

Sung-Han Lin

sunghlin@gmail.com | +1-213-810-2139 | <https://sunghlin.github.io/about>

EDUCATION

Ph.D. in Computer Science - *University of Southern California*, Los Angeles, California, U.S.A. 2010 – 2017

- Dissertation: **Distributed Resource Management for QoS-Aware Service Provision**

- Teaching Assistant, **CSci 402: Operating Systems**, Fall 2012 - Summer 2017

B.S. and M.S. in Computer Science (CSIE) - *National Taiwan University*, Taipei, Taiwan 2002 – 2008

SELECTED WORK EXPERIENCE

Performance Analysis Engineer & Data Scientist - *NetApp Inc.* October 2017 - Present

- **ONTAP AI: Improving GPU Data Pipeline**: [*Python, TensorFlow, Horovod, Container, MPI, Kubernetes, S3*]

- Designed and Deployed the ML/DL reference architecture combining NVIDIA DGX systems and NetApp storage systems for distributed TensorFlow training and inference over the RoCE connections.

- Identified the bottleneck of the end-to-end data pipelines; improved the GPU utilization, data feeding rate, and training speed via better configuring the CPU workload and multithreading, the GPU memory usage, and network interconnections, mainly on image recognition training and inference (*ResNet, VGG, Inception*).

- Adopted NVIDIA NGC and libraries to build our reference architecture; implemented and optimized ML/DL training models, and published joined Technical Report with NVIDIA on various verticals - autonomous vehicle (*Mask-R-CNN*), financial (*Semi-supervised GAN, Auto-Encoder*), and conversational AI (*BERT*).

- **Data Operation Failover**: Built the automation and investigated the executing path of failover, where the secondary cluster takes over the write operations from the primary cluster; helped reduce the I/O resumption latency via improving the logic of triggering failover and reducing the function overhead. [*C/C++, Python*]

- **Synchronous Replication**: Built the automation and investigated the latency of replicating write operations to the remote cluster synchronously; helped reduce the overhead of transmission via identifying the bottleneck and improving the concurrency of parallel operations. [*C/C++, Tcl, Python*]

- **Advanced Block Fetching**: Investigated the block prefetching strategies to reduce the read overhead and latency caused by the new checksum mechanism for disks; improved the read throughput by 20%. [*C/C++*]

Intern - *Teradata* June 2014 - August 2014

- **Improving SQL-MapReduce Execution Engine**: Redesigned and implemented the data buffering and transmission mechanisms to improve the query throughput by 15% via reducing the I/O latency and function overhead between user defined functions and databases. [*C/C++, Hadoop, SQL*]

SELECTED PUBLICATIONS

- **A Model-based Approach to Streamlining Distributed Training for Asynchronous SGD**, by *Sung-Han Lin*, Marco Paolieri, Cheng-Fu Chou, Leana Golubchik, IEEE MASCOTS, 2018

- **Performance Driven Resource Sharing Markets for the Small Cloud**, by *Sung-Han Lin*, Ranjan Pal, Marco Paolieri, Leana Golubchik, IEEE ICDCS, 2017

- **Sustaining Ad-Driven P2P Streaming Ecosystems A Market-Based Approach**, by *Sung-Han Lin*, Ranjan Pal, Bo-Chun Wang, Leana Golubchik, IEEE/ACM IWQoS, 2015

SELECTED RESEARCH PROJECTS

- **Throughput Maximization for Large-scale Deep Learning**: Analyzed the traffic pattern of the Asynchronous SGD training in the parameter-server architecture, and built a queueing network model to estimate the training speed (examples processed per second) of distributed TensorFlow GPU training jobs; leveraged the estimation model to address a problem of scheduling heterogeneous distributed DNN training jobs in a shared cluster.

TALK

- **Performance Considerations for AI and ML Deployments**, NetApp INSIGHT Las Vegas, October 2019