

Tzu-Heng (Brian) Huang

 thuang273@wisc.edu



zihengh1.github.io



zihengh1



zihengh1

Education

2021 – Jul. 2026 (Expected)

■ **Ph.D. in Computer Science. University of Wisconsin-Madison.**

Advised by Frederic Sala.

2016 – 2020

■ **B.S. in Computer Science. National Chengchi University.**

Advised by Man-Kwan Shan and Ling-Jyh Chen. Major GPA: 3.96/4.00.

Research Summary

My research focuses on ***data-centric AI for multimodal models***, with the goal of enabling systems to learn more from less but higher-quality supervision. Several ***data lifecycle*** projects I have worked on, including **(i)** visual RL with *sample-specific, synthetic data verifiers*, **(ii)** *online data mixing* for overhead-reduced multimodal training, **(iii)** *fine-grained data selection* for efficient pretraining (**ICML'25 DataWorld Oral**), **(iv)** *data curation* via ensemble and objective detection (**1st place on the DataComp'23 leaderboard**), and **(v)** a 500x cheaper *automated data labeling* system as an alternative to LLM annotators (**NeurIPS'24 Spotlight**).

Research Experience

Oct. 2025 – Jan. 2026

■ **AIML Research Intern. Apple Inc.**

advised by Manjot Bilkhu and Javier Movellan.

— *Rubric-based RL for Dense Captioning Tasks*.

May. 2025 – Sep. 2025

■ **Research Scientist Intern. Meta GenAI (now MSL).**

advised by Ernie Chang, Sang Michael Xie, Yiting Lu, and David Kant.

— *Learnability-aware Synthetic Data Generation*.

May. 2024 – Dec. 2024

■ **AIML Research Intern. Apple Inc.**

advised by Javier Movellan and Manjot Bilkhu.

— *Automated Model-aware Data Selection for Efficient Pretraining*.

— *Optimizing Domain Mixtures for MLLM Pretraining*.

Aug. 2021 – Present

■ **Graduate Research Student. UW-Madison.**

advised by Frederic Sala.

— *Data-centric AI for Foundation Models: Auto-labeling and Data Curation*.

— *Parameter Marketplace: Through Model Merging and Auction Agents*.

May. 2023 – May. 2024

■ **Founder. Awan.AI LLC (integrated with TechTCM).**

— *Tongue Syndrome Diagnosis and LLM for Traditional Chinese Medicine*.

— *Automating TCM Diagnosis: Herbal-based Recommendation System*.

Jun. 2019 – Sep. 2019

■ **Research Intern. Argonne National Laboratory.**

advised by Charlie Catlett and Rajesh Sankaran.

— *Ensemble-based Time Series Calibration for Low-cost Sensors*.

Sep. 2018 – Aug. 2021

■ **Research Assistant. National Chengchi University.**

advised by Man-Kwan Shan.

— *Spatio-temporal Modeling in Large-scale Sensor Networks*.

Feb. 2018 – Jul. 2020

■ **Research Intern. Academia Sinica.**

advised by Ling-Jyh Chen.

— *Large-scale Air Quality Sensor Networks (AirBox)*.

Publications and Preprints

- 1 T.-H. Huang, S. Salekin, J. Movellan, F. Sala, and M. Bilkh, "RubiCap: Synthesized Rubrics as Rewards to Guide RL for Dense Captioning," in *submission*, 2026.
- 2 J. Zhao, C. Shin, T.-H. Huang, S. S. S. Namburi, and F. Sala, "CARE: Confounder-Aware Aggregation for Reliable LLM Evaluation," in *submission*, 2026.
- 3 J. Saad-Falcon, E. K. Buchanan, M. F. Chen, T.-H. Huang, B. McLaughlin, T. Bhathal, S. Zhu, B. Athiwaratkun, F. Sala, S. Linderman, A. Mirhoseini, and C. Re, "Shrinking the Generation-Verification Gap by Scaling Compute for Verification," in *Neural Information Processing Systems (NeurIPS), ICML Workshop: Efficient Systems for Foundation Models (ES-FoMo III), and ICML Workshop: Multi-Agent Systems in the Era of Foundation Models: Opportunities, Challenges and Futures (MAS)*, 2025. ⚡ URL: <https://arxiv.org/abs/2506.18203>.
- 4 T.-H. Huang, H. Vishwakarma, and F. Sala, "Time to Impeach LLM-as-a-Judge: Programs are the Future of Evaluation," in *ICML Workshop: Programmatic Representations for Agent Learning (PRAL)*, 2025. ⚡ URL: <https://arxiv.org/abs/2506.10403>.
- 5 J. Zhao, C. Shin, T.-H. Huang, S. S. S. Namburi, and F. Sala, "From Many Voices to One: A Statistically Principled Aggregation of LLM Judges," in *NeurIPS Workshop: Evaluating the Evolving LLM Lifecycle: Benchmarks, Emergent Abilities, and Scaling*, and *NeurIPS Workshop: Reliable ML from Unreliable Data*, 2025. ⚡ URL: <https://openreview.net/pdf?id=0u53DNvjx7>.
- 6 A. Ge, T.-H. Huang, J. Cooper, A. Trost, Z. Chu, S. S. S. Namburi, Z. Cai, K. Park, N. Roberts, and F. Sala, "R&B: Domain Regrouping and Data Mixture Balancing for Efficient Foundation Model Training," in *ICML Workshop: Unifying Data Curation Frameworks Across Domains (DataWorld)*, and *ICML Workshop: Data in Generative Models (The Bad, the Ugly, and the Greats) (DIG-BUGS)*, 2025. ⚡ URL: <https://arxiv.org/abs/2505.00358>.
- 7 T.-H. Huang, M. Bilkh, J. Cooper, F. Sala, and J. Movellan, "Evaluating Sample Utility for Efficient Data Selection by Mimicking Model Weights," in *ICML Workshop: Unifying Data Curation Frameworks Across Domains (DataWorld) [Oral Paper]*, 2025. ⚡ URL: <https://arxiv.org/abs/2501.06708>.
- 8 T.-H. Huang, C. Cao, V. Bhargava, and F. Sala, "The ALCHEmist: Automated Labeling 500x CHEaper than LLM Data Annotators," in *Neural Information Processing Systems (NeurIPS) [Spotlight Paper (Top 2.08%)]*, 2024. ⚡ URL: <https://arxiv.org/abs/2407.11004>.
- 9 W. Tan, N. Roberts, T.-H. Huang, J. Zhao, J. Cooper, S. Guo, C. Duan, and F. Sala, "MoRe Fine-Tuning with 10x Fewer Parameters," in *ICML Workshop: Efficient Systems for Foundation Models (ES-FoMo)*, and *ICML Workshop: Foundation Models in the Wild.*, 2024. ⚡ URL: <https://arxiv.org/abs/2408.17383>.
- 10 N. Roberts, X. Li, D. Adila, S. Crompt, T.-H. Huang, J. Zhao, and F. Sala, "Geometry-Aware Adaptation for Pretrained Models," in *Neural Information Processing Systems (NeurIPS)*, 2023. ⚡ URL: <https://arxiv.org/abs/2307.12226>.
- 11 T.-H. Huang, H. Vishwakarma, and F. Sala, "Train 'n Trade: Foundations of Parameter Markets," in *Neural Information Processing Systems (NeurIPS)*, 2023. ⚡ URL: <https://arxiv.org/abs/2312.04740>.
- 12 T.-H. Huang, C. Shin, S. J. Tay, D. Adila, and F. Sala, "Multimodal Data Curation via Object Detection and Filter Ensembles," in *ICCV Workshop: Towards the Next Generation of Computer Vision Datasets (TNGCV) [1st place on the Datacomp leaderboard (small-scale filtering track)]*, 2023. ⚡ URL: <https://arxiv.org/abs/2401.12225>.
- 13 T.-H. Huang, C. Cao, S. Schoenberg, H. Vishwakarma, N. Roberts, and F. Sala, "ScriptoriumWS: A Code Generation Assistant for Weak Supervision," in *ICLR Workshop: Deep Learning For Code (DL4C)*, 2023. ⚡ URL: <https://arxiv.org/abs/2502.12366>.
- 14 N. Roberts, X. Li, T.-H. Huang, D. Adila, S. Schoenberg, C.-Y. Liu, L. Pick, H. Ma, A. Albarghouthi, and F. Sala, "AutoWS-Bench-101: Benchmarking Automated Weak Supervision with 100 Labels," in *Neural Information Processing Systems (NeurIPS)*, 2022. ⚡ URL: <https://arxiv.org/abs/2208.14362>.

15

- T.-H. Huang**, C.-H. Tsai, and M.-K. Shan, “Key Sensor Discovery for Quality Audit of Air Sensor Networks,” in *ACM International Conference on Mobile Systems, Applications, and Services (MobiSys)*, 2020.
- 🔗 URL: <https://dl.acm.org/doi/abs/10.1145/3386901.3396606>.

Miscellaneous

Awards

- 2025 ┣ **Oral Paper: Grad-Mimic (Over 100+)**, selected by ICML'25 DataWorld Workshop.
- 2024 ┣ **Spotlight Paper (Top 2.08%): The Alchemist**, selected by NeurIPS'24.
- 2023 ┣ **ICCV Datacomp Competition**, won the first place in the small-scale filtering track.
- ┣ **Scholar Award**, granted by NeurIPS'23.
- 2021 ┣ **First-year Departmental Scholarship**, granted by UW-Madison.
- 2020 ┣ **Research Intern Scholarship**, granted by National Chengchi University.
- ┣ **Undergrad Research Scholarship**, granted by Ministry of Science and Technology.

Invited Talks

- Jan. 2026 ┣ **Data Recipes: Labeling, Selection, and Verification**, invited by National Tsing Hua University (NTHU).
- ┣ **Data Recipes: Labeling, Selection, and Verification**, invited by National Taiwan University (NTU).
- ┣ **Data Recipes: Labeling, Selection, and Verification**, invited by Taiwan Semiconductor Manufacturing Company (TSMC).
- ┣ **Data Recipes: Labeling, Selection, and Verification**, invited by National Yang Ming Chiao Tung University (NYCU).
- Aug. 2025 ┣ **Data Recipes: Automated Labeling and Efficient Selection**, invited by Scaling Intelligence Lab (Azalia Mirhoseini's group) in Stanford University.
- Apr. 2025 ┣ **Spatio-temporal Modeling for Underwater Sensor Networks**, invited by National Taipei University of Technology (NTUT).
- Dec. 2019 ┣ **Air Quality Sensor Network Developments**, invited by Nangang High School (Taiwan).
- Sep. 2019 ┣ **Internship Research Talk**, invited by National Chengchi University.
- Jul. 2019 ┣ **LASS Conference: International Session**, invited by Academia Sinica.
- Mar. 2019 ┣ **Techbang Magazine: PiM25 Project**, invited by Techbang Magazine.
- ┣ **Raspberry Pi Jam: PiM25 Project**, invited by Raspberry Pi Foundation (Taiwan).
- Jan. 2019 ┣ **Raspberry Pi Meetup: PiM25 Project**, invited by Raspberry Pi Foundation (Taiwan).

Academic Services

- 2026 – Present ┣ **Organizer**, Data Foundations of AI (data-focus online seminars).
- 2021 – Present ┣ **Active Paper Reviewer**, NeurIPS, ICLR, and ICML.
- ┣ **Co-organizer**, AutoML Cup in AutoML Conference.
- 2022 – 2023 ┣ **President of Student Association of Taiwan**, UW-Madison.
- 2021 – 2022 ┣ **Vice President of Student Association of Taiwan**, UW-Madison.

Skills

- | | |
|-----------------------|--|
| Programming Languages | ┣ Python, R, C++/C, SQL, L ^A T _E X, and Shell Programming. |
| Technologies | ┣ (Distributed) PyTorch, Tensorflow, Keras, PostgreSQL, and Vim. |