

ASTRA-sim2.0: Modeling Hierarchical Networks and Disaggregated Systems for Large-model Training at Scale



Georgia Tech College of Computing
Center for Research into
Novel Computing Hierarchies

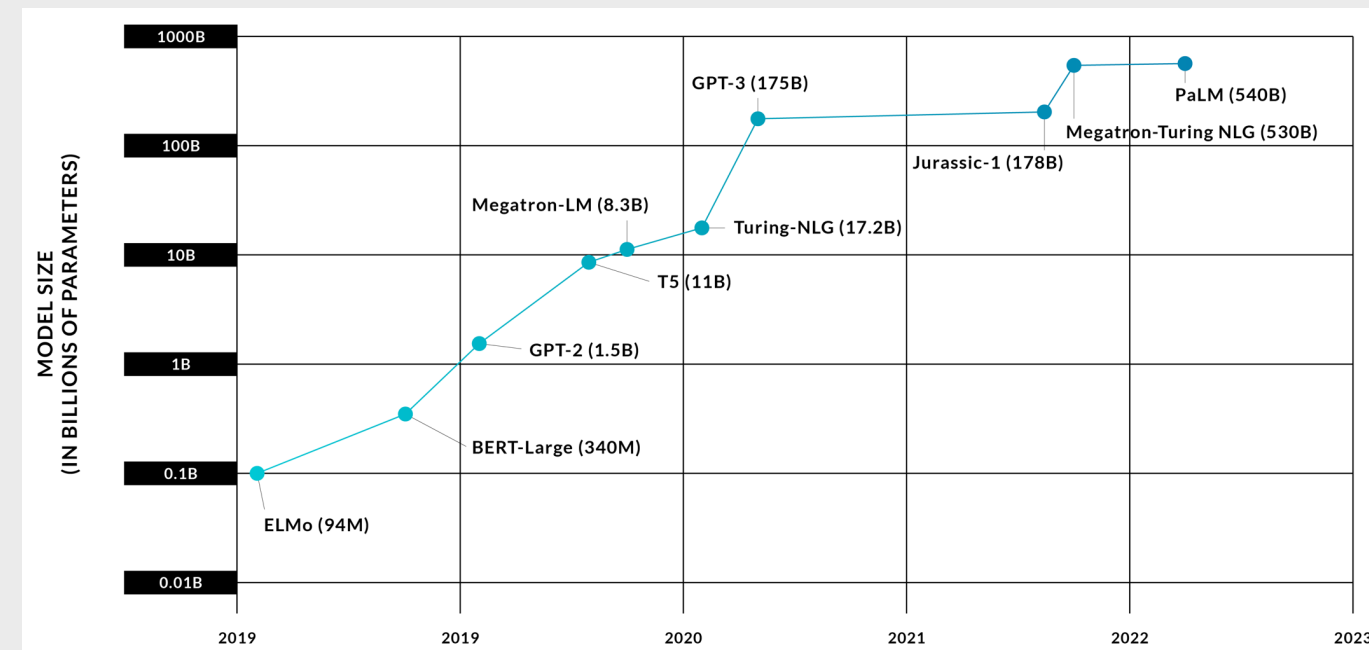
William Won¹, Taekyung Heo², Saeed Rashidi³, Srinivas Sridharan², Sudarshan Srinivasan⁴, Tushar Krishna¹

¹Georgia Institute of Technology, ²NVIDIA, ³HP Labs, ⁴Intel

Distributed ML

Machine Learning (ML) models and training data are scaling quickly

- 2× / 3.4 months

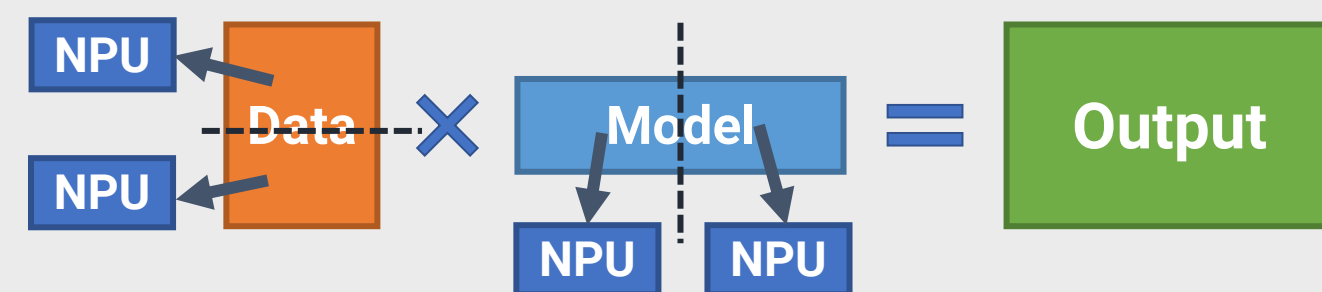


Single-NPU execution is impractical

- 355 years to train GPT-3 on NVIDIA V100

Distributed ML is necessitated

- Model Parallel: Shard model to fit on NPU
- Data Parallel: Shard input data to speedup



State-of-the-art ML Clusters

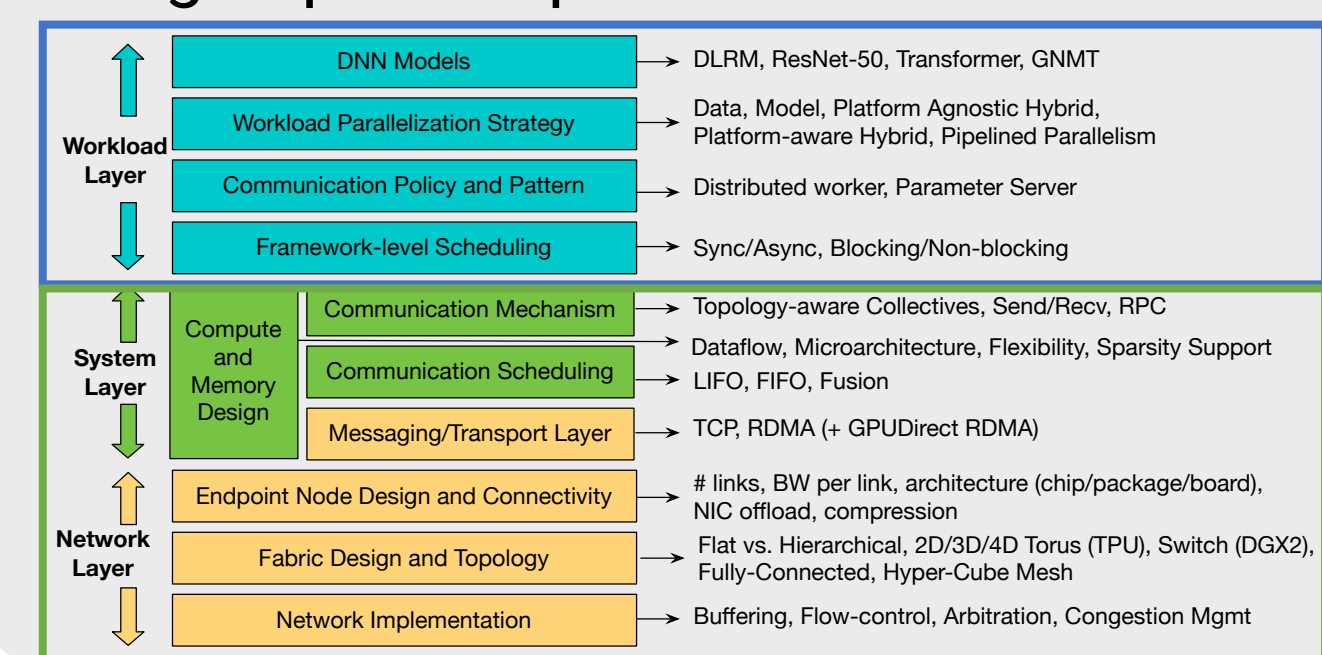


Intel Aurora Google TPUv4 NVIDIA HGX-H100

Design Space

Design-space of Distributed ML is vast and complex

- Workload, System, and Network Layers
- Needs a methodology for swift design-space exploration



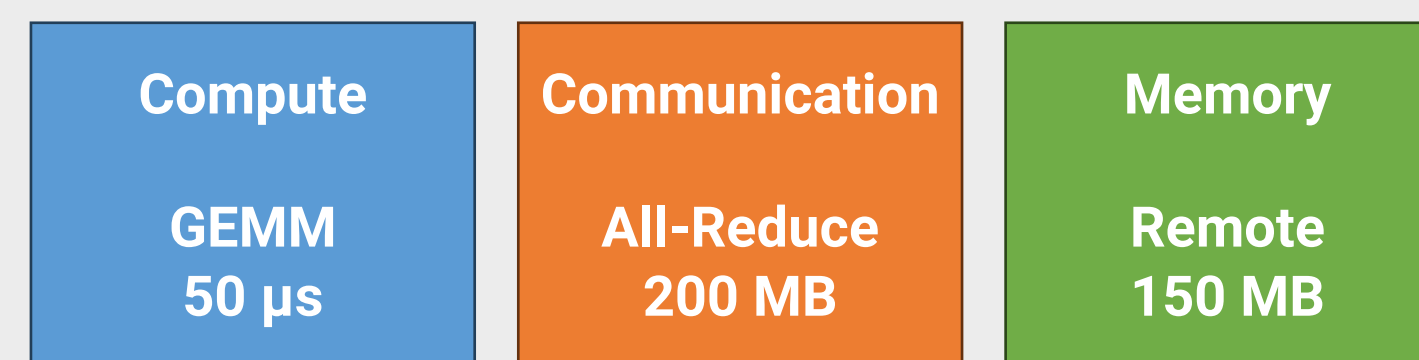
Chakra Execution Trace

Chakra ET: Represents arbitrary distributed ML workloads

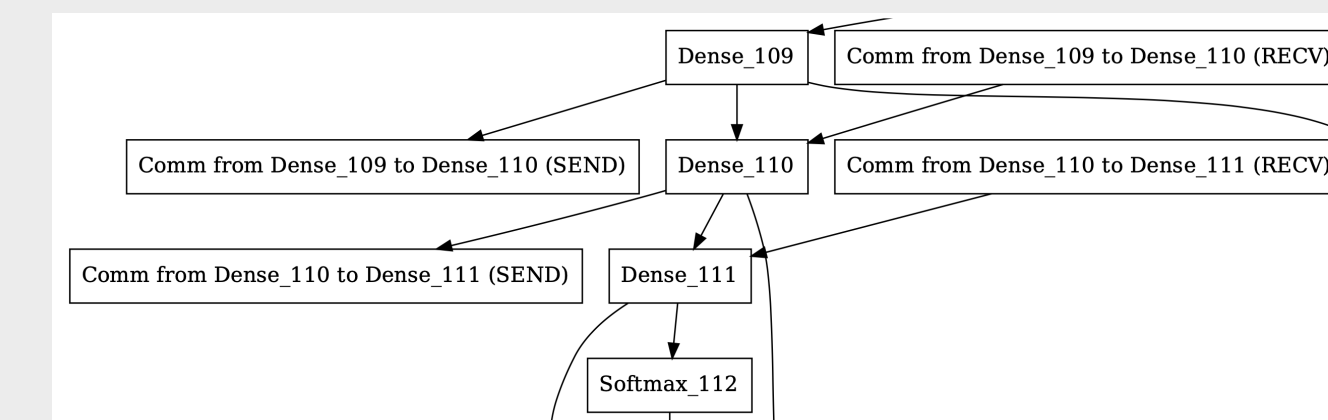
- In Directed Acyclic Graph (DAG) format

Three types of Chakra nodes

- Compute: Operand size or runtime
- Communication: Collective or Peer-to-peer
- Memory: Local or Remote memory access



Example Chakra Execution Trace



Chakra MLCommons



Chakra Paper

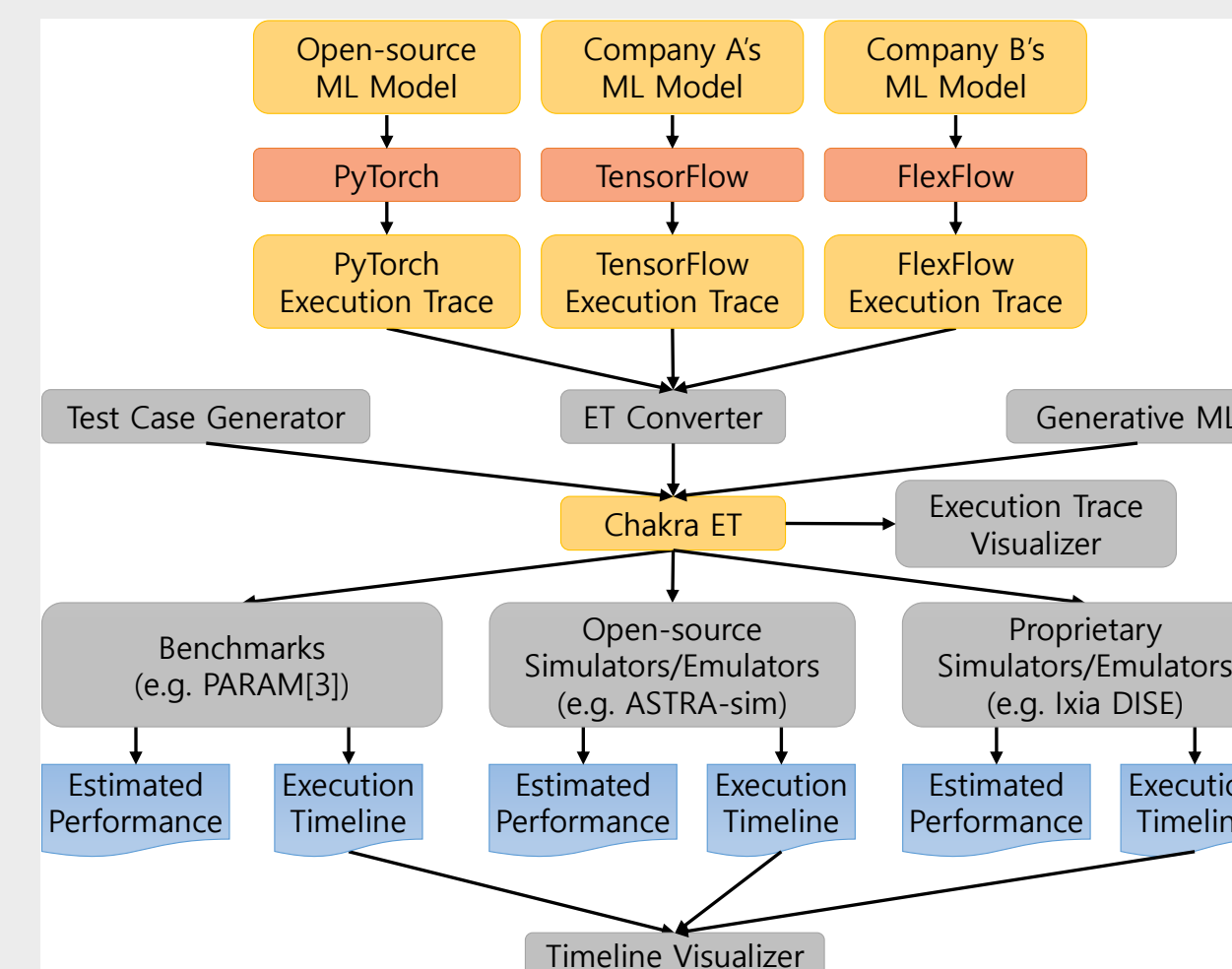


ASTRA-sim Website



ASTRA-sim2.0 Paper

Workflow of Chakra Execution Trace



Advantages of Chakra Execution Trace

- Flexible to capture arbitrary ML schemes
- Suitable for benchmark, replay, simulation
- Obfuscates key proprietary workload data
- Facilitates the exchange of ML traces
- Being standardized via MLCommons

Case Studies

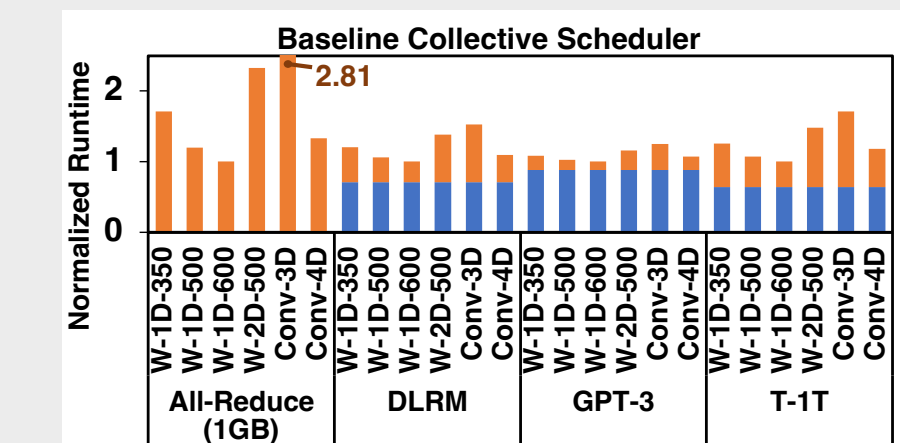
Case 1: Conventional vs. Wafer-scale

- Conventional: High-dimensional (scale-out)
- Wafer-scale: Low-dimensional (scale-up)

| Topology | Shape | NPU Size | BW (GB/s) |
|----------|---------------------|----------|----------------|
| W-1D | Switch | 512 | 350, 500, 600 |
| W-2D | Switch_Switch | 32×16 | 250_250 |
| Conv-3D | Ring_FC_Switch | 16×8×4 | 200_100_50 |
| Conv-4D | Ring_FC_Ring_Switch | 2×8×8×4 | 250_200_100_50 |

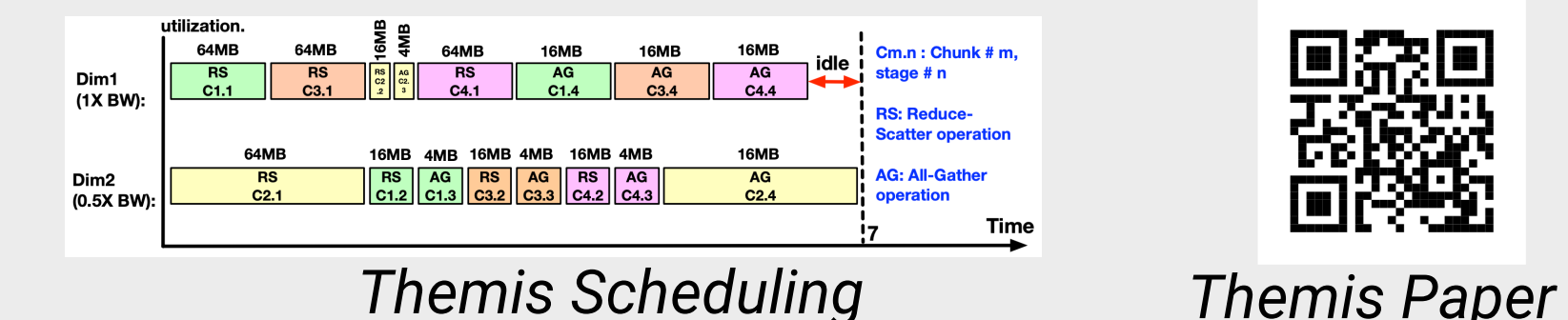
Case 1 Result

- Wafer-scale system perform better
- Due to the multi-dimensional network overhead



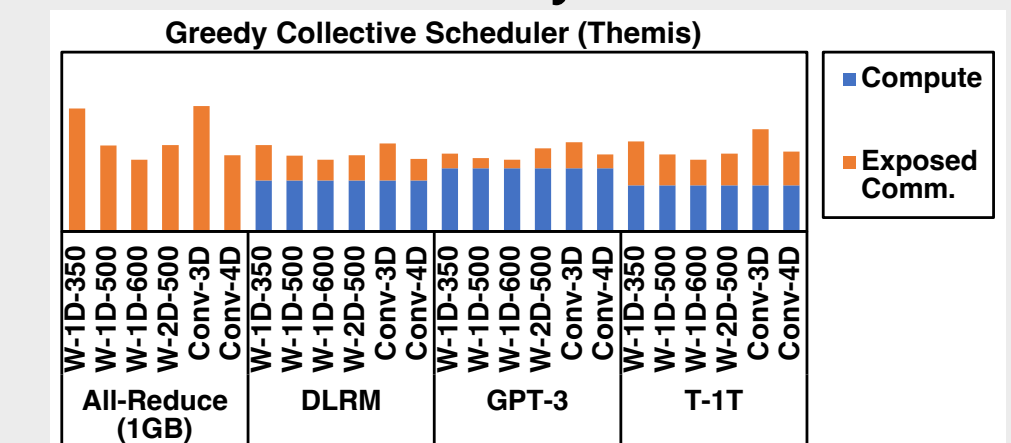
Case 2: Smart Chunk Scheduling

- Themis (ISCA '22): greedy-based chunk scheduler for multi-dimensional networks



Case 2 Result

- Themis mitigates the overhead of conventional scale-out systems



CONCLUSIONS

Needs for Distributed ML

- Models and training data are scaling
- Single-NPU execution is impractical

Needs for simulation methodology

- For swift design-space exploration
- Of emerging ML platforms

Chakra ET and ASTRA-sim2.0

- Chakra ET: standardized, graph-based distributed ML workload representation
- ASTRA-sim2.0: simulating emerging ML systems