

Weitao Wang

wtwang@rice.edu • (+1) 281-236-7550 • weitaowang.site/about

Education

Rice University, Houston, Texas

Aug. 2018 – Apr. 2024

- Ph.D. Candidate in Computer Science
- **Research Interests:** Decentralized algorithm/protocol design, application-infrastructure co-design, data center network, network-on-chip, cluster scheduling, programmable hardware, systems for AI

Shanghai Jiao Tong University, Shanghai, China

Sept. 2014 – Jun. 2018

- B.S. in Engineering
- IEEE Honor Class
- Major GPA: 90.45 / 100 (Rank: 4 / 69)

Industry Experiences

High Precision Datacenter-scale Clock Synchronization | Google Cloud May 2022 – Present

- A new large-scale clock synchronization system designed to achieve nanosecond-level precision.
- Design a decentralized sync algorithm to reduce both systematic and random errors at scale.
- Leverage redundancy design to achieve agile and reliable clock synchronization services.
- Reduce the synchronization error from 100s of nanoseconds to single-digit nanoseconds at large scale.

Data Center Congestion Control Via Deployable INT | Google Cloud May 2021 – May 2022

- Poseidon: a new congestion control for next-generation network with the in-network telemetry (INT)
- Achieve ultra-low latency ($< 50 \mu s$), high utilization ($> 96\%$), and max-min fairness.
- Reduce the average job completion time by 42.4% and tail completion time by 99.1% for RPC workload.
- Validate on both [simulator](#) and testbed in the production environment, wide deployment expected in 2024.

Minimize the Precision Loss in WCMP Deployment | Google Cloud May 2020 – Aug. 2020

- Explore ILP approaches to minimize the WCMP precision loss within the data center network.
- Build a simulator with Integer Linear Programming solvers based on SCIP + Cpp.
- Reduce the precision loss up to 60% compared to the current strategy.
- [Project repository](#) with 6000+ lines of code (Cpp + Proto3 + CMake).

Selected Research Experiences

A Decentralized Host-based Weighted Bandwidth Allocation Algorithm Mar. 2023 – Present

- Provide a bandwidth allocation system that each flow could arbitrarily get proportionally higher bandwidth.
- Without a central controller or smart switches, the weight can be changed directly by flows on the end-host.
- Achieves an agile weight update for all flows and an accurate weighted bandwidth allocation.

A Decentralized Zero-queue Congestion Control with Max-min Fair May 2022 – Mar. 2023

- A congestion control algorithm that achieves zero-queuing despite traffic patterns and topology.
- Precisely monitor and maintain the bandwidth usage of links for $< 100\%$ with in-network telemetry.
- Achieve zero-queuing, fast convergence, max-min fair, and maintain network utilization at $\geq 90\%$.

A Decentralized Task Scheduler with Resource Sharing Knowledge Apr. 2021 – Mar. 2022

- Schedule the tasks in MXDAG (a cluster APP abstraction) precisely by understanding the resource sharing.
- Parse the code to obtain MXDAG and communicate with the cluster for precise resource allocation.
- Reduce the job completion time and minimize the resource usage for all the cloud service users.

A Default-Off Network Diagnostic System with Programmable Switches Mar. 2020 – May 2021

- Achieve network-wide monitoring with zero overhead and reactive diagnosis with low latency.
- Implement the prototypes on both **Barefoot Tofino** Switches and **NS3-Bmv2** simulator.
- Reduce the memory overhead by more than **99%** compared to a record-all monitoring solution.

A Reconfigurable Pod-Centric Data Center Network Architecture

Oct. 2018 – May 2020

- Optimize the data center network traffic with the dynamically reconfigurable network topology.
- Deploy a **Hadoop/MPI/Memcached** datacenter prototype with 16 servers and 5 **openflow** switches.
- Achieve an average path length **35%** shorter and improve the job completion times by **1.1-2.7x**.

Support High Throughput Low Delay Multicast with Optical Network Apr. 2019 – Mar. 2020

- Support multicast traffic with optimal bandwidth provisioning using a dedicated network.
- Implemented a **Hadoop/MPI** network prototype with 16 servers and 5 OpenFlow switches.
- Speed up raw broadcast **2.6x** and improve end-to-end application performance by up to **23%**.

Publications

- [\[NSDI'23\]](#) Poseidon: Efficient, Robust, and Practical Datacenter CC via Deployable INT
Weitao Wang, Masoud Moshref, Yuliang Li, Gautam Kumar, T. S. Eugene Ng, Neal Cardwell, Nandita Dukkipati
- [\[NSDI'22\]](#) RDC: Relieving Data Center Network Congestion with Topological Reconfigurability at the Edge
Weitao Wang, Dingming Wu, Sushovan Das, Afsaneh Rahbar, Ang Chen, T. S. Eugene Ng
- [\[NSDI'22\]](#) SpiderMon: Harnessing Wait-For Relations for Performance Debugging with Programmable Switches
Weitao Wang, Xinyu Crystal Wu, Praveen Tamma, Ang Chen, T. S. Eugene Ng
- [\[HotNets'21\]](#) MXDAG: A Hybrid Abstraction for Cluster Applications
Weitao Wang, Sushovan Das, Xinyu Crystal Wu, Zhuang Wang, Ang Chen, T. S. Eugene Ng
- [\[OptSys'21\]](#) Abstractions for Reconfigurable Hybrid Network Update and A Consistent Update Approach
Weitao Wang, Sushovan Das, T. S. Eugene Ng
- [\[SoSR'20\]](#) Grasp the Root Causes in the Data Plane: Diagnosing Latency Problems with SpiderMon
Weitao Wang, Praveen Tamma, Ang Chen, T. S. Eugene Ng
- [\[SIGCOMM'23\]](#) Augmented Queue: A Scalable In-Network Abstraction for Data Center Network Sharing
Xinyu Wu, Zhuang Wang, **Weitao Wang**, T. S. Eugene Ng
- [\[WORDS'23\]](#) Aurelia: CXL Fabric with Tentacle
Shu-Ting Wang, **Weitao Wang**
- [\[NeurIPS'23\]](#) Scissorhands: Exploiting the Persistence of Importance Hypothesis for LLM KV Cache Compression at Test Time
Zichang Liu, Aditya Desai, Fangshuo Liao, **Weitao Wang**, Victor Xie, Zhaozhuo Xu, Anastasios Kyrillidis, Anshumali Shrivastava
- [\[ToN'22\]](#) Shufflecast: An Optical, Data-rate Agnostic and Low-Power Multicast Architecture for Next-Generation Compute Clusters
Sushovan Das, Afsaneh Rahbar, Xinyu Wu, Zhuang Wang, **Weitao Wang**, Ang Chen, T. S. Eugene Ng
- [\[OptSys'21\]](#) Towards All-optical Circuit-switched Network Cores: Mitigating Traffic Skewness at the Edge
Sushovan Das, **Weitao Wang**, T. S. Eugene Ng
- [\[SoSR'19\]](#) Say No to Rack Boundaries: Towards A Reconfigurable Pod-Centric DCN Architecture
Dingming Wu, **Weitao Wang**, Ang Chen, T. S. Eugene Ng

Ongoing Submissions

- [\[Under Submission\]](#) Roptopia: A New Approach for Path-Aware Max-Min Fairness for Datacenter Networks
Weitao Wang, Liangcheng Yu, Vincent Liu, T. S. Eugene Ng
- [\[Under Submission\]](#) Zero: A New Congestion Control Algorithm with Zero Queuing
Weitao Wang, Xinyu Crystal Wu, Sushovan Das, Ang Chen, T. S. Eugene Ng
- [\[Draft\]](#) Soze: Zero-Coordination Weighted Bandwidth Allocation for Datacenter Traffic
Weitao Wang, Sushovan Das, Ang Chen, T. S. Eugene Ng

Patents

- [\[Patent Link\]](#) Congestion Control for Networks Using Deployable INT
Masoud Moshref Javadi, **Weitao Wang**, Yuliang Li, Gautam Kumar, Nandita Dukkupati, Neal Douglas Cardwell

Services

ACM SIGCOMM 2022 Artifact Evaluation Committee Jun. 2022

Invited Talks

Decentralized Weighted Bandwidth Allocation With Zero-coordination Nov. 2023

- Invited talk by [Duke Systems Group](#), Virtual event

Poseidon: Efficient, Robust, and Practical Datacenter CC via Deployable INT Aug. 2022

- Invited talk by [FlexNet](#), Houston, TX

Poseidon: Efficient, Robust, and Practical Datacenter CC via Deployable INT May 2022

- Invited talk by [Google Cloud](#), Sunnyvale, CA

MXDAG: A Hybrid Abstraction for Emerging Applications Feb. 2022

- Invited talk by [Google Networking Research Summit](#), Virtual event

An Important In-network Signal for Achieving Consistent Application Performance Jul. 2021

- Invited talk by [Intel](#), Virtual event

Teaching Experience

Mentor | Research Experiences for Undergrads Program May 2023 – Aug. 2023

- Advised 2 research projects, one with a sophomore student and one with a senior student.
- Developed a decentralized priority scheduling algorithm for modern data centers.
- Developed a defense system for the CC algorithm against spoofing attacks and man-in-the-middle attacks.

Teaching Assistant | Approximate Computing System For Big Data Aug. 2019 – Dec. 2019

- Taught bi-weekly labs for over 100 students.
- Co-designed the midterm and final exams with the course instructors.
- Led the team discussions for the final project and designed an APP for patient EEG monitoring.