

Lin Ma

Carnegie Mellon University
Department of Computer Science
Gates-Hillman Center 9215
Pittsburgh, PA 15213-3891 USA

Voice: +1-412-519-7097
E-mail: lin.ma@cs.cmu.edu
Web: <http://www.cs.cmu.edu/~malin199>
GitHub: <https://github.com/malin1993ml>

EDUCATION

Ph.D., Computer Science
Carnegie Mellon University
Advisor: [Andy Pavlo](#)

2015-PRESENT
Pittsburgh, PA USA

M.Sc., Computer Science
Carnegie Mellon University
Advisor: [Andy Pavlo](#)

2015-2018
Pittsburgh, PA USA

B.Sc., Computer Science
Peking University
Advisor: [Bin Cui](#)

2011-2015
Beijing, China

RESEARCH EXPERIENCE

Graduate Research Assistant
Carnegie Mellon University

2015-PRESENT

- **NoisePage** - <https://noise.page/>

Working with students at CMU to develop a relational database management system that is designed for autonomous operation. My research is focused on the following areas:

Workload Forecasting: A framework that allows the DBMS to identify the trends and patterns of the queries in the workload and optimize the system based on the predicted workload in the future. It applies an ensemble method with linear regression models and recurrent neural networks trained with PyTorch and TensorFlow to improve the forecasting accuracy. It uses query templatzation and an on-line clustering technique to reduce the storage and computation overhead.

Self-Driving DBMS: A system architecture that leverages advancements in deep neural networks, improved hardware, and adaptive system designs to remove the human capital impediments of deploying DBMSs. We are working on an architecture that autonomously characterizes and predicts the workload of the DBMS, plans optimizations to improve its performance with receding-horizon control models, and applies reinforcement learning techniques to learn the feedback from these optimizations.

Peloton Query Optimizer: This is an objected-oriented Cascade-style framework written in C++ for query optimization. I implemented the original framework and co-advised two master students to add support for additional features. This optimizer is later imported into NoisePage.

- **H-Store** - <http://hstore.cs.brown.edu>

Modern Storage Hardware for In-Memory DBMS: This is an evaluation of the design decisions on managing cold data in H-Store. We analyzed the issues of these decisions against different storage devices, including NVMs, SSDs, HDDs, and SMRs. We also proposed combinations of techniques that are optimized for the target hardware, including tuple retrieval strategies, merging threshold estimates using a Count-Min Sketch, and data access methods.

Research Intern
Data Management, Exploration and Mining Group
Microsoft Research, Redmond

SUMMER 2018

- **Auto-Indexing in Cloud** - <https://www.microsoft.com/en-us/research/project/autoadmin>

Active Learning for ML Enhanced Database Systems: This is a strategy that leverages B-Instances in the cloud to actively collect more training data for machine learning models that enhance the DBMS performance (e.g., predicting the runtime of query plans). We explored active learning based methods to improve the accuracy of the models with limited budget efficiently.

Undergraduate Research Intern

2013-2015

Peking University

- **Graph Computation**

Social Network Analysis System: A C++ prototype built from scratch that answers four types of social queries, such as finding large interest communities or the most central people.

PSgL: A parallel subgraph listing framework that iteratively divides the problem into partial tasks and balances the loads between concurrent workers. The prototype is written in Java on Apache Giraph.

Page: A partition aware graph computation engine that uses the online statistics of graph partitioning results to guide the resource allocation for parallel processing. The prototype is written in Java on Apache Giraph.

PUBLICATIONS

- [1] **Lin Ma**, Bailu Ding, Sudipto Das, and Adith Swaminathan. Active learning for ml enhanced database systems. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, pages 175–191, 2020.
- [2] Andrew Pavlo, Matthew Butrovich, Ananya Joshi, **Lin Ma**, Prashanth Menon, Dana Van Aken, Lisa Lee, and Ruslan Salakhutdinov. External vs. internal: an essay on machine learning agents for autonomous database management systems. *IEEE bulletin*, 42(2), 2019.
- [3] **Lin Ma**, Dana Van Aken, Ahmed Hefny, Gustavo Mezerhane, Andrew Pavlo, and Geoffrey J. Gordon. Query-based workload forecasting for self-driving database management systems. In *Proceedings of the 2018 ACM International Conference on Management of Data*, SIGMOD '18, 2018.
- [4] Andrew Pavlo, Gustavo Angulo, Joy Arulraj, Haibin Lin, Jiexi Lin, **Lin Ma**, Prashanth Menon, Todd C Mowry, Matthew Perron, Ian Quah, et al. Self-driving database management systems. In *CIDR*, 2017.
- [5] **Lin Ma**, Joy Arulraj, Sam Zhao, Andrew Pavlo, Subramanya R Dulloor, Michael J Giardino, Jeff Parkhurst, Jason L Gardner, Kshitij Doshi, and Stanley Zdonik. Larger-than-memory data management on modern storage hardware for in-memory oltp database systems. In *Proceedings of the 12th International Workshop on Data Management on New Hardware*, page 9. ACM, 2016.
- [6] Huanchen Zhang, David G Andersen, Andrew Pavlo, Michael Kaminsky, **Lin Ma**, and Rui Shen. Reducing the storage overhead of main-memory oltp databases with hybrid indexes. In *Proceedings of the 2016 International Conference on Management of Data*, pages 1567–1581. ACM, 2016.
- [7] Yingxia Shao, Bin Cui, and **Lin Ma**. Page: a partition aware engine for parallel graph computation. *IEEE Transactions on Knowledge and Data Engineering*, 27(2):518–530, 2015.
- [8] Yingxia Shao, Bin Cui, Lei Chen, **Lin Ma**, Junjie Yao, and Ning Xu. Parallel subgraph listing in a large-scale graph. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*, pages 625–636. ACM, 2014.
- [9] Yingxia Shao, Junjie Yao, Bin Cui, and **Lin Ma**. Page: A partition aware graph computation engine. In *Proceedings of the 22nd ACM International Conference on Information and Knowledge Management*, pages 823–828. ACM, 2013.

ACADEMIC TALKS

- **Active Learning for ML Enhanced Database Systems**
SIGMOD, June 2020.
- **Self-Driving Databases: It All Starts with Workload Forecasting**
Percona Live, May 2019.
- **Efficiently Leveraging B-Instances for Query Plan Predictions**
Microsoft Research, August 2018.
- **Query-based Workload Forecasting for Self-Driving DBMSs**
SIGMOD, June 2018.
- **Query-based Workload Forecasting for Self-Driving DBMSs**
Microsoft Research, May 2018.
- **Query-based Workload Forecasting for Self-Driving DBMSs**
PDL Retreat, October 2017.
- **Larger-than-Memory Data Management on Modern Storage Hardware for In-Memory OLTP Database Systems**
SIGMOD, June 2016.
- **The Self-Driving DBMS**
PDL Retreat, October 2016.
- **Multi-Level Anti-Caching for NVM+SSD in H-Store**
PDL Retreat, October 2015.
- **Finalist Presentation of Programming Contest**
SIGMOD, June 2014.
- **Using Less to Do More With Anti-Caching in OLTP Database Systems**
Carnegie Mellon University, August 2014.

AWARDS AND SCHOLARSHIPS

- The China Computer Federation Outstanding Undergraduate Award - 2015
- China National Scholarship - 2014
- SIGMOD Programming Contest Finalist - 2014
- SIGMOD Travel Award - 2014

SERVICE

- Program Committee - AIDB 2020
- External Reviewer - DAPD 2019, SIGMOD Demo 2017
- CSD Faculty Search Committee – Carnegie Mellon University, 2020
- CSD MS Admissions Committee – Carnegie Mellon University, 2018
- Graduate Student Recruitment (Open House) Committee – Carnegie Mellon University, 2018