

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

ASTRA-sim Tutorial
@HotI 2024
Aug 23, 2024

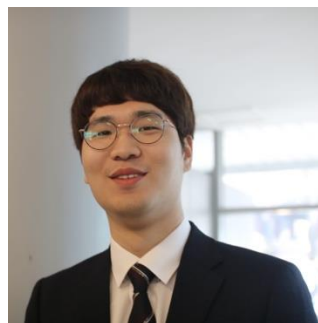
ASTRA-sim and Chakra Tutorial: *Network Layer*

Will Won

Ph.D. Student

School CS, Georgia Institute of Technology

william.won@gatech.edu



ASTRA-sim Tutorial - Agenda

Time (PDT)	Topic	Presenter
3:00 – 3:30 pm	Introduction to Distributed ML	Tushar Krishna
3:30 – 3:45 pm	Overview of Chakra and ASTRA-sim	Tushar Krishna
3:45 – 4:35 pm	Deeper Dive into Chakra and ASTRA-sim	Will Won
	Workload, System, and Network Layers	
4:35 – 4:45 pm	Demo	Will Won
4:45 – 5:00 pm	Closing Remarks	Tushar Krishna

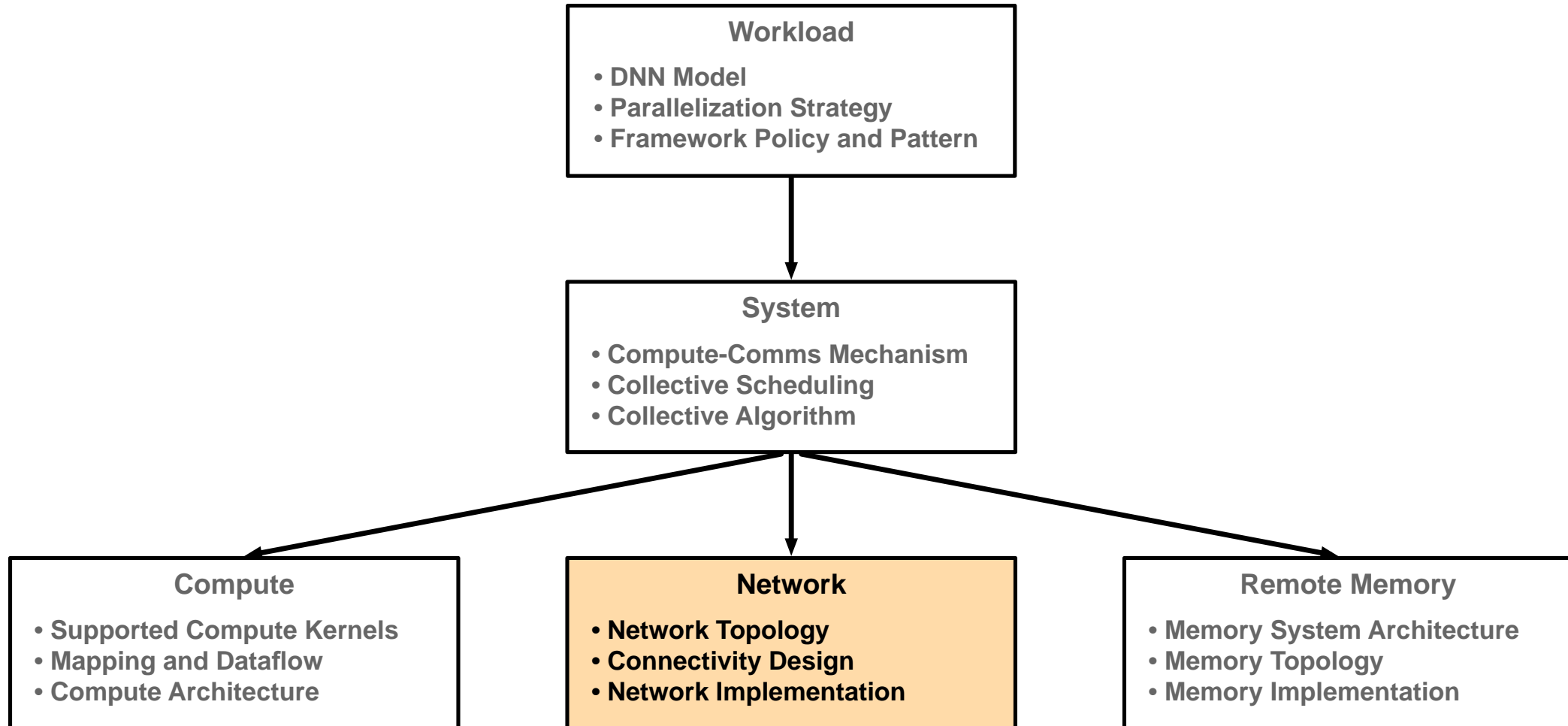
Tutorial Website

includes agenda, slides, ASTRA-sim installation instructions (via source + docker image)

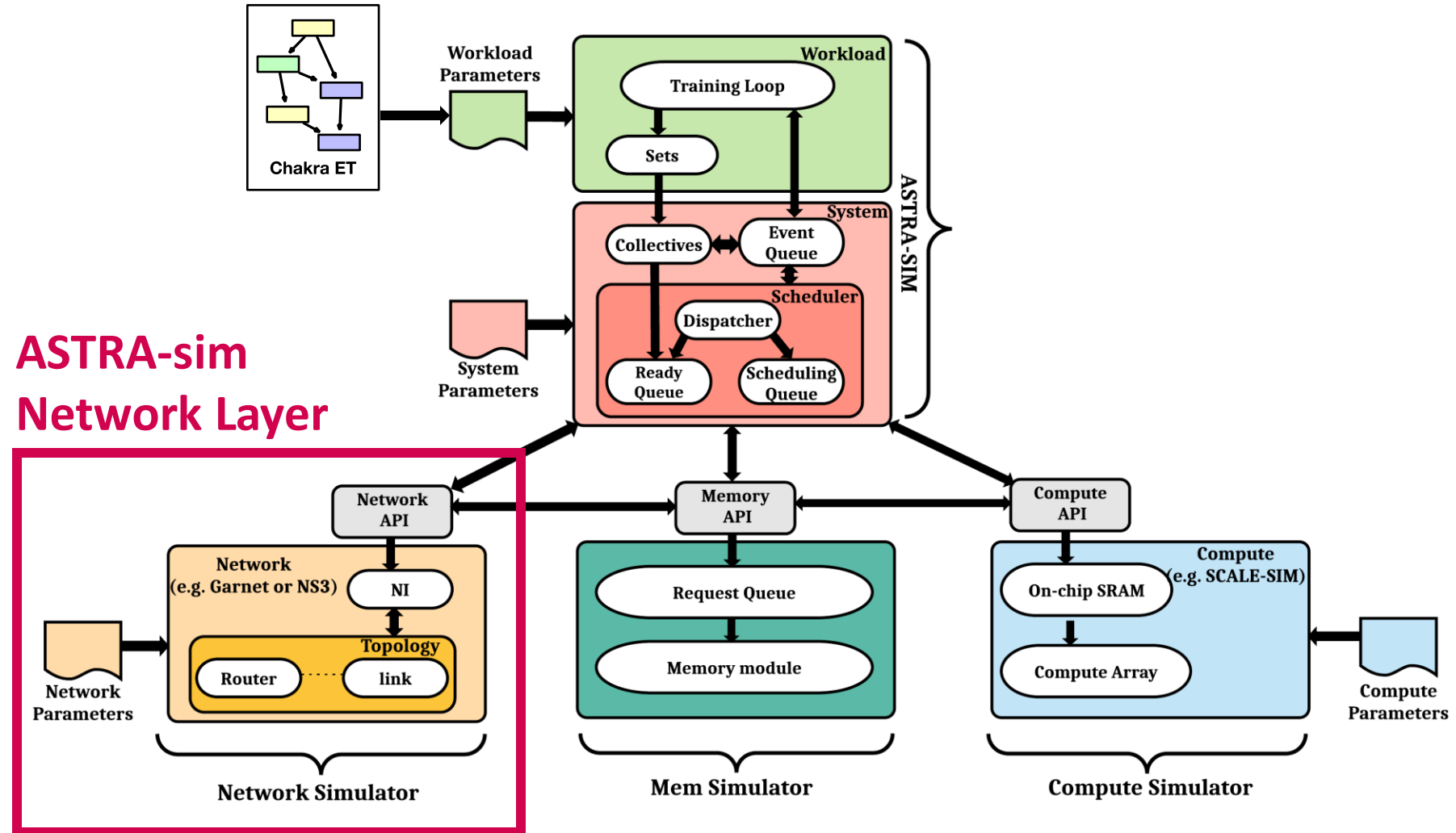
<https://astra-sim.github.io/tutorials/hoti-2024>

Attention: Tutorial is being recorded

Design Space: Network



ASTRA-sim: Network Layer



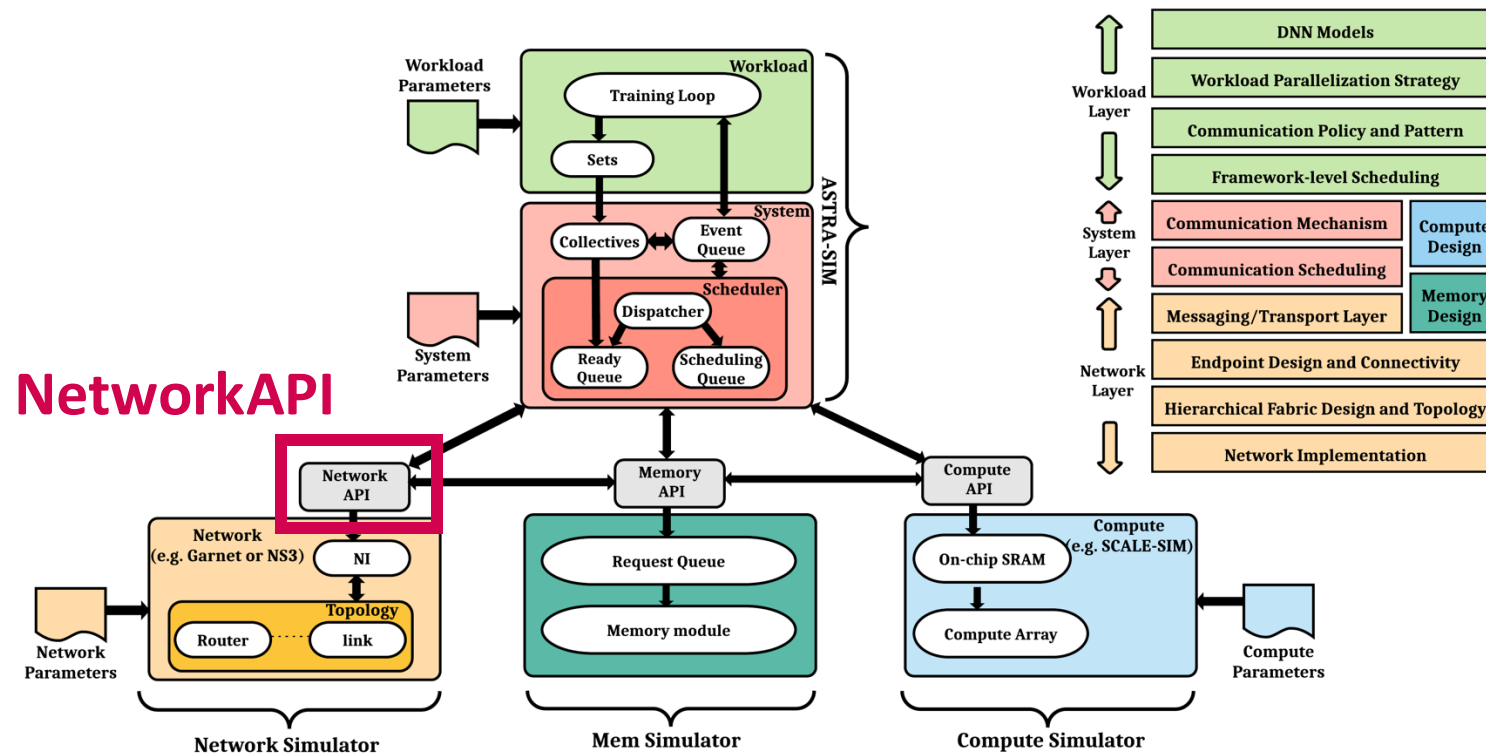
Network Layer

- Network layer simulates **actual network behaviors**
 - Communication protocols (TCP, RDMA, etc.)
 - Network topology
 - BW/latency per link
 - In-network collective communication
 - NIC offloading
 - Compression
 - Buffering, Arbitration

- Through easy **plug-and-play** of any **network simulators**
 - Enabled via **NetworkAPI**

NetworkAPI

- Interface between System layer and Network backend
- Any network simulator implementing the NetworkAPI could be used as ASTRA-sim backend



(HOTI '20) Scalable Distributed Training of Recommendation Models: An ASTRA-SIM + NS3 case-study with TCP/IP transport

Example NetworkAPIs

- **sim_send(msg_size, src, dest, callback)**
 - Simulate sending a message of size msg_size from src through dest and **invoke callback function** once transmission has finished
- **sim_recv(msg_size, src, dest, callback)**
 - Simulate receiving a message of size msg_size from src through dest and **invoke callback function** once transmission has finished
- **sim_schedule(delta, callback)**
 - Invoke callback function after delta time
- **sim_get_time()**
 - Return current time of simulation to the frontend

NetworkAPI at System Layer

- Ring All-Reduce algorithm implementation

```
bool Ring::ready() {  
    (...)  
    stream->owner->sim_send(0, Sys::dummy_data, msg_size, UINT8, packet.preferred_dest, stream->stream_id,  
    &snd_req, &Sys::handleEvent, nullptr);  
  
    (...)  
    stream->owner->sim_rcv(0, Sys::dummy_data, msg_size, UINT8, packet.preferred_src, stream->stream_id,  
    &rcv_req, &Sys::handleEvent, ehd);  
  
    reduce();  
    return true;  
}
```

Send a chunk

Receive a chunk

NetworkAPI Implementation: Example

- Ring All-Reduce algorithm implementation

```
int CongestionAwareNetworkApi::sim_send(...) {  
    (...)  
    // create chunk  
    auto chunk_arrival_arg = std::tuple(tag, src, dst, count, chunk_id);  
    auto arg = std::make_unique<decltype(chunk_arrival_arg)>(chunk_arrival_arg);  
    const auto arg_ptr = static_cast<void*>(arg.release());  
    const auto route = topology->route(src, dst);  
    auto chunk = std::make_unique<Chunk>();  
  
    // initiate transmission from src -> dst.  
    topology->send(std::move(chunk));  
}
```

Select a static route

Trigger actual network simulation

NetworkAPI Implementation varies by network simulation backend

Available Network Backends

- Network backends are maintained separately and are imported as **submodule**.
- We currently have **4 network backends** which implement NetworkAPI

Backend	Purpose	Notable Feature
analytical/analytical	analytical equation-based simulation	fast simulation, hierarchical topologies
analytical/congestion	congestion-aware analytical simulation	first-order congestion (queueing) modeling
Garnet	on-chip/scale-up network simulation	packetization, flow control, congestion
ns-3	inter-network simulation	large parallel GPU clusters

Caveat: Garnet currently only works with ASTRA-sim 1.0 and should be updated

Analytical Backend

- Leverages **analytical equation** to estimate communication delay

- $\text{delay}(\text{msg_size}, \text{src}, \text{dest}) =$
 $\underbrace{(\#hops) \times (\text{link latency})}_{\text{link delay}} + \underbrace{(\text{msg_size}) / (\text{link BW})}_{\text{serialization delay}}$

- `sim_send(msg_size, src, dest, callback)`
 - Estimate communication delay
 - Assign callback to event queue after delay
- No congestion modeling
 - Appropriate for topology-aware collectives without network congestion
- **Fast simulation** for large-scale systems

(ISPASS '23) ASTRA-sim2.0: Modeling Hierarchical Networks and Disaggregated Systems for Large-model Training at Scale

Congestion-aware Analytical Backend

- First-order congestion modeling by per-link queueing
- Per-link delay is calculated using analytical equation
- e.g., `send(msg_size: 1 MB, route: [1, 2, 3, 4, 5])`
 - `send(1 MB, 1 → 2)`
 - `send(1 MB, 2 → 3)`
 - `send(1 MB, 3 → 4)`
 - `send(1 MB, 4 → 5)`
 - each send can be queued per each link
 - link processes pending chunks in-order
- **Fast simulation** for large-scale systems **with network congestion**

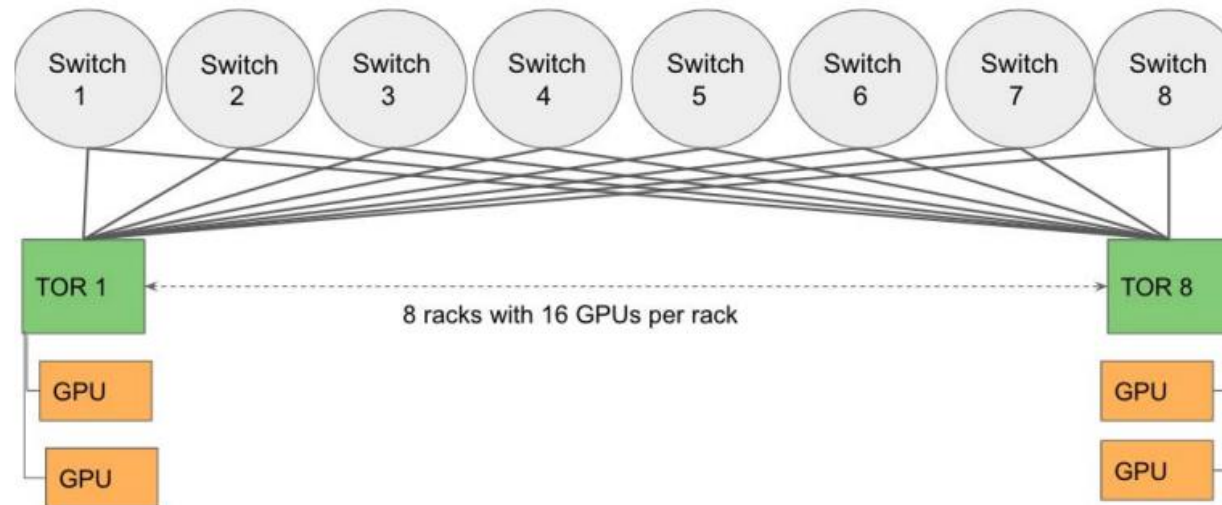
Garnet Backend

- Leverages **Garnet (interconnection network) simulator** as backend
- Appropriate for on-chip/scale-up networks
- Simulates interconnection network behaviors:
 - Message Packetization
 - Credit-based flow control
 - Congestion modeling
 - etc.
- Slower than analytical backend for large systems/models
- Supports switch-based/torus-based topologies

Caveat: Garnet currently only works with ASTRA-sim 1.0 and should be updated

ns-3 Backend

- Network simulator for **internet (inter-node) communication**
- Used to model ML training in **largely parallel GPU clusters**
- NPUs connected with ToR/spine switch, etc.



(HOTI '22) Current RoCE congestion control methods have little impact on ML training workloads

Slide courtesy: Jinsun Yoo <jinsun@gatech.edu>

ns-3 Network Configurations

- Detailed internetwork behavior modeling/simulation

PACKET_PAYLOAD_SIZE	packet size
CC_MODE	Congestion control algorithm
BUFFER_SIZE	switch buffer size
ACK_HIGH_PRIO	0: ACK has same priority with data packet 1: prioritize ACK
RATE_BOUND	Bound rate to a limited rate
ENABLE_QCN	Whether QCN (Quantized Congestion Notification) is enabled
L2_BACK_TO_ZERO	(Go-Back-N protocol) Layer 2 go back to zero transmission
L2_CHUNK_SIZE	(Go-Back-N protocol) Layer 2 chunk size
L2_ACK_INTERVAL	(Go-Back-N protocol) Layer 2 Ack intervals
HAS_WIN	Whether to use a window
GLOBAL_T	0: different server pairs use their own RTT as T 1: use the max base RTT as the global T
VAR_WIN	Whether the window size is variable
RATE_BOUND	Use rate limiter
ACK_HIGH_PRIO	Prioritize acknowledgement packets
KMAX_MAP	a map from link bandwidth to ECN threshold kmax
KMIN_MAP	a map from link bandwidth to ECN threshold kmin
PMAX_MAP	a map from link bandwidth to ECN threshold pmax
RATE_AI	Rate increment unit in AI period
RATE_HAI	Rate increment unit in hyperactive AI period
MIN_RATE	Minimum rate of a throttled flow