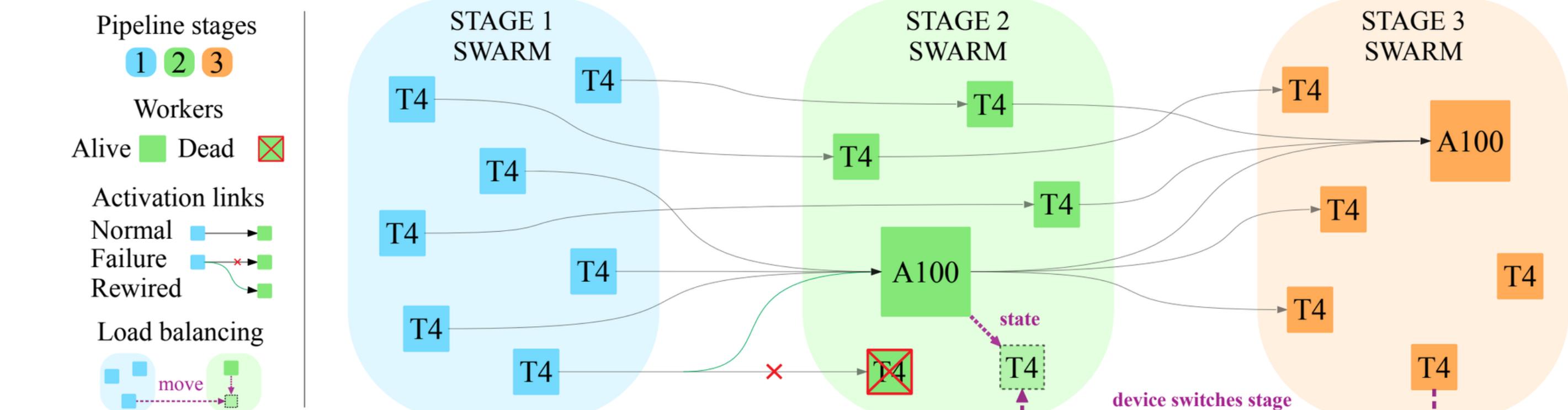
↳ Infini-k Model
Data Scale

↳ On-Off
Ramping



Decentralized SWARM

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



LET'S BENCHMARK

TASK

Model

Objective

Data

Algo

Baseline 1

Baseline 2

ALGORITHM

SETTING

Worker

Network

Throughput (T/s)

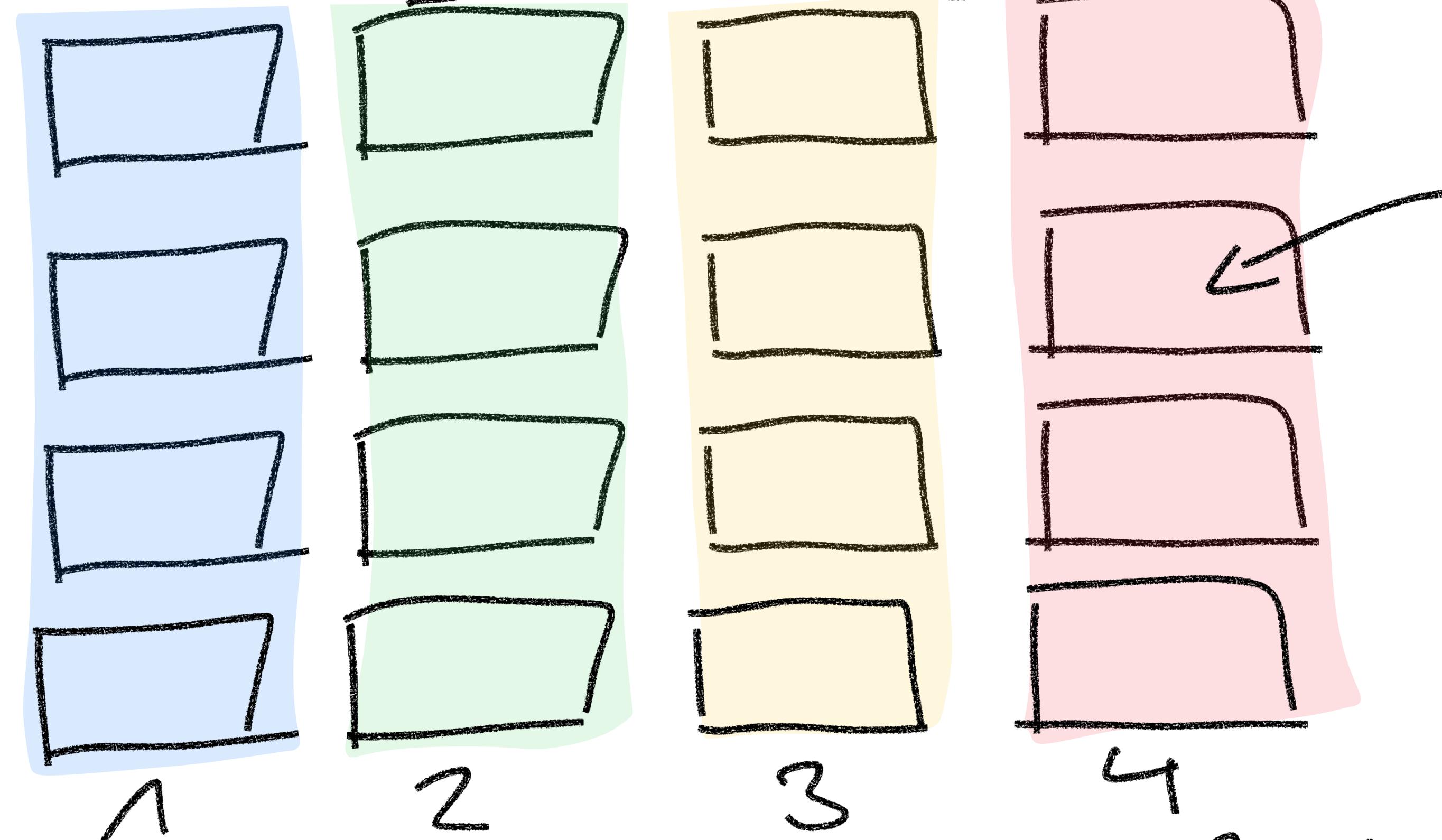
Utilisation (%)

Covergence
(h/step)

METRICS

SWARM Experiment 1

No latency / $100 \pm 50\text{ms}$



V100 (28GB)
holding 4 layers

Metrics

Throughput (Time / 6250 ex)

AllReduce (Time)

Baselines

GPipe

Zero offload

Models

16 Transformer layers w/

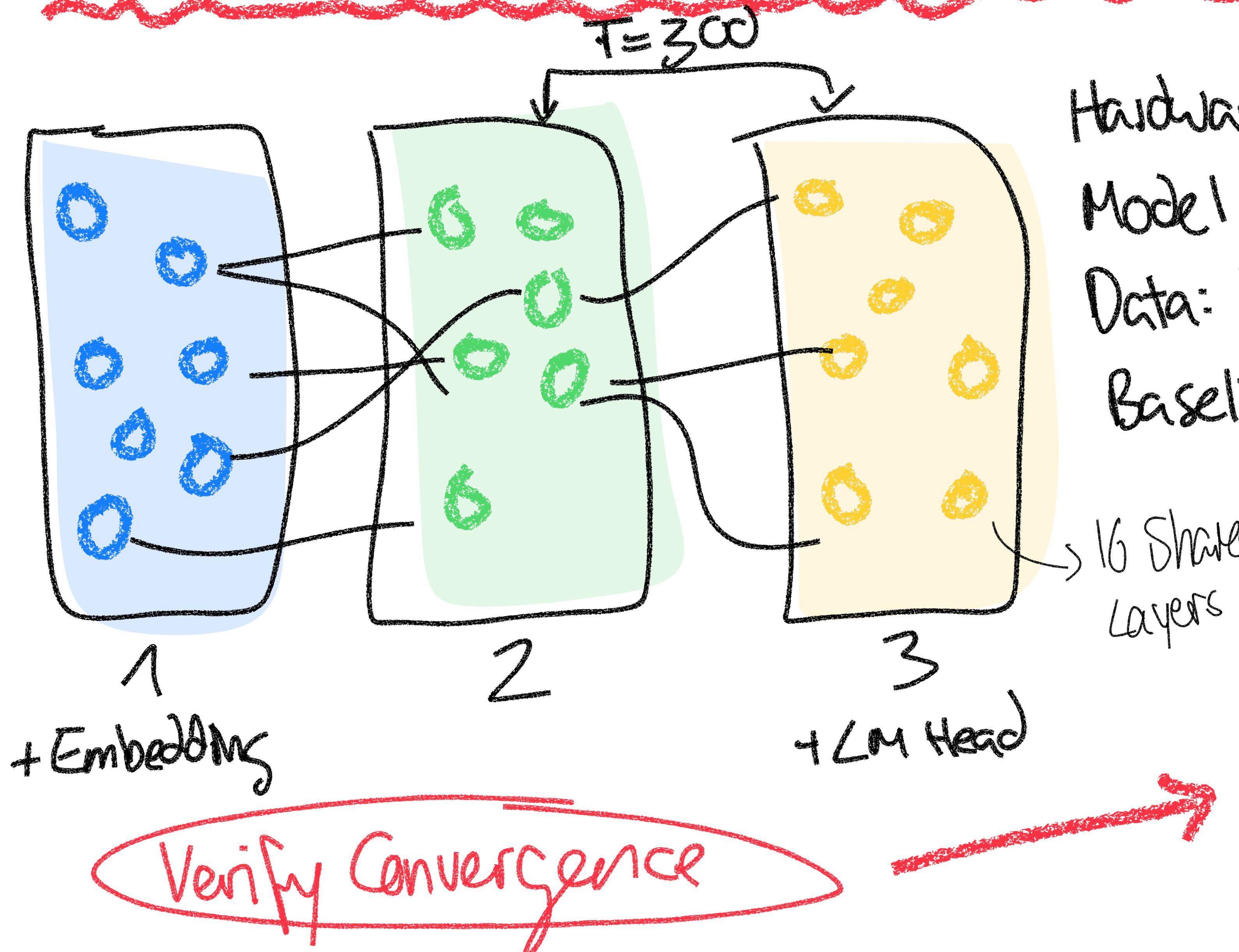
① $d=4096, h=32, b=4$

② $d=12288, h=96, b=1$

③ $d=4096, h=32, b=1$

+ Scale of
Comm. Efficiency
Experiment

SWARM Experiment 2



Hardware: 400+14

Model: 1.1B (effective 13B)

Data: Pile

Baseline: DDP + offloading on 128 A100

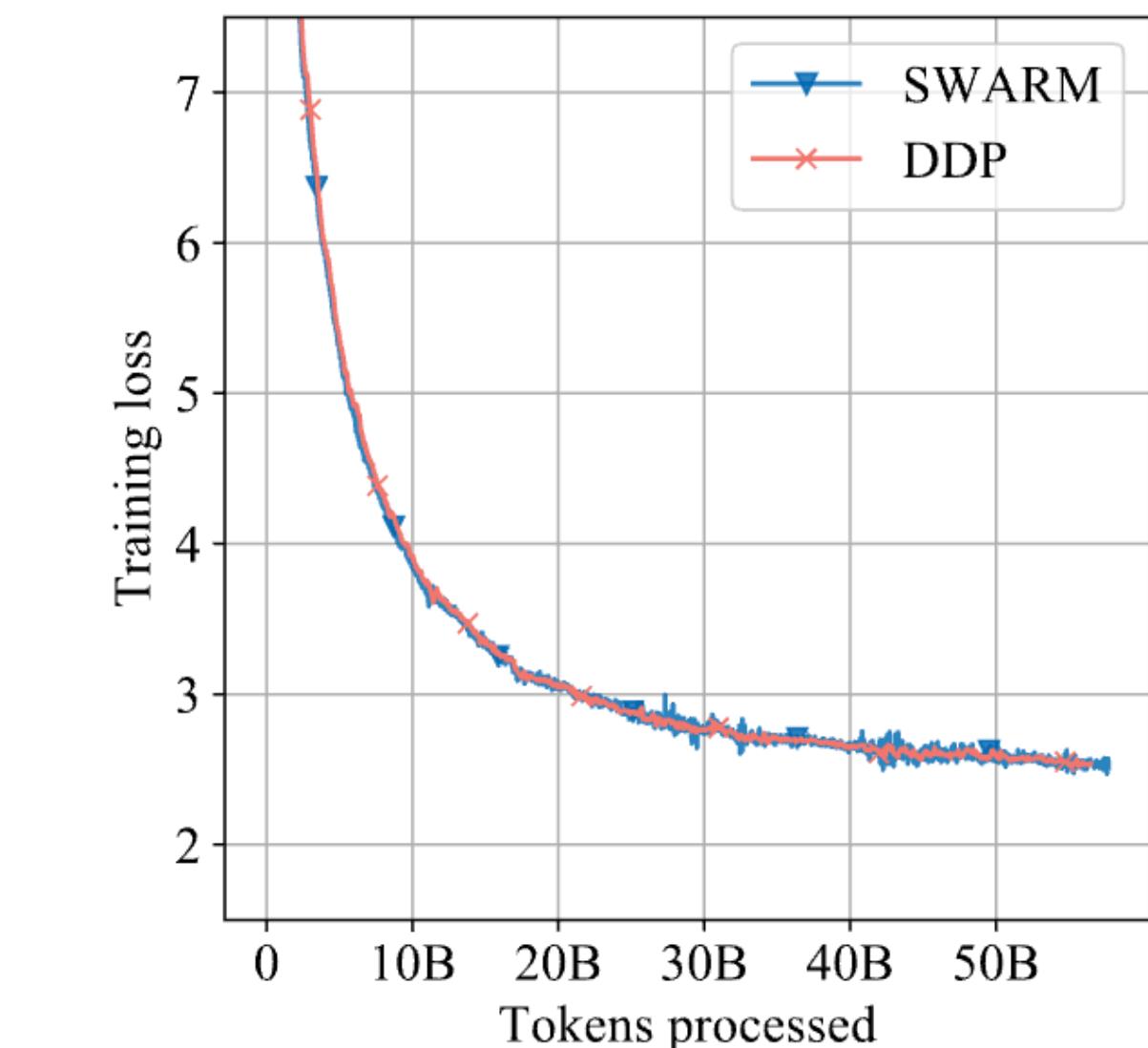
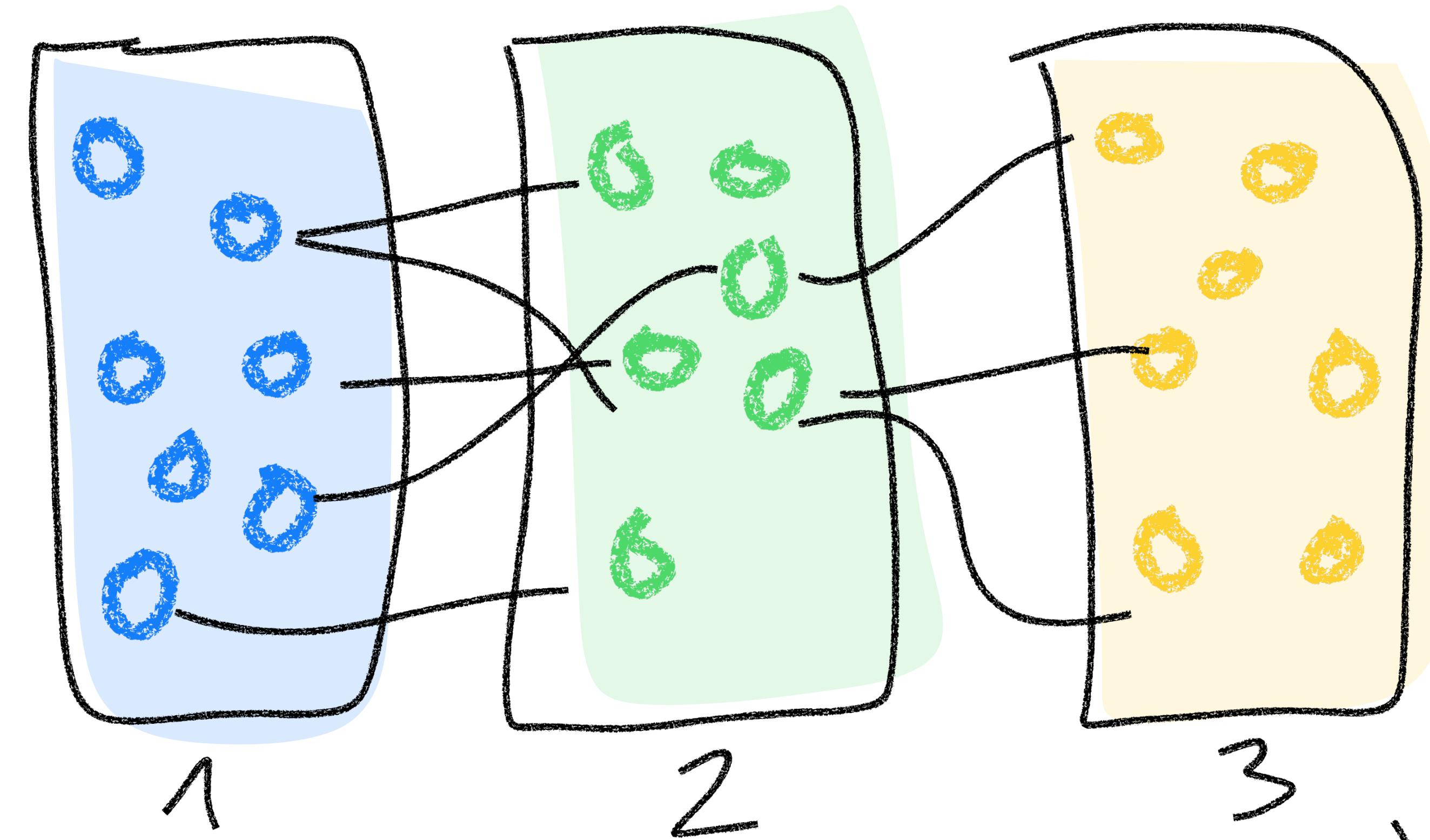


Figure 4: Training convergence comparison.

SWARM Experiment 3



1
+ Embedding

2

3
+ LM Head

Hardware: 400T4 / 7x8A100

Model:

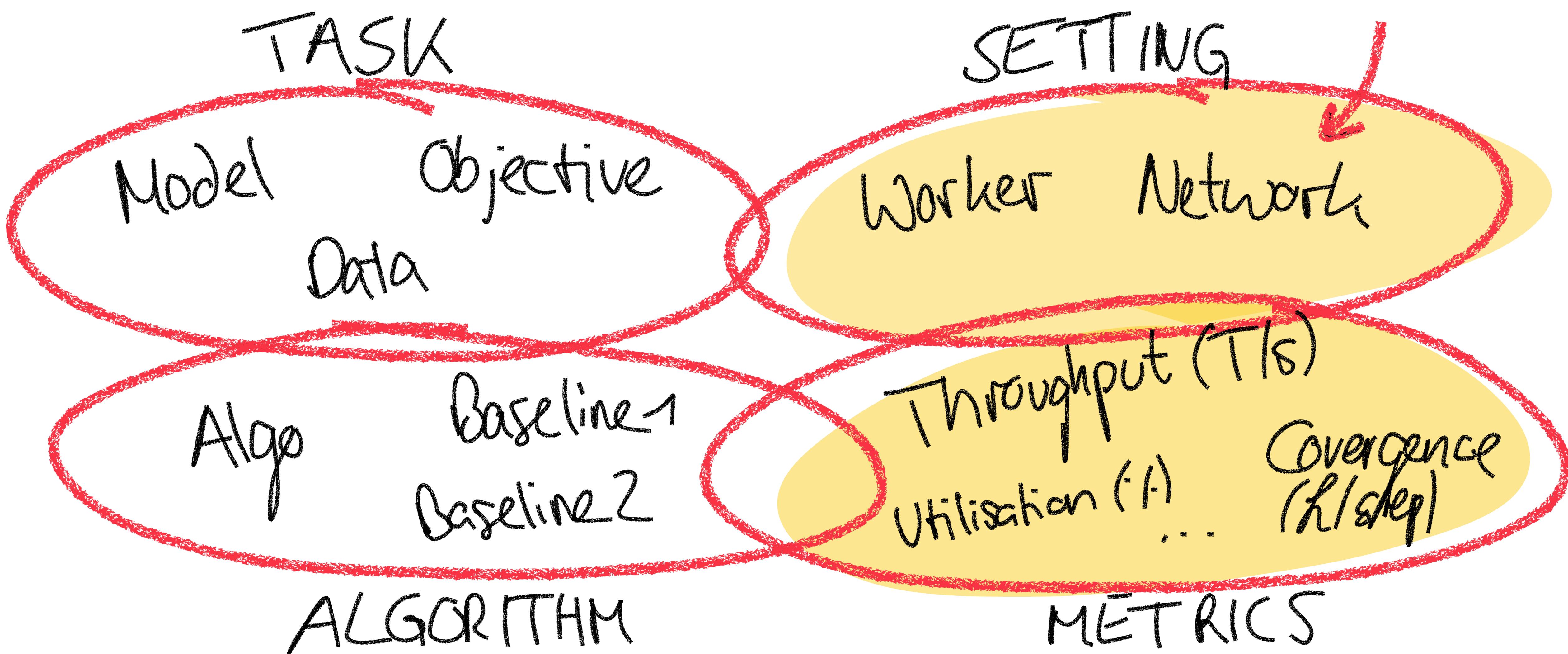
Data: Pile

Baseline: Ideal Throughput

(Avg. Processed Tokens/Stage ×
3 × Num Nodes)

+ Experiment
on Rebalancing

Focus



Next Steps

Hivemind / Qtouch / zeroboard

→ Home

- ① White-up Experiment Setup → CY
- ② Implementation

baselines ↗ T4 / A100 /
H100

↳ Existing Frameworks / Code?

↳ MVP Sandbox Implementation

↳ Distribute step-by-step

- # Questions
- ① SWARM
 - When to sync what inter-layer?
 - How to sync intra-layer?
 - Scheduling of backward path?
 - ② Tooling
 - Hivemind
 - OpenDILoCo, Petals
 - ③ Resources
 - Prime Compute
 - Local EPFL cluster?
 - ④ Metrics
 - What is relevant?