

YUFAN ZHANG

New York, NY | 914-720-6892 | yz2894@cornell.edu | [GitHub](#) | [LinkedIn](#) | [Website](#)

EDUCATION

- Cornell Tech (Cornell University)**, New York, NY Aug 2023 – May 2025
M.S. in Information Systems | GPA: 4.0/4.0 | Merit Scholarship Recipient
- Duke University / Duke Kunshan University**, Kunshan, China Aug 2019 – May 2023
B.S. in Data Science | GPA: 3.7/4.0 | Dean's Lists

TECHNICAL SKILLS

- Coding Languages:** Python, SQL, HTML/CSS/JavaScript
Machine Learning & AI: PyTorch, PyTorch Lightning, TensorFlow, Keras, Scikit-Learn, HuggingFace, OpenCV, NLTK
Data Handling & Analysis: Pandas, NumPy, Spark, Matplotlib, Plotly, Tableau, PowerBI
Miscellaneous Skills: Git, Jupyter Notebooks, AWS, Azure, GCP, Linux, Wandb, Firebase, Jira

EXPERIENCE

- insitro | Data Science and Machine Learning Intern**, South San Francisco, CA May 2024 – Aug 2024
- Onboarded a ViT-based cell segmentation model, improving the mIoU metric from **0.79** to **0.86** on the internal evaluation dataset by adopting an **DETR** object detector to predict bounding boxes, which are used for prompting a **Segment Anything** Model for mask generation.
 - Migrated the previous deployed cell segmentation model, Cellpose, to **PyTorch Lightning**, enhancing scalability of the model implementation.
 - Monitored and organized over **100+** training jobs using Weights & Biases (Wandb), ensuring efficient tracking and reproducibility of experiments.
 - Developed helper functions to provide a unified interface for data access on **AWS S3** and **Azure Blob Storage**, facilitating the company's transition from AWS to Azure for deep learning model training.
- eBay | Product Manager Intern, Cloud Data & Storage Team**, Shanghai, China Mar 2023 – Jun 2023
- Leveraged the knowledge of cloud computing, including **Apache Kafka** and **Flink**, to achieve feature extensions and performance optimization of eBay's internal data streaming platform, resulting in **10** new use case onboardings and a **13%** user satisfaction score improvement.
 - Conducted data analysis and visualization with **PowerBI** on user feedback surveys, deriving actionable insights for the engineering team, which resulted in the platform's integration of **4** new data storage connectors, significantly enhancing the platform's data accessibility.
- Duke Kunshan University | Research Assistant**, Kunshan, China Jul 2022 – Nov 2022
- Leveraged **Python**, **NetworkX** and **Pandas** to build **2000+** daily transaction networks, extracting key graph features and detecting graph core-periphery with cpnet library, which unveiled the decline in inclusiveness in **4** major DeFi banks.
 - Developed and optimized **SQL** scripts to extract **over 2 million** Ethereum transaction records on **Google BigQuery**.
 - Visualized and presented key metrics' temporal trends with **Plotly** to Liquity's founder, influencing their strategy of DeFi protocol designs.
 - First-authored** a research paper for this cross-sectional comparative study framework on DeFi banks; received **16+ citations**. [\[Paper\]](#) [\[GitHub\]](#)

PROJECTS

- Data-drive Restaurant Recommendation System for Yelp**, (Python, Spark) [\[GitHub\]](#) Spring 2024
- Engineered restaurant recommendation systems with **PySpark**, achieving an RMSE of **1.081** using a hybrid recommendation approach, which is a **13.5%** improvement over the content-based filtering approach and a **51.7%** improvement over the ALS-based collaborative filtering method.
 - Conducted data cleaning and **feature engineering** on 2 million users and 150,000 businesses data with **Python** and **Pandas**, including scaling numerical attributes, one-hot encoding for categorical attributes, and ordinal encoding with custom scoring functions.
 - Experimented with various types of hybrid recommendation systems, identifying ridge regression with feature combination as the optimal model with an R^2 of **0.396** over ensemble methods and other predictive machine learning algorithms (e.g., linear regression, random forest regression).
- Text-to-SQL Translation with Modern Natural Language Processing (NLP) Practices**, (Python, PyTorch) [\[GitHub\]](#) Spring 2024
- Experimented with various NLP techniques for translating natural language instructions into SQL queries using **PyTorch**, including prompt engineering with a large language model (**LLM**), training language models, and fine-tuning pre-trained language models.
 - Fine-tuned a pre-trained T5 language model using **Hugging Face Transformers**, achieving the best F1 score of **0.627** over other approaches.
 - Implemented various prompt engineering techniques for LLMs, experimenting with **zero-shot** and **few-shot prompting**, demonstrating that the best prompt design improved F1 score by **37.3%** compared to the baseline prompting method.
- miniTorch: Python Re-implementation of the PyTorch API**, (Python) [\[GitHub\]](#) Fall 2023
- Engineered a **Python**-based alternative library to the **Torch** API, resulting in **100%** compatibility with native **PyTorch** code.
 - Architected a custom Tensor data structure pivotal for deep learning model training and evaluation, supporting tensor backend operations including **broadcasting**, mathematical operation **overloads**, **auto-differentiation**, and **backpropagation**.
 - Achieved a **10x** speedup in training by implementing parallel computations with **Numba JIT** for essential tensor operations (map, zip & reduce).
- GAN-based Font Image Style Transfer**, (Python, PyTorch) [\[Paper\]](#) [\[GitHub\]](#) Summer 2022
- Designed an end-to-end **GAN**-based image generation model with **PyTorch**, improving SSIM by **17.3%** over the previous SOTA models
 - Adopted the **self-attention** mechanism in the image style encoder to capture both the local and global font style and an adaptive **skip connection** mechanism to improve content fidelity, evidenced by a **14%** and **12%** improvement in SSIM respectively through **ablation studies**.
 - First-authored** a research paper at the top multimedia computing conference, **ACM Multimedia 2022**; received **5+ citations**.