

Shuangqi Li

📍 Lausanne, Switzerland ✉ shuangqi.li@epfl.ch 🌐 lishuangqi.com

ABOUT

PhD candidate at EPFL, advised by [Mathieu Salzmann](#) [🔗](#). My research investigates **how training data shapes model behavior**—through data attribution, curation, synthesis—to understand and improve AI training.

Strong engineering foundation: Python, PyTorch, CUDA, with a competitive programming background.

EDUCATION

EPFL (Swiss Federal Institute of Technology in Lausanne) Lausanne, Switzerland
PhD in Machine Learning Sept 2022 – 2027

EPFL (Swiss Federal Institute of Technology in Lausanne) Lausanne, Switzerland
Master in Data Science Sept 2020 – July 2022

University of California, San Diego Remote
Master in Computer Science (Quit due to visa issues and COVID-19.) Sept 2019 – June 2020

University of Electronic Science and Technology of China Chengdu, China
Bachelor in Microelectronic Science and Engineering Sept 2015 – June 2019

WORK EXPERIENCE

Research Intern Zurich, Switzerland
Oracle Labs July 2021 – Sept 2021

- Developed a time series model that detects anomalous Linux sessions in the cloud servers.

Algorithm Engineering Intern Beijing, China
DiDi (China's largest taxi-hailing platform) Oct 2018 – Feb 2019

- Developed an algorithm for learning road segment weights from historical ride data, significantly improving route planning quality for ride-hailing services in production environment.

PROJECTS

Dense Credit Assignment for RL via Token-Level Data Attribution Jan 2026 – present
Ongoing research

- Proposed a novel data attribution framework for reinforcement learning to estimate the marginal contribution of individual tokens to total rewards for GRPO/DAPO-style algorithms.
- Achieved fine-grained dense credit assignment, effectively mitigating the reward sparsity limitations in reasoning and agentic RL training.

Scalable Training Data Attribution for Large Language Models Aug 2025 – Jan 2026
Under review for ICML 2026. [arXiv](#) [🔗](#)

- Developed a novel, highly scalable method for training data attribution in large-scale models by exploiting the low-rank properties of gradients, cutting storage cost and query latency 20×.
- Enabled, for the first time, the ability to efficiently trace the output of a 70-billion-parameter LLM back to individual examples in their SFT training data.

Learning to Weight Parameters for Training Data Attribution Feb 2025 – Sept 2025
ICLR 2026. [arXiv](#) [🔗](#)

- Identified the heterogeneity of attribution signal across parameters/layers in diffusion models and LLMs.
- Proposed a method to re-weight layers, boosting attribution accuracy and enabling interpretable attribution.

LLM Development from Scratch

Mar 2025 – Jul 2025

Collaborative project with 25 PhD students to build a large language model from scratch.

- Engineered the pre-training pipeline, including environment setup and investigating optimal data mixing recipes for the training corpus.
- Implemented and validated the evaluation suite by reproducing the SmolLM2 benchmark to establish a robust performance baseline.

Enhancing Text-to-Image Generation with Reliable Random Seeds

Jan 2024 – Oct 2024

ICLR 2025 Spotlight. [arXiv](#) [🔗](#)

- Identified the significant role of initial noise in text-to-image inconsistencies for diffusion models.
- Proposed a method that identifies reliable random seeds to improve text-to-image generation. Leveraging reliable seeds to synthesize high-quality data for fine-tuning diffusion models.

Controlling the Fidelity and Diversity of Deep Generative Models

Feb 2023 – Mar 2024

TMLR 2024 (poster presentation at ICLR 2025). [arXiv](#) [🔗](#)

- Proposed an approach to bias generative models towards generating data with either enhanced fidelity or increased diversity.
- Enabled model training with data of better fidelity or diversity.

Interlock-Free Multi-Aspect Rationalization for Text Classification

Sept 2021 – Feb 2022

Semester project. [arXiv](#) [🔗](#)

- Proposed a multi-stage training method to alleviate the interlocking issue in training interpretable models.

TEACHING & SUPERVISION

- **Improving Waste Detection and Sorting @ WasteFlow – Project Supervisor** Fall 2025
- **CS-233: Introduction to Machine Learning – Head Teaching Assistant** Spring 2025
- **CS-401: Applied Data Analysis – Teaching Assistant** Fall 2024
- **CS-233: Introduction to Machine Learning – Teaching Assistant** Spring 2024
- **COM-407: TCP/IP Networking – Teaching Assistant** Fall 2023
- **COM-112: Object-Oriented Programming (in C++) – Teaching Assistant** Spring 2023
- **CS-456: Deep Reinforcement Learning – Student Teaching Assistant** Spring 2022

HONORS & AWARDS

- **National Scholarship** Sept 2018
- **China Collegiate Programming Contest – GOLD MEDAL** May 2018
- **China Collegiate Computing Contest – FIRST PRIZE** Mar 2018
- **First-class People's Scholarship** Dec 2017
- **ACM ICPC (Asia Regional) – BRONZE MEDAL** Oct 2017
- **China Collegiate Computing Contest – FIRST PRIZE** Apr 2017
- **First-class People's Scholarship** Dec 2016

SKILLS

Programming: Python, C++, CUDA, Coding competition

Frameworks & Tools: PyTorch, Docker, Git, Linux, PySpark, Cursor, Claude Code

Languages: Chinese (native), English (fluent), French (basic)

PUBLICATIONS

Low-Rank Influence Functions for Scalable Training Data Attribution

Shuangqi Li, Hieu Le, Jingyi Xu, and Mathieu Salzmann

[arXiv:2601.21929](#) 

Submitted to ICML 2026

Learning to Weight Parameters for Training Data Attribution

Shuangqi Li, Hieu Le, Jingyi Xu, and Mathieu Salzmann

[ICLR 2026](#) 

Enhancing Compositional Text-to-Image Generation with Reliable Random Seeds

Shuangqi Li, Hieu Le, Jingyi Xu, and Mathieu Salzmann

[ICLR 2025](#) 

Spotlight (top 4%)

Controlling the Fidelity and Diversity of Deep Generative Models via Pseudo Density

Shuangqi Li, Chen Liu, Tong Zhang, Hieu Le, Sabine Süssstrunk, and Mathieu Salzmann

[TMLR 2024](#) 

Presented at ICLR 2025

Interlock-Free Multi-Aspect Rationalization for Text Classification

Shuangqi Li, Diego Antognini, and Boi Faltings

[arXiv:2205.06756](#) 