

# Yibai Meng

Mountain View, CA

✉ [yibai@meng.engineer](mailto:yibai@meng.engineer)

🌐 [linkedin.com/in/yibai-meng](https://www.linkedin.com/in/yibai-meng)

☎ +1 6692009099

## SKILLS

---

**Programming Languages:** Python, C/C++, CUDA, Go, JavaScript/HTML

**Machine Learning:** PyTorch, Tensorflow, XLA, MLIR, Triton, quantization

**Development tools:** Git, Vim, VSCode, Bazel, CMake, Linux utilities

## INDUSTRY EXPERIENCES

---

### Waymo

Jan 2023 - Present

*Software Engineer*

*Mountain View, California*

- Working in ML Runtime team under ML Infrastructure.
- **Model optimization:** Improved the latency and stability of on-vehicle ML models using techniques such as quantization and operation fusion. Modernized graph manipulation workflow with technologies like MLIR.
- **Accelerator enablement:** Integrated GPU kernels into on-vehicle models. Tested, benchmarked and modified multi-head attention kernels, using both Cutlass and Triton libraries.
- **Foundation model support:** Unblocked foundation model adoption by overcoming longstanding limitations due to tensorflow implementation details. Made numerous contributions to tensorflow codebase.
- **Developer tooling:** Created tool for inspecting models and providing optimization suggestions. Made tool to verify numerics after graph manipulation, leveraging existing integration testing infrastructure.

### TikTok

May 2022 - Aug 2022

*Software Engineer Intern*

*Mountain View, California*

- As a member of the Software Defined Network team, implemented a novel data plane network verification algorithm in C++ from scratch, based on a recent academic paper. Network verification models network topology and packet forwarding, and checks for the existence of unwanted behavior, such as loops and blackholes.
- Drastically increased the performance of network topology modelling and invariant checks (loop) by 500+ times, compared with existing Python tools, enabling real-time verification of TikTok's expansive global network.
- Implemented incremental rule checking, allowing network configuration to be updated on the fly.
- Made preparations to deploy the verification tool as a internal service, to be used by TikTok's network engineers.

## EDUCATION

---

### University of California, Berkeley

Aug 2021 - Dec 2022

*Master of Engineering in Electrical Engineering and Computer Science*

*Berkeley, California*

### Peking University

Sep 2016 - May 2020

*Bachelor of Science in Electronics and Information Science and Technology*

*Beijing, China*

## ACADEMIC EXPERIENCES

---

### Center for Energy-Efficient Computing and Applications, Peking University

July 2020 – June 2021

*Research Assistant*

*Beijing, China*

- GPU acceleration of elfPlace using CUDA, an algorithm for the placement phase of FPGA physical synthesis. Shortened runtime by 7 times on average.
- Ported the algorithm to an PyTorch based framework, framing this nonlinear nonconvex optimization problem as training a neural network. Used C++ & CUDA extension on Python to speed up critical segments while maintaining low code complexity. Refactored the individual functionalities into “operators”, following the paradigm of high-cohesion low-coupling.
- Extended the algorithm to consider clock network routing resource constraints with a quadratic penalty, with consideration of both global placement convergence and design legality.
- Extended the algorithm to include new cell types, allowing it to process a non-academic real-life industry architecture and substantially improved its performance.
- Resulted in two academic publications in top journals.

## PUBLICATIONS

---

- **elfPlace: Electrostatics-based Placement for Large-Scale Heterogeneous FPGAs:** Yibai Meng, Wuxi Li, Yibo Lin and David Z. Pan. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 2021
- **Multi-Electrostatic FPGA Placement Considering SLICEL-SLICEM Heterogeneity and Clock Feasibility:** Jing Mai, Yibai Meng, Zhixiong Di and Yibo Lin. *Design Automation Conference (DAC)*, 2022