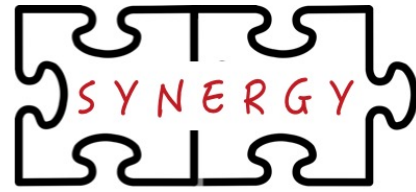




Georgia Tech School of Electrical and
Computer Engineering
College of Engineering



<http://synergy.ece.gatech.edu>



<https://astra-sim.github.io>



<https://github.com/mlcommons/chakra>

Chakra + ASTRA-sim: An ecosystem for advancing benchmarking and co-design for diverse AI supercomputing platforms

Tushar Krishna

Associate Professor, School of ECE
Georgia Institute of Technology

tushar@ece.gatech.edu

<https://tusharkrishna.ece.gatech.edu/>

ModSim Workshop
August 13, 2025

Acknowledgments



Will Won



Jinsun Yoo



Changhai Man



Joongun Park



Saeed Rashidi



Brad Beckmann



Taekyung Heo



Srinivas Sridharan



OPEN
Compute Project®



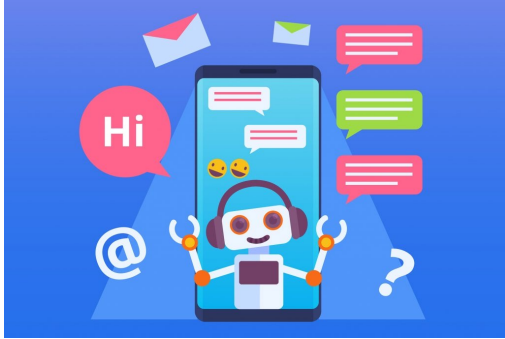
**Semiconductor
Research
Corporation**



+ many more

AI is pervasive today!

Chatbots

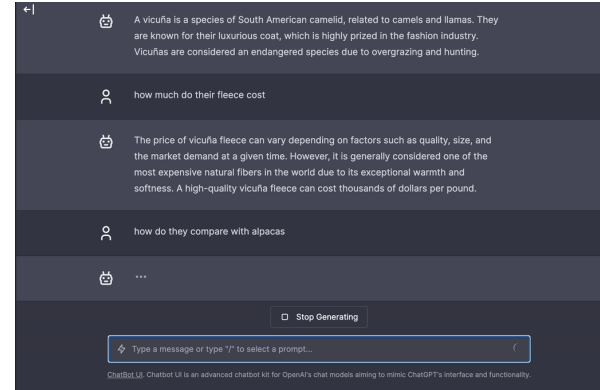


Code Generation

```
1 segmentation.rb
2
3 def segmentation(items, separator)
4   curr = []
5   segments = []
6   items.each do |item|
7     ...
8   end
9
10  segmentation([1,2,4,0,2,5,0,3,0], 0).each do |segment|
11    puts(segment.join(", "))
12  end
13
```

"25% of code at google in last quarter was AI generated"

Text Generation



Language Translation

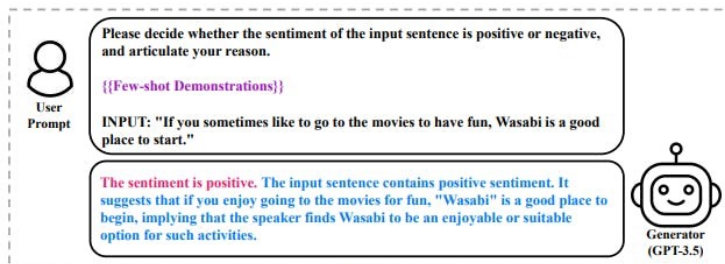
[Instruction]: Translate the following sentences from English to Chinese.
[Input]: Did you see it go?

LLM

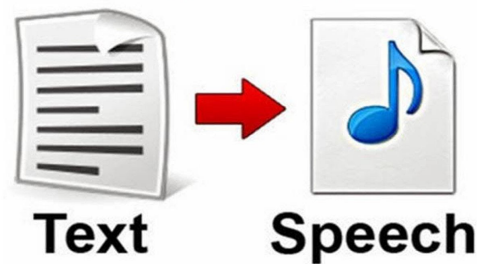


[Output]: 看清楚了吗?

Sentiment Analysis



Text to Speech

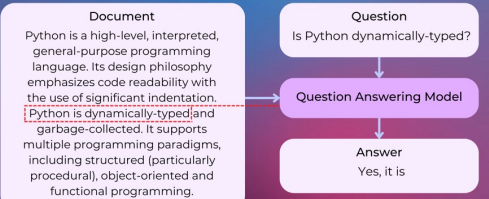


Recommendations



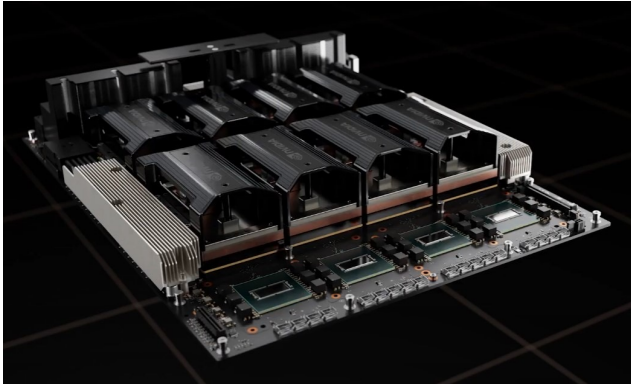
Question Answering

Generative Question Answering



➡ Algorithmic view of AI (Datasets and Models)

Computer Architect's view of AI



NVIDIA
HGX



Google
Cloud TPU



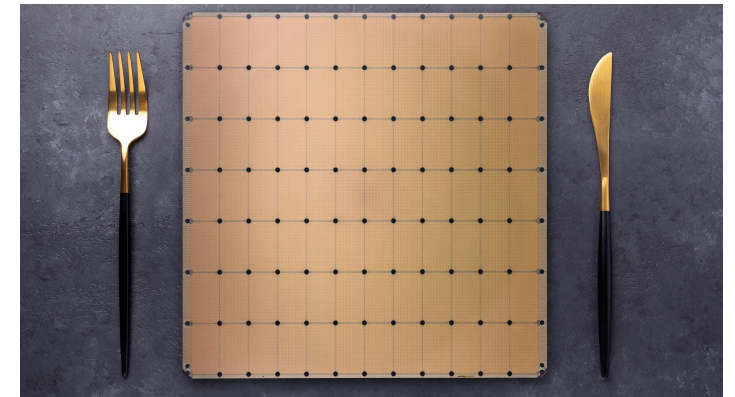
SambaNova
SN40L



AMD
Instinct Platforms

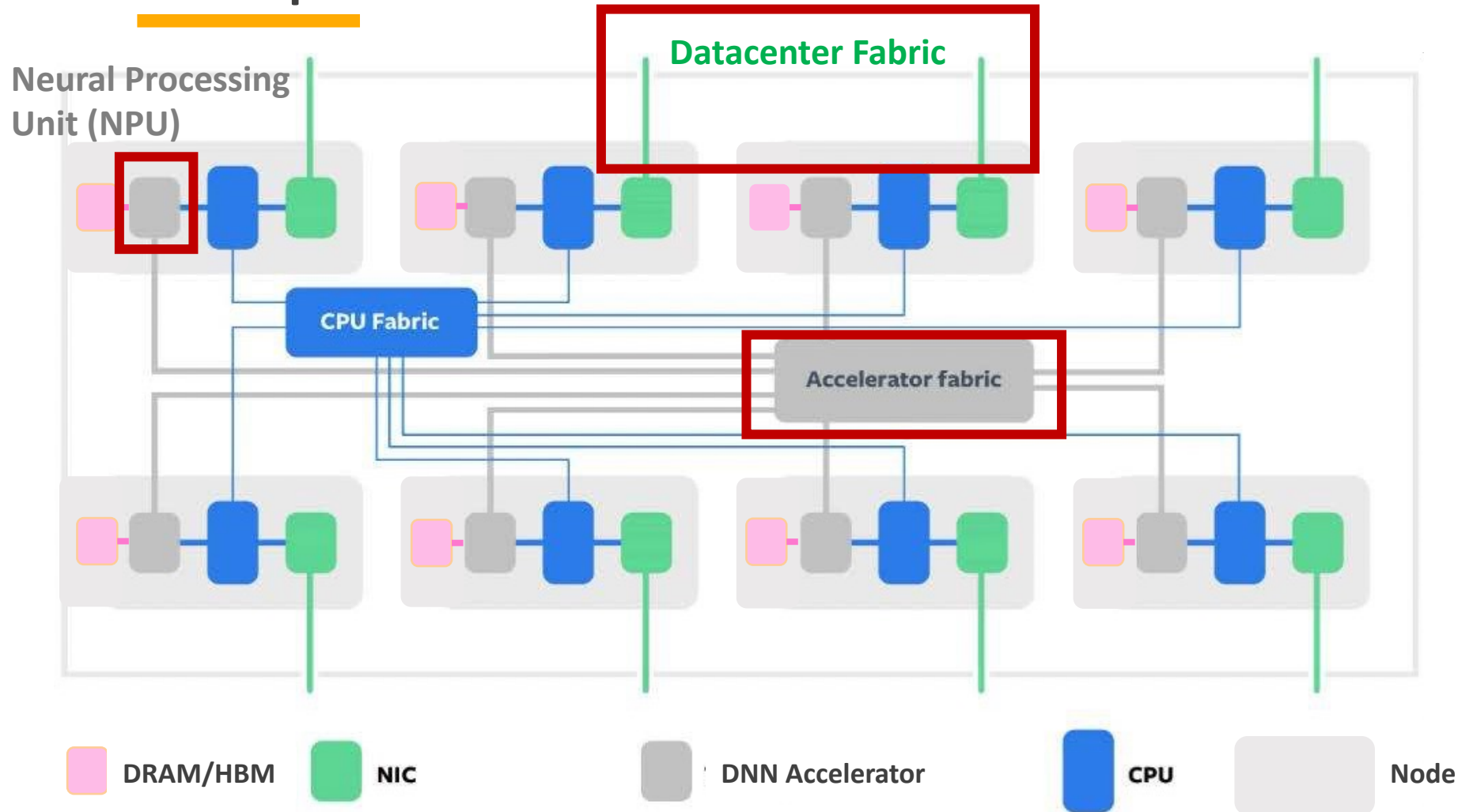


Intel
Gaudi



Cerebras
Andromeda

Computer Architect's view of AI




- ✓ Customized accelerators for AI (aka NPUs)
- ✓ Customized network fabrics to scale the AI task across multiple accelerators

Why?

Figure modified from "Zion: Facebook Next- Generation Large Memory Training Platform", Misha Smelyanskiy, Hot Chips 31"

“Large” Language Models

Hundreds of ZettaFLOPs of compute



Model (Company)	Company	#Parameters [Billion]	Model Footprint (Assuming 2B/Param)	Training Footprint (Assuming 16B/Param*)
Clause 3 Opus	Anthropic	2,000	4.00 TB	32.00 TB
GPT-4	OpenAI	1,760	3.52 TB	28.16 TB
Gemini 1.5 Pro	Google	1,500	3.00 TB	24.00 TB
Samba-1	SambaNova	1,400	2.80 TB	22.40 TB
Cerebras-1T	Cerebras	1,000	2.00 TB	16.00 TB
Grok-3	xAI	928	1.86 TB	14.85 TB
DeepSeep-R1	DeepSeek-AI	685	1.37 TB	10.96 TB
PaLM	Google	540	1.08 TB	8.64 TB

<https://lifearchitected.ai/models/>

*Assuming Mixed-precision Adam Optimizer (See Microsoft ZeRO)

AI is a distributed systems problem!

Consider GPT-4

- **Total parameters** — ~1.8 trillion (over 10x more than GPT-3)
- **Architecture** — Uses a mixture of experts (MoE) model to improve scalability
- **Training compute** — Trained on ~25,000 Nvidia A100 GPUs over 90-100 days
- **Training data** — Trained on a dataset of ~13 trillion tokens
- **Inference compute** — Runs on clusters of 128 A100 GPUs for efficient deployment
- **Context length** — Supports up to 32,000 tokens of context

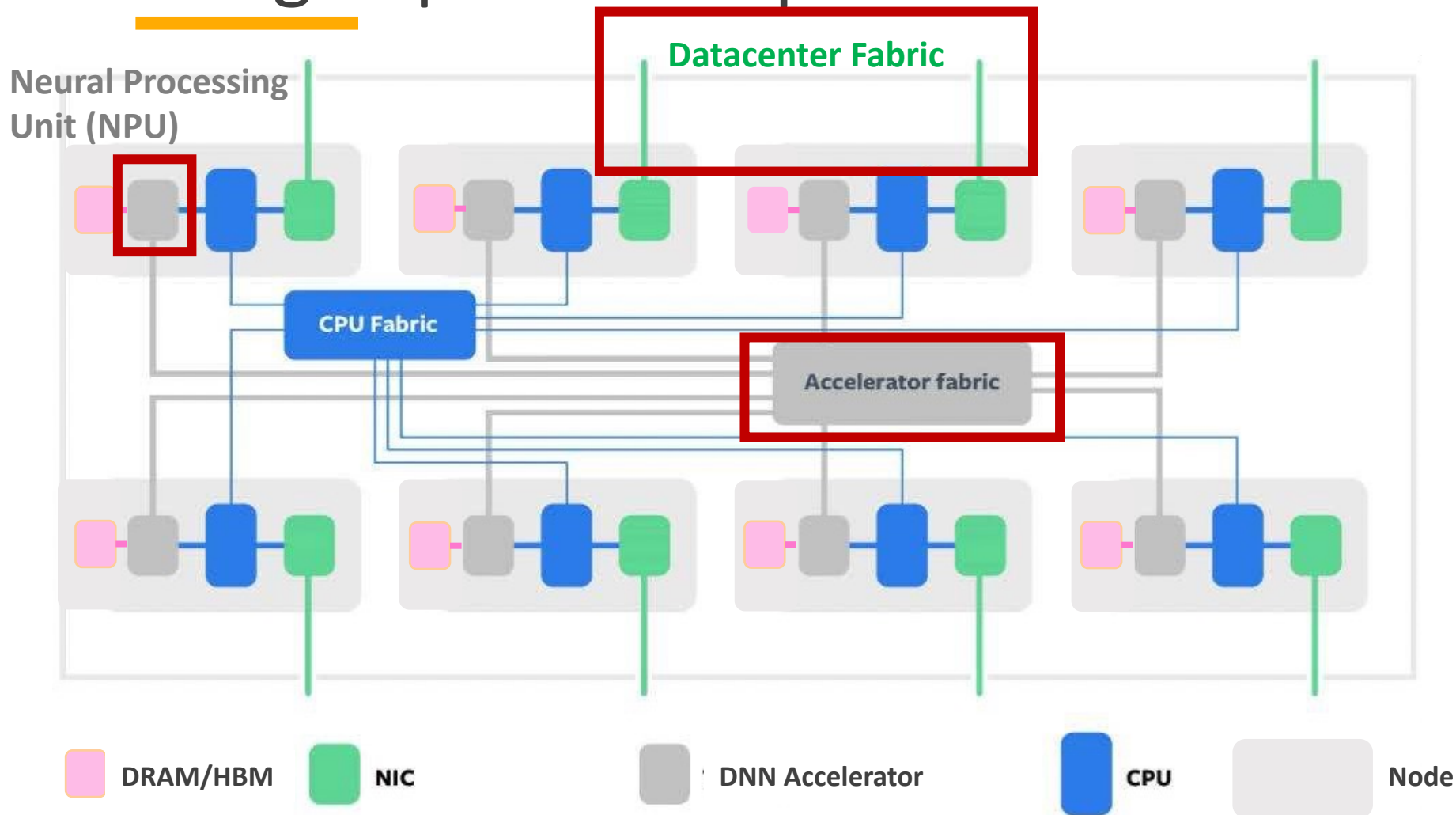
Outline

- Design Space of AI Platforms
- Chakra Workload Execution Traces
- ASTRA-sim Layers and APIs
- External Engagements
- Conclusion

Outline

- **Design Space of AI Platforms**
- Chakra Workload Execution Traces
- ASTRA-sim Layers and APIs
- External Engagements
- Conclusion

Design Space of AI platforms

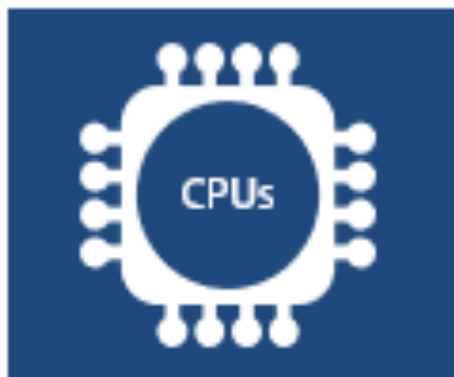


- ✓ Customized accelerators for AI (aka NPUs)
- ✓ Customized network fabrics to scale the AI task across multiple accelerators

Figure modified from “Zion: Facebook Next- Generation Large Memory Training Platform”, Misha Smelyanskiy, Hot Chips 31”

Diverse Compute

Intel
AMD



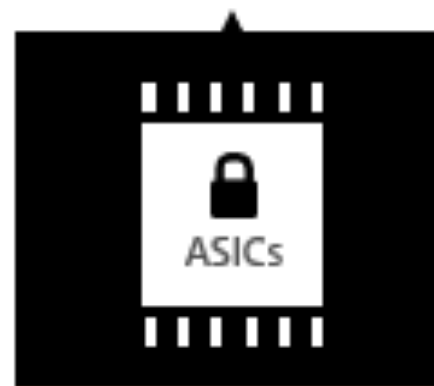
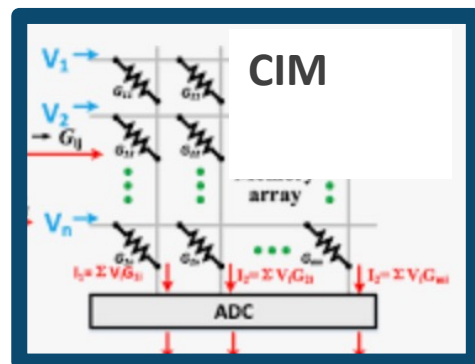
NVIDIA
AMD



Microsoft



EnChargeAI

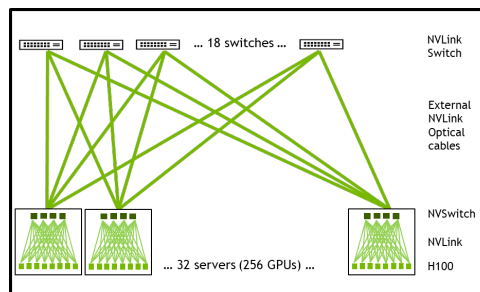


Google
Amazon
Meta
Microsoft
Groq
Cerebras
Rebellions

Diverse Networks

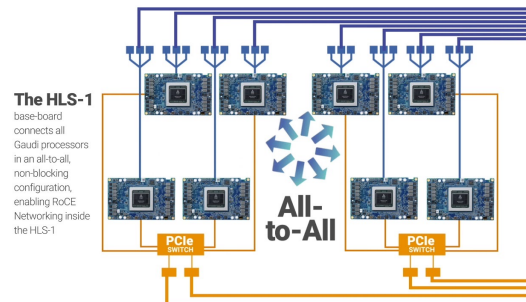
Scale-up → Scale-out

NVIDIA



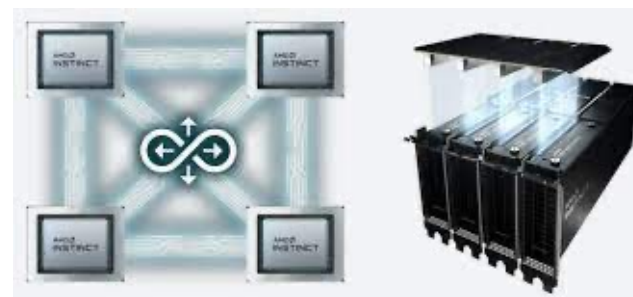
NVswitch → Infiniband

Intel



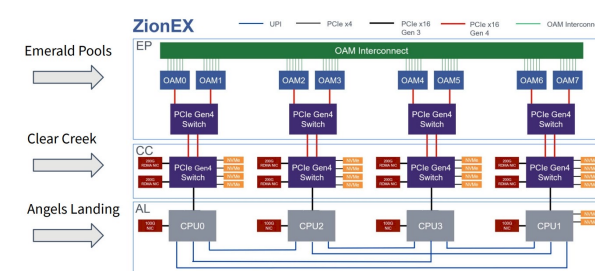
Custom NICs → RoCE

AMD



Infiniti → Infiniti

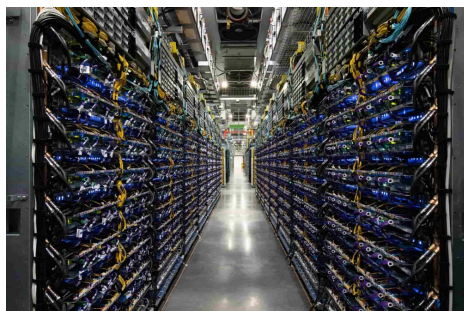
Meta



NVlink → RoCE

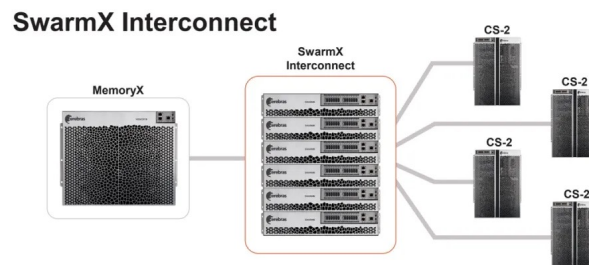
“Classic” Technologies

Google



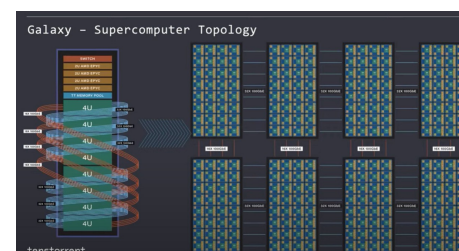
3D Electrical Torus → Optical

Cerebras



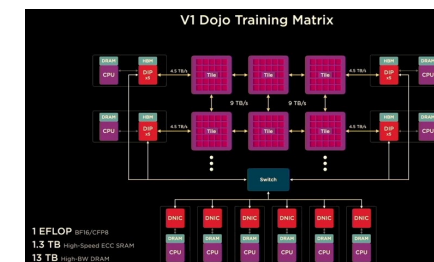
Wafer-scale → SwarmX Tree

Tensorrent



On-package Mesh → off-chip mesh

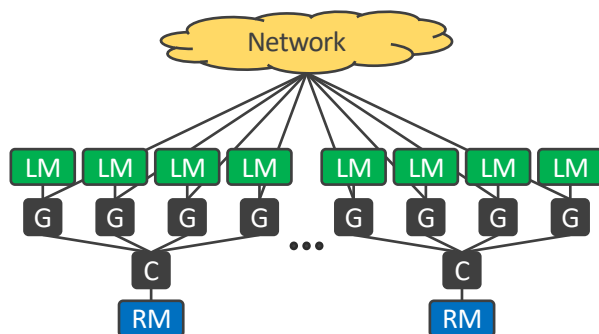
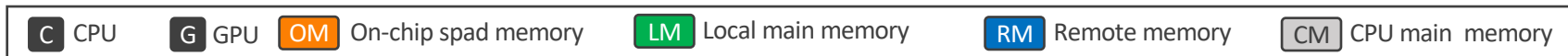
Tesla



On-package Mesh → Ethernet

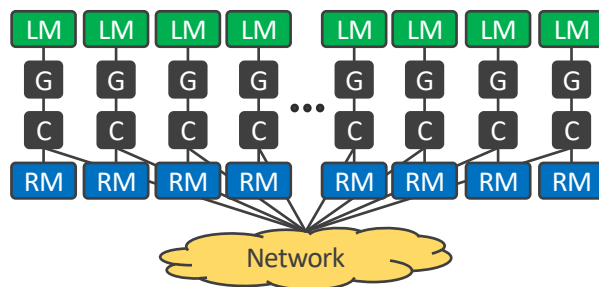
“Emerging” Technologies

Diverse Memory Systems

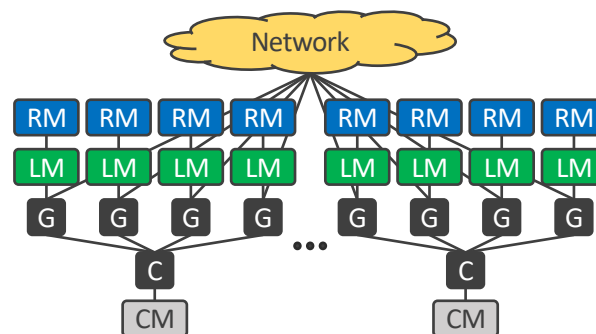


① Per-node memory expansion

Intel Habana, NVIDIA DGX2

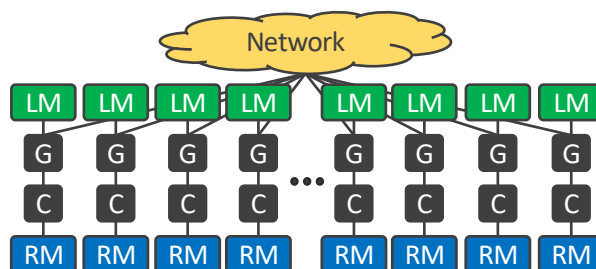


③ Integ. chip mem. exp with CPU network



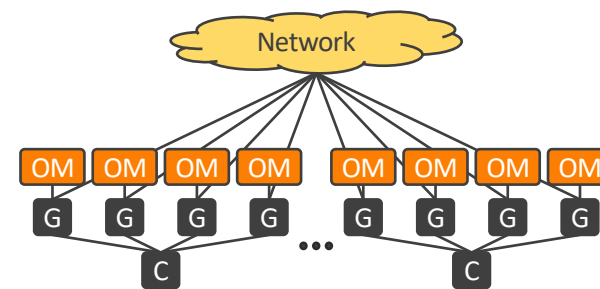
② Per-GPU memory expansion

Astera Labs



④ Integ. chip mem. exp with GPU network

NVIDIA Grace Hopper

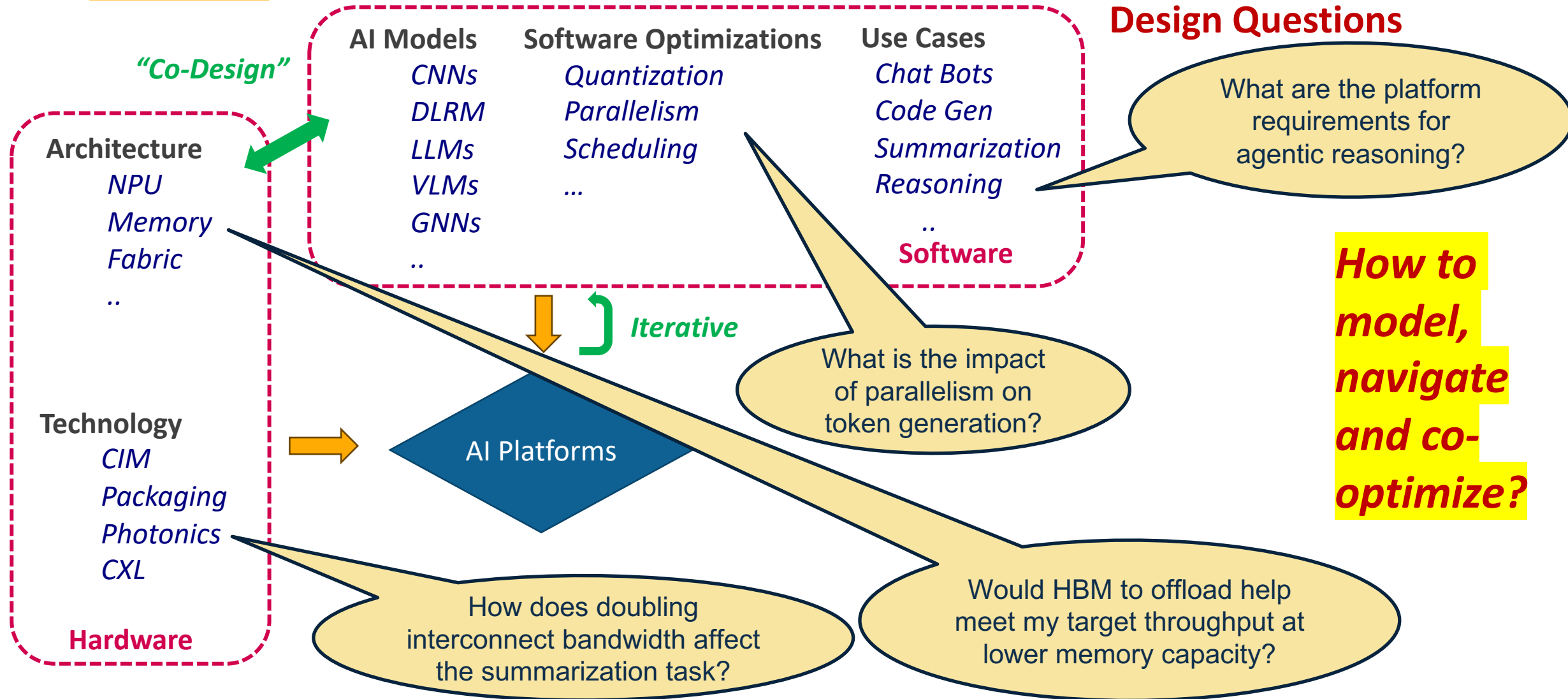


(5) On-chip (SRAM) memory

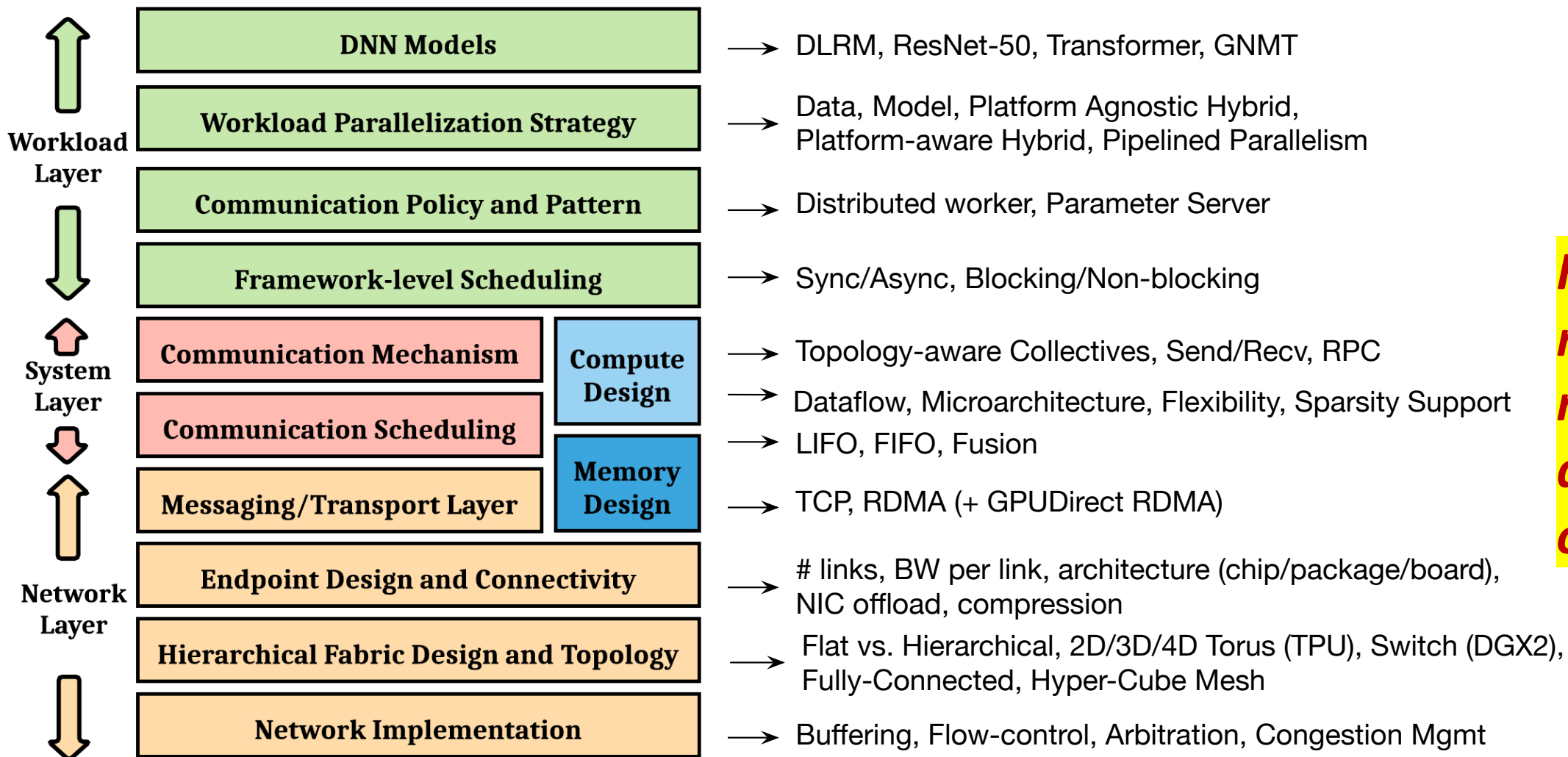
Groq

Cerebras CS-3

Cross-coupled Design-Space



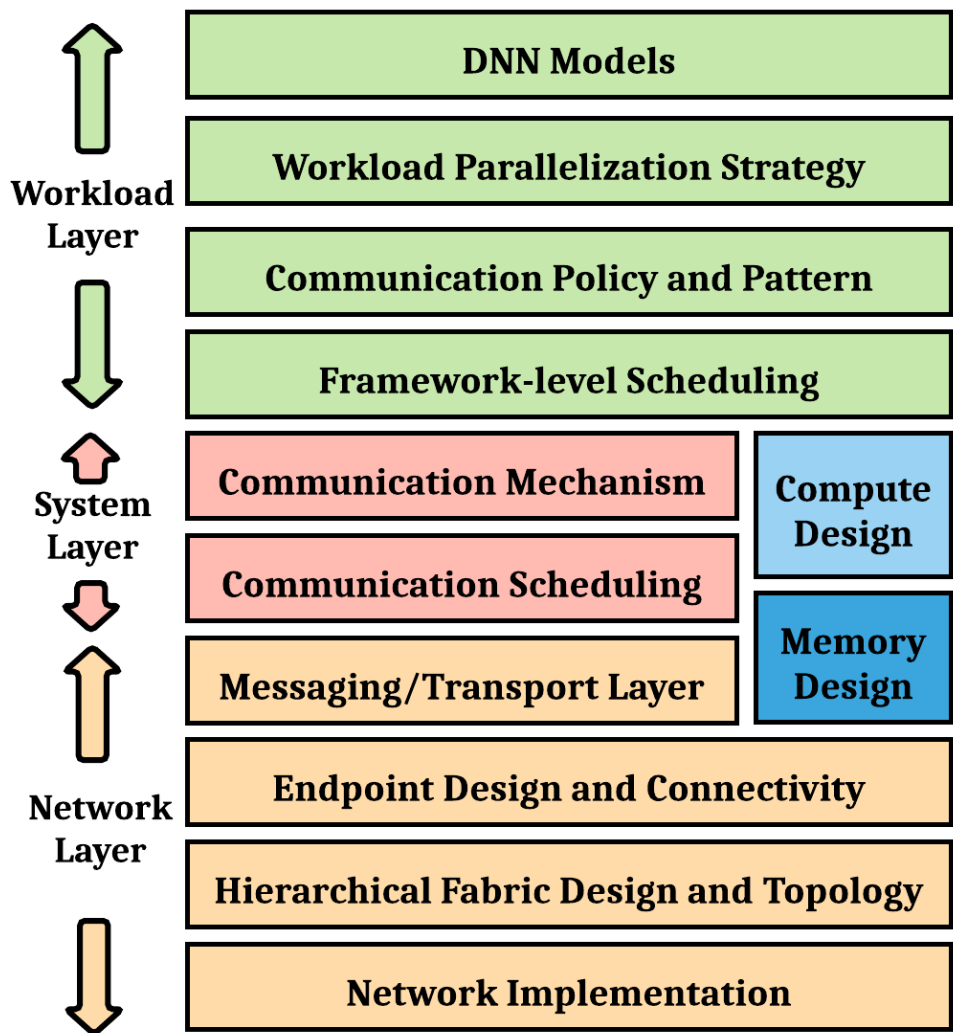
HW/SW Co-Design Space



*How to
model,
navigate
and co-
optimize?*

Figure Courtesy: Srinivas Sridharan (NVIDIA)

Introducing Chakra and ASTRA-sim



Chakra Execution Graph/Trace: an open graph-based representation of AI/ML workload execution



ASTRA-sim: Distributed AI system simulator



Figure Courtesy: Srinivas Sridharan (NVIDIA)

Outline

- Design Space of AI Platforms
- **Chakra Workload Execution Traces**
- ASTRA-sim Layers and APIs
- External Engagements
- Conclusion

HW/SW Co-Design Space

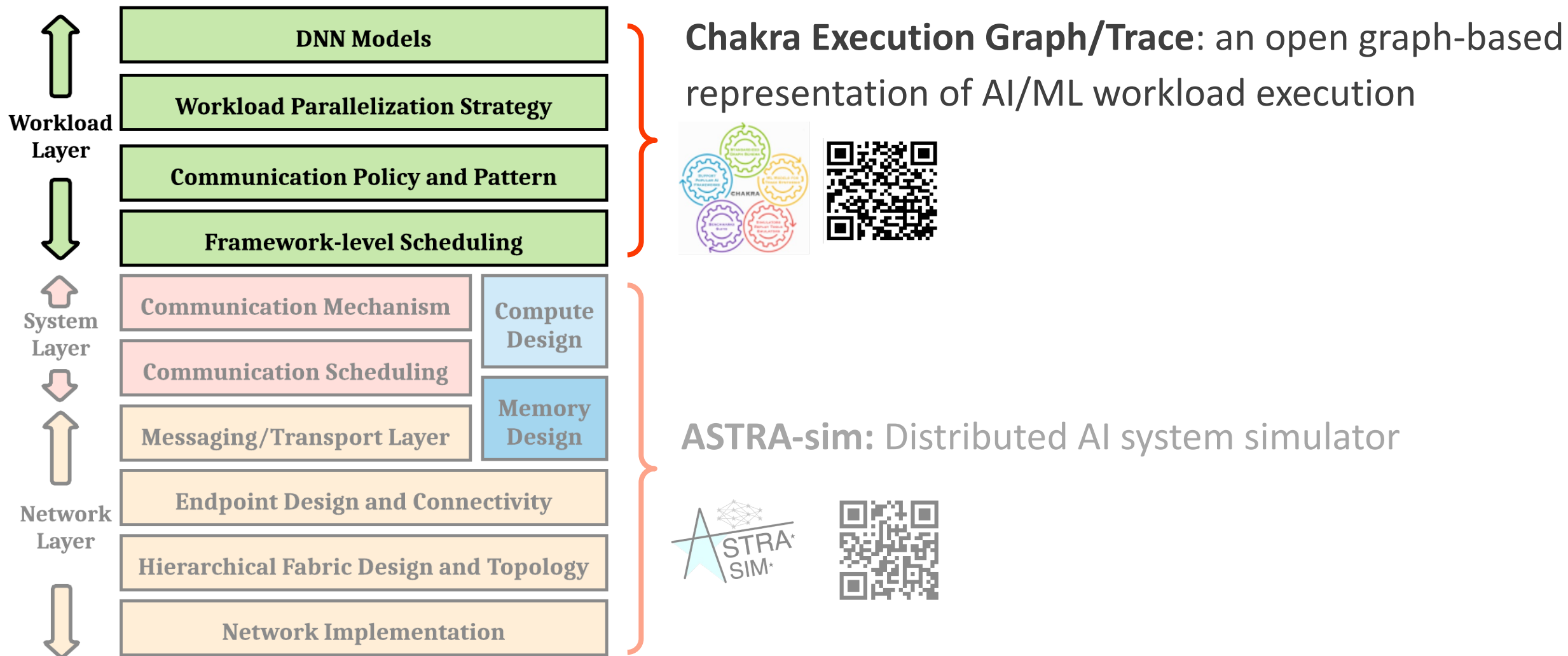
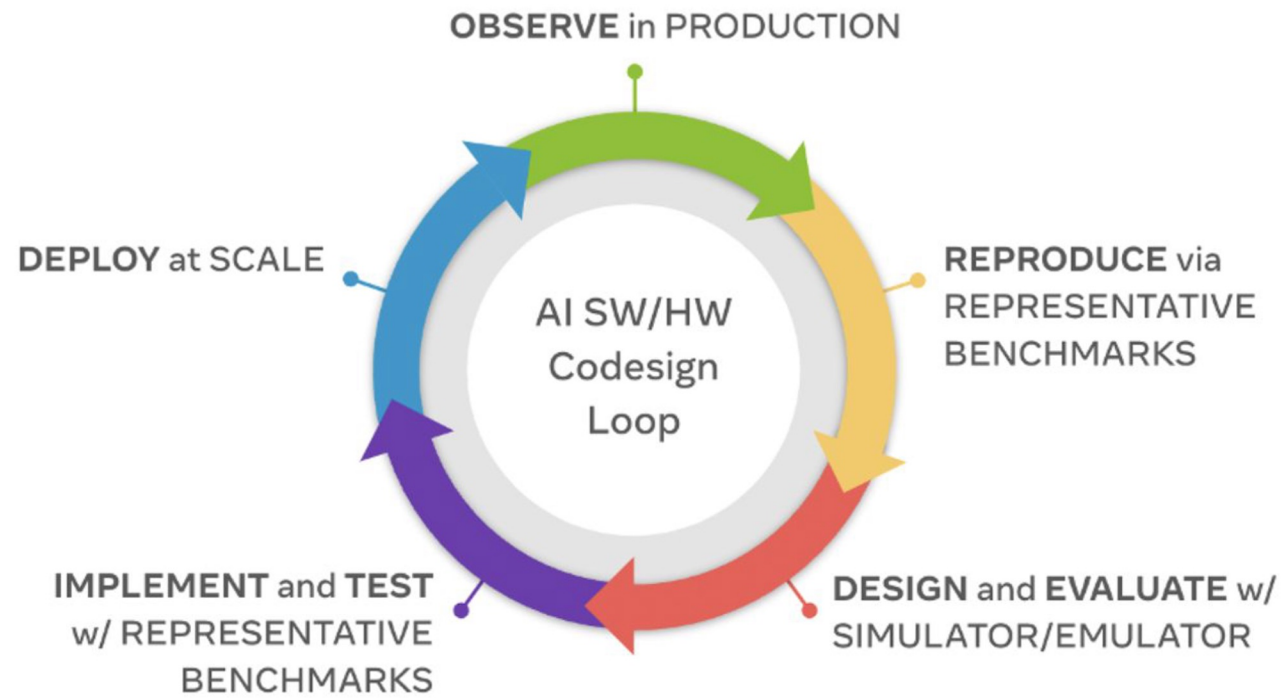


Figure Courtesy: Srinivas Sridharan (NVIDIA)

Chakra: Motivation



Motivation

- High-cost of running full workload benchmarks
- Requires cross-domain full-stack expertise
- Difficult to isolate specific HW/SW bottlenecks
- Difficult to isolate compute, memory, network behavior
- Cannot keep up with the pace of AI innovation
- Hard to obfuscate proprietary AI model details
- Hard to reproduce without support infrastructure

Chakra Execution Graphs

• Hierarchical DAG

• Nodes

- Primitive operators: compute, comms, memory
- Tensor objects: shape, size, device (local/remote)
- Timing and resource constraints

• Edges

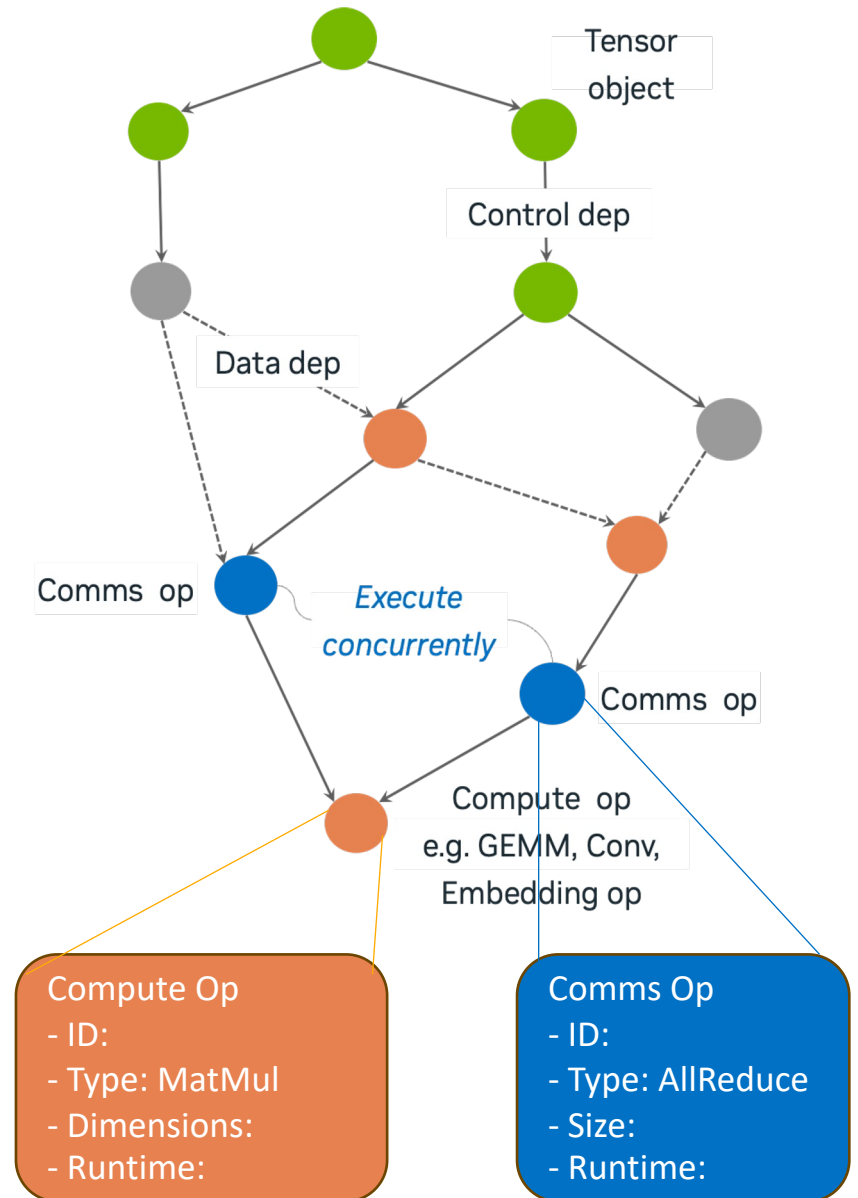
- Data dependency
- Control dependency (e.g. call stack)

• Higher-level abstractions (e.g., components)

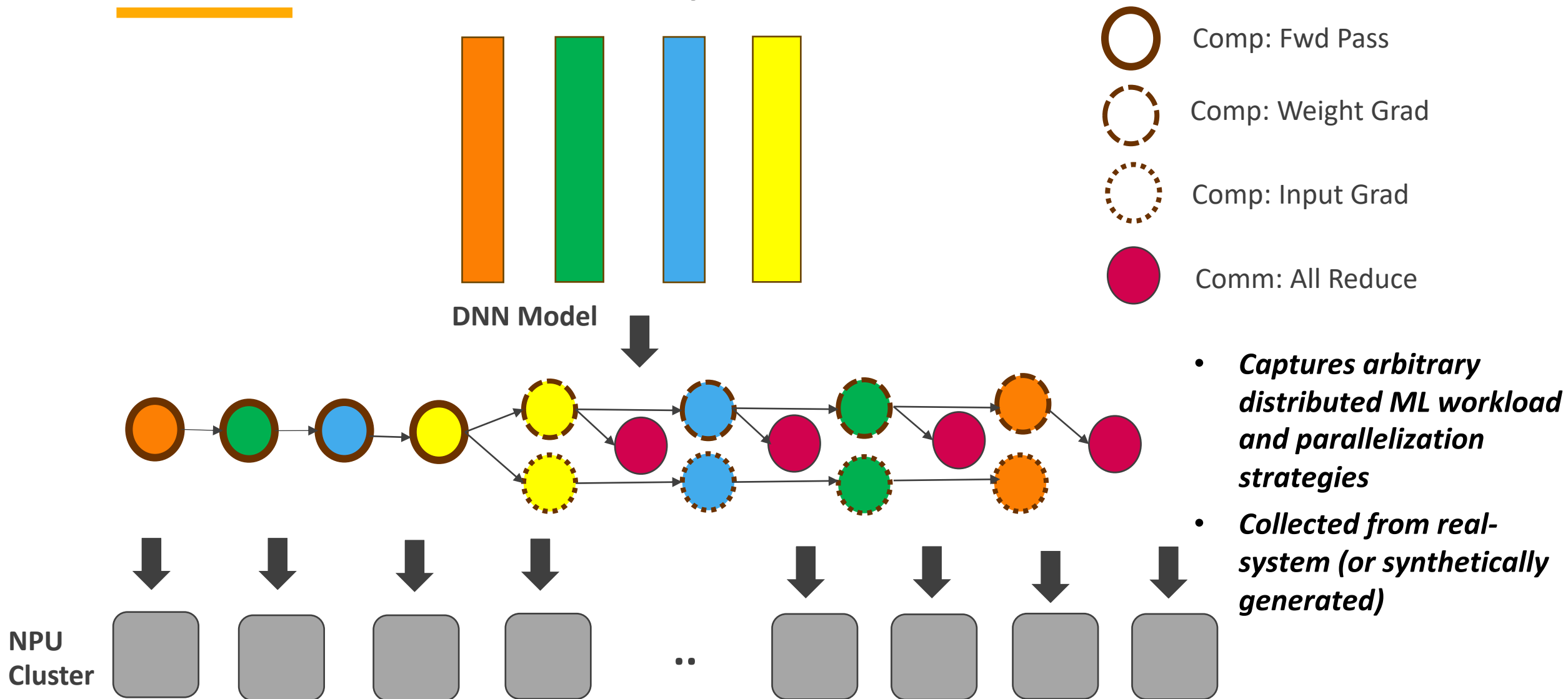
- Comprises of other components or primitive ops

• Benefits

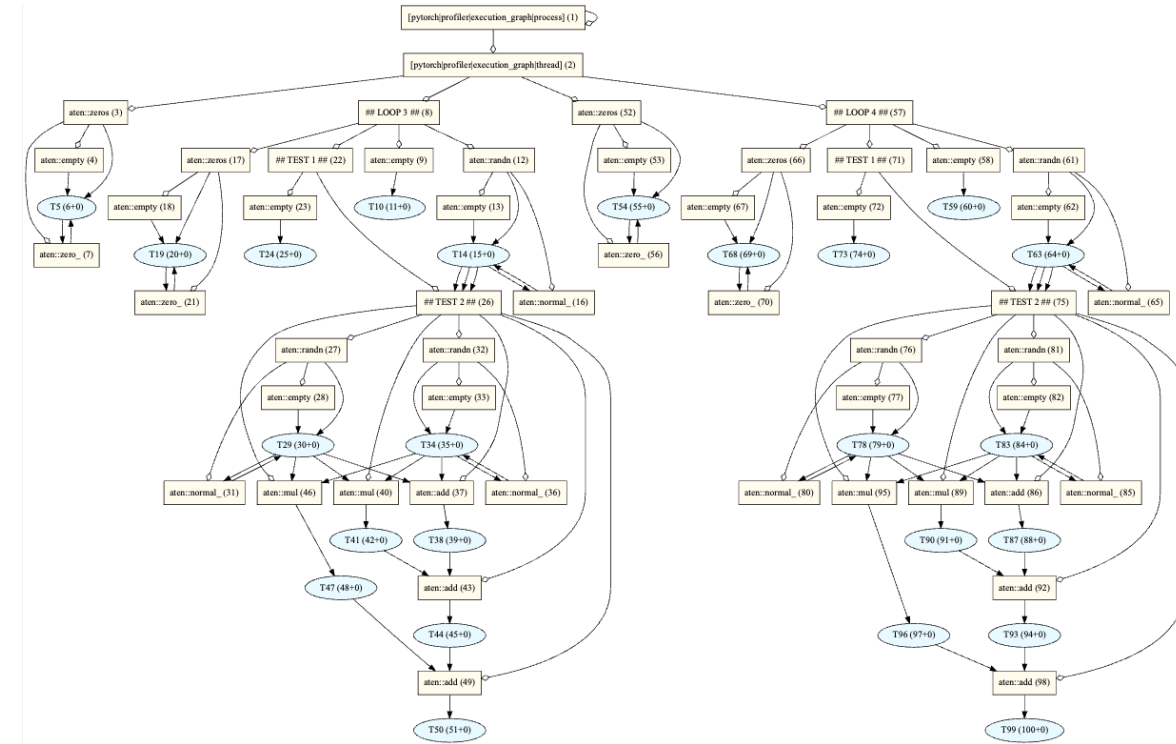
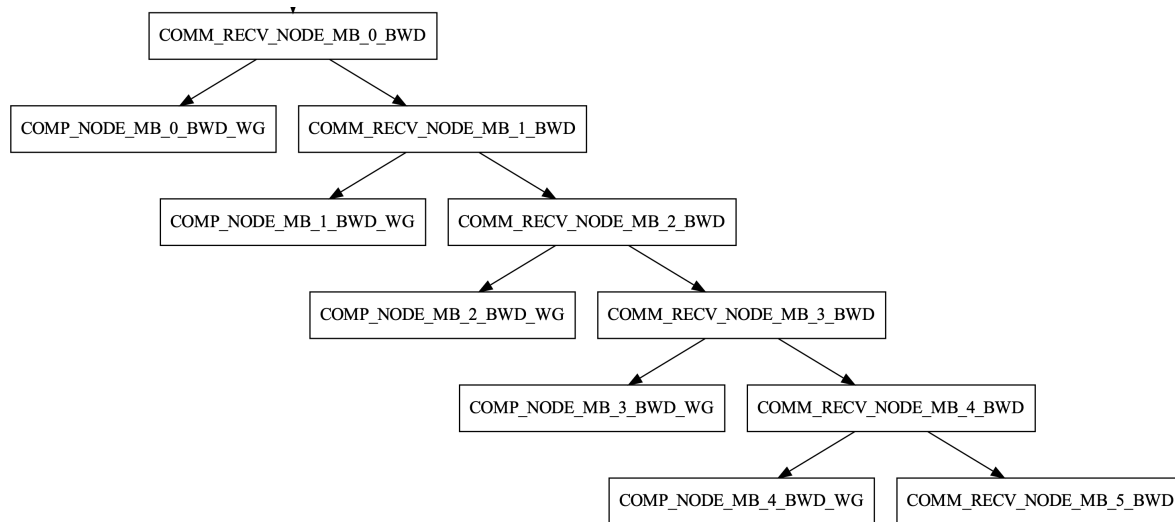
- Isolate comms and compute operators
- Operator, dependencies, and timing for replay, simulation, and analysis
- Flexible to represent both workloads and collective implementations
- Graph transformations to obscure sensitive IP



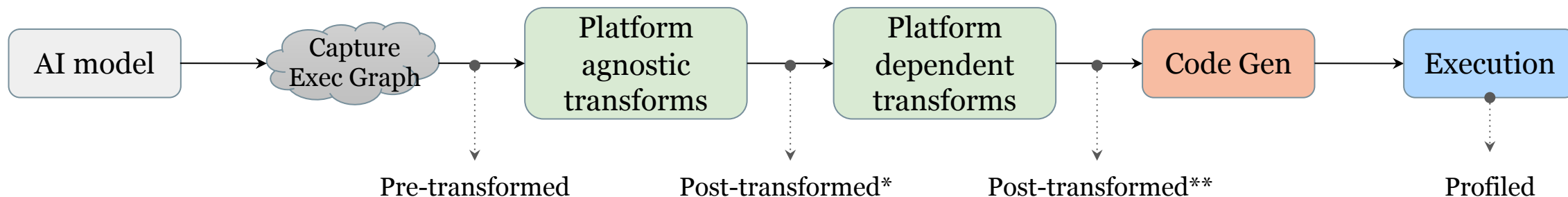
Chakra Execution Graphs



Sample Chakra Traces from Real systems



Chakra Traces: Collection and Synthesis



Type of Execution Traces

1. **Pre-transformed:** original model
2. **Post-transformed#:** optimized graph (may or may not be platform dependent) *through PyTorch2.0 FXgraphs*
3. **Profiled Traces:** graph executed on a specific platform
4. **Synthesized:** via analytical or statistical models

https://github.com/astra-sim/symbolic_tensor_graph

ICLR 2025 paper (GT, Meta, NV): <https://arxiv.org/abs/2411.02322>

[PyTorch] Integrate Execution Graph Observer into PyTorch Profiler #75358

❗ Closed louisfeng wants to merge 1 commit into pytorch:master from louisfeng:export-035342394

```

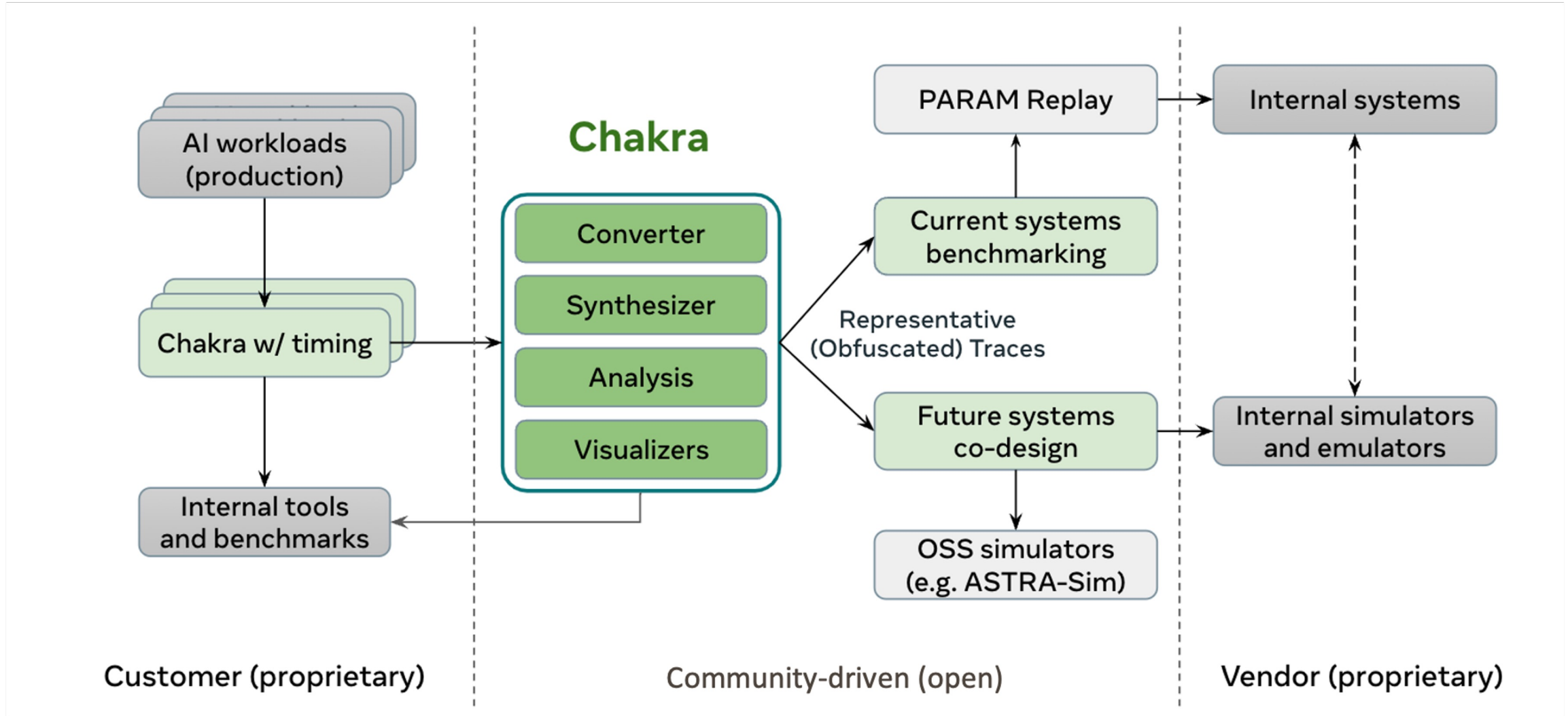
eg = None
if args.eg:
    eg_file = f"{out_file_prefix}_eg.json"
    eg = ExecutionGraphObserver()
    eg.register_callback(eg_file)
    eg.start()

with torch.autograd.profiler.profile(
    args.profile, use_cuda=use_cuda, use_kineto=True, record_shapes=False
) as prof:
    with record_function(f"[param]{run_options['device']}"):
        benchmark.run()

if eg:
    eg.stop()
    eg.unregister_callback()
    logger.info(f"exection graph: {eg_file}")
  
```

Code modifications

Chakra Ecosystem and End-to-End Flow



Outline

- Design Space of AI Platforms
- Chakra Workload Execution Traces
- **ASTRA-sim Layers and APIs**
- External Engagements
- Conclusion

HW/SW Co-Design Space

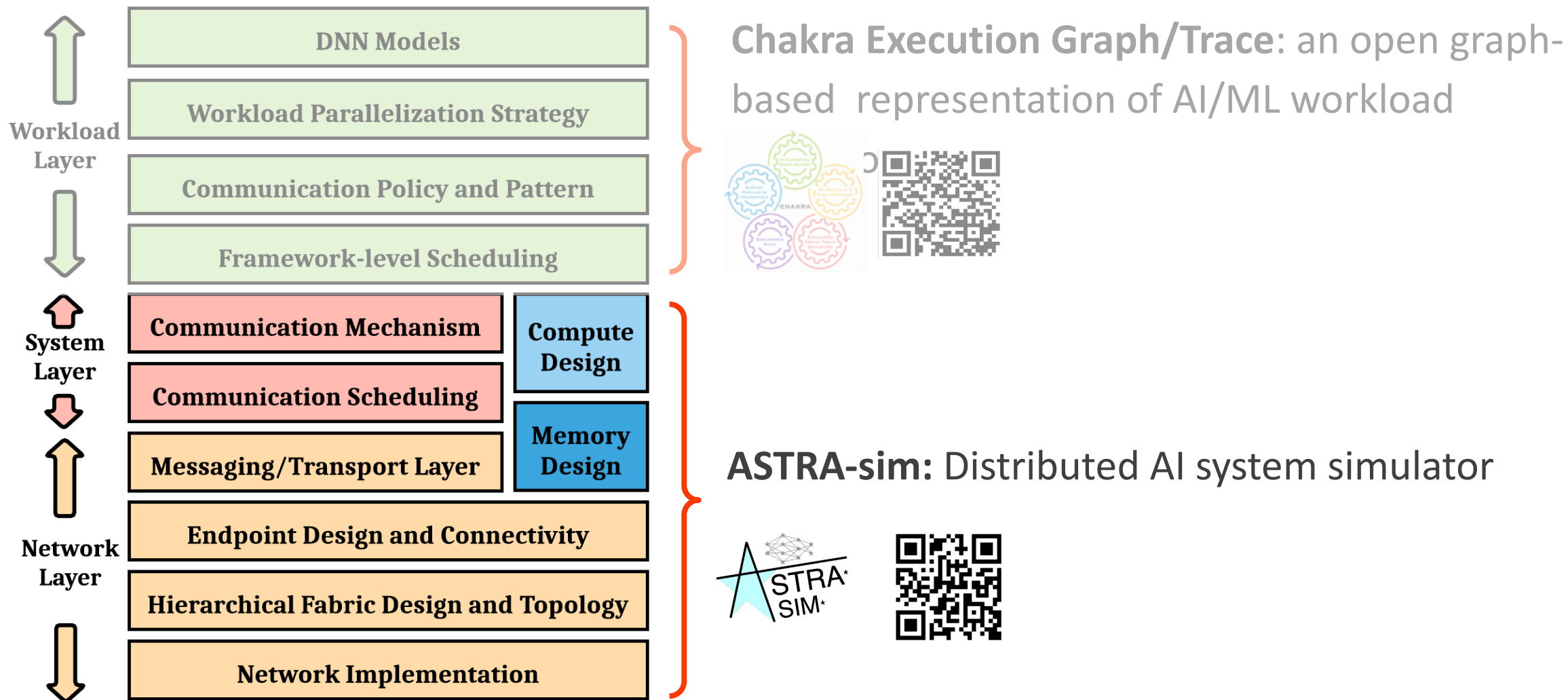
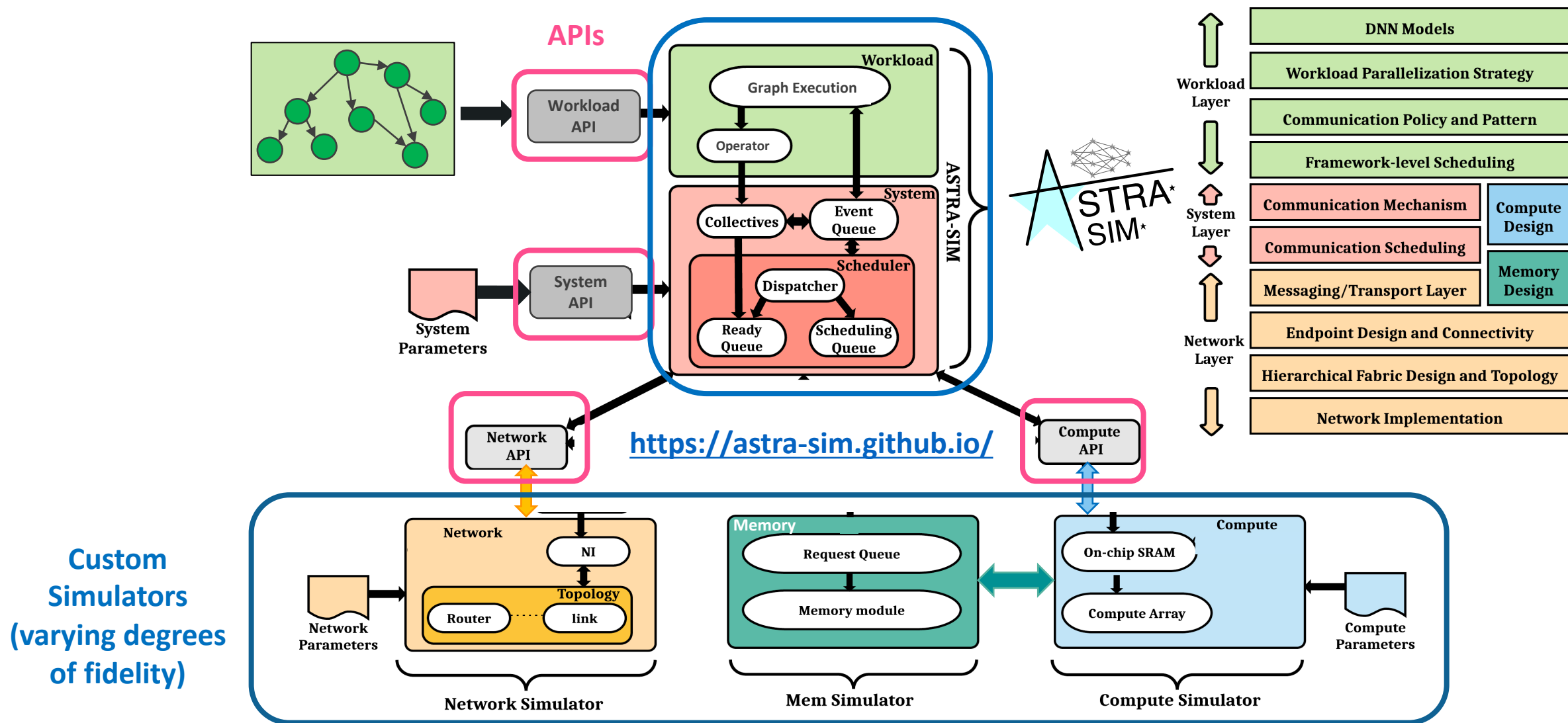


Figure Courtesy: Srinivas Sridharan (NVIDIA)

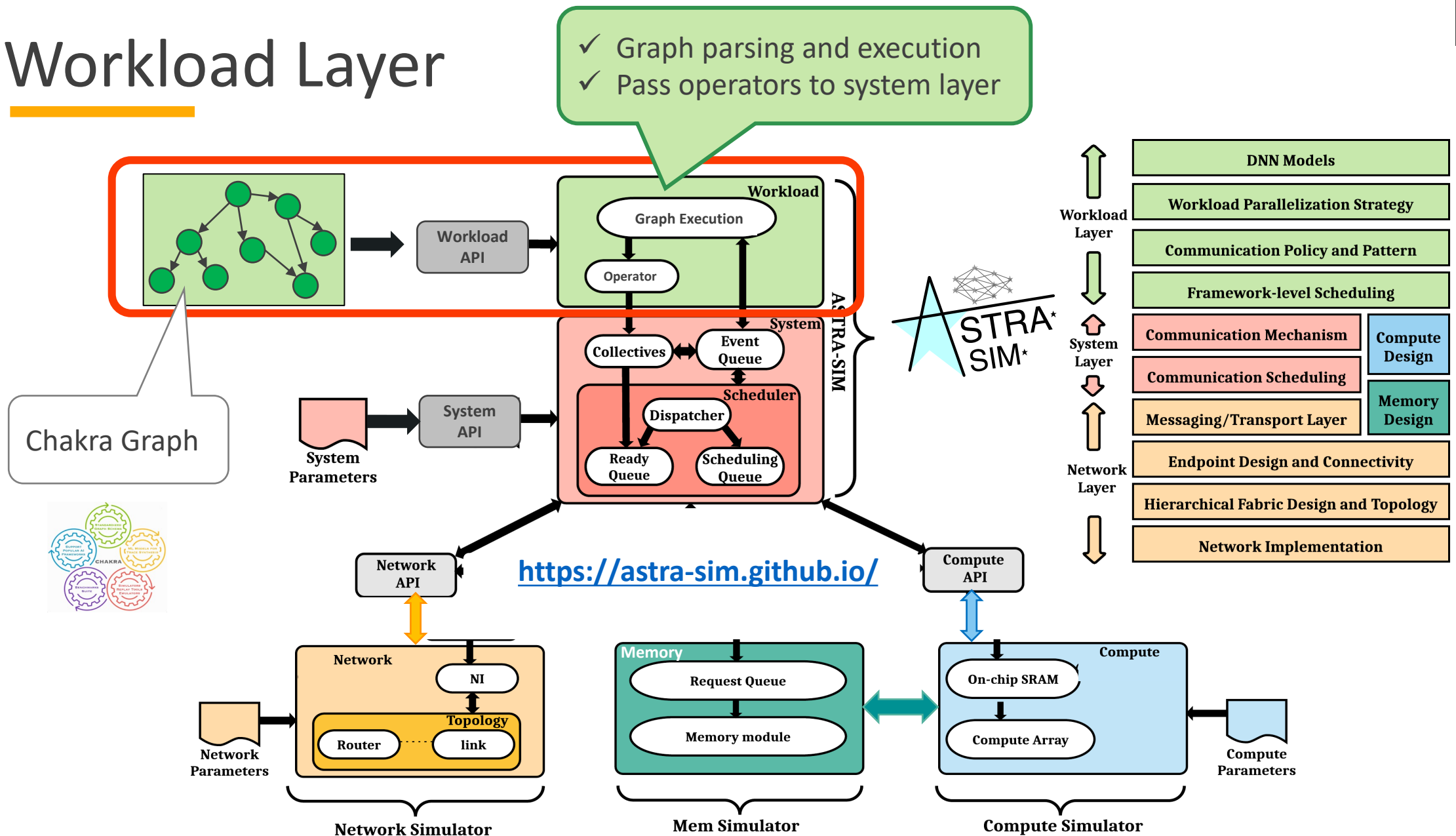
Overview of ASTRA-sim



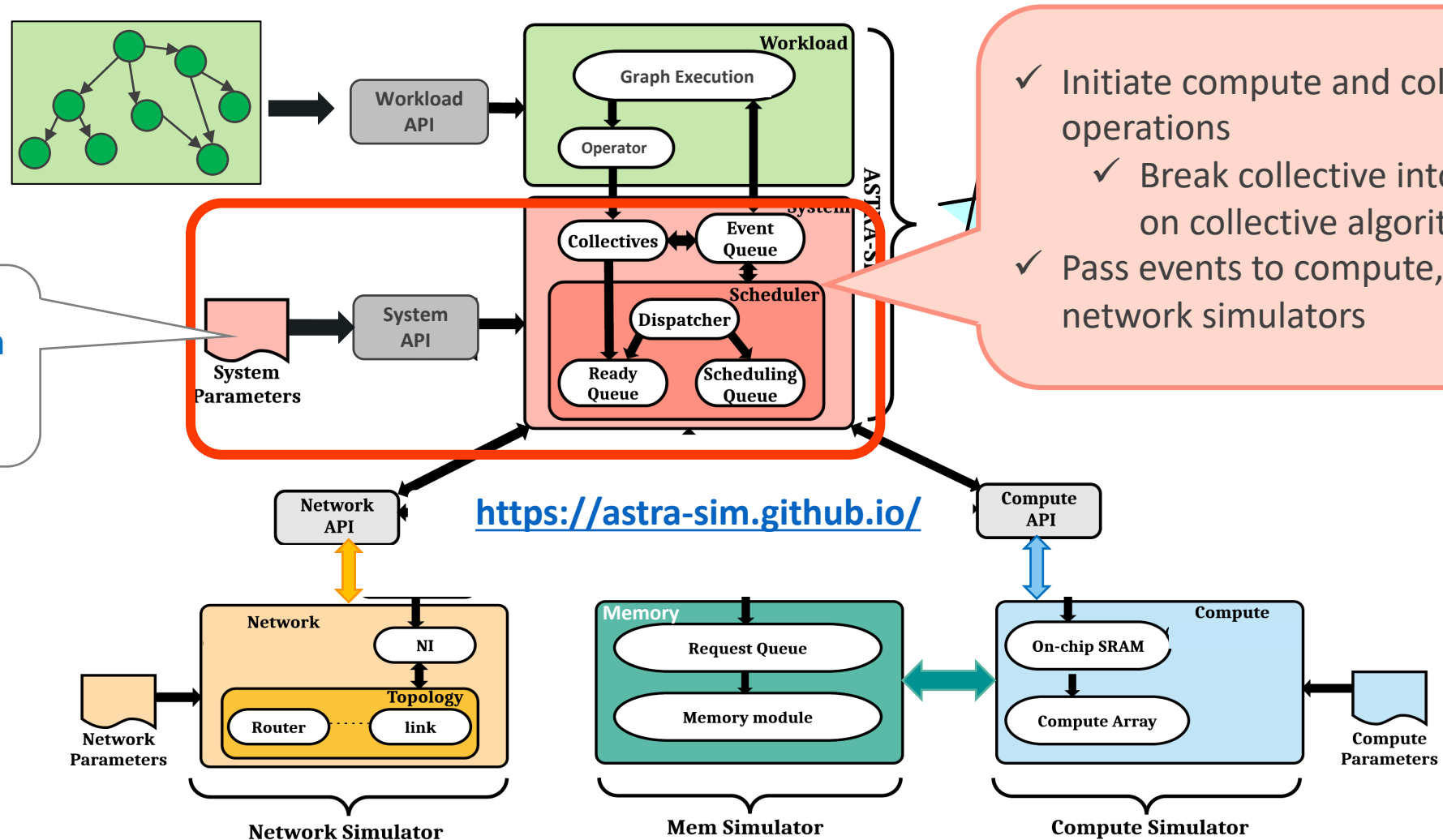
ASTRA-sim: Design Principles

- User determines the model/simulator for compute/network/memory depending on the level of detail and simulation time they want
- **Key enabler:** APIs for plugging in diverse external open/proprietary tools (i.e., composable simulators)
- **Reference Implementation:** <http://github.com/astra-sim/astra-sim>
- **Website:** <https://astra-sim.github.io/>
- **Tutorials:** <https://astra-sim.github.io/tutorials>

Workload Layer



System Layer



Collective Algorithms

**Collective
Algorithm
Synthesizer**

e.g., TACCL, TACOS

<https://github.com/astra-sim/tacos>

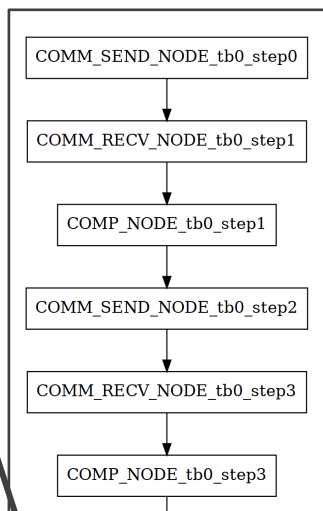
**MSCCLang
DSL**

<https://github.com/microsoft/msccl>

Custom DSLs

**XML
Representation
(MSCCL-IR)**

CollectiveAPI



Workload Layer

**System
Layer**

Collective Pattern

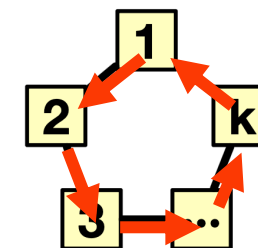
**Collective
Algorithm**

**Custom
Algorithm**

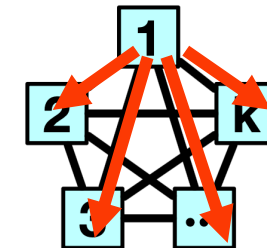
**ASTRA-sim
Native**

NetworkAPI

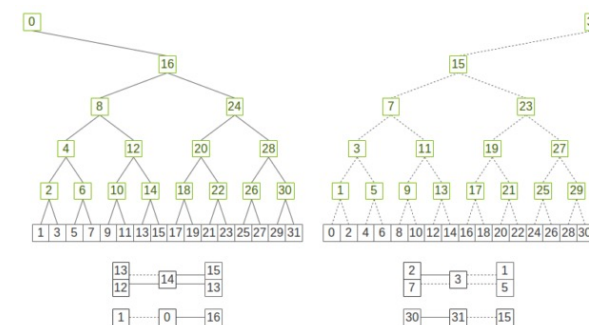
Network Simulator



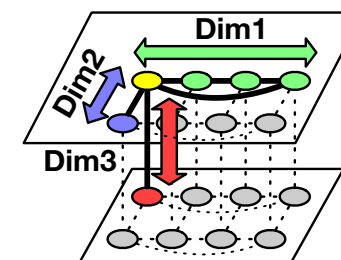
Ring



Direct



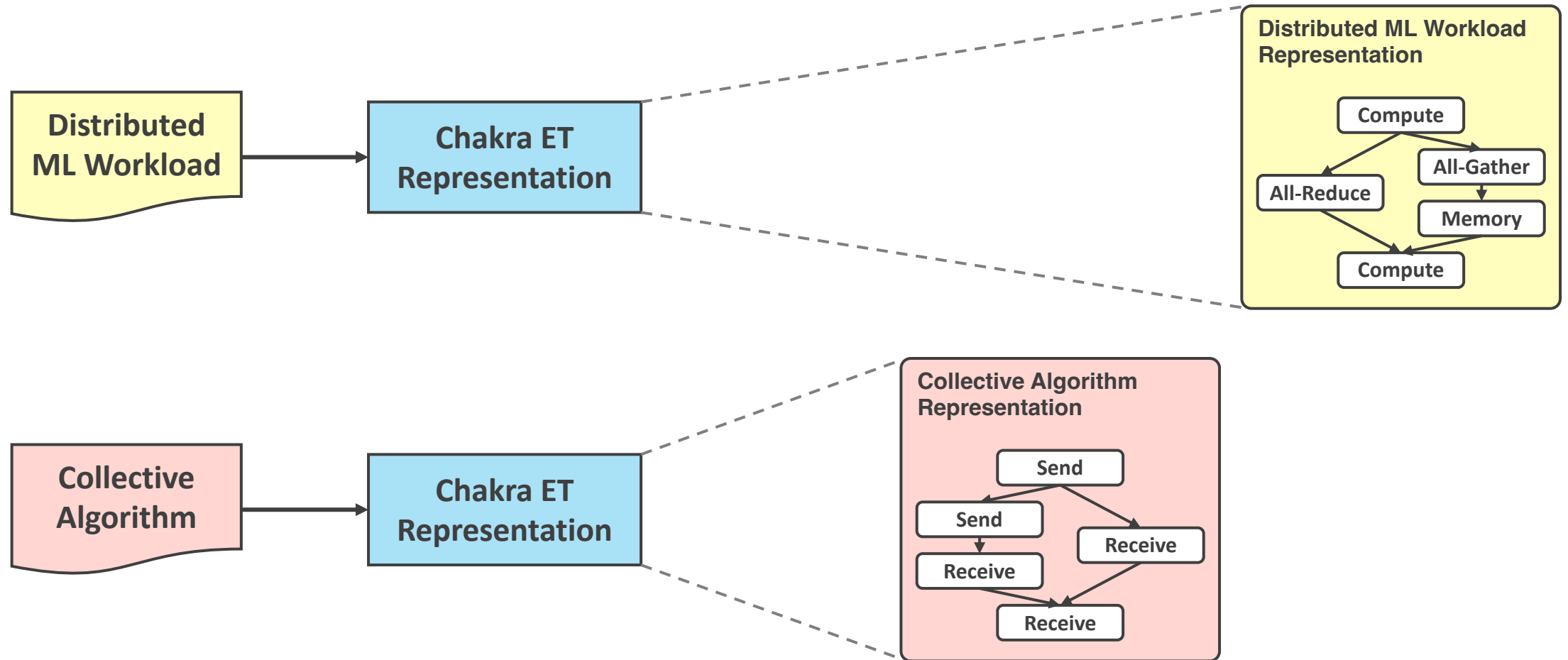
Double Binary Tree



Hierarchical

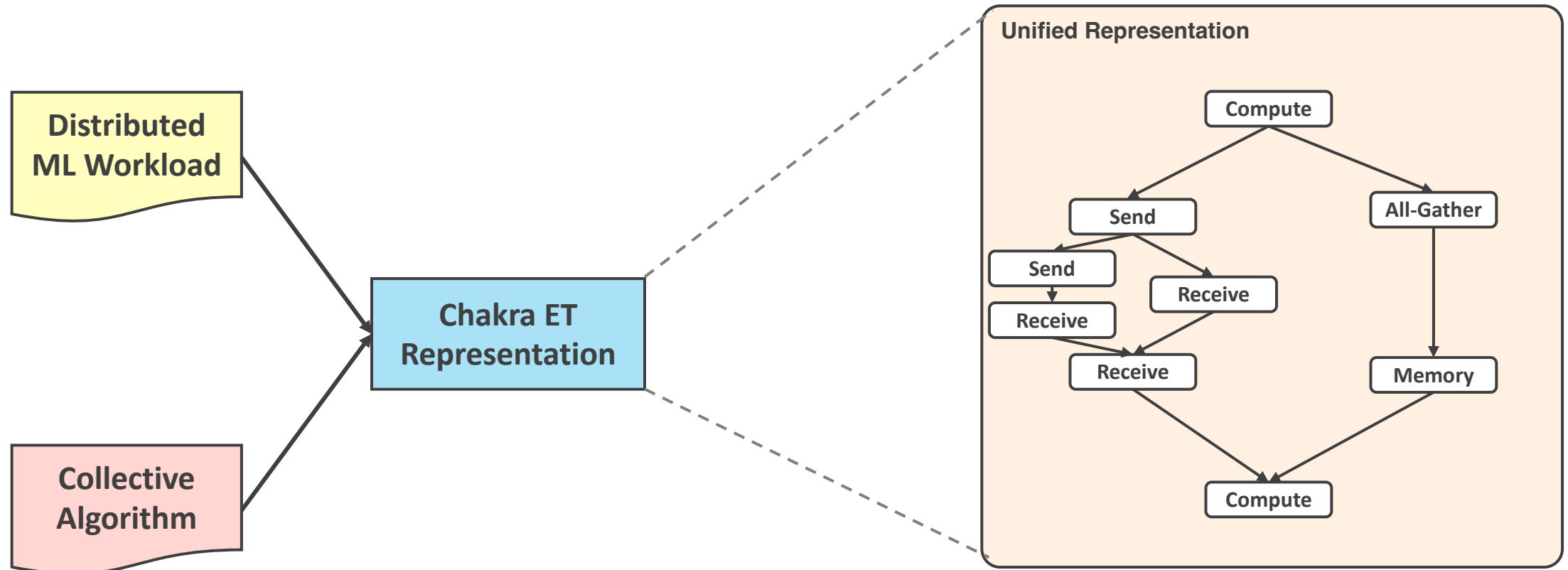
- (1) Dim 1: Reduce-Scatter
- (2) Dim 2: Reduce-Scatter
- (3) Dim 3: Reduce-Scatter
- (4) Dim 3: All-Gather
- (5) Dim 2: All-Gather
- (6) Dim 1: All-Gather

Collective API (uses Chakra format)



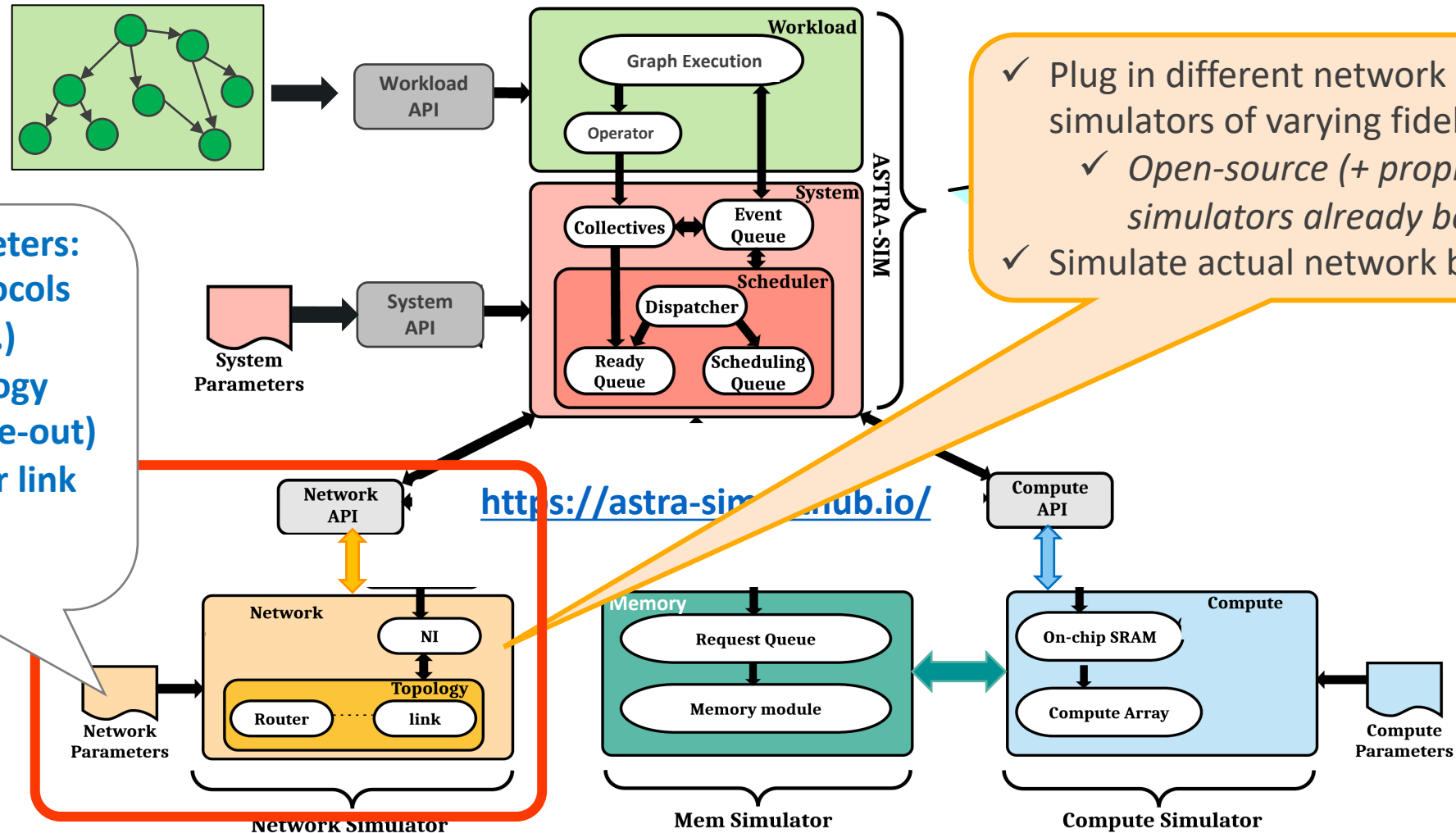
Jinsun Yoo, William Won, Meghan Cowan, Nan Jiang, Benjamin Klenk, Srinivas Sridharan, and Tushar Krishna, **“Towards a Standardized Representation for Deep Learning Collective Algorithms”**, *In Proc. of the 31st IEEE Hot Interconnects Symposium (HotI)*, Aug 2024
<https://arxiv.org/abs/2408.11008>

Vision: Compute-Collective Co-Optimization



Common representation using **Chakra ET** format
→ enable fine-grained co-optimization

Network Layer



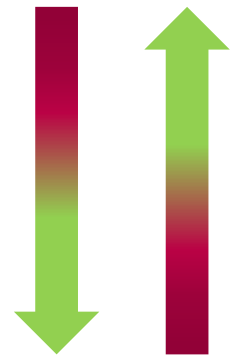
Network parameters:

- Transport protocols (TCP, RDMA, etc.)
- Network topology (scale-up → scale-out)
- BW/latency per link
- Buffering and Arbitration

- ✓ Plug in different network models / simulators of varying fidelity
- ✓ Open-source (+ proprietary) simulators already being used
- ✓ Simulate actual network behavior

Network Layer

- Simulates **actual network behavior** (send/recv)
- Supports **multiple** network models/simulators through NetworkAPI
 - Enabling the simulation of various scales/fidelity
 - We currently support **4 network simulators** implementing NetworkAPI

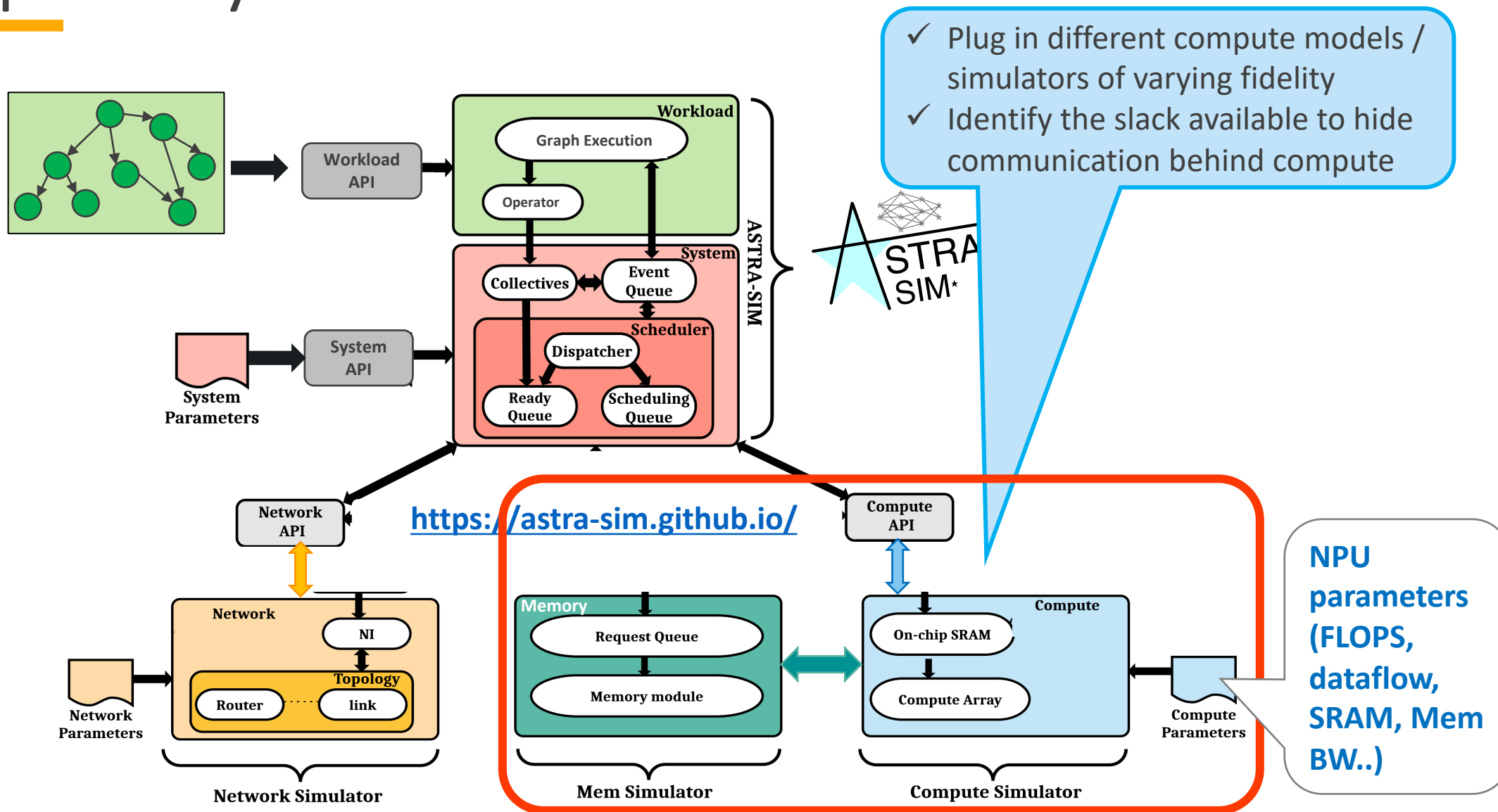
Model/Simulator	Purpose	Notable Feature	
analytical	analytical equation-based simulation	fast simulation hierarchical topologies	 Fidelity Speed
congestion-aware	congestion-aware analytical simulation	+ congestion (queueing) modeling	
gem5-Garnet	on-chip/scale-up network simulation	packetization, flow control, congestion	
HTsim	inter-cluster simulation	TCP-IP, congestion-control, UEC use cases	
ns-3	inter-cluster simulation	RDMA, congestion-control	
Real Network*	Transmit packets through real network	Measured network performance	

Depracated

*In progress

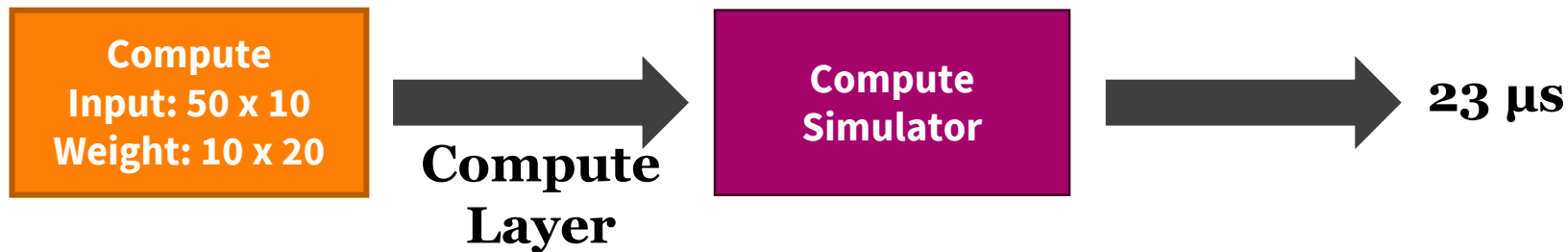
Additional proprietary network simulators implementing the NetworkAPI being used by AMD, Alibaba, HPE Labs, Keysight, some startups ..

Compute Layer



Compute Layer

- Simulates compute times required for Chakra ET's compute node

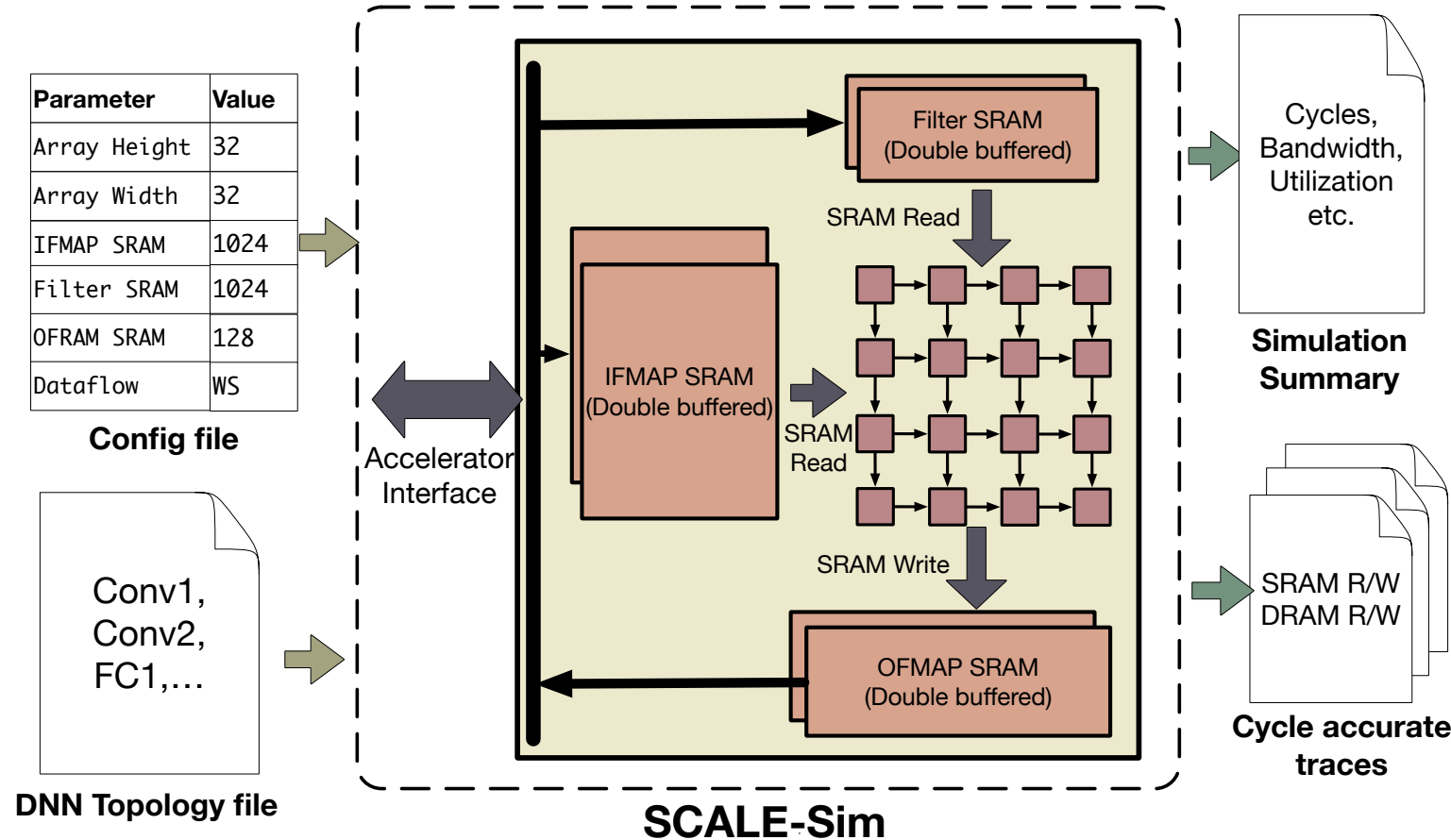


Model	Purpose	Notable Feature	
Roofline	Analytical: First-order roofline	Fast analysis for compute vs memory boundness	
SCALE-sim	Cycle-accurate: systolic array and memory	Models Google-TPUv5-like SoC	
Accel-sim*	GPU simulator	Can run CUDA code	
Real GPU*	Run compute operator on real GPU	Measured runtime	

Fidelity Speed

**In progress*

Sneak Peak: SCALE-Sim NPU Simulator



Cycle accurate systolic array based accelerator simulator

Inputs

1. Architecture configuration
 - Array dimensions
 - Buffer sizes
 - Dataflow (OS, WS, IS)

2. Workload hyper parameters

SCALE sim generates,

1. Runtime in cycles
2. Cycle accurate memory traces
3. Interface bandwidth requirements

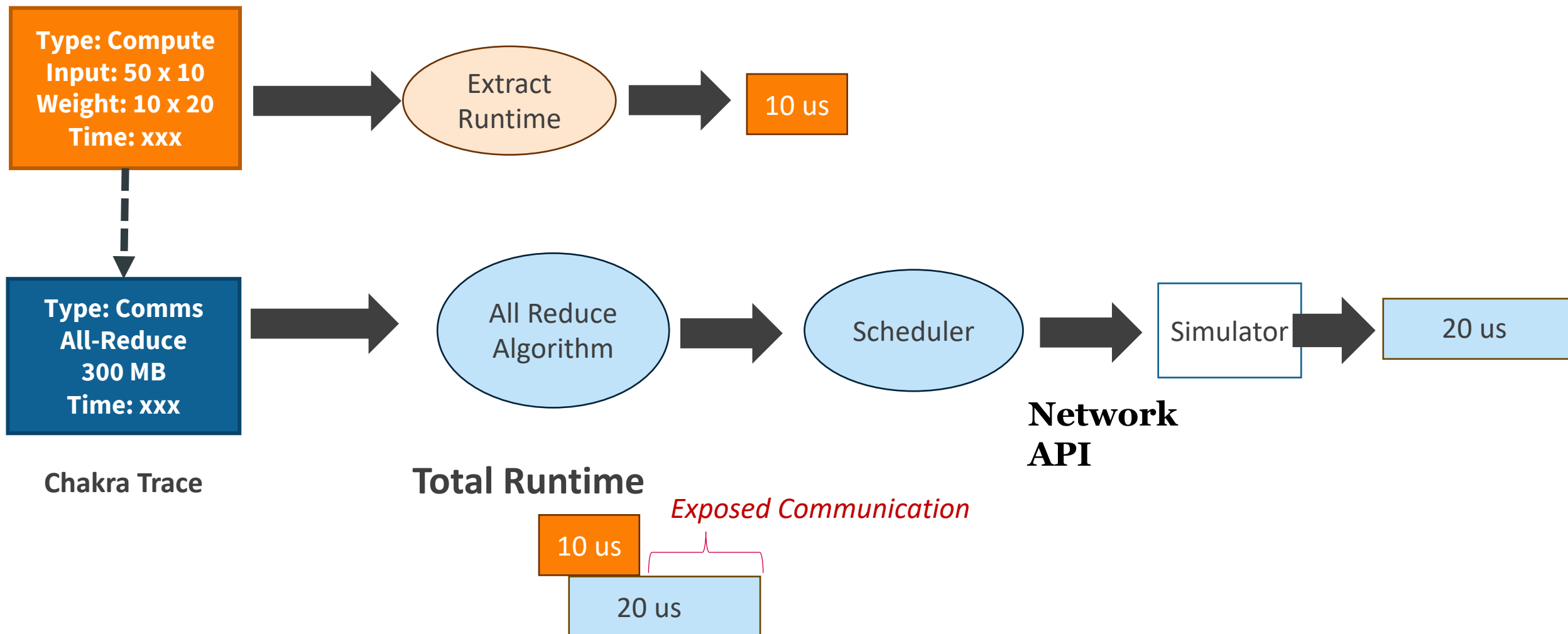
**RAPID-FIRE Talk Tomorrow
at 3:25 PM!**

<https://github.com/scalesim-project/scale-sim>

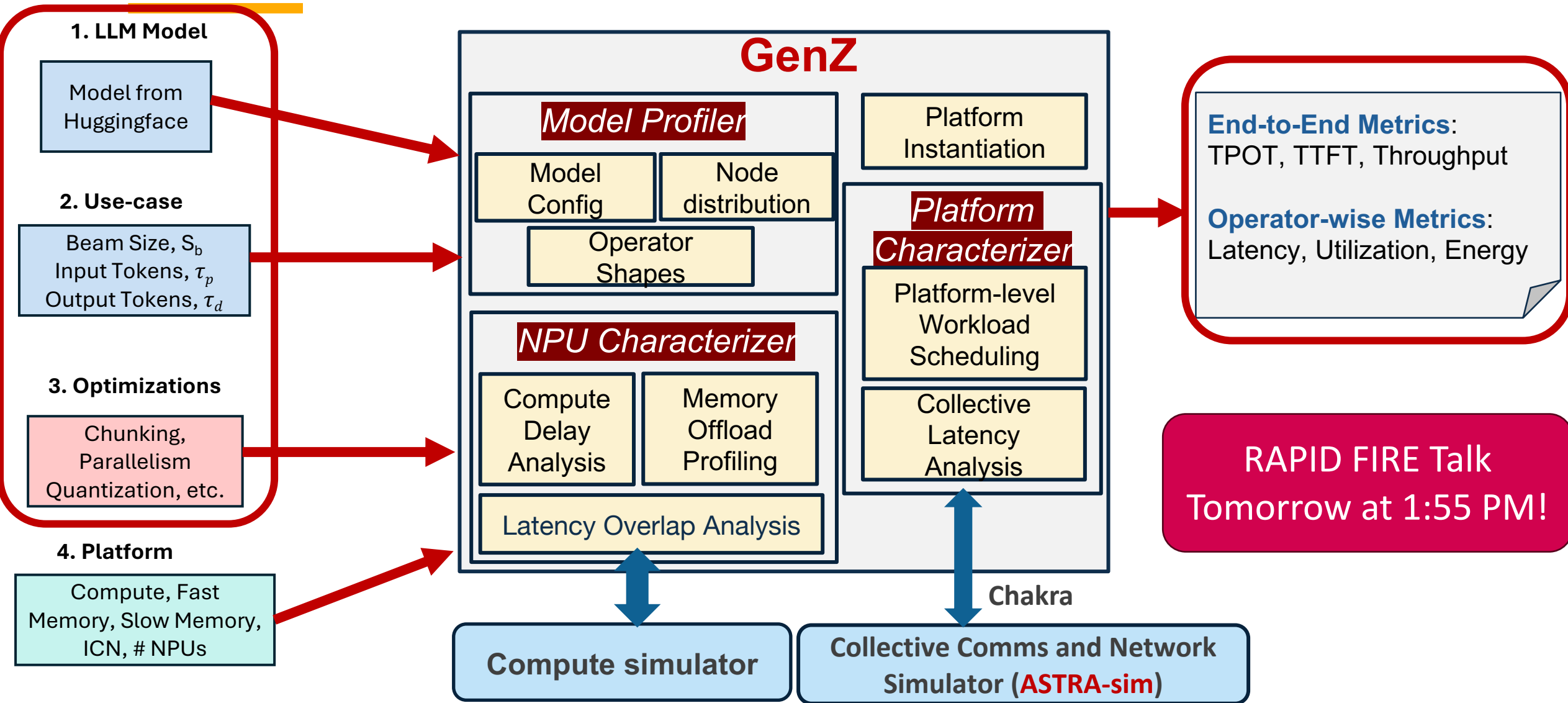
R Raj, S Banerjee, N Chandra, Z Wan, J Tong, A Samajdar, T Krishna, "SCALE-Sim v3: A modular cycle-accurate systolic accelerator simulator for end-to-end system analysis", ISPASS 2025

Example of ASTRA-sim in action

Simulating a system with new network fabric but same GPU



Sneak Peak: Performance Analysis of LLM Inference



Outline

- Design Space of AI Platforms
- Chakra Workload Execution Traces
- ASTRA-sim Layers and APIs
- **External Engagements**
- Conclusion

Chakra has been adopted by MLCommons

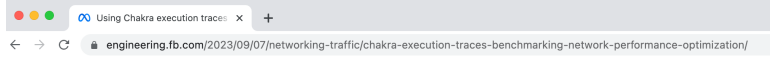
ML Commons



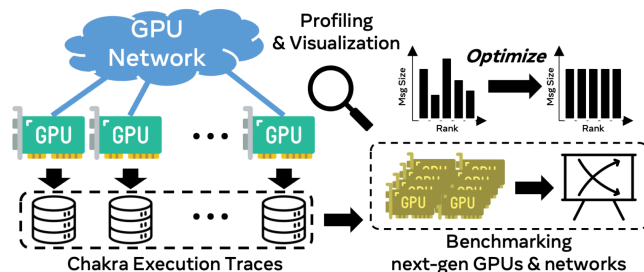
07.31.2023 — San Francisco, CA

Chakra: Advancing Benchmarking and Co-design for Future AI Systems

Announcing Chakra, execution traces and benchmarks working group



Using Chakra execution traces for benchmarking and network performance optimization



Blog Article from Meta

By Louis Feng, Shengbao Zheng, Zhaodan Wang, Wenyin Fu, James Hongyi Zeng

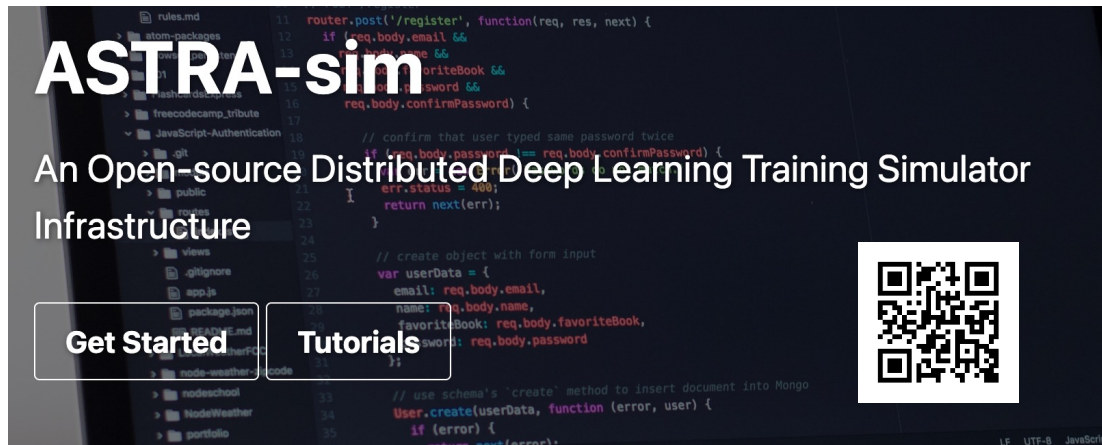


- **MLCommons** helps define industry-standard ML benchmarks
 - Maintainer of MLPerf
- **Chakra Working Group in MLCommons** **Please join!**
 - **Chairs:** Tushar Krishna (GT) + Srinivas Sridharan (NVIDIA)
 - **Active Members (and Users):** NVIDIA, AMD, Intel, Meta, ByteDance, HPE, Juniper, Keysight, Marvell, (+startups), ..

<https://mlcommons.org/working-groups/research/chakra/>

- **Goals of Working Group**
 - Trace Format Standardization
 - Enable easier sharing between hyperscaler/cloud and vendors (with/without NDA)
 - Vendors can focus on different components (compute/memory/network)
 - Trace Collection support (PyTorch/TensorFlow/JAX)
 - Trace Replay and Simulation support tools
 - Trace Benchmark Suite Creation

ASTRA-sim Open-Source Community



Wiki

- <https://astra-sim.github.io/astra-sim-docs/index.html>

Github

- <https://github.com/astra-sim/astra-sim>

Tutorials

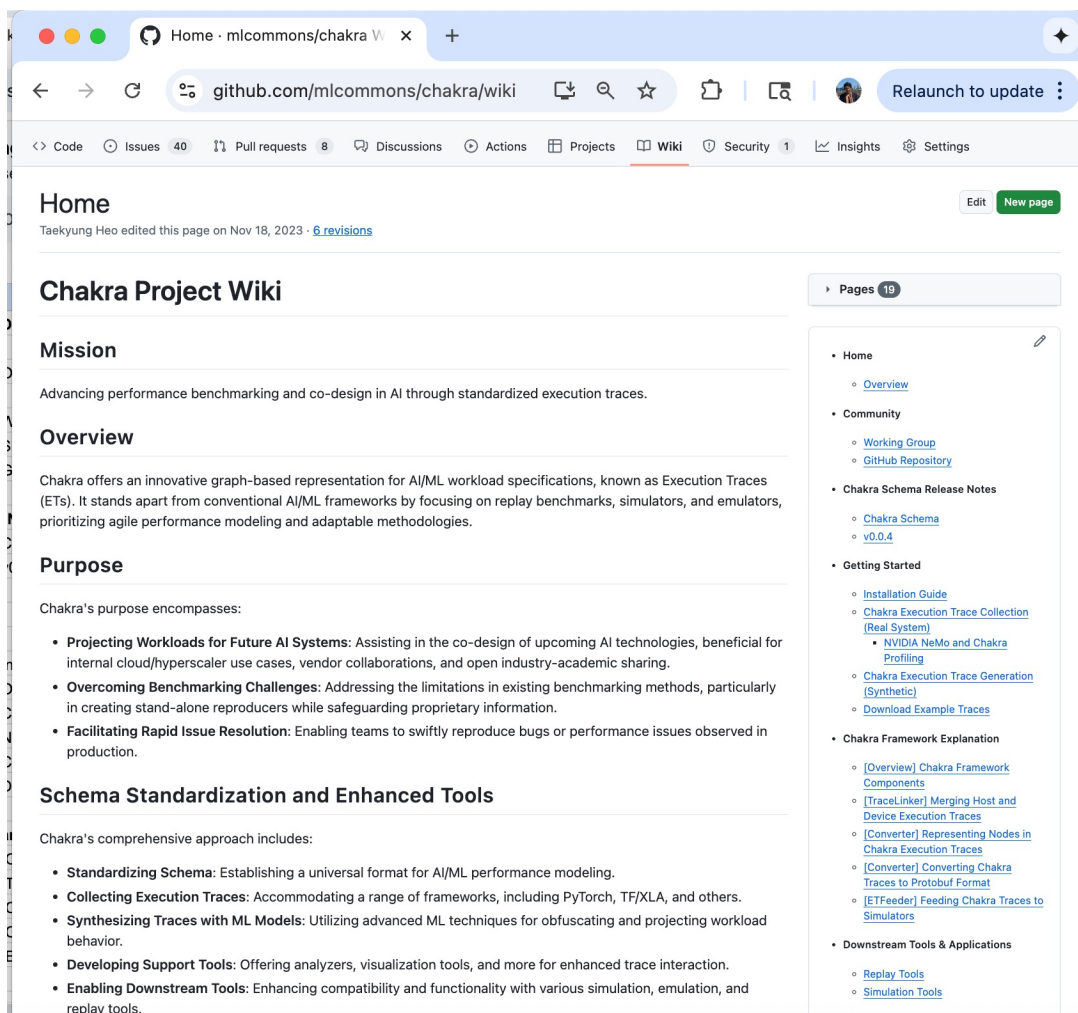
- ASPLOS (2022, 2023), ISCA (2022), MLSys (2022), HotI (2024), MICRO (2024)
- <https://astra-sim.github.io/tutorials>

Growing Userbase (410 stars on github, 150 forks)

- Industry
 - Google, Meta, HPE, Alibaba, ByteDance
 - **AMD***, NVIDIA, Intel, **IMEC***, Qualcomm
 - **Keysight***, Micron, **Marvell***, Samsung
 - ...
- Several startups
- Open Compute Project (OCP) WG on Co-Design
- ML Commons Chakra WG
- Many universities

**Contributing new features*

Resources



Home · mlcommons/chakra Wiki

github.com/mlcommons/chakra/wiki

Code Issues 40 Pull requests 8 Discussions Actions Projects Wiki Security 1 Insights Settings

Home

Taekyung Heo edited this page on Nov 18, 2023 · 6 revisions

Chakra Project Wiki

Mission

Advancing performance benchmarking and co-design in AI through standardized execution traces.

Overview

Chakra offers an innovative graph-based representation for AI/ML workload specifications, known as Execution Traces (ETs). It stands apart from conventional AI/ML frameworks by focusing on replay benchmarks, simulators, and emulators, prioritizing agile performance modeling and adaptable methodologies.

Purpose

Chakra's purpose encompasses:

- Projecting Workloads for Future AI Systems:** Assisting in the co-design of upcoming AI technologies, beneficial for internal cloud/hyperscaler use cases, vendor collaborations, and open industry-academic sharing.
- Overcoming Benchmarking Challenges:** Addressing the limitations in existing benchmarking methods, particularly in creating stand-alone reproducers while safeguarding proprietary information.
- Facilitating Rapid Issue Resolution:** Enabling teams to swiftly reproduce bugs or performance issues observed in production.

Schema Standardization and Enhanced Tools

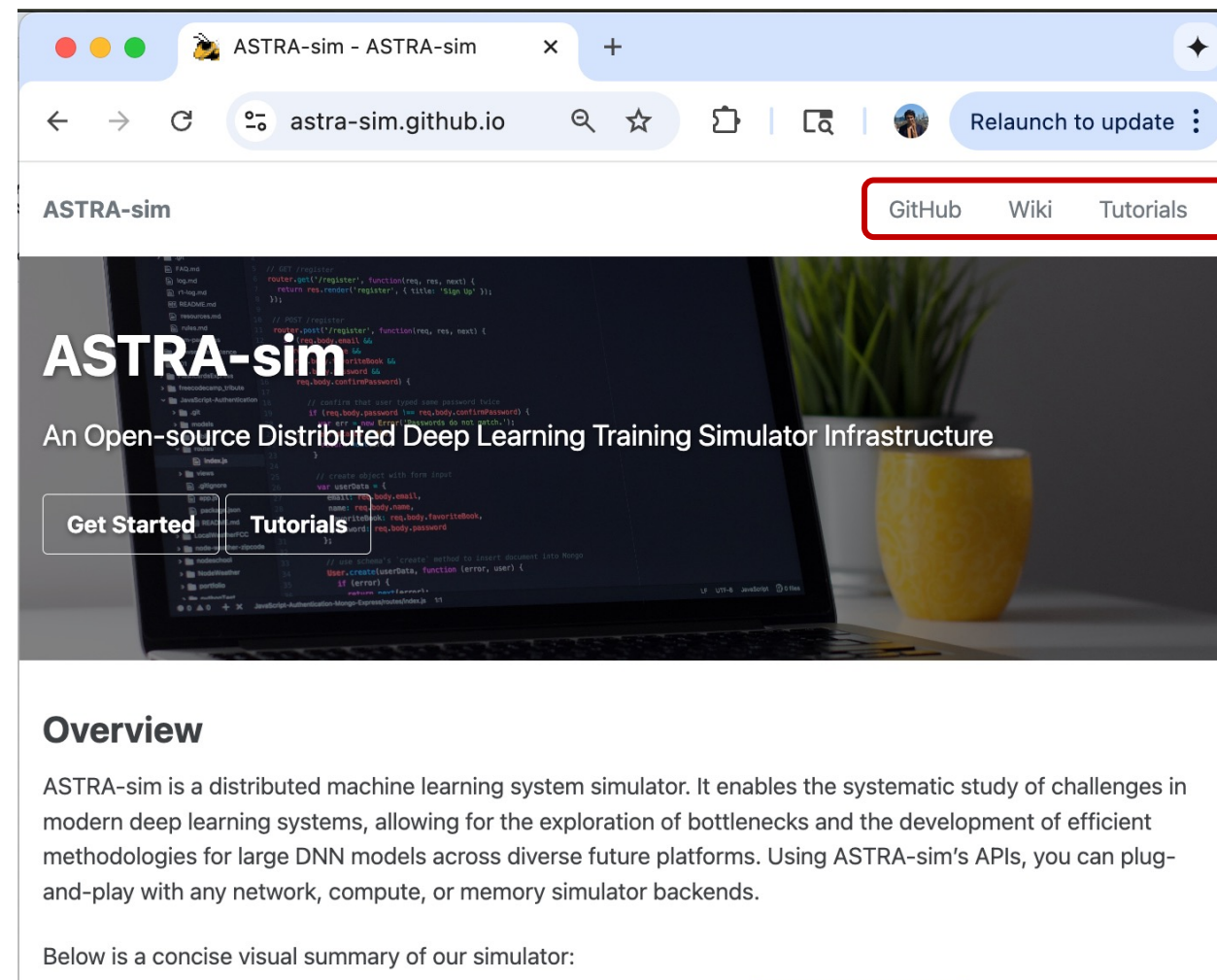
Chakra's comprehensive approach includes:

- Standardizing Schema:** Establishing a universal format for AI/ML performance modeling.
- Collecting Execution Traces:** Accommodating a range of frameworks, including PyTorch, TF/XLA, and others.
- Synthesizing Traces with ML Models:** Utilizing advanced ML techniques for obfuscating and projecting workload behavior.
- Developing Support Tools:** Offering analyzers, visualization tools, and more for enhanced trace interaction.
- Enabling Downstream Tools:** Enhancing compatibility and functionality with various simulation, emulation, and replay tools.

Pages 19

- Home
 - Overview
- Community
 - Working Group
 - GitHub Repository
- Chakra Schema Release Notes
 - Chakra Schema
 - v0.0.4
- Getting Started
 - Installation Guide
 - Chakra Execution Trace Collection (Real System)
 - NVIDIA NeMo and Chakra Profiling
 - Chakra Execution Trace Generation (Synthetic)
 - Download Example Traces
- Chakra Framework Explanation
 - [Overview] Chakra Framework Components
 - [TraceLinker] Merging Host and Device Execution Traces
 - [Converter] Representing Nodes in Chakra Execution Traces
 - [Converter] Converting Chakra Traces to Protobuf Format
 - [ETFeeder] Feeding Chakra Traces to Simulators
- Downstream Tools & Applications
 - Replay Tools
 - Simulation Tools

<https://github.com/mlcommons/chakra/wiki>



ASTRA-sim - ASTRA-sim

astra-sim.github.io

Relaunch to update

ASTRA-sim

GitHub Wiki Tutorials

ASTRA-sim

An Open-source Distributed Deep Learning Training Simulator Infrastructure

Get Started Tutorials

Overview

ASTRA-sim is a distributed machine learning system simulator. It enables the systematic study of challenges in modern deep learning systems, allowing for the exploration of bottlenecks and the development of efficient methodologies for large DNN models across diverse future platforms. Using ASTRA-sim's APIs, you can plug-and-play with any network, compute, or memory simulator backends.

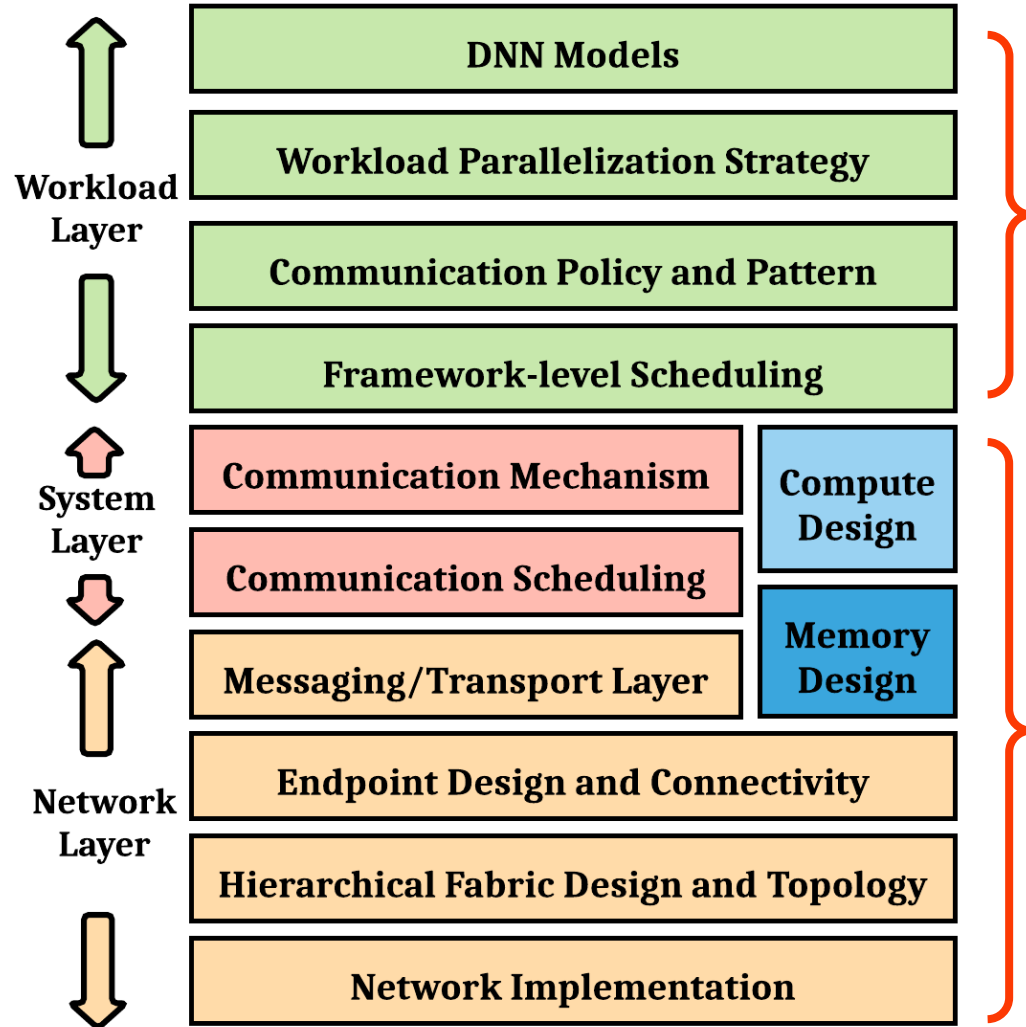
Below is a concise visual summary of our simulator:

<https://astra-sim.github.io/>

Outline

- Design Space of AI Platforms
- Chakra Workload Execution Traces
- ASTRA-sim Layers and APIs
- External Engagements
- **Conclusion**

Summary: Chakra and ASTRA-sim



Chakra Execution Graph/Trace: an open graph-based representation of AI/ML workload execution

- enables isolation and optimization of compute, memory, communication behavior
- an ecosystem for benchmarking, performance analysis, and performance projection

ASTRA-sim: Distributed AI system simulator

- models distributed AI system co-design stack
- allows mix-and-match of performance models for compute, memory and network (API-based)

Thank you!

Figure Courtesy: Srinivas Sridharan (NVIDIA)