

Wannan (Winnie) Yang

📍 San Francisco, USA ✉ winnieyangwn96@gmail.com 🌐 [personal website](#) 🎓 [Google Scholar](#) 🐙 [github](#)

EDUCATION

Ph.D. candidate. New York University, Buzsáki Lab 🔗	Graduating in Spring 2026
Visiting Student. MIT, Tye Lab 🔗	2018.6 – 2019.6
B.S. in Computational Neuroscience. University of Edinburgh. <i>GPA: 4.0 (USA equivalent)</i>	2014.9 – 2018.6

RESEARCH EXPERIENCE

Self-improving AI 2025.6 – present

Student Researcher at Meta, Collaborator: [Jenny Zhang](#)

- Contributing to the research project that develops a self-improve system to automatically produce agents for solving downstream tasks across various domains.
- In charge of enabling self-improvement in automatic reward function design, which enables RL agents to perform tasks like locomotion, robotic arm control in Embodied AI/Robotics settings (with the Genesis simulator).

LLM Alignment and Post Training ([MI NeurIPS](#) [🔗](#)) 2025.3 – 2025.09

Research Scientist Internship at Meta GenAI

- Developed a novel training algorithm to mitigate hallucination in LLMs.
- The new algorithm delivers ~ **30x higher compute efficiency (FLOPs per token)** and requires ~**20x less training data** comparing to competitive baselines like SFT and DPO.
- Developed a modality-agnostic training pipeline, effectively mitigating hallucination in both text-only and multimodal models.
- It is *the first* ever steering-based training framework with general applicability to **both dense and Mixture-of-Experts (MoE) models**.

LLM Alignment and Interpretability ([preprint](#) [🔗](#)) 2024.3 – 2024.12

Collaborator: [Chen Sun](#), Google DeepMind

- Designed and conducted experiments to study *deception* in LLMs.
- Implemented a battery of interpretability tools including contrastive activation steering, activation patching and sparse auto-encoders (SAEs) to understand and monitor the internal activity of LLMs.
- Built a pipeline ([github](#) [🔗](#)) that allow easy hypothesis testing and fast experiments to evaluate, analyze and steer 25+ large language models from different model families (Gemma, Llama, Qwen and Yi) of different sizes (from 1.5 billion to 405 billion parameters).

Memory Representation and Consolidation ([Science](#) [🔗](#), [Nature](#) [🔗](#), [NeurIPS](#) [🔗](#)) 2022.9 – 2024.3

Mentor: [György Buzsáki](#), NYU

- Led a project to study a key mechanism for selective memory consolidation in the brain. This novel discovery has led to a publication in *Science* (leading author).
- Developed a novel latent-space based decoding method and applied various ML tools (including Bayesian decoding) to decode the content of memory reactivations ('replays') from neural population activity during learning and sleep.
- Cultivated research-oriented software engineering skills. Created a pipeline for decoding large-scale (50TB) electrophysiology data.
- Implemented variants of the method to different datasets and projects, which enabled further key publications, including a collaboration project recently accepted at *Nature* (in press) and a first author paper at *NeurIPS* Symmetry and Geometry in Neural Representations Workshop.
- Open-sourced [demo codes](#) [🔗](#) and [tutorials](#) [🔗](#). The neural data processing and decoding pipeline has been widely used by lab members and colleagues from other research labs.

Brain-inspired Deep Reinforcement Learning ([NeurIPS](#) [🔗](#)) 2021.9 – 2022.9

Collaborator: [Chen Sun](#), Google DeepMind

- Co-developed a brain-inspired (memory consolidation and reflection) framework to build a novel deep RL algorithm.
- The resulting simple and scalable algorithm greatly improved long-term credit assignments in a diverse set of RL tasks (including grid-world, Montezuma's Revenge and other Atari games).

PUBLICATIONS

- W. Yang, X. Qiu, L. Yu, Y. Zhang, O. A. Yang, N. Kokhlikyan, N. Cancedda, D. Garcia-Olano. (2025). **Hallucination Reduction with CASAL: Contrastive Activation Steering for Amortized Learning** *MI NeurIPS* [↗](#).
- W. Yang, Z. Yang, C. Sun, G. Buzsáki. (2025). **How Large Language Models Lie: Rotation of the Truth Direction as a Universal Motif.** *preprint* [↗](#).
- I. Zutshi, A. Apostolelli, W. Yang, Z. Zheng, T. Dohi, E. Balzani, A. H. Williams, C. Savin, G. Buzsáki. (2025). **Hippocampal neuronal activity is aligned with action plans.** *Nature* [↗](#).
- W. Yang, C. Sun, G. Buzsáki. (2024). **Interpretability for Safe AI: LLM Lying as a case study.** *NeurIPS* [↗](#) (SafeGenAi Workshop).
- W. Yang, C. Sun, G. Buzsáki. (2024). **Interpretability for Safe AI: Jailbreak as a case study.** *In preparation*.
- W. Yang, C. Sun, R. Huszár, T. Hainmueller, K. Kiselev, G. Buzsáki. (2024). **Selection of experience for memory by hippocampal sharp wave ripple.** *Science* **383**, 1478-1483. [↗](#)
- C. Sun, W. Yang, T. Jiralerspong, D. Malenfant, B. Alsbury-Nealy, Y. Bengio, B. Richards. (2023). **Contrastive Retrospection: honing in on critical steps for rapid learning and generalization in RL.** *NeurIPS*. [↗](#)
- W. Yang, C. Sun, R. Huszár, G. Buzsáki. (2023). **Changes in the geometry of hippocampal representations across brain states.** Symmetry and Geometry in Neural Representations Workshop *NeurIPS*. [↗](#)
- E. Y. Kimchi, A. Burgos-Robles, G. A. Matthews, T. Chakoma, M. Patarino, J. Weddington, C. A. Siciliano, W. Yang, S. Foutch, R. Simons, M. Fong, M. Jing, Y. Li, D. B. Polley, Kay M. Tye. (2023). **Reward contingency gates selective cholinergic suppression of amygdala neurons.** *eLife* [↗](#)
- S. Tennant, I. Hawes, H. Clark, W. Tam, J. Hua, W. Yang, K. Gerlei, E. Wood, M. Nolan. (2022). **Analogue representation of a spatial memory by ramp-like neural activity in retrohippocampal cortex.** *Current Biology* [↗](#)
- C. Sun, W. Yang, J. Martin, S. Tonegawa. (2020). **Hippocampal neurons represent events as transferable units of experience.** *Nature Neuroscience* [↗](#).

SKILLS

LLM Training: SFT, DPO, PPO

LLM Interpretability and Alignment: transformer-lens, Contrastive Activation Steering, Activation Patching, SAE Steering

Programming: Python, MATLAB, HTML, LaTeX

Statistical Data Analysis: Large-scale High-dimensional Data Analysis, Signal processing, Linear and Nonlinear Dimensionality Reduction, Time Series Data Analysis, Neural Data Decoding, Multimodal Data Analysis