

## EDUCATION

---

- |   |                     |
|---|---------------------|
| Stanford University   | Sep 2018 - Jun 2022 |
| <ul style="list-style-type: none"> <li>◦ B.S., Computer Science (AI Specialization) • GPA: 4.18 / 4.0 • Graduated with Distinction (Summa Cum Laude)</li> <li>◦ <a href="#">Endorsed</a> by Fields Medalist Prof. <a href="#">Ngo Bao Chau</a> for exceptional mathematical talent and research potential.</li> </ul> |                     |

## AWARDS

---

- |  |   |
|--|---|
| International Mathematical Olympiad: <b>Bronze Medal</b> | American Mathematics Competition: Gold Medal    |
| Vietnam Math Olympiad: <b>Gold Medal (Rank #1)</b>       | International Math Local Tournament: Gold Medal |
| Hanoi Open Math Competition: Gold Medal                  | Regional Math Competitions: 5+ Gold Medals      |

## WORK EXPERIENCE

---

- |  |                                  |                     |
|--|----------------------------------|---------------------|
| Meta Platforms   | Senior Machine Learning Engineer | July 2022 - Present |
| <ul style="list-style-type: none"> <li>◦ <b>#1 impact driver</b> to organizational topline metrics, directly contributing to <b>90-120% of revenue and session targets</b> and <b>40-50% of DAU goals</b> through end-to-end ownership of high-leverage ML solutions. Recognized with Exceeded Expectations+ ratings over multiple halves.</li> <li>◦ <b>Finetuned and productionized a compact LLM</b> to enhance content generation and candidate ranking.</li> <li>◦ Pushed the frontier of recommendation modeling by pioneering advanced architectures and optimization strategies. Developed novel techniques including <b>reinforcement learning, contextual and long-horizon representation learning, dynamic user interest modeling, causal inference, knowledge distillation, pairwise ranking, and value-aware multi-task learning</b> to capture latent interests and maximize user engagement.</li> <li>◦ Spearheaded <b>large-scale, cross-organizational initiatives spanning 20+ engineers</b> across infrastructure, product, and modeling teams. Independently defined strategies, aligned stakeholders, and executed system-wide upgrades that materially advanced Meta's notifications recommendation system.</li> <li>◦ <b>Scaled team capability by 3x</b> by proactively scoping and launching new technical roadmaps, <b>mentoring and onboarding 10+ IC4/IC5 engineers</b>, and enabling sustained growth and research-grade innovation.</li> </ul> |                                  |                     |
| Meta Platforms   | Software Engineering Intern      | Jun 2021 - Sep 2021 |
| <ul style="list-style-type: none"> <li>◦ Developed large-scale ML models in PyTorch to generate contextual embeddings of users and ads, enhancing <b>semantic understanding</b> and <b>fine-grained personalization</b> in ad retrieval and ranking systems.</li> <li>◦ Designed and deployed representation learning frameworks for user-ad matching across multiple surfaces, significantly improving relevance and driving <b>ad revenue growth</b>. Used Python, SQL, C++, PHP, and Spark.</li> </ul>  |                                  |                     |
| Meta Platforms   | Software Engineering Intern      | Jun 2020 - Sep 2020 |
| <ul style="list-style-type: none"> <li>◦ Enhanced latency profiling tools with module-level debugging. Achieved a <b>5x increase in runtime efficiency</b>.</li> <li>◦ Accelerated Conv1D and channel shuffle operations via low-level kernel optimizations for the <b>PyTorch</b> framework, delivering up to <b>10x operator-level speedup</b> for on-device speech and NLP models. Used Python and C++.</li> </ul>  |                                  |                     |

## RESEARCH EXPERIENCE

---

- |   |                          |                     |
|---|--------------------------|---------------------|
| Stanford AI Lab   | Undergraduate Researcher | Apr 2020 - Sep 2021 |
| <ul style="list-style-type: none"> <li>◦ Built large-scale data pipelines in Python, NumPy, and Pandas to ingest, align, and preprocess satellite imagery and global forest loss driver labels. Implemented advanced <b>data augmentation</b> and <b>stratified sampling</b> techniques to improve signal diversity and model robustness across diverse geographies.</li> <li>◦ Designed and trained deep learning models in PyTorch, including <b>CNNs, LSTMs, and multimodal fusion architectures</b>, to classify forest loss drivers from multi-temporal satellite imagery, achieving <b>80% classification accuracy</b> and supporting scalable environment monitoring.</li> </ul> |                          |                     |
| Stanford InfoLab  | Undergraduate Researcher | Jan 2020 - Apr 2020 |
| <ul style="list-style-type: none"> <li>◦ Engineered a high-throughput input pipeline for loading and preprocessing underwater video data, including <b>image normalization, spatial augmentation, and temporal slicing</b>, to support robust model training.</li> <li>◦ Developed <b>Mask R-CNN</b> and <b>U-Net</b> models in TensorFlow to detect, localize, and temporally track coral structures in underwater environments, enabling fine-grained analysis and monitoring of reef health over time.</li> </ul>  |                          |                     |
| Computer Science Lab  | Undergraduate Researcher | Sep 2019 - Dec 2019 |
| <ul style="list-style-type: none"> <li>◦ Researched model compression techniques, including <b>structured pruning, regularization-based sparsity, and weight quantization</b>, to reduce the computational footprint of deep convolutional networks.</li> <li>◦ Applied regularization and pruning strategies to ResNet, achieving a <b>15% reduction in FLOPs</b> with a <b>2% improvement in accuracy</b>, demonstrating efficient compression without performance trade-offs.</li> </ul>   |                          |                     |

## PUBLICATIONS

---

FlapAI Bird: Training an Agent to Play Flappy Bird Using Reinforcement Learning Techniques	
Tai Vu, Leon Tran	[paper] [github]
How Not to Give a FLOP: Combining Regularization and Pruning for Efficient Inference	
Tai Vu, Emily Wen, Roy Nehoran	[paper] [github]
Privacy Preserving Inference of Personalized Content for Out of Matrix Users	
Michael Sun, Tai Vu, Andrew Wang	[paper] [github]
GANime: Generating Anime and Manga Character Drawings from Sketches with Deep Learning	
Tai Vu, Robert Yang	[paper] [github]
BERT-VQA: Visual Question Answering on Plots	
Tai Vu, Robert Yang	[paper] [github]
Pixel-Perfect Piloting: Superhuman Control of Pixelcopter via Reinforcement Learning	
Tai Vu, Brad Nikkel, Jenny Yang	[paper] [github]
Beyond the Panels: A Deep Neural Network Approach for Manga Object Detection	
Tai Vu, Robert Yang	[paper] [github]
Amplifying Emotional Signals: Data-Efficient Deep Learning for Robust Speech Emotion Recognition	
Tai Vu	[paper] [github]
From Bayes to BERT: A Comprehensive Benchmark for State-of-the-Art Intent Detection	
Tai Vu, Robert Yang	[paper] [github]

## SKILLS

---

**Languages:** Python • C • C++ • SQL • JavaScript • TypeScript • Java • R • PHP

**Technologies:** PyTorch • TensorFlow • HuggingFace • NumPy • Pandas • vLLM • Accelerate • FSDP • TRL  
OpenRLHF • Scikit-learn • AWS • Google Cloud • Azure • Spark • HTML • CSS • React • Node.js • MongoDB  
Express.js • Next.js • React Native • Django Heroku • Sass • Git • Linux

**Coursework:** Machine Learning • Deep Learning • LLMs • NLP • Reinforcement Learning • Computer Vision  
Information Retrieval • Spoken Language Processing • Graph ML • Parallel Computing • Data Structures &  
Algorithms • Systems • Web Applications • Statistics • Linear Algebra • Convex Optimization • Game Theory

## PROJECTS

---

FlapAI Bird: Training AI Agents to Play Flappy Bird

- Designed a Flappy Bird agent in PyTorch and NumPy that achieved superhuman performance with **2,000+ scores** by implementing and fine-tuning **reinforcement learning** algorithms, including **SARSA**, **Q-learning**, **function approximation**, and **deep Q networks**. Recognized by [Dong Nguyen](#), creator of Flappy Bird.

Privacy Preserving Inference of Personalized Content for Out of Matrix Users

- Architected a novel graph-based recommender system using **BERT-powered embeddings** and **graph neutral networks** to deliver personalized, privacy-preserving recommendations for cold-start users, outperforming WMF by **7x** and DropoutNet by **1.5x** in user recall. Used PyTorch, HuggingFace, NumPy, Pandas, and Azure.

Pixelcopter-RL: Building Reinforcement Learning Algorithms for Playing Pixelcopter

- Developed an AI agent in PyTorch and NumPy to master Pixelcopter using **reinforcement learning** techniques (**Q-learning**, **SARSA**, **function approximation**, and **policy gradients**), achieving a **peak score of 379**.

ConnAIsseur: An AI-driven Recipe Recommendation Website

- Engineered a **web application** with a **BERT-based recommendation engine** that infers latent taste profiles from unstructured user preference data to provide highly personalized recipe suggestions. Used Python, Javascript, PyTorch, HuggingFace, NumPy, Pandas, Django, React, Node.js, HTML, CSS, and AWS.

GANime: Generating Anime Character Drawings from Sketches

- Developed generative models in TensorFlow and AWS, including **neural style transfer**, **Pix2Pix**, and **CycleGAN**, to automate the colorization of anime sketches, attaining **state-of-the-art scores of 220.5 FID and 0.76 SSIM**.

MangaNet: Building an Object Detection System for Mangas

- Designed advanced object detection model architectures, including **Faster R-CNN**, **RetinaNet**, and **YOLOv3**, for the manga domain, achieving a **state-of-the-art mAP of 71.0**. Used PyTorch, NumPy, and Pandas.