

DUY NGUYEN

✉ duyng@cs.unc.edu  [duykhongnguyen.github.io](https://github.com/duykhongnguyen)  github.com/duykhongnguyen

Research Interests

My research focuses on post-training methods for enhancing the capabilities and continually updating the knowledge of (multimodal) LLMs. Additionally, I am interested in mechanistic interpretability and inference-time interventions for interpreting and monitoring the model behaviors.

Education

The University of North Carolina at Chapel Hill

Aug. 2024 – Aug. 2029 (expected)

Ph.D. in Computer Science

Chapel Hill, NC, USA

- Advisor: Prof. Mohit Bansal

Hanoi University of Science and Technology

Sep. 2018 – Sep. 2022

B.S. in Computer Science

Hanoi, Vietnam

- GPA: 3.65/4.00, graduated with Excellent Degree

Publications

Duy Nguyen, Archiki Prasad, Elias Stengel-Eskin, Mohit Bansal. GrAIInS: Gradient-based Attribution for Inference-Time Steering of LLMs and VLMs. In *Association for Computational Linguistics (ACL)*, 2026. [\[pdf\]](#)

Duy Nguyen*, Archiki Prasad*, Elias Stengel-Eskin, Mohit Bansal. LAsER: Learning to Adaptively Select Reward Models with Multi-Arm Bandits. In *Neural Information Processing Systems (NeurIPS)*, 2025. [\[pdf\]](#)

Duy Nguyen, Archiki Prasad, Elias Stengel-Eskin, Mohit Bansal. Multi-Attribute Steering of Language Models via Targeted Intervention. In *Association for Computational Linguistics (ACL)*, 2025. [\[pdf\]](#)

Bao Nguyen, Binh Nguyen, **Duy Nguyen**, Viet Anh Nguyen. Risk-Aware Distributional Intervention Policies for Language Models. In *Findings of Empirical Methods in Natural Language Processing (Findings of EMNLP)*, 2025. [\[pdf\]](#)

Ngoc Bui, **Duy Nguyen**, Man-Chung Yue, Viet Anh Nguyen. Coverage-Validity-Aware Algorithmic Recourse. In *Operations Research*, 2024. [\[pdf\]](#)

Hieu Nguyen, **Duy Nguyen**, Khoa Doan, Viet Anh Nguyen. Cold-start Recommendation by Personalized Embedding Region Elicitation. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2024. [\[pdf\]](#)

Duy Nguyen, Ngoc Bui, Viet Anh Nguyen. Distributionally Robust Recourse Action. In *International Conference on Learning Representations (ICLR)*, 2023. [\[pdf\]](#)

Duy Nguyen, Ngoc Bui, Viet Anh Nguyen. Feasible Recourse Plan via Diverse Interpolation. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2023. [\[pdf\]](#)

Ngoc Bui, **Duy Nguyen**, Viet Anh Nguyen. Counterfactual Plans under Distributional Ambiguity. In *International Conference on Learning Representations (ICLR)*, 2022. [\[pdf\]](#)

Tuan-Duy H. Nguyen, Ngoc Bui, **Duy Nguyen**, Man-Chung Yue, Viet Anh Nguyen. Robust Bayesian Recourse. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2022. [\[pdf\]](#)

(*) denotes equal contribution

Preprints

Duy Nguyen, Hanqi Xiao, Archiki Prasad, Elias Stengel-Eskin, Hyunji Lee, Mohit Bansal.
Conflict-Resolving and Sharpness-Aware Minimization for Generalized Knowledge Editing with Multiple Updates. *Under Review*. [\[pdf\]](#)

Duy Nguyen, Bao Nguyen, Viet Anh Nguyen. Cost Adaptive Recourse Recommendation by Adaptive Preference Elicitation. *Forever Preprint*. [\[pdf\]](#)

Experience

Amazon Science **May 2025 – Aug. 2025**
Applied Scientist Intern *Seattle, WA, USA*

- Research topics: LLM reasoning for tabular data

VinAI Research (now Qualcomm AI Research) **Aug. 2022 – Aug. 2024**
Research Resident *Hanoi, Vietnam*

- Advisor: Prof. Viet Anh Nguyen
- Research topics: LLM safety, interpretability and explainability

Honors and Awards

NeurIPS 2025 Scholar Award **Oct. 2025**
NeurIPS

Honorable Mention - Undergraduate Operations Research Prize **Oct. 2022**
INFORMS

Best thesis presentation award **Aug. 2022**
School of Information and Communication Technology, HUST

Excellence Scholarship for the academic year **Sep. 2019**
School of Information and Communication Technology, HUST

Professional Academic Services

Reviewer at ICLR (2025, 2026), ICML (2025, 2026), NeurIPS (2023–2025), ACL Rolling Review (2025, 2026).