

# Tra My (Chiffon) Nguyen

San Francisco, CA, USA | [hi@mychifffonn.com](mailto:hi@mychifffonn.com) | [github.com/mychifffonn](https://github.com/mychifffonn) | [mychifffonn.com](https://mychifffonn.com)

## RESEARCH INTERESTS

Current and future AI systems aligned with and representative of diverse human experiences: **AI safety** (pluralistic alignment, evaluation, cooperative AI, AI control), **multilingual and multicultural AI**, and **socially responsible AI**.

## EDUCATION

Minerva University, College of Computational Sciences <i>B.Sc in Computational Sciences (Machine Learning and Statistics)</i> , GPA: 3.7/4.0	Sep 2021 — May 2025 San Francisco, CA, USA
<ul style="list-style-type: none"><li>• <b>Relevant Coursework:</b> Machine Learning (A), AI Ethics, Bayesian Modeling (A), Statistical Modeling and Causal Inference (A), Optimization Methods (A), Probability and Statistics (A-), Software Engineering</li><li>• <b>Global Rotation:</b> Seoul (South Korea), Chinese Taipei, Hyderabad (India), Buenos Aires (Argentina), Berlin (Germany)</li><li>• <b>Self-study:</b> <a href="#">Technical AI Alignment</a> (ARENA, 2025), <a href="#">Introduction to AI Alignment</a> (Bluedot Impact, 2025)</li></ul>	
<hr/>	

## RESEARCH EXPERIENCE

AI Research Mentee (Multilingual Agentic Evaluation) SEACrowd, SEACrowd 2026 Research Apprenticeship	Feb 2026 — Present Remote
<ul style="list-style-type: none"><li>• Project: Extending Tau2 Bench to low-resource languages, new cultural domain, and visual modality</li><li>• Mentors: <a href="#">Samuel Cahyawijaya</a> (Cohere Labs, HKUST PhD) and <a href="#">Patomporn Payoungkhamdee</a> (VISTEC PhD Student)</li></ul>	
AI Research Fellow (Chain-of-thought Monitorability) Algoverse AI Research, AI Research Program Fall 2025 (Mentor: <a href="#">Yeonwoo Jang (MATS 8 Scholar)</a> )	Oct 2025 — Present Remote / US
<ul style="list-style-type: none"><li>• Project: Chain-of-thought monitorability in relation to target-monitor capability gaps.</li><li>• <u>Description:</u> Evaluating relationships between success rates (ROC-AUC) of monitoring <b>sandbagging</b> and capability gap of monitor-target pairs within Qwen3 model family (8B–480B) (and cross-family), using <b>Inspect AI</b></li></ul>	[Code]
<hr/>	
Qualitative Researcher & Data Scientist Intern Human Resources Department, T-Hub (Telangana's innovation hub)	Feb 2023 — Apr 2023 Hyderabad, Telengana, India
<ul style="list-style-type: none"><li>• Designed preliminary mental health framework serving 1,500 employees across 650+ startups by analyzing 30+ hours of interviews and employee databases in Excel, identifying 5 key stress factors to inform T-Hub's pilot program</li><li>• Led primary research in a team of four, conducting 30 structured interviews and distributing surveys to 50 employees, generating actionable insights that shaped T-Hub's mental health strategy for India's largest startup ecosystem</li></ul>	

## TEACHING & MENTORING EXPERIENCE

Curious Cardinals, <i>Passion Project</i> , Mentor	Nov 2025 — Present
<ul style="list-style-type: none"><li>• Computational neuroscience: Association between HEMA genes and Parkinson disease</li><li>• <a href="#">CCIR</a> project: Robustness of Fact-Checking Language Models under Evidence Corruption</li></ul>	
Minerva University, <i>PR51 Programming with Python</i> , Lead Peer Tutor and Data Analyst	Spring 2025
<ul style="list-style-type: none"><li>• Taught 40+ first-year students from 20+ countries in <b>weekly hands-on programming labs</b> for 11 weeks, covering Python, OOP, debugging, security, and computing fundamentals</li><li>• Extracted <b>40 data-driven pedagogical insights</b> using Google Drive API, Google Sheet trackers, student and tutor surveys, improving hands-on learning and student engagement for the next class iteration by 15%</li></ul>	
Minerva University, <i>FA50/FA51 Logic, Probability &amp; Statistics</i> , Lead Teaching Assistant	Fall 2023 — Spring 2024
<ul style="list-style-type: none"><li>• Guided <b>150+ students each semester</b> for four semesters in formal logic, probability and statistics, algorithmic thinking, and simulation, through weekly office hours</li><li>• Provided <b>formative assessment on 25 quizzes</b> for 50 students to correct and shape their learning</li><li>• Assisted professors in <b>grading</b> three math and programming assignments per semester</li></ul>	

## LEADERSHIP

SEACrowd Communications Associate & Design Engineer	Aug 2025 — Present
---	--------------------

## SELECTED PROJECTS

---

More projects on [mychifffonn.com/projects](http://mychifffonn.com/projects) and [github.com/mychifffonn](https://github.com/mychifffonn)

### Replication: Unsupervised Elicitation of Language Models ([github.com/mychifffonn/icm](https://github.com/mychifffonn/icm))

Dec 2025

Replicated [Wen et. al \(2025\)](#)'s Internal Coherence Maximization, which elicits human concepts from base language models by maximizing mutual predictability and local consistency among concept-related examples.

### SportConnect: Connect Through Local Sport Events ([github.com/mychifffonn/sport-connect](https://github.com/mychifffonn/sport-connect))

Nov 2025

Scalable web application connecting users to local recreational sports events, featuring secure multi-provider OAuth authentication, comprehensive event management, easy event discovery and filtering, and real-time RSVP tracking. Two-person full-stack project with TypeScript, React, Express, PostgreSQL, TailwindCSS, DaisyUI, and BetterAuth.

### Academic Portfolio Theme for Researchers ([github.com/mychifffonn/website](https://github.com/mychifffonn/website))

Aug — Oct 2025

High-performance academic theme, enabling fast, accessible, and multilingual publishing of publications and technical blogs while ensuring top-tier SEO (99) and Lighthouse (100) scores using Astro, TailwindCSS, and shadcn/ui

### Mini-LLM2 PyTorch Implementation ([github.com/mychifffonn/cmu-advanced-nlp-minllama](https://github.com/mychifffonn/cmu-advanced-nlp-minllama))

May 2024

- Implemented the core architecture of Llama-2 from scratch in PyTorch, including critical components including Rotary Positional Embeddings (RoPE), RMSNorm, and SwiGLU activation functions
- Developed a custom training loop with AdamW optimization to pretrain the model on a small corpus and fine-tune it for sentiment classification tasks (SST-5), resulting in convergence and coherent text generation

### Replication: Synthetic Control (Causal Inference) ([github.com/mychifffonn/synthetic-control-rep](https://github.com/mychifffonn/synthetic-control-rep))

Dec 2023

- Replicated and extended [Chrisinger \(2021\)](#)'s analysis of Philadelphia's SNAP benefit redemption in R, analyzing policy impacts across 4 counties and 50+ months of longitudinal data; identified critical limitations in dataset reliability and magnitude discrepancies between original and replicated results
- Conducted new leave-one-out robustness analysis on synthetic control models, showing model instability

## CERTIFICATES

---

- **Advanced Web Development**, CodePath ([drive.google.com/file/d/1n4dHj4TFM8HWlDXMTt9ZGjEXVIpkP-F-](https://drive.google.com/file/d/1n4dHj4TFM8HWlDXMTt9ZGjEXVIpkP-F-/))
- **Natural Language Specialization**, deeplearning.ai ([coursera.org/verify/specialization/3FJ3W7QJX8GK](https://coursera.org/verify/specialization/3FJ3W7QJX8GK)) Nov 2023
- **Applied Data Science**, World Quant University ([creddly.com/badges/2e1e6902-aae4-47c4-97e2-0ad9265e5561](https://creddly.com/badges/2e1e6902-aae4-47c4-97e2-0ad9265e5561)) Aug 2023
- **Machine Learning Specialization**, deeplearning.ai ([coursera.org/verify/specialization/G9898XKB9EAV](https://coursera.org/verify/specialization/G9898XKB9EAV)) Jun 2022

## SKILLS

---

- **Programming Languages:** Python, TypeScript, SQL, R, Bash
- **Machine Learning:** PyTorch, Inspect AI, unsloth, trl, scikit-learn, LangGraph, LlamaIndex
- **Web/App Development:** Astro, React, FastAPI, Flask, Express.js, PostgreSQL, TailwindCSS, shadcn/ui
- **Tools & DevOps:** Git, Docker, Python tooling (uv, ruff, ty), Render, Netlify, LaTeX, Zotero
- **Languages:** Vietnamese (native), English (fluent/C2), Mandarin Chinese (lower-intermediate/band 4)

Last Updated: Feb 19, 2026