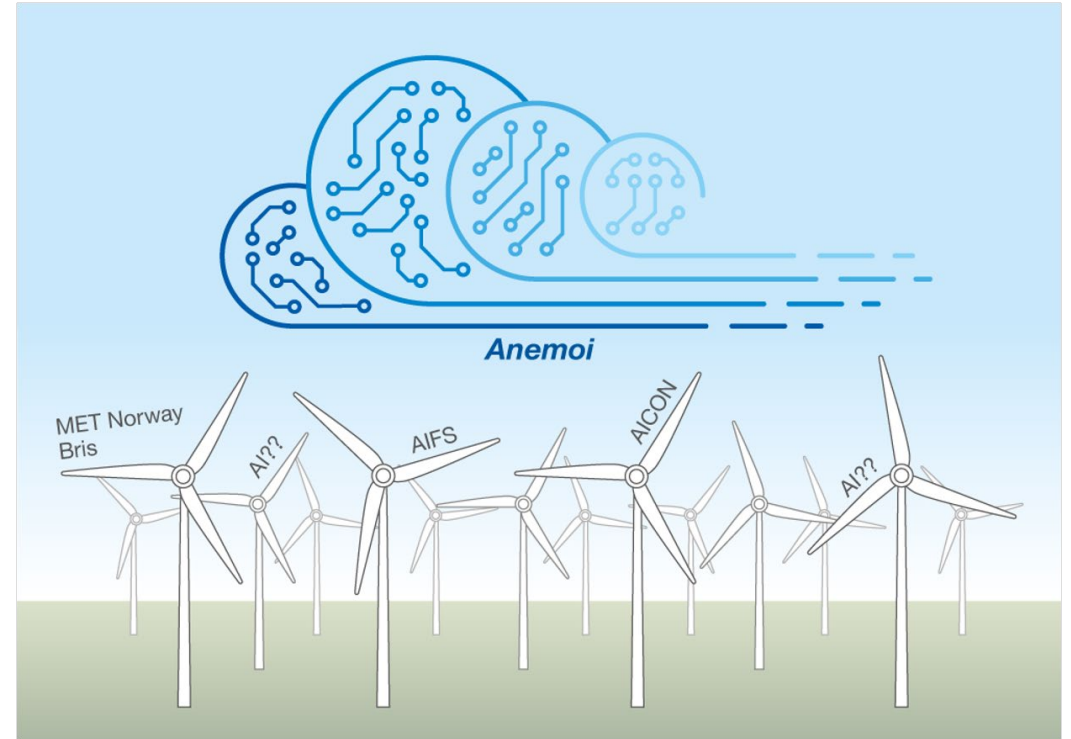# Anemoi scalability improvements

21st HPC workshop
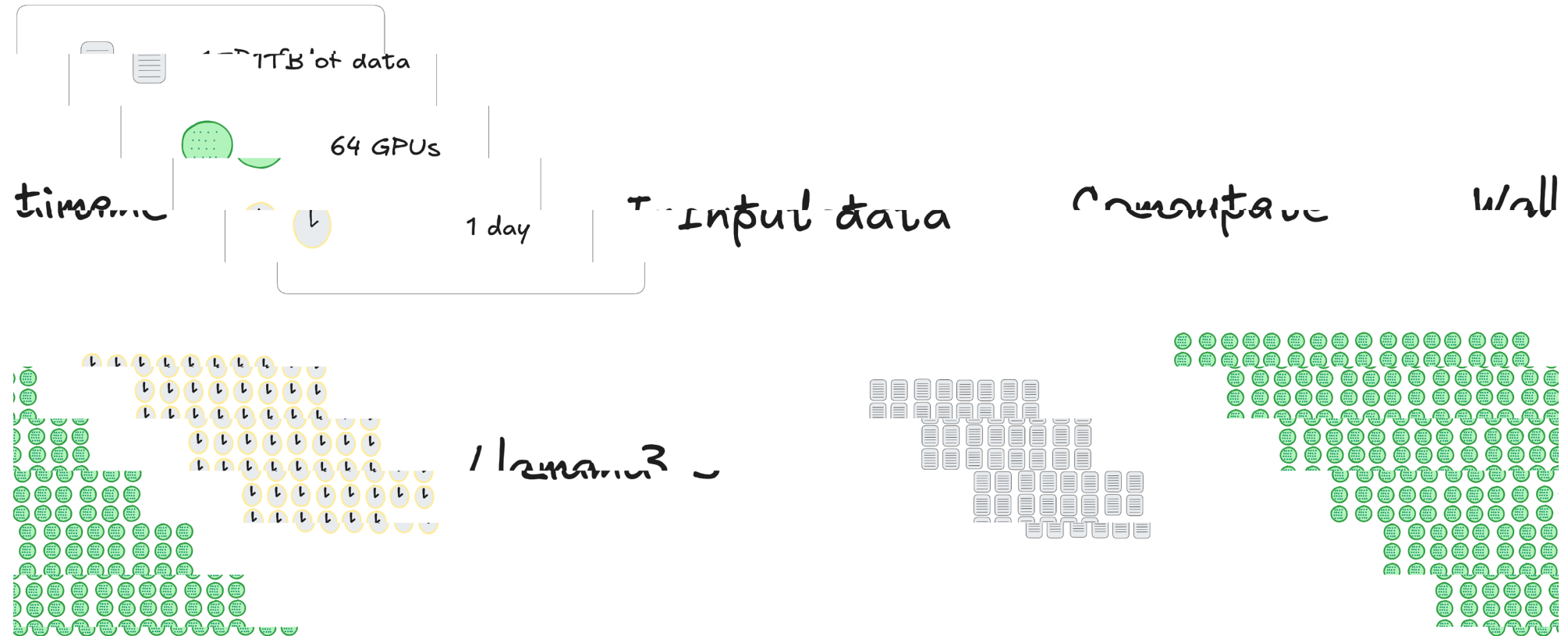
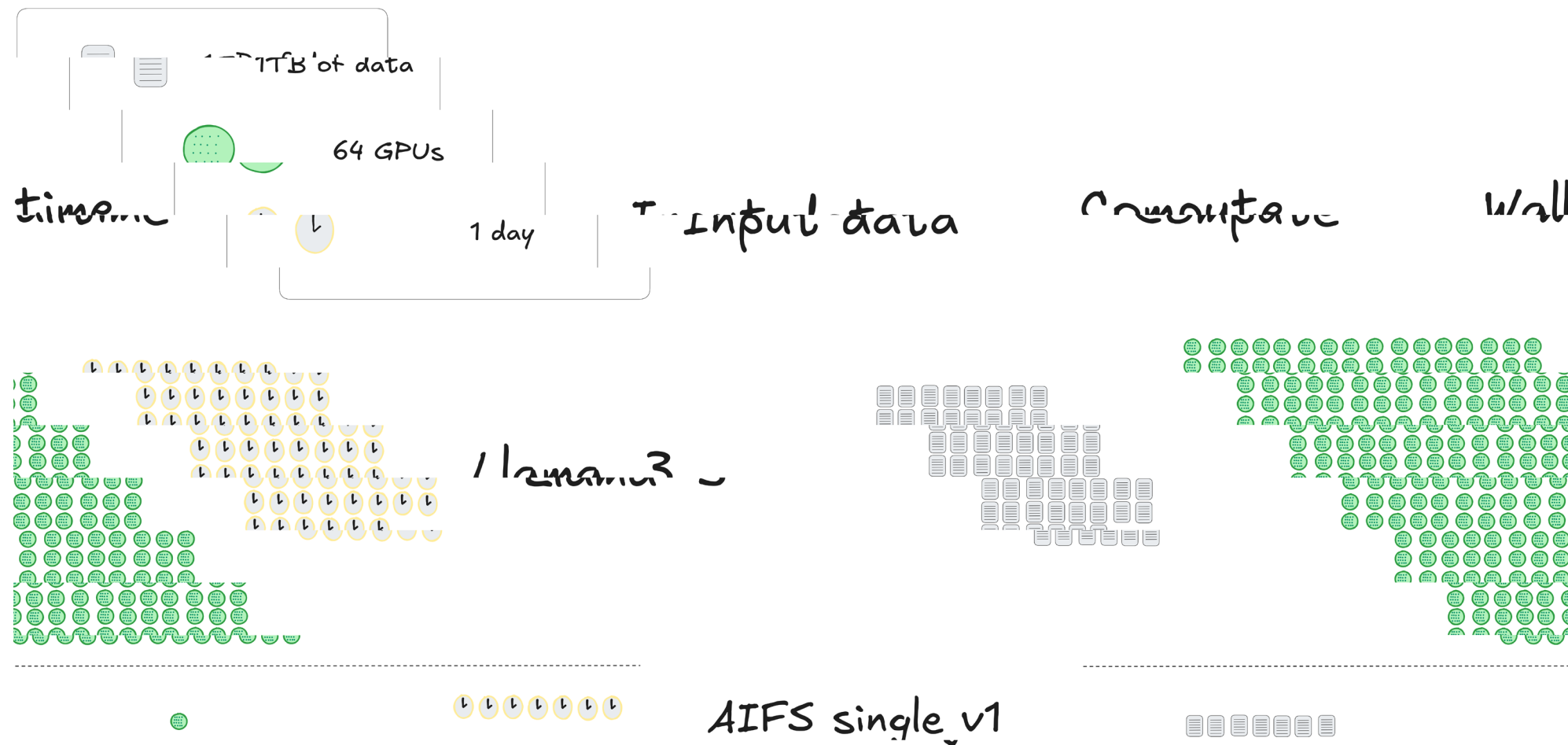Jan Polster and Cathal O'Brien

**ECMWF**

# What is Anemoi

- A framework for developing data-driven weather forecasting models

- Entire framework from datasets to inference

- Multiple use-cases
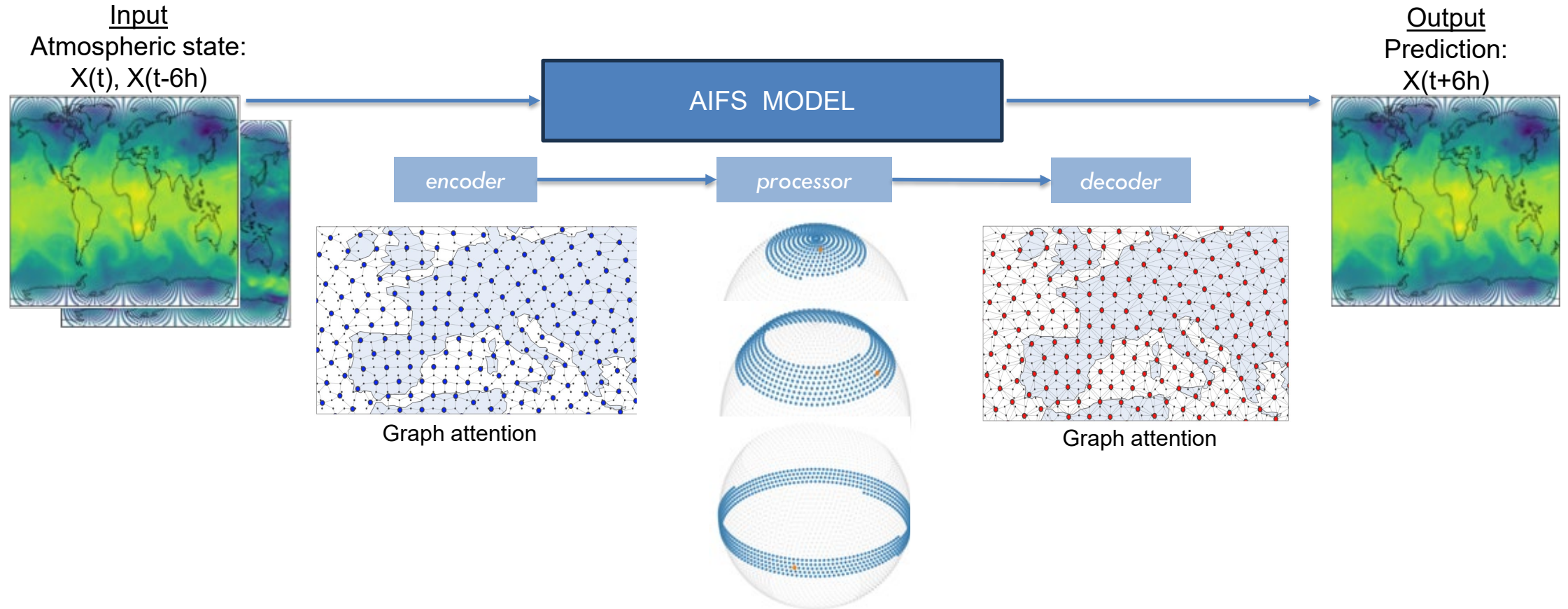    - Global, LAM, stretched grid, ensembles

# 2024: A tale of two models



1TB of data

64 GPUs

1 day

time

input data

compute

wall

# 2024: A tale of two models



1TB of data

64 GPUs

1 day

time...     Input data     Compute     Wall

AIFS single_v1

# AIFS: Data-driven Weather Forecasting Systems

OPERATIONAL AIFS SINGLE

Input
Atmospheric state:
X(t), X(t-6h)

AIFS MODEL

Output
Prediction:
X(t+6h)

encoder → processor → decoder

Graph attention

Graph attention

Transformer blocks and windowed
attention (attention across regional bands).
On a coarser grid than input grid
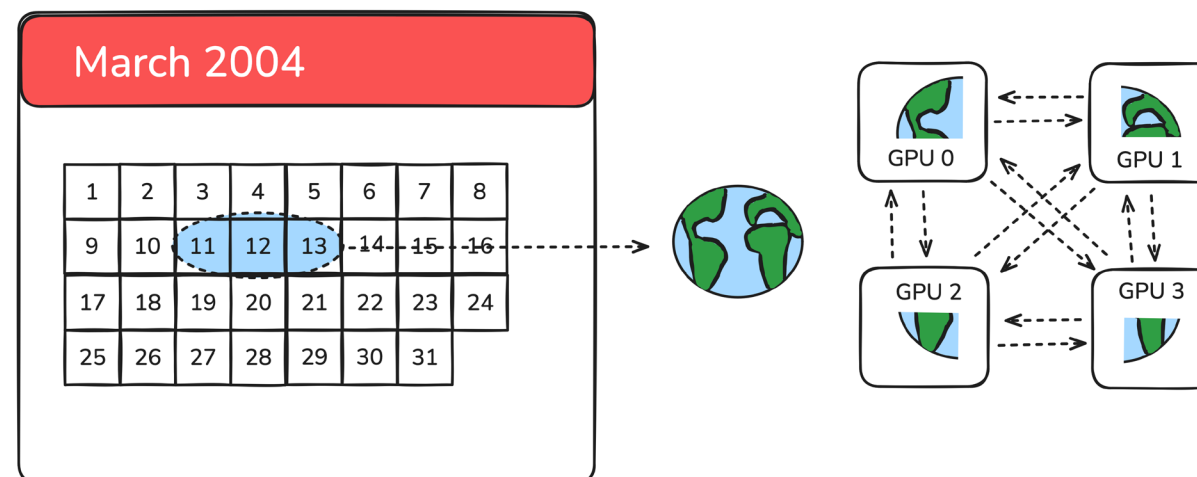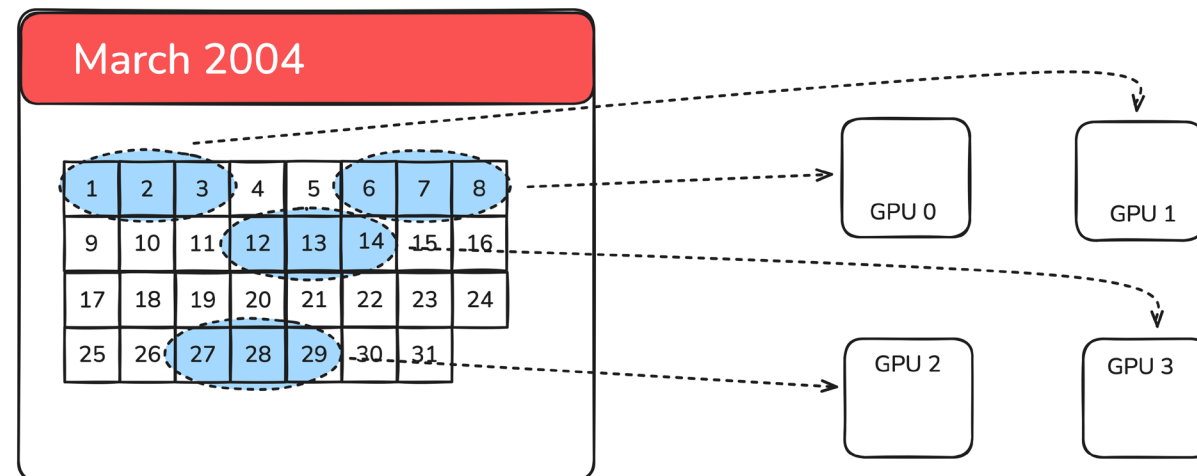
Implemented in

anemoi

# Parallelism in Anemoi

- Data Parallelism: DistributedDataParallel (PyTorch)

  o Distribute training batch across model replicas

  o Aggregate gradients via all-reduce, good scaling :)

  o Limited by batch size :(

- Model parallelism: domain-specific sharding

  o Distribute input data and activations across GPUs

  o Collective communication to handle synchronisation

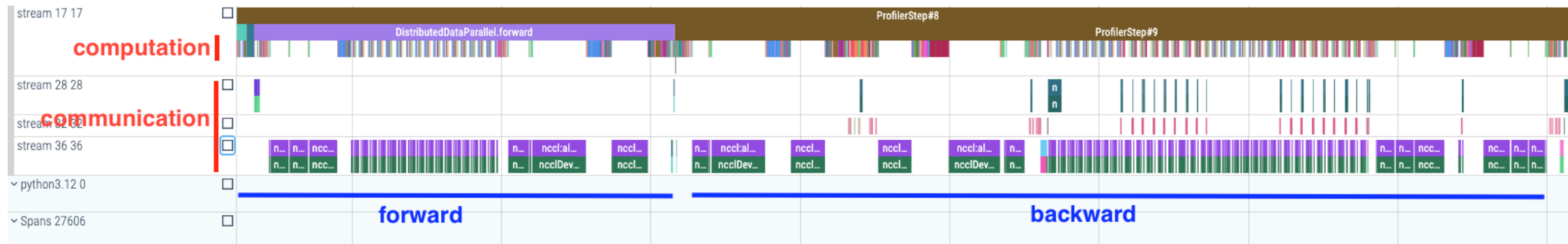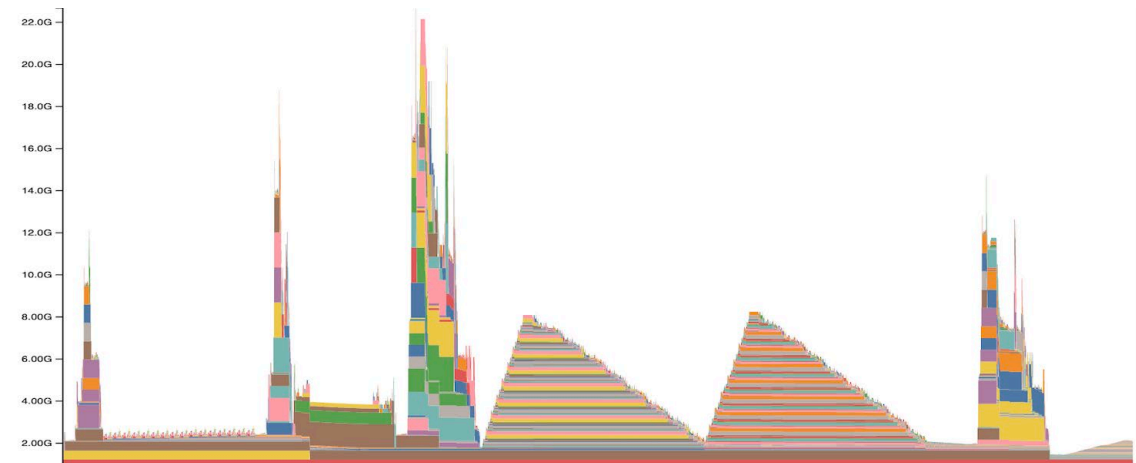  o Limited by communication overheads

# Scaling the input resolution

| Resolution | Batch Size |
|---|---|
| O96 (100km) | 0.05 GB |
| N320 (32km) | 0.65 GB |
| O1280 (9km) | 7.40 GB |
| O2560 (4km) | ~30 GB |

- **Bottlenecks**:
  1. Loading/storing full training batch on device
  2. Encoder/Decoder: communication and memory overhead
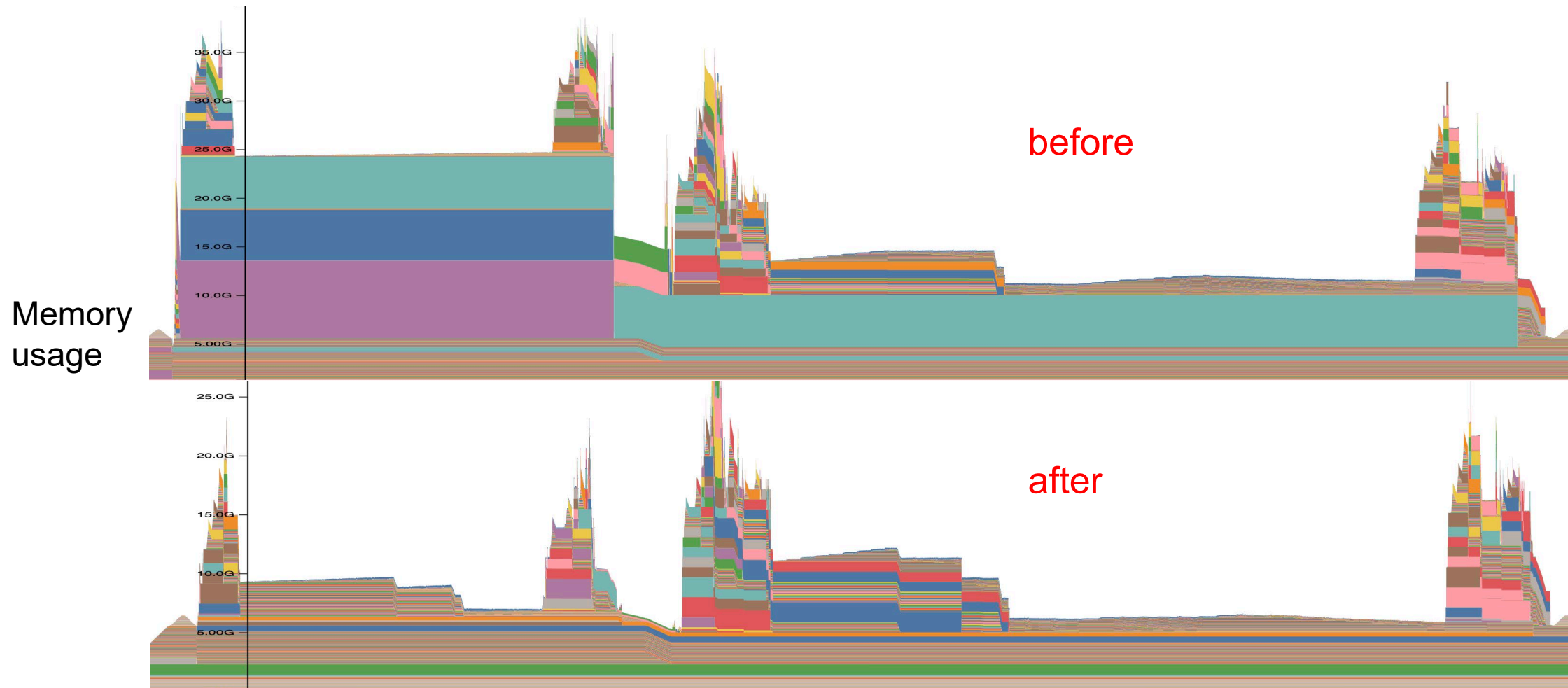


- **Improvements:**
  1. Keep batch sharded
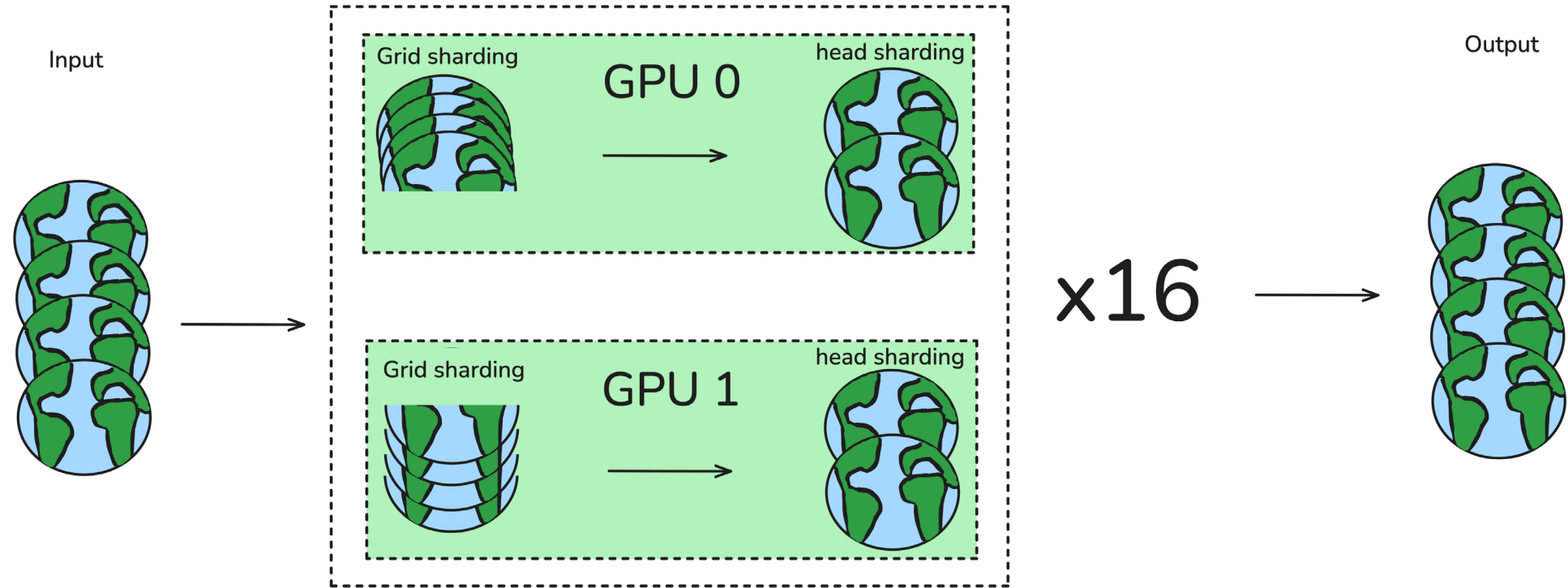  2.1. Mapper edge sharding
  2.2. Mapper chunking

# Keep batch sharded

- Avoid materialising full input/output grid in memory

  - Load batches in shards

  - Compute loss locally + all-reduce for global loss

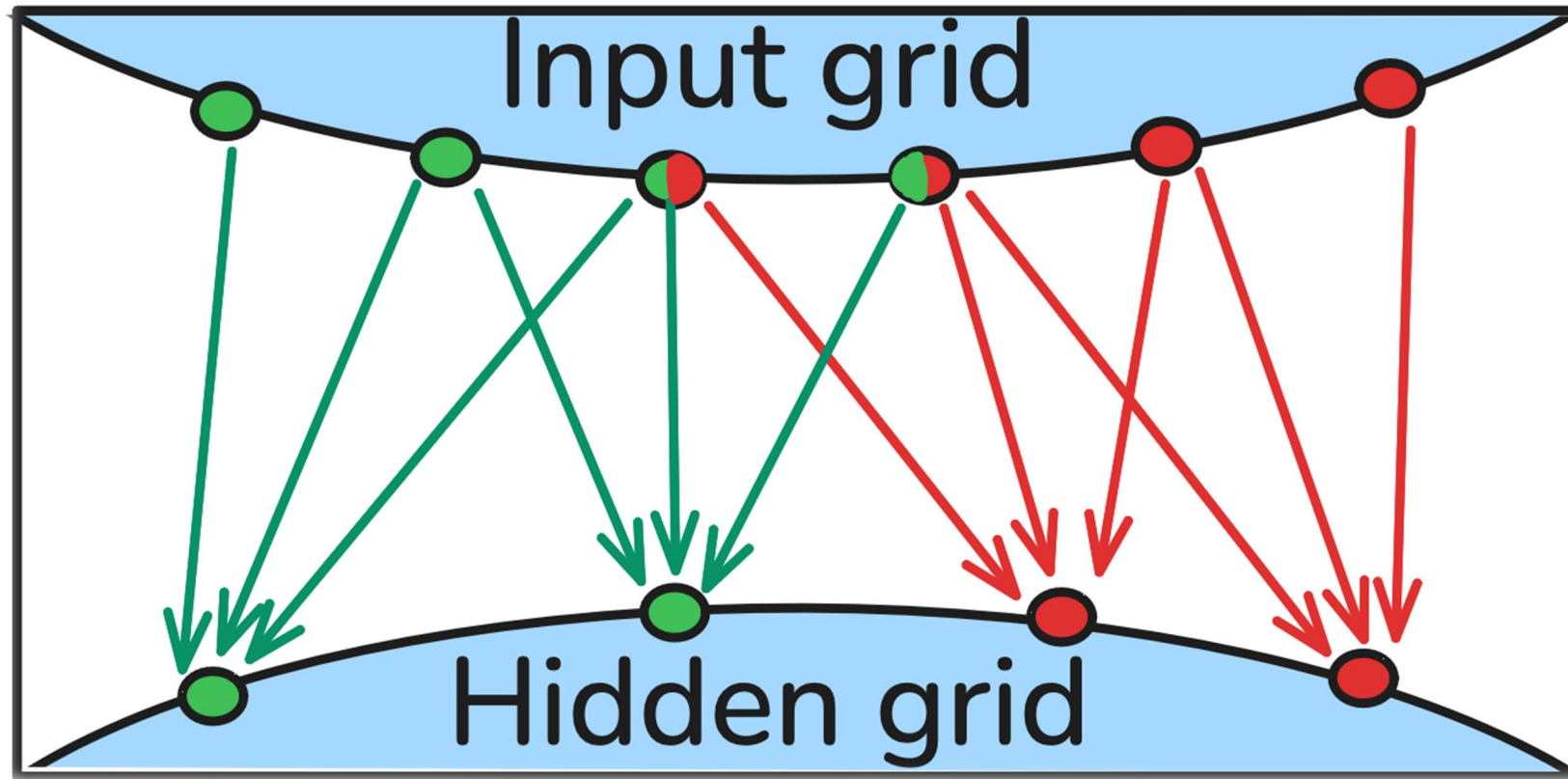| Resolution | Batch Size |
|---|---|
| O96 (100km) | 0.05 GB |
| N320 (32km) | 0.65 GB |
| O1280 (9km) | 7.40 GB |
| O2560 (4km) | ~30 GB |

Memory
usage

before

after

# Sharding Attention



- 4x all-to-all communication per fwd/bwd
- But: we can do better by exploiting sparsity + locality of graph attention
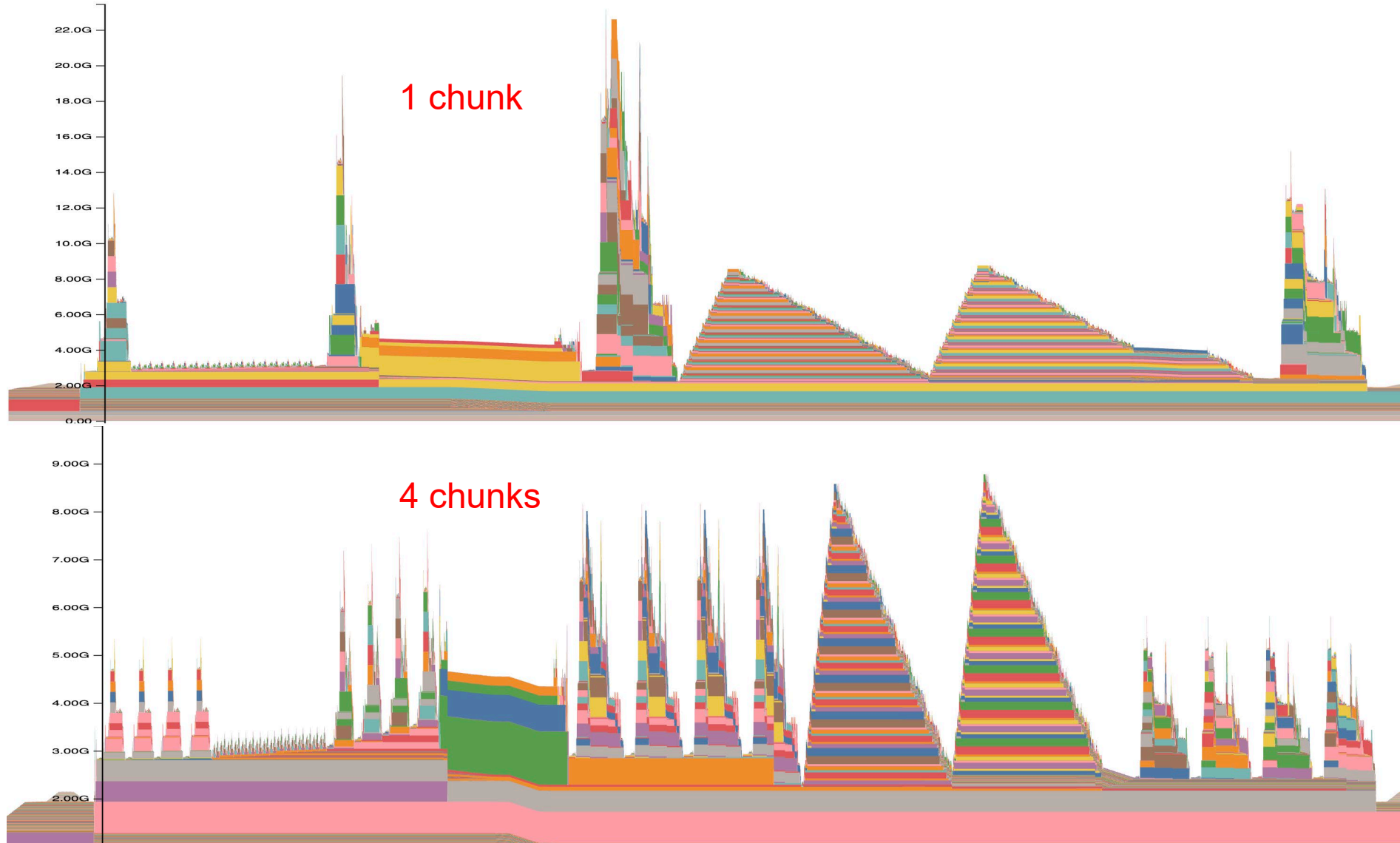
# Mapper edge sharding



- 1x all-gather in fwd + reduce-scatter in bwd
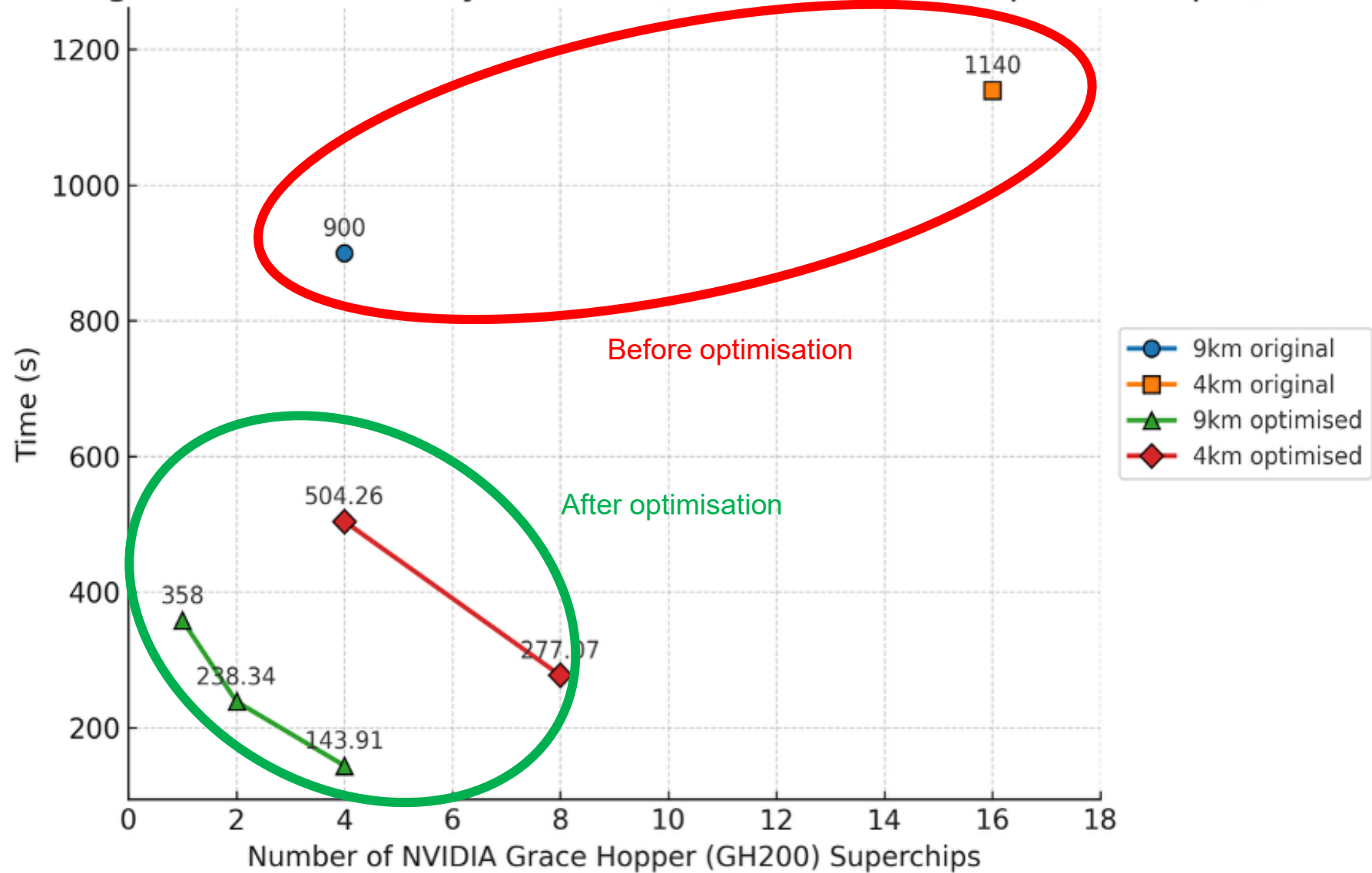- + independent subgraphs on each device (with small overlap)

# Mapper chunking

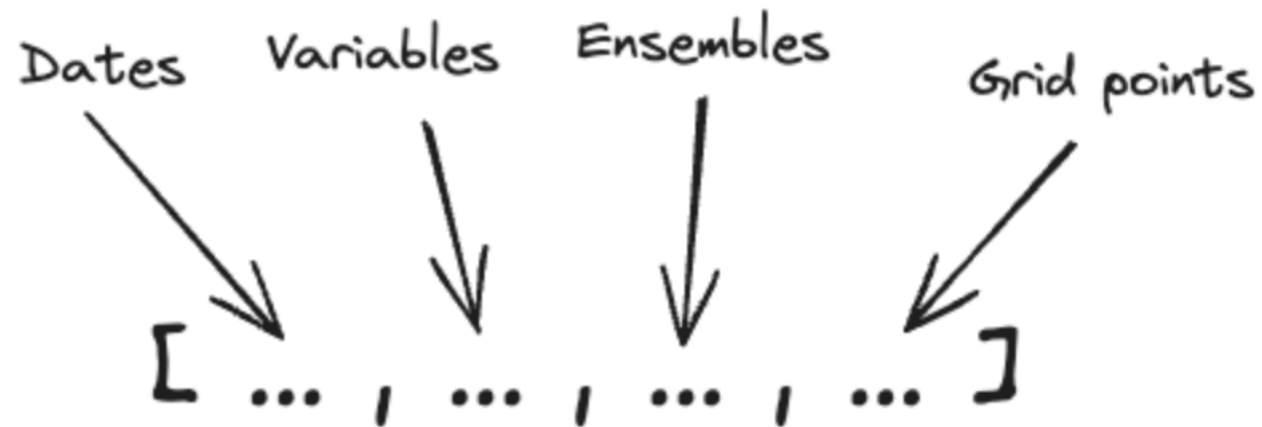- Same idea as edge sharding but sequentially



Memory usage

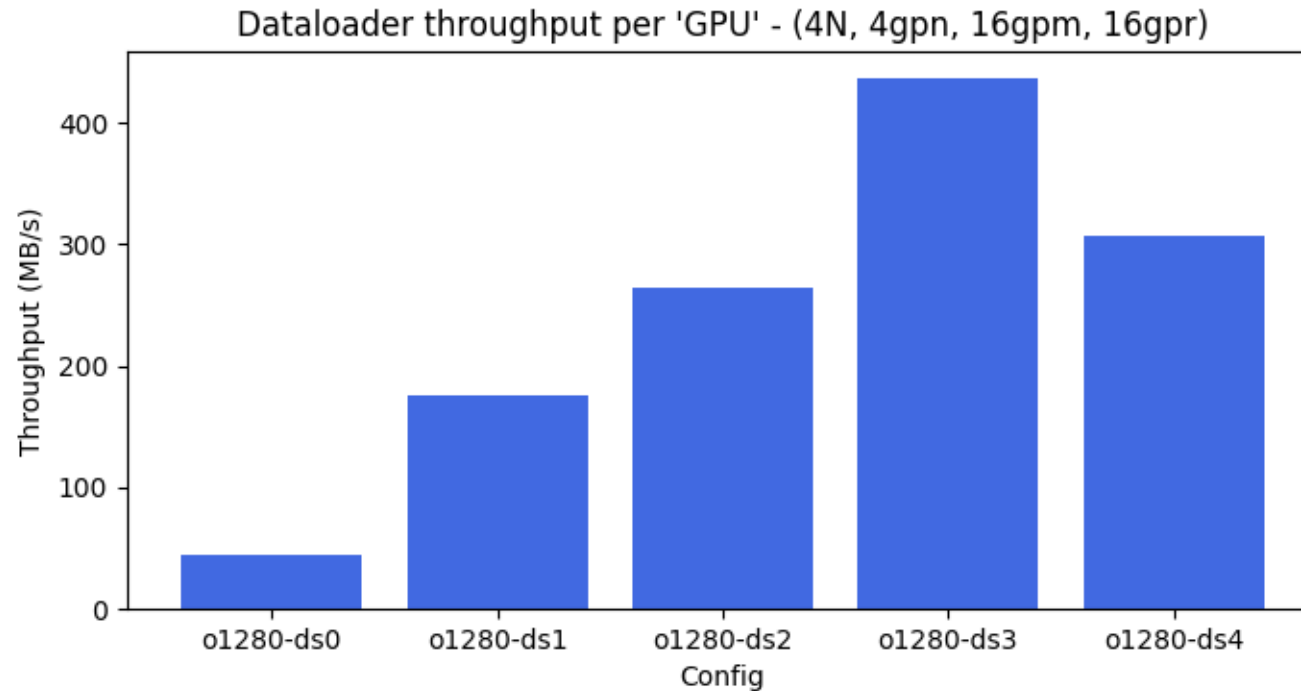AIFS single inference 15-day forecast (100 fields, 150 steps, no output)

# Anemoi datasets

| Resolution | files per date |
|---|---|
| O96 (100km) | 1 |
| N320 (32km) | 1 |
| O1280 (9km) | 4 |
| O2560 (4km) | 16 |

- Anemoi-datasets built on top of Zarr
  - Offers an array-like view of a collection of files
  - Each date is at least 1 file
  - Larger resolutions are split across multiple files

Dates    Variables    Ensembles    Grid points

[ ... , ... , ... , ... ]

# Reformatted datasets



Dataloader throughput per 'GPU' - (4N, 4gpn, 16gpm, 16gpr)
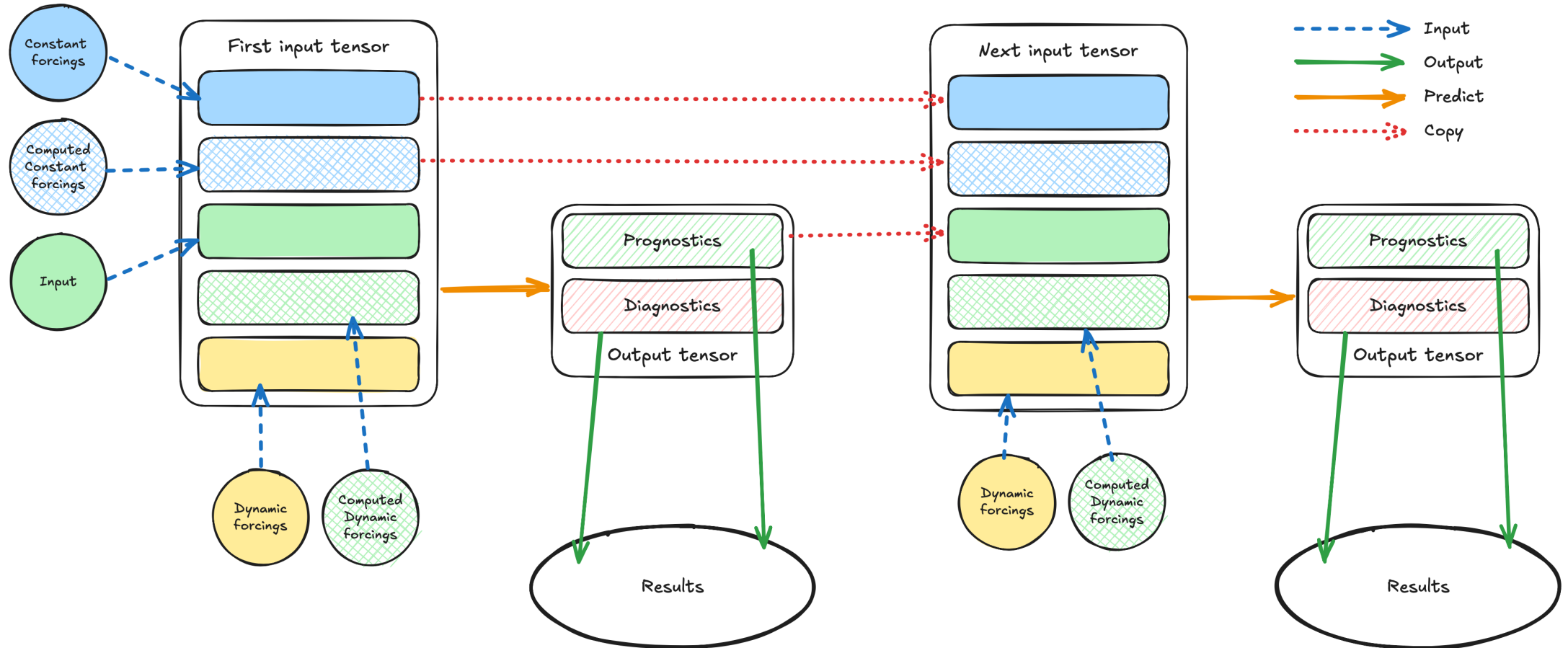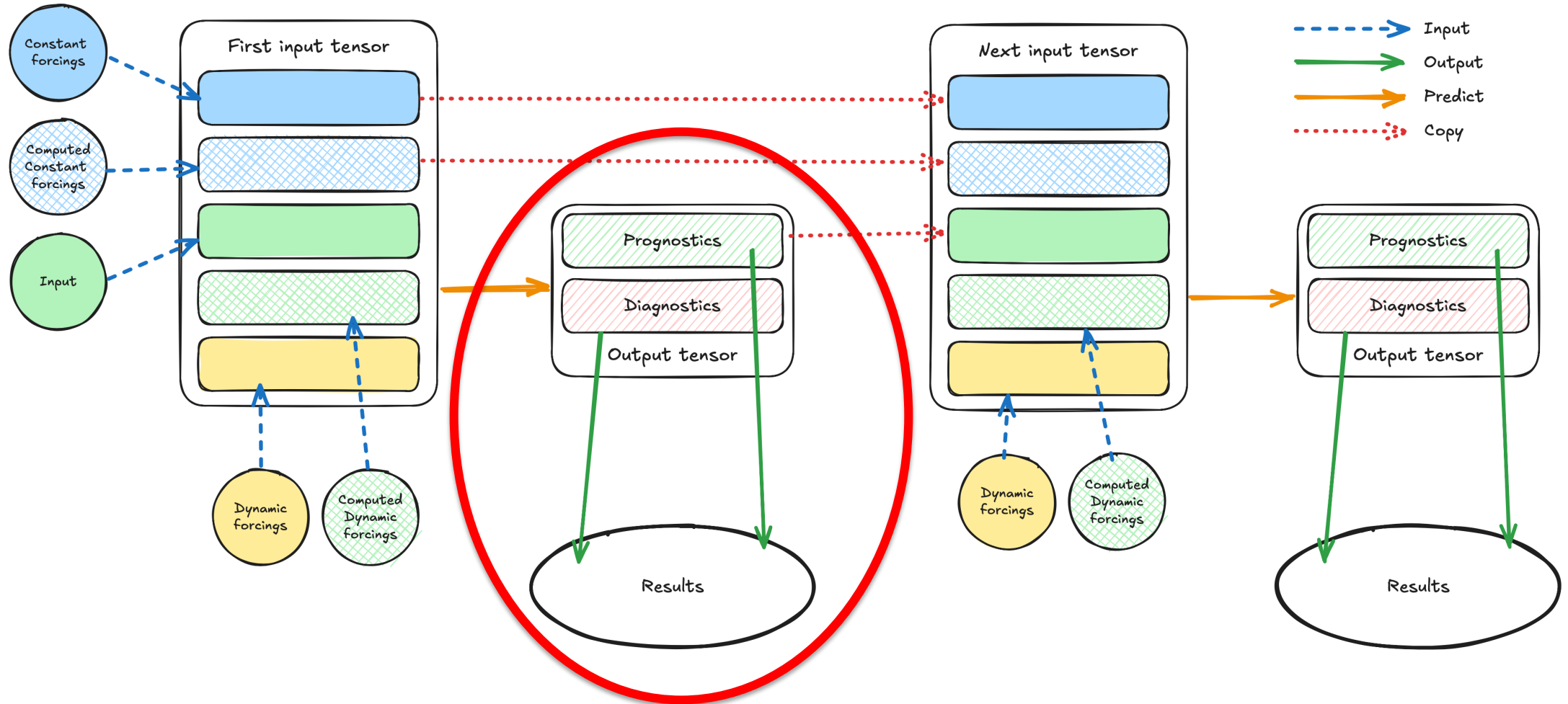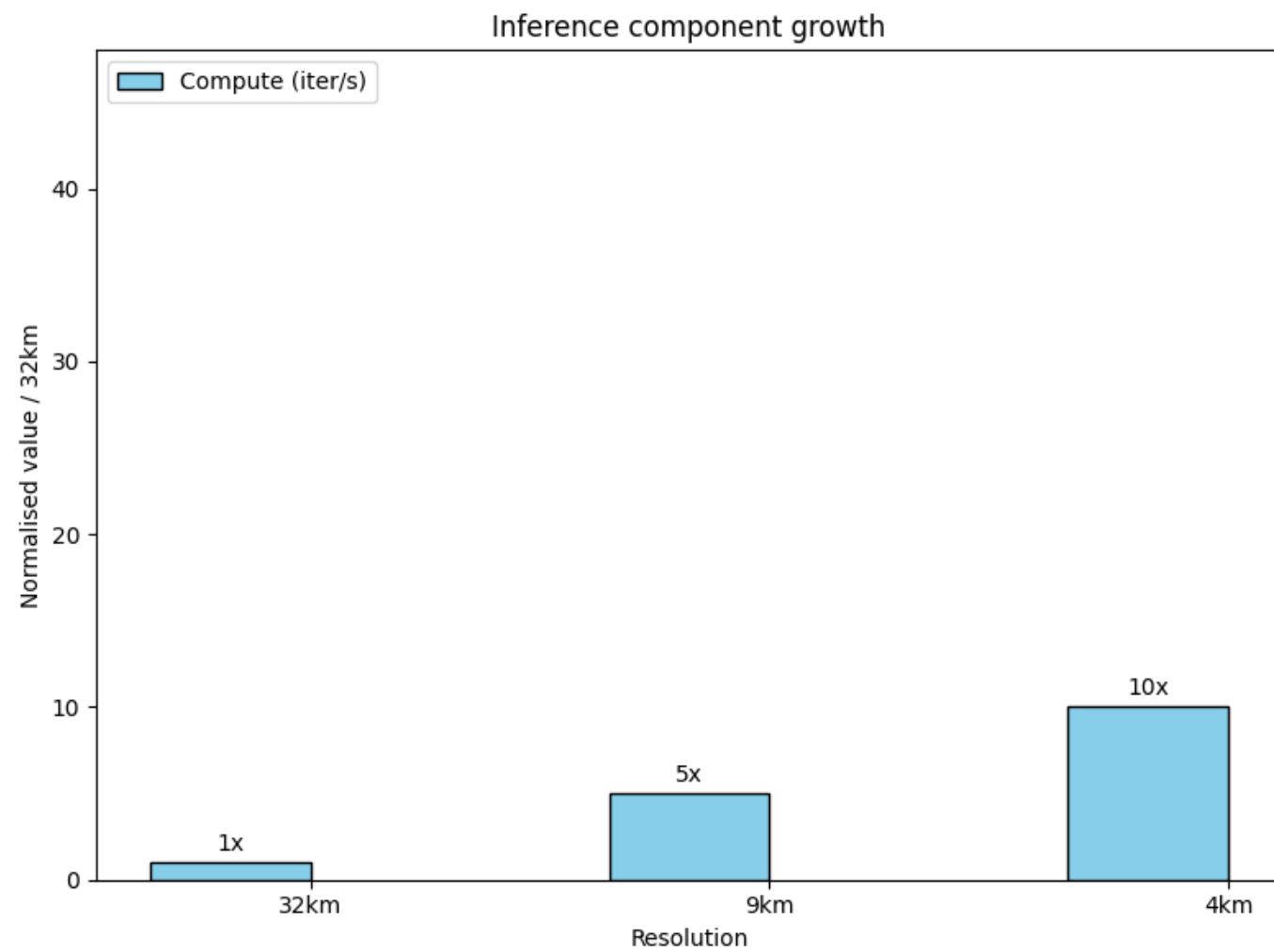
o1280-ds0=/home/mlx/ai-ml/datasets/aifs-od-an-oper-0001-mars-o1280-2023-2023-6h-v1-one-month.zarr
o1280-ds1=/ec/res4/scratch/naco/aifs/inputs/custom/aifs-od-an-oper-0001-mars-o1280-2023-2023-6h-v1-one-month-4gridchunks.zarr
o1280-ds2=/ec/res4/scratch/naco/aifs/inputs/custom/aifs-od-an-oper-0001-mars-o1280-2023-2023-6h-v1-one-month-8gridchunks.zarr
o1280-ds3=/ec/res4/scratch/naco/aifs/inputs/custom/aifs-od-an-oper-0001-mars-o1280-2023-2023-6h-v1-one-month-16gridchunks.zarr
o1280-ds4=/ec/res4/scratch/naco/aifs/inputs/custom/aifs-od-an-oper-0001-mars-o1280-2023-2023-6h-v1-one-month-32gridchunks.zarr
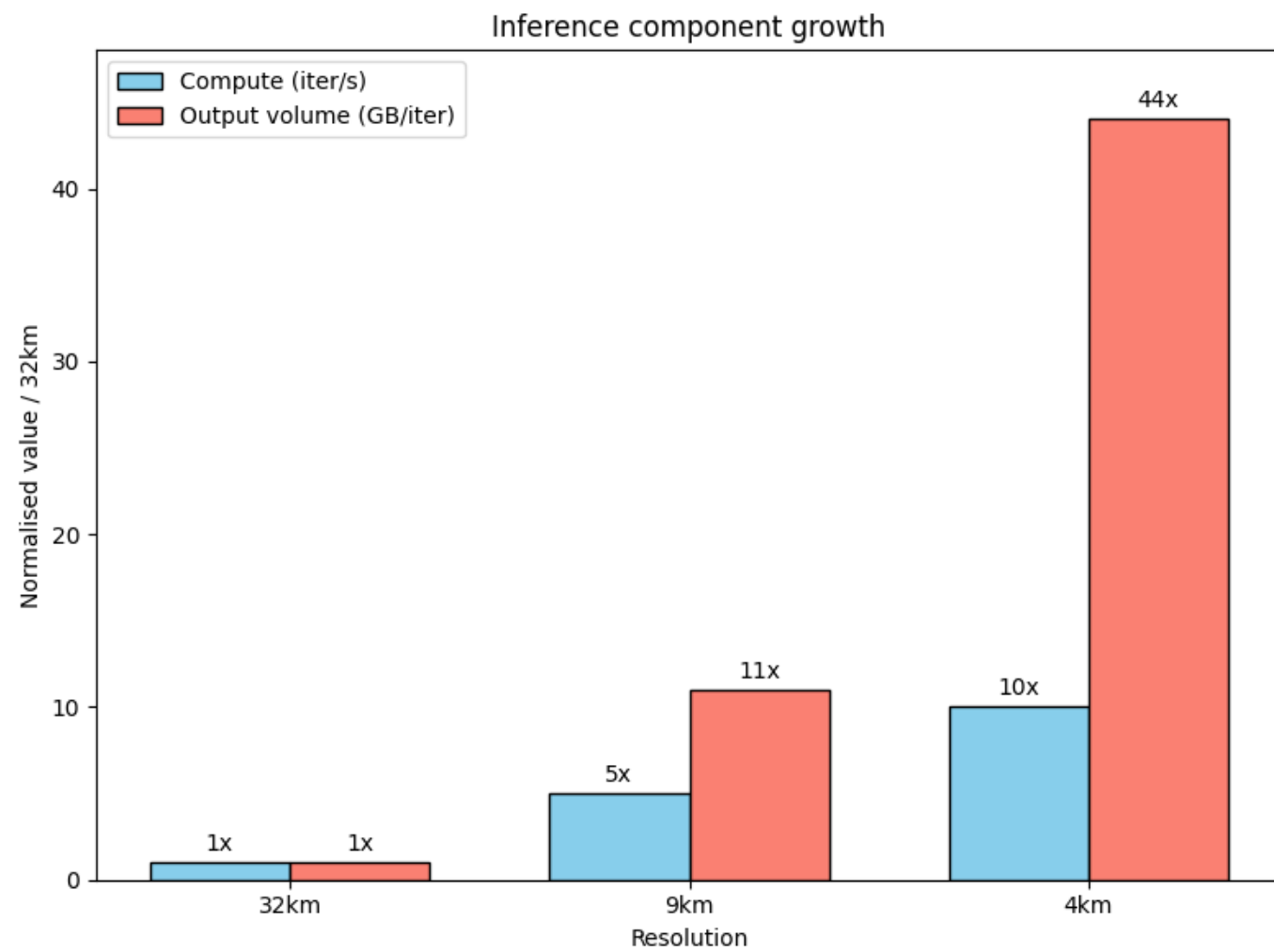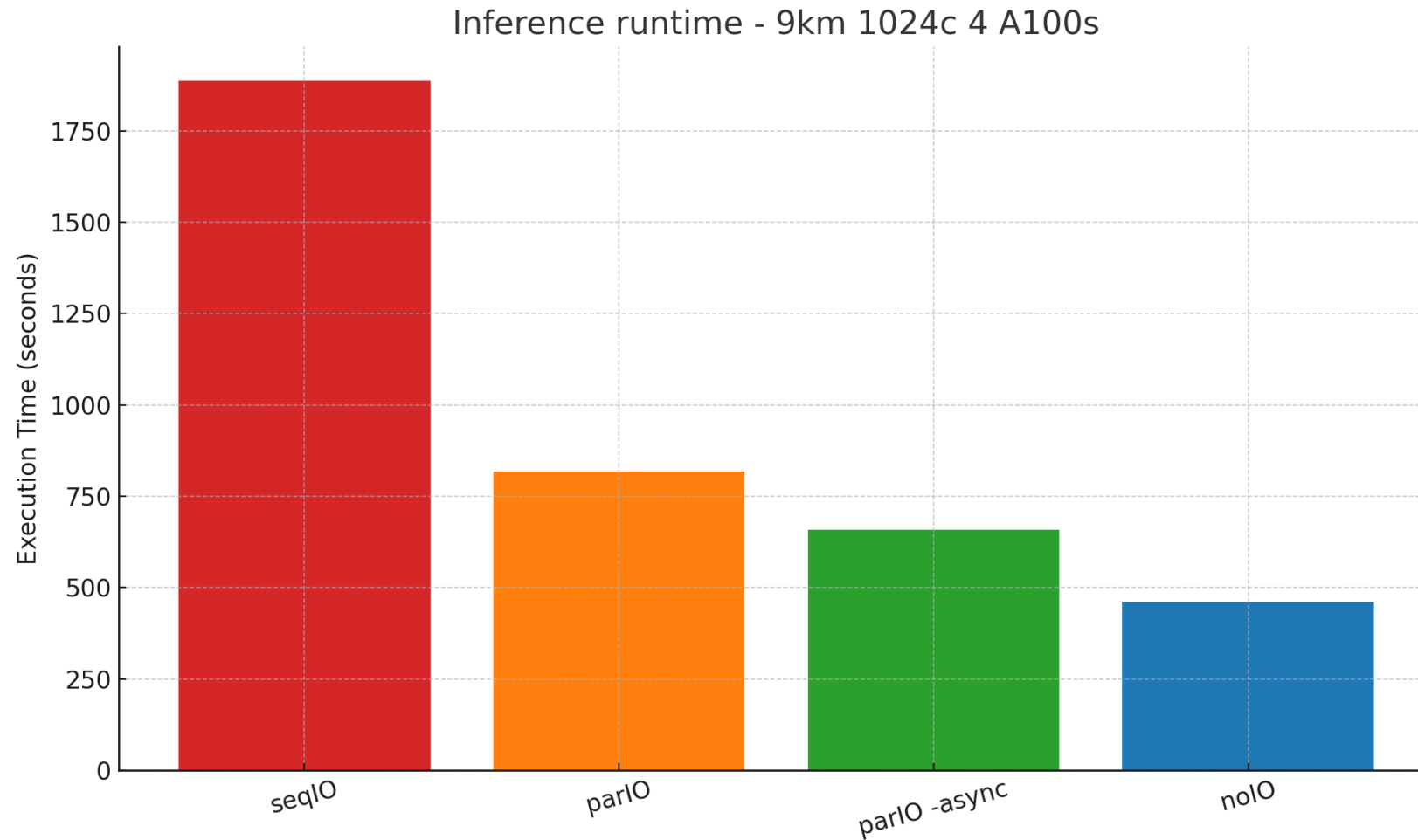
# Anemoi inference

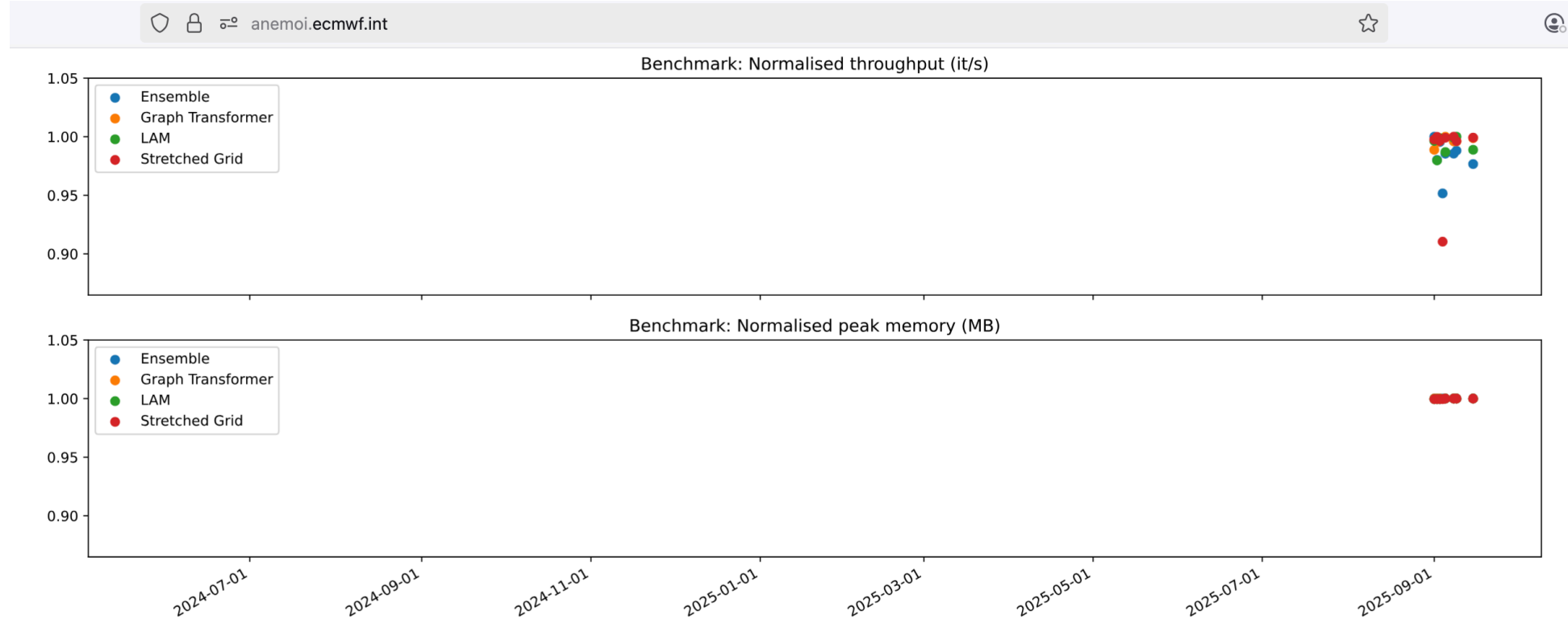**ECMWF** EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS

# Anemoi inference

Inference component growth

**ECMWF** EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS

# Parallel output



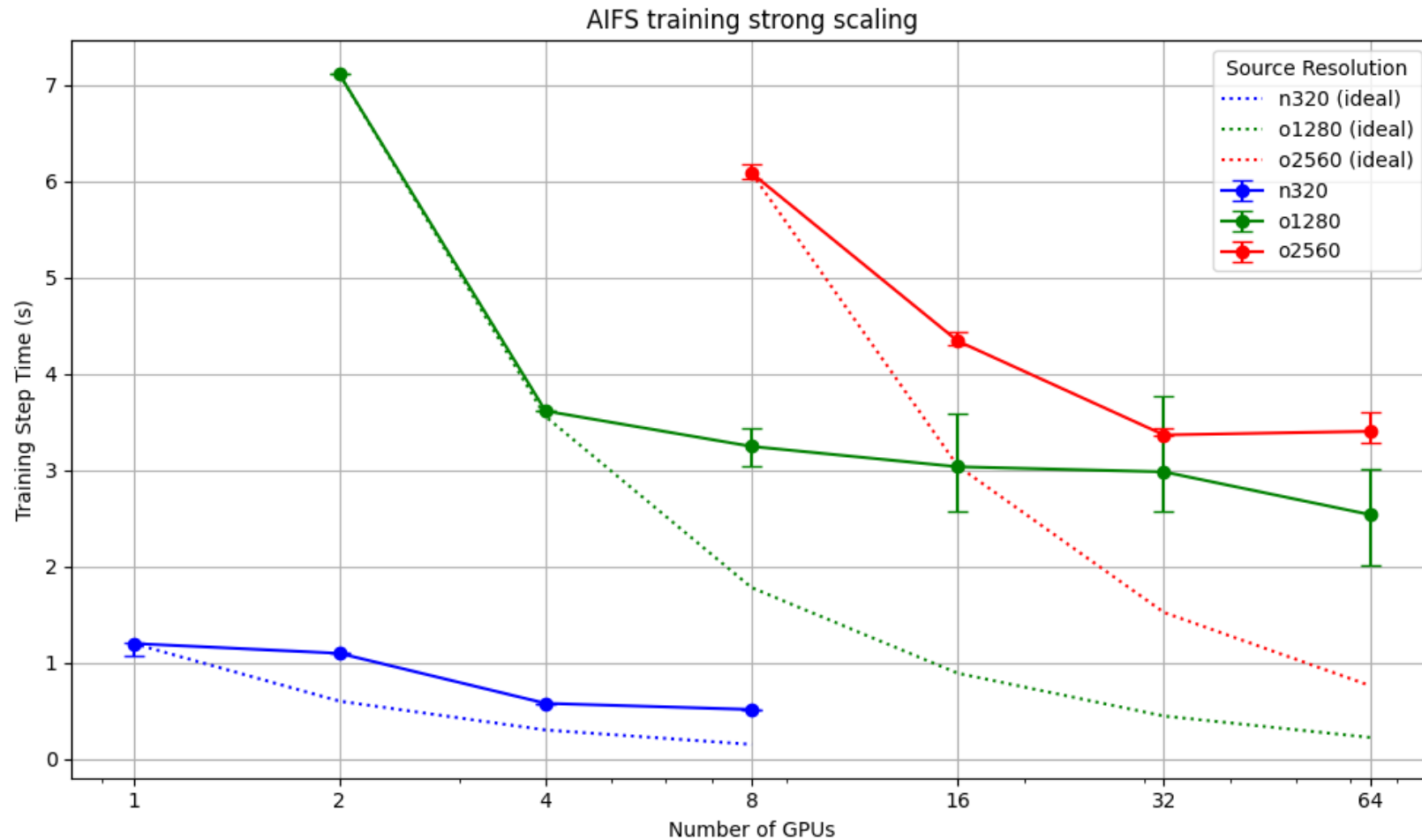Inference runtime - 9km 1024c 4 A100s

# Performance CI testing



- Runs nightly

  – A variety of use cases tested: Global, LAM, Stretched grid and ensemble

  – Throughput and memory usage monitored

  – Model split over 2 GPUs

**ECMWF**  EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS

# Next steps: pushing model sharding

# Thank you



Open-source development driven by
European Meteorological Centers and ECMWF