

Nupoor Gandhi

phone: (408) 893-3764 | email: nmgandhi@cs.cmu.edu | web: nupoorgandhi.github.io

RESEARCH STATEMENT

My research focuses on **low-cost structured information extraction** through the development of **efficient annotation methods** in combination with **task decomposition**. I am particularly interested in low-resource settings where large-scale expert annotation is prohibitively expensive.

EDUCATION

Carnegie Mellon University

Ph.D., Machine Learning and Public Policy

Advisor: Emma Strubell

Thesis: Center-based Decomposition for Efficient Annotation of Structured Representations

Expected Graduation 2026

Pittsburgh, PA

University of Illinois at Urbana-Champaign

Bachelor of Science, Computer Science; Minor in Mathematics

Advisor: ChengXiang Zhai

Senior Thesis: Forecasting Public Health Outcomes with Social Media Data

August 2016 – May 2020

Urbana, IL

PEER-REVIEWED PUBLICATIONS

1. **Decomposing Unitization and Typing for Efficient and Consistent Span-Bound Concept Annotation.**

N. Gandhi, M. Bada, E. Strubell

Findings of the Association for Computational Linguistics (ACL), 2026.

2. **SynthTextEval: Synthetic Text Data Generation and Evaluation for High-Stakes Domains.**

K. Ramesh, D. Smolyak, Z. Zhao, N. Gandhi, R. Agarwal, M. Bjarnadóttir, A. Field

Proceedings of EMNLP – Demo Track, 2025.

3. **Beyond Text: Characterizing Domain Expert Needs in Document Research.**

S. Gururaja, N. Gandhi, J. Milbauer, E. Strubell

Findings of the Association for Computational Linguistics (ACL), 2025.

4. **Evaluating Differentially Private Synthetic Data Generation in High-Stakes Domains.**

K. Ramesh, N. Gandhi, P. Madaan, L. Bauer, C. Peris, A. Field

Findings of Empirical Methods in Natural Language Processing (EMNLP), 2024.

5. **Challenges in End-to-End Policy Extraction from Climate Action Plans.**

N. Gandhi, T. Corringham, E. Strubell

Proceedings of the ClimateNLP Workshop at ACL, 2024.

6. **Mention Annotations Alone Enable Efficient Domain Adaptation for Coreference Resolution.**

N. Gandhi, A. Field, E. Strubell (*Selected for Oral Presentation*)

Proceedings of the Association for Computational Linguistics (ACL), 2023.

7. **Examining Risks of Racial Biases in NLP Tools for Child Protective Services.**

A. Field, A. Coston, N. Gandhi, A. Chouldechova, E. Putnam-Hornstein, D. Steier, Y. Tsvetkov

Proceedings of ACM FAccT, 2023.

8. **Improving Span Representation for Domain-adapted Coreference Resolution.**

N. Gandhi, A. Field, Y. Tsvetkov

Proceedings of the EMNLP Workshop on Computational Models of Reference, Anaphora and Coreference (CRAC), 2021.

9. **Predicting Opioid Abuse with Text-based Twitter Features.**

N. Gandhi, A. Morales, S. Chan, D. Albarracin, C. Zhai

Proceedings of the AAAI Conference on Artificial Intelligence – Student Abstract, 2020.

- Multi-Attribute Topic Feature Construction for Social Media-based Prediction.**
A. Morales, **N. Gandhi**, S. Chan
IEEE International Conference on Big Data (Big Data), 2018.

NON-ARCHIVAL WORKSHOPS AND PREPRINTS

- Task Decomposition for Efficient Annotation.**
N. Gandhi, E. Strubell
Under Review, 2026.
- Extracting Structured Policy Information from Climate Action Plans**
T. Corringham, **N. Gandhi**, B. Flores, and E. Strubell, S. Gururaja, T. Romanov, J. Dunafon
NeurIPS 2025 Workshop on Tackling Climate Change with Machine Learning, 2025.
- LLMs for Few-shot Information Extraction with Climate Action Plans**
N. Gandhi, T. Corringham, and E. Strubell
Text-as-Data Workshop 2023

INVITED TALKS

- Centering as a Critical Challenge for Efficient Annotation of Structured Representations**
Poster Presentation at Midwest Speech and Language Days 2025
- Challenges in End-to-End Policy Extraction from Climate Action Plans**
Poster Presentation at ClimateNLP Workshop, ACL 2024
Oral Presentation at Climate NLP Reading Group 2024
- LLMs for Few-shot Information Extraction with Climate Action Plans**
Poster Presentation at Text-as-Data Conference 2023
- Mention Annotations Alone Enable Efficient Domain Adaptation for Coreference Resolution**
Oral Presentation at ACL 2023
- Risk Prediction and Information Extraction for Child Welfare in Allegheny County**
Guest Lecture for Computational Ethics in NLP 2022
- Improving Span Representation for Domain-adapted Coreference Resolution**
Oral Presentation at EMNLP CRAC Workshop 2021
- Predicting Opioid Abuse with Text-based Twitter Features**
Poster Presentation at AAAI Conference on Artificial Intelligence 2020
- Reducing Response Times to Citizen Legal Questions in Uganda**
Lightning Talk at KDD – Social Impact Session 2020
Oral Presentation at Data Science for Social Good: DataFest 2019
- Natural Selection and the Bystander Effect**
Poster Presentation at Undergraduate Research Symposium 2019
- Multi-dimensional Features for Prediction with Tweets**
Poster Presentation at ACM Tapia Conference 2018
- Predicting Sexually Transmitted Infections with Twitter Data**
Rapid-Fire Talk at Society of Women Engineers Conference 2018
Poster Presentation at Illinois Scholars in Undergraduate Research Expo 2018

TEACHING EXPERIENCE

- Probabilistic Graphical Models (10-708)** with *Andrej Risteski* Spring 2026
Machine Learning Department, CMU
- Computational Ethics in NLP (11-830)** with *Emma Strubell and Alan Black* Spring 2022
Language Technologies Institute, CMU

FELLOWSHIPS & AWARDS

AAAI Travel Grant	2020
Knights of St. Patrick , <i>highest award for leadership by College of Engineering, UIUC</i>	2020
Outstanding Poster Presentation - Science & Math , <i>Undergraduate Research Week, UIUC</i>	2019
Gold Collegiate Section Award , <i>Society of Women Engineers Conference</i>	2018
Rapid Fire Research Finalist , <i>Society of Women Engineers Conference</i>	2018
Illinois Scholars Research Grant , <i>awarded for excellence in research, UIUC</i>	2017

RESEARCH EXPERIENCE

Graduate Student Researcher <i>Carnegie Mellon University, Structure in(g) LAnguage LAB</i> Advisor: Emma Strubell	January 2022 – Present <i>Pittsburgh, PA</i>
<ul style="list-style-type: none">• Methods for data-efficient domain transfer with limited annotation schemes.• Cross-domain alignment of structured knowledge using shallow semantic structures.	
Graduate Student Researcher <i>Carnegie Mellon University, CMU TsvetShop</i> Advisor: Yulia Tsvetkov and Alexandra Chouldechova	September 2020 – December 2021 <i>Pittsburgh, PA</i>
<ul style="list-style-type: none">• Adapting end-to-end coreference resolution model for clinical notes from child welfare (Allegheny County Dept. Human Services) and medical domains.• Developing new methods to use external knowledge to adapt coreference models to new domains.	
Undergraduate Student Researcher <i>University of Illinois at Urbana-Champaign, Text Information Management Analysis Group</i> Advisor: ChengXiang Zhai	August 2016 – September 2020 <i>Urbana, IL</i>
<ul style="list-style-type: none">• Predicted outbreaks of sexually-transmitted infections and opioid mortality rates using LDA topic features over a Twitter corpus.• Constructed features to represent regional colloquialisms following geolocation of tweets.	
Data Science for Social Good Fellow <i>Imperial College London</i> Advisor: Rayid Ghani, Leid Zejnilovic	June 2019 – August 2019 <i>London, UK</i>
<ul style="list-style-type: none">• Built an Information Retrieval system to assist lawyers at nonprofit Barefoot Law reach citizens 60% faster.	

INDUSTRY EXPERIENCE

Software Engineering Intern <i>Qualcomm</i>	May 2018 – August 2018 <i>San Diego, CA</i>
<ul style="list-style-type: none">• Achieved over 10% reduction in memory consumption for transceiver chip.• Developed tool to eliminate redundancies in radiofrequency settings data for 5G modem.• Modeled embedded software with Abstract Syntax Trees to automate optimization of memory layout.	
Software Engineering Intern <i>Synopsys</i>	May 2017 – August 2017 <i>Mountain View, CA</i>
<ul style="list-style-type: none">• Built ML API to analyze runtime performance of EDA tools.• Developed web-based push notification system for SQL server.	

SERVICE

Reviewing	
ACL Rolling Review	2026
ClimateNLP Workshop at ACL	2024–2026
Journal of Comparative Policy Analysis	2024
EMNLP	2023
Academic	
Application Reviewer, Data Science for Social Good Fellowship	2022, 2026
Mentor, CMU AI Mentoring Program	2021 – Present
Committee Member, Machine Learning Department Masters Admissions	2022

GSA Representative, <i>Graduate Student Assembly, CMU</i>	2023 – Present
Student Leadership Committee Representative, <i>CMU Senate</i>	2025–2026
Student Volunteer, ACL Conference	2023
Volunteer, <i>Scopeathon Data Science for Social Good, CMU</i>	2022, 2025
President, <i>Society of Women Engineers, UIUC</i>	2019–2020
External Vice President, <i>Society of Women Engineers, UIUC</i>	2017–2018

Community

Canvasser, <i>Daeja Baker Shaler County School Board Campaign</i>	2025
Canvasser, <i>Summer Lee U.S. House Congressional Campaign (Primary)</i>	2022, 2024
Canvasser, <i>Kamala Harris U.S. Presidential Campaign</i>	2024
Organizer, <i>Pittsburgh Public Parks Coalition</i> (campaign to plug abandoned oil wells in SW PA)	2025
Volunteer Data Scientist, <i>Illinois Legal Aid Online</i>	2020

SKILLS

Python, Java, C, TensorFlow, PyTorch, NumPy/SciPy/pandas, Transformers, Huggingface, NLTK, R, L^AT_EX