# Vijay Viswanathan

**CONTACT INFORMATION**

| | |
|---|---|
| **Homepage:** | https://www.cs.cmu.edu/~vijayv/ |
| **E-mail:** | vijayv@andrew.cmu.edu |
| **Github**: | github.com/viswavi |

**RESEARCH STATEMENT**

I do research on making AI models more reliable at specialized tasks. My research primarily investigates the use of synthetic data to achieve this goal.

**EDUCATION**

**Carnegie Mellon University**, School of Computer Science, Pittsburgh, PA  2022–2027 (expected)

- Ph.D in Language and Information Technologies                    GPA: 3.92/4.30
- *Advisors*: Sherry Tongshuang Wu and Graham Neubig

**Carnegie Mellon University**, School of Computer Science, Pittsburgh, PA  2020–2022 (expected)

- M.S. in Intelligent Information Systems                    GPA: 4.05/4.30
- *Advisors*: Prof. Graham Neubig and Dr. Pengfei Liu
- *Relevant Courses*: Algorithms for NLP, Probabilistic Graphical Models, Question Answering, Convex Optimization, Search Engines

**Carnegie Mellon University**, Mellon College of Science, Pittsburgh, PA                    2012–2016

- B.S. in Mathematics, with minors in Computer Science and Language Technologies
- *Undegradute Research Advisor*: Prof. Eric Nyberg
- *Relevant Courses*: Real Analysis I and II, Intro to Theoretical Computer Science, Algebraic Structures, Basic Logic, Probability

**FELLOWSHIPS AND AWARDS**

- Oustanding Demo Paper, ACL 2022 (top 3 papers out of 75 system demo papers presented)
- NEC Student Research Fellow, 2022-2023
- Co-PI, Microsoft "Accelerate Foundation Models Research" Grant, 2023 (*$20,000, from Microsoft Research*)
- Johnson & Johnson Undergraduate Research Award (presented at CMU's "Meeting of the Minds"), 2015
- Thiel "20 Under 20" Fellowship Finalist, 2012

**PUBLICATIONS**

1. Jaixin Ge, Xueying Jia, **Vijay Viswanathan**, Hongyin Luo, Graham Neubig, Graham Neubig. 2024. " Training Task Experts through Retrieval Based Distillation". In *arXiv:2407.12874*.

2. Ian Wu, Sravan Jayanthi, **Vijay Viswanathan**, Simon Rosenberg, Sina Pakazad, Tongshuang Wu, Graham Neubig. 2024. "Synthetic Multimodal Question Generation". In *Findings of the Association for Computational Linguistics: EMNLP 2024 (EMNLP Findings 2024)* (to appear).

3. Chenyang Zhao, Xueying Jia, **Vijay Viswanathan**, Tongshuang Wu, Graham Neubig. 2024. "Large Language Models Enable Few-Shot Clustering". In *1st Conference on Language Modeling (COLM) 2024*.

4. Saumya Gandhi*, Ritu Gala*, **Vijay Viswanathan**, Tongshuang Wu, Graham Neubig. 2024. "Better Synthetic Data by Retrieving and Transforming Existing Datasets". In *Findings of the Association for Computational Linguistics: ACL 2024 (ACL Findings 2024)*.

5. **Vijay Viswanathan**, Kiril Gashteovski, Carolin Lawrence, Tongshuang Wu, Graham Neubig. 2024. "Large Language Models Enable Few-Shot Clustering". In *Transactions of the Association for Computational Linguistics (TACL)*.

6. **Vijay Viswanathan**\*, Chenyang Zhao*, Amanda Bertsch, Tongshuang Wu, Graham Neubig. 2023. "Prompt2Model: Generating Deployable Models from Natural Language Instructions". In *Conference on Empirical Methods in Natural Language Processing Demo Track (EMNLP Demo Track 2023)*.

7. **Vijay Viswanathan**, Luyu Gao, Tongshuang Wu, Pengfei Liu, Graham Neubig. 2023. "DataFinder: Scientific Dataset Recommendation from Natural Language Descriptions". In *Annual Conference of the Association for Computational Linguistics (ACL 2023)*.

8. Aryeh Tiktinsky*, **Vijay Viswanathan***, Danna Niezni, Dana Meron Azagury, Yosi Shamay, Hillel Taub-Tabib, Tom Hope, Yoav Goldberg. 2021. "An *N*-ary Relation Extraction Dataset for Drug Combinations". In *Proceedings of the 2022 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2022)*. **Oral**.

9. Yang Xiao, Jinlan Fu, Weizhe Yuan, **Vijay Viswanathan**, Zhoumianze Liu, Yixin Liu, Graham Neubig, Pengfei Liu. 2022. "'DataLab: A Platform for Data Analysis and Intervention". In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics, System Demonstrations Track (ACL Demo Track 2022)*. **\*Outstanding Demo Paper\***.

10. **Vijay Viswanathan**, Graham Neubig, Pengfei Liu. 2021. "CitationIE: Leveraging the Citation Graph for Scientific Information Extraction". In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers) (ACL-IJCNLP 2021)*. **Oral**.

11. Siddhant Arora*, Alissa Ostapenko*, **Vijay Viswanathan***, Siddharth Dalmia*, Florian Metze, Shinji Watanabe, Alan W Black. 2021. "Rethinking End-to-End Evaluation of Decomposable Tasks: A Case Study on Spoken Language Understanding". In *22nd Annual Conference of the International Speech Communication Association (Interspeech 2021)*.

   *\* authors contributed equally*

PRESS
COVERAGE

- **2023**
  - "The push to make big AI small" — Axios
  - "CMU & Tsinghua U's Prompt2Model Generates Deployable Models Following Natural Language Instructions" — SyncedReview / Jiqi Zhixin
  - "Researchers from CMU and Tsinghua University Propose Prompt2Model: A General Purpose Method that Generates Deployable AI Models from Natural Language Instructions" — MarkTechPost

MENTORED
STUDENTS

- Jiaxin Ge (Peking Uni.)                                                                 Dec 2023 - Now
  Co-advised with Graham Neubig and Hongyin Luo (MIT).

- Xueying Jia (CMU)                                                                       Oct 2023 - Now
  Collaborating with Chenyang Zhao and co-advised with Graham Neubig and Tongshuang Wu. Published a paper (as co-first author) at COLM 2024.

- Ritu Gala (CMU)                                                                    Aug 2023 - May 2024
  Collaborating with Saumya Gandhi and co-advised with Graham Neubig and Tongshuang Wu. Published their first paper (as co-first author) at ACL Findings 2024.

- Saumya Gandhi (CMU)                                                               Aug 2023 - May 2024
  Collaborating with Ritu Gala and co-advised with Graham Neubig and Tongshuang Wu. Published a paper as co-first author at ACL Findings 2024.

- Vanya Bannihatti Kumar (CMU)                                                            Aug 2023 - Now
  Co-advised with Graham Neubig and Tongshuang Wu.

- Chenyang Zhao (Tsinghua Uni.)                                                     Mar 2023 - Aug 2024
  Had first paper accepted to EMNLP Demo 2023 (co-led with Vijay Viswanathan), and had a second paper published as co-first author at COLM 2024. Co-advised with Graham Neubig and Tongshuang Wu.

- Yuanchen Bai (CMU)                                                                May 2022 - Oct 2023
  Had first paper accepted to a workshop at AACL in Nov 2023 (co-led with Raoyi Huang). Co-advised with Tzu-Sheng Kuo, and Tongshuang Wu.

- Raoyi Huang (CMU)                                                                 May 2022 - Oct 2023
  Had first paper accepted to a workshop at AACL in Nov 2023 (co-led with Yuanchen Bai). Co-advised with Tzu-Sheng Kuo, and Tongshuang Wu.

| | |
|---|---|
| **Research**<br>**Experience** | • **Cohere**     May 2024–Sep 2024 |

RESEARCH
EXPERIENCE

- **Cohere** — May 2024–Sep 2024
  *Intern of Technical Staff*, with Pat Verga on the *Agents & RAG Team* — Remote
  - Generated a large-scale synthetic dataset for training web navigation agents.
  - Leveraged weak supervision from a web search dataset to guide the data synthesis process.

- **Carnegie Mellon University, Language Technologies Institute** — Aug 2022–*present*
  *Doctoral Researcher* with Profs. Sherry Tongshuang Wu and Graham Neubig — Pittsburgh, PA
  - Studying applications of data generation for NLP as an NEC Student Research Fellow.
  - Designing new algorithms for semi-supervised clustering for diverse tasks such as open knowledge base canonicalization and document clustering.

- **Carnegie Mellon University, Language Technologies Institute** — Aug 2020–May 2022
  *Research Assistant* for Prof. Graham Neubig and Dr. Pengfei Liu — Pittsburgh, PA
  - Studied extracting critical structured information (e.g. methods and metrics) from full scientific texts, using the *SciREX* dataset
  - Proposed using the citation graph in a neural multi-task information extraction system, giving state-of-the-art results (oral presentation at the *ACL-IJCNLP 2021* conference)
  - Working on automatically recommending datasets to use for a given natural language system description (submitted to *ARR Dec '22*, currently in revision)

- **Allen Institute for Artificial Intelligence (AI2)** — Summer 2021
  *Research Intern* at AI2 Israel, with Prof. Yoav Goldberg — Remote
  - Worked on extracting a database of drug interactions from text in biomedical abstracts
  - Supported the construction of a novel annotated dataset for discovering drug interactions by implementing an efficient modeling framework, implementing new relation extraction metrics, and benchmarking state-of-the-art models for this new task (submitted to *ARR Nov '22*)
  - Developed novel task-specific pretraining scheme to leverage unlabeled data

- **Carnegie Mellon University, Language Technologies Institute** — Jan 2015 - May 2016
  *Research Assistant* for Prof. Eric Nyberg — Pittsburgh, PA
  - Created a query-to-question conversion system to convert a search engine query to its intent
  - Won $1000 Johnson & Johnson Award at CMU's undergraduate research symposium
  - Developed method for using an ensemble of adaptive sampling methods in active learning

WORK
EXPERIENCE

- **Uber ATG** — Apr 2019 - Aug 2020
  *Software Engineer*, Prediction team (under Dr. Micol Marchetti-Bowick) — Pittsburgh, PA
  - Developed model enabling autonomous cars to jointly predict spatial paths and vehicle-vehicle interactions, considerably improving the comfort of autonomous driving in traffic
  - Owned system to measure the autonomous car's predictions of other vehicles via simulation
  - Built application to crowdsource labels for road interactions in dense traffic from onboard sensor logs

- **Scaled Inference** — Mar 2017 - Mar 2019
  *Inference Engineer*, Statistical Modeling team — Palo Alto, CA
  - Implemented algorithms for contextual, adaptive A/B optimization. Wrote first production version of our policy search algorithm and efficient Bayesian inference from scratch in Go
  - Responsible for debugging and improving statistical models and RL algorithms in production

- **Yik Yak** — Jun 2016 - Dec 2016
  *Software Engineer*, Machine Intelligence Team (under Dr. Marsal Gavalda) — Atlanta, GA
  - Designed and implemented the ranking component of a user recommendation system for Yik Yak's location-based social network
  - Trained deep models on large user-generated text and image data (from >100M messages)
  - Built user interface to crowdsource user relevance data to tune our ranking model
  - Deployed ranker to production via Docker and Kubernetes

SERVICE

- Student Member, CMU LTI Diversity, Equity, and Inclusion (DEI) Committee, 2022-2023

- Reviewer, Controllable Generative Modeling in Language and Vision (CtrlGen) Workshop at NeurIPS, 2021

- Chair, Student Committee for CMU LTI Faculty Hiring, 2022

- Member, Student Committee for CMU LTI Faculty Hiring, 2021

- Teaching Assistant for "Multilingual Natural Language Processing" (11-737), Spring 2022

STANDARDIZED
TEST SCORES

- GRE: 170/170 (Quantitative), 167/170 (Verbal), 5.5/6.0 (Writing) taken on 10/2/2021

REFERENCES

Dr. **Graham Neubig**
Associate Professor at Carnegie Mellon University
E-mail: gneubig@andrew.cmu.edu

Dr. **Sherry Tongshuang Wu**
Assistant Professor at Carnegie Mellon University
E-mail: sherryw@cs.cmu.edu

Dr. **Yoav Goldberg**
Research Director at Allen Institute for Artificial Intelligence, Israel and Professor at Bar Ilan University
E-mail: yoavg@allenai.org

Dr. **Kiril Gashteovski**
Senior Research Scientist at NEC Laboratories Europe
E-mail: kiril.gashteovski@neclab.eu

Dr. **Micol Marchetti-Bowick**
Tech Lead Manager at Aurora Innovation
(formerly, Manager at Uber ATG)
E-mail: micol.mb@gmail.com