

# Swabha Swayamdipta

GABILAN ASSISTANT PROFESSOR OF COMPUTER SCIENCE, UNIVERSITY OF SOUTHERN CALIFORNIA

✉ swabhas@usc.edu | 🌐 <https://swabhs.com> | 📺 [swabhs](#) | 🐦 [@swabhz](#)

## Appointments

---

### University of Southern California

GABILAN ASSISTANT PROFESSOR OF COMPUTER SCIENCE

- Co-Associate Director of the Center for AI and Society

Los Angeles, CA, USA

Aug 2022 - PRESENT

### Amazon AGI

AMAZON SCHOLAR

- Team: Responsible AI

Los Angeles, CA, USA

Jan 2025 - PRESENT

### Allen Institute for AI

POSTDOCTORAL INVESTIGATOR

- Team: MOSAIC led by Yejin Choi

Seattle, WA, USA

Aug 2019 - June 2022

### Oracle Server Technologies

MEMBER TECHNICAL STAFF

- Team: Unified messaging for Fusion Middleware

Bangalore, India

2010 - 2011

## Education

---

### Carnegie Mellon University

PHD IN LANGUAGE AND INFORMATION TECHNOLOGIES

- **Advisors:** Noah A. Smith, Chris Dyer
- **Thesis:** Syntactic Inductive Biases in NLP

Pittsburgh, PA, USA

2013 - 2019

### Columbia University

MASTERS IN COMPUTER SCIENCE

- **Advisors:** Owen Rambow, Michael Collins

City of New York, NY, USA

2011 - 2012

### National Institute of Technology

BACHELORS OF TECHNOLOGY IN COMPUTER SCIENCE AND ENGINEERING

Calicut, Kerala, India

2006 - 2010

## Awards and Grants

---

2025	<b>Simons Institute Research Fellowship</b> , Special Program on LLMs and Transformers, Spring 2025	UC Berkeley
2024	<b>Outstanding Paper Award</b> , OATH-Frames: Characterizing Online Attitudes Towards Homelessness with LLM Assistants	EMNLP
2024	<b>NSF CISE Core Medium</b> , RII: Hard Data to the Model	NSF
2024	<b>Zumberge Fellowship</b> , DEI in Research	USC OORI
2023	<b>Intel Rising Stars Award</b> , LLM Inference for Scaling Data Creation	Intel Labs
2022	<b>AI2 Young Investigators Award</b> , Synthesizing Data Towards Commonsense and Causal Reasoning	Allen Institute for AI
2022	<b>Gabilan Assistant Professor Fellowship</b> ,	USC WISE
2022	<b>Outstanding Paper Award</b> , Understanding Dataset Difficulty with V-Usable Information	ICML
2021	<b>Outstanding Paper Award</b> , MAUVE: Measuring the Gap Between Neural Text and Human Text using Divergence Frontiers	NeurIPS
2020	<b>Best Paper: Honorable Mention</b> , Don't Stop Pretraining: Adapt Language Models to Domains and Tasks	ACL
2019	<b>Rising Star in EECS</b> , UIUC	
2014	<b>Best Student Presentation</b> , Diversity in Dependency Parsing	CMU LTI Student Research Symposium
2009	<b>Sun Campus Ambassador Scholarship</b> , Bangalore, India	Sun Microsystems
2003	<b>National Talent Search Scholarship</b> , Rourkela, India	NCERT, India

# Publications

---

## CONFERENCE AND JOURNAL PAPERS

Better Language Model Inversion by Compactly Representing Next-Token Distributions	<i>NeurIPS 2025</i>
• M. Nazir, M. Finlayson, J. X. Morris, X. Ren, <b>S. Swayamdipta</b>	
Teaching Models to Understand (but not Generate) High-risk Data	<i>CoLM 2025</i>
• R. Wang, M. Finlayson, L. Soldaini, <b>S. Swayamdipta</b> , R. Jia	
Improving LLM Personas via Rationalization with Psychological Scaffolds	<i>EMNLP 2025</i>
• B. Joshi, X. Ren, <b>S. Swayamdipta</b> , R. Koncel-Kedziorski, T. Paek	
Evaluating Evaluation Metrics – The Mirage of Hallucination Detection	<i>EMNLP-Findings 2025</i>
• A. Kulkarni, Y. Zhang, J. R. A. Moniz, X. Ge, B.-H. Tseng, D. Piraviperumal, <b>S. Swayamdipta</b> , H. Yu	
Robust Data Watermarking in Language Models by Injecting Fictitious Knowledge	<i>ACL-Findings 2025</i>
• X. Cui, J. Wei, <b>S. Swayamdipta</b> , R. Jia,	
ELI-Why: Evaluating the Pedagogical Utility of LLM Explanations	<i>ACL-Findings 2025</i>
• B. Joshi, K. He, S. Ramnath, S. Sabouri, K. Zhou, S. Chattopadhyay, <b>S. Swayamdipta</b> , X. Ren,	
Compare without Despair: Reliable Preference Evaluation with Generation Separability	<i>EMNLP 2024</i>
• S. Ghosh, T. Srinivasan, <b>S. Swayamdipta</b> ,	
Out-of-Distribution Detection through Soft Clustering with Non-Negative Kernel Regression	<i>EMNLP 2024</i>
• A. Gulati, X. Dong, C. Hurtado, S. Shekkizhar, <b>S. Swayamdipta</b> , A. Ortega	
OATH-Frames: Characterizing Online Attitudes towards Homelessness via LLM Assistants	<i>EMNLP 2024</i>
• J. Ranjit, B. Joshi, R. Dorn, L. Petry, O. Koumoundouros, J. Bottarini, P. Liu, E. Rice, and <b>S. Swayamdipta</b>	
• <b>Outstanding Paper Award @ EMNLP 2024 &amp; Best Poster @ ShowCAIS 2024</b>	
Logits of API-Protected LLMs Leak Proprietary Information	<i>COLM 2024</i>
• M. Finlayson, X. Ren, <b>S. Swayamdipta</b> ,	
Crowd-Calibrator: Can Annotator Disagreement Inform Calibration in Subjective Tasks?	<i>COLM 2024</i>
• U. Khurana, E. Nalisnick, A. Fokkens, <b>S. Swayamdipta</b> ,	
Annotating FrameNet via Structure-Conditioned Language Generation	<i>ACL 2024</i>
• X. Cui, and <b>S. Swayamdipta</b>	
Closing the Curious Case of Neural Text Degeneration	<i>ICLR 2024</i>
• M. Finlayson, J. Hewitt, A. Koller, <b>S. Swayamdipta</b> , A. Sabharwal	
Does Video Summarization Require Videos? Quantifying the Effectiveness of Language in Video Summarization	<i>ICASSP 2024</i>
• Y. Nam, A. Lehavi, D. Yang, D. Bose, <b>S. Swayamdipta</b> , and S. Narayanan	
NeuroComparatives: Neuro-Symbolic Distillation of Comparative Knowledge	<i>NAACL-Findings 2024</i>
• P. Howard, J. Wang, V. Lal, G. Singer, Y. Choi, and <b>S. Swayamdipta</b>	
Generative Explanations for Program Synthesizers	<i>ICSE + VL/HCC 2024</i>
• A. Nazari, S. Chattopadhyay, <b>S. Swayamdipta</b> , M. Raghothaman	
We're Afraid Language Models Aren't Modeling Ambiguity	<i>EMNLP 2023</i>
• A. Liu, Z. Wu, J. Michael, A. Suhr, P. West, A. Koller, <b>S. Swayamdipta</b> , N. A. Smith, and Y. Choi	
MAUVE Scores for Generative Models: Theory and Practice	<i>JMLR 2023</i>
• K. Pillutla, L. Liu, J. Thickstun, S. Welleck, <b>S. Swayamdipta</b> , R. Zellers, S. Oh, Y. Choi, and Z. Harchaoui	
I2D2: Inductive Knowledge Distillation with NeuroLogic and Self-Imitation	<i>ACL 2023</i>
• C. Bhagavatula, J. D. Hwang, D. Downey, R. Le Bras, X. Lu, K. Sakaguchi, <b>S. Swayamdipta</b> , P. West, and Y. Choi	
REV: Information-Theoretic Evaluation of Free-Text Rationales	<i>ACL 2023</i>
• H. Chen, F. Brahman, X. Ren, Y. Ji, Y. Choi, and <b>S. Swayamdipta</b>	
COBRA Frames: Contextual Reasoning about Effects and Harms of Offensive Statements	<i>ACL-Findings 2023</i>
• X. Zhou, H. Zhu, A. Yerukola, T. Davidson, J. D. Hwang, <b>S. Swayamdipta</b> , and M. Sap	

Investigating the Benefits of Free-Form Rationales J. Sun, <b>S. Swayamdipta</b> , J. May, and X. Ma	<i>EMNLP-Findings 2022</i>
NeuroCounterfactuals: Beyond Minimal-Edit Counterfactuals for Richer Data Augmentation P. Howard, G. Singer, V. Lal, Y. Choi, and <b>S. Swayamdipta</b>	<i>EMNLP-Findings 2022</i>
WaNLI: Worker and AI Collaboration for Natural Language Inference Dataset Creation A. Liu, <b>S. Swayamdipta</b> , N. A. Smith, and Y. Choi	<i>EMNLP-Findings 2022</i>
Reframing Human-AI Collaboration for Generating Free-Text Explanations S. Wiegreffe, J. Hessel, <b>S. Swayamdipta</b> , M. Riedl, and Y. Choi	<i>NAACL 2022</i>
Annotators with Attitudes: How Annotator Beliefs And Identities Bias Toxic Language Detection M. Sap, <b>S. Swayamdipta</b> , L. Vianna, X. Zhou, Y. Choi, and N. A. Smith	<i>NAACL 2022</i>
Understanding Dataset Difficulty with V-Usable Information • K. Ethayarajh, Y. Choi, and <b>S. Swayamdipta</b> • <b>Outstanding Paper Award</b>	<i>ICML 2022</i>
MAUVE: Measuring the Gap Between Neural Text and Human Text using Divergence Frontiers • K. Pillutla, <b>S. Swayamdipta</b> , R. Zellers, J. Thickstun, S. Wellecks, Y. Choi, and Z. Harchaoui • <b>Outstanding Paper Award</b>	<i>NeurIPS 2021</i>
Contrastive Explanations for Model Interpretability A. Jacovi, <b>S. Swayamdipta</b> , S. Ravfogel, Y. Elazar, Y. Choi, and Y. Goldberg	<i>EMNLP 2021</i>
On-the-Fly Controlled Text Generation with Experts and Anti-Experts A. Liu, M. Sap, X. Lu, <b>S. Swayamdipta</b> , C. Bhagavatula, N. A. Smith, and Y. Choi	<i>ACL 2021</i>
Challenges in Automated Debiasing for Toxic Language Detection • X. Zhou, M. Sap, <b>S. Swayamdipta</b> , N. A. Smith, and Y. Choi	<i>EACL 2021</i>
Dataset Cartography: Mapping and Diagnosing Datasets with Training Dynamics • <b>S. Swayamdipta</b> , R. Schwartz, N. Lourie, Y. Wang, H. Hajishirzi, N. A. Smith, Y. Choi	<i>EMNLP 2020</i>
Generative Data Augmentation for Commonsense Reasoning • Y. Yang, C. Malaviya, J. Fernandez, <b>S. Swayamdipta</b> , R. LeBras, J. Wang, C. Bhagavatula, Y. Choi, and D. Downey	<i>EMNLP-Findings 2020</i>
Adversarial Filters of Dataset Biases • R. LeBras, <b>S. Swayamdipta</b> , C. Bhagavatula, R. Zellers, M. E. Peters, A. Sabharwal, and Y. Choi	<i>ICML 2020</i>
The Right Tool for the Job: Matching Model and Instance Complexities • R. Schwartz, G. Stanovsky, <b>S. Swayamdipta</b> , J. Dodge, and N. A. Smith	<i>ACL 2020</i>
Don't Stop Pretraining: Adapt Language Models to Domains and Tasks • S. Gururangan, A. Marasović, <b>S. Swayamdipta</b> , K. Lo, I. Beltagy, D. Downey, and N. A. Smith • <b>Best Paper Honorable Mention</b>	<i>ACL 2020</i>
Syntactic Scaffolds for Semantic Structures • <b>S. Swayamdipta</b> , S. Thomson, K. Lee, L. Zettlemoyer, C. Dyer, N. A. Smith.	<i>EMNLP 2018</i>
Learning Joint Semantic Parsers from Disjoint Data • H. Peng, S. Thomson, <b>S. Swayamdipta</b> , N. A. Smith	<i>NAACL 2018</i>
Annotation Artifacts in Natural Language Inference Data • S. Gururangan*, <b>S. Swayamdipta</b> *, O. Levy, R. Schwartz, S. Bowman, and N. A. Smith • * equal contribution	<i>NAACL 2018</i>
Polyglot Semantic Role Labeling • P. Mulcaire, <b>S. Swayamdipta</b> , and N. A. Smith	<i>ACL 2018</i>
Multi-Mention Learning for Reading Comprehension with Neural Cascades • <b>S. Swayamdipta</b> , A. Parikh, T. Kwiatkowski	<i>ICLR 2018</i>
Greedy, Joint Syntactic and Semantic Parsing with Stack LSTMs • <b>S. Swayamdipta</b> , M. Ballesteros, C. Dyer, N. A. Smith	<i>CoNLL 2016</i>
A Dependency Parser for Tweets • L. Kong, N. Schneider, <b>S. Swayamdipta</b> , A. Bhatia, C. Dyer, N. A. Smith	<i>EMNLP 2014</i>

The Pursuit of Power and its Manifestation in Written Dialog

ICSC 2012

- **S. Swayamdipta**, O. Rambow

## WORKSHOP PAPERS

Sample, Align, Synthesize: Graph-Based Response Synthesis with CONGRS

*ScalR@COLM / ER@NeurIPS 2025*

- S. Ghosh, S. S. Warrach, D. Tarsadiya, G. Yauney, **S. Swayamdipta**

Evaluation Under Imperfect Benchmarks and Ratings: A Case Study in Text Simplification

*LLM-Evals@NeurIPS 2025*

- J. Liu, Y. Nam, X. Cui, **S. Swayamdipta**

Uncovering Intervention Opportunities for Suicide Prevention with Language Model Assistants

*EAAMO / GenAI for Health @NeurIPS 2025*

- J. Ranjit, H. J. Cho, C. J. Smerdon, Y. Nam, M. Phung, J. May, J. R. Blosnich, **S. Swayamdipta**

ChEmREF: Evaluating Language Model Readiness for Chemical Emergency Response Assistance

*LLM-Evals@NeurIPS 2025*

- R. Surana, Q. Ye, **S. Swayamdipta**

Sister Help: Data Augmentation for Frame-Semantic Role Labeling

*LAW-DMR @ EMNLP 2021*

- A. Pancholy, **S. Swayamdipta**, M. R. L. Petrucc

Multi-Task Learning for Incremental Parsing using Stack LSTMs

*WiML @ NeurIPS 2016*

- **S. Swayamdipta**, M. Ballesteros, C. Dyer, N. A. Smith

CMU: Arc-Factored, Discriminative Semantic Dependency Parsing

*SemEval 2014*

- S. Thomson, D. Bamman, J. Dodge, **S. Swayamdipta**, N. Schneider, C. Dyer, N. A. Smith

The CMU Machine Translation Systems

*WMT 2014*

- A. Matthews, C. Dyer, A. Lavie, G. Hanneman, W. Ammar, A. Bhatia, **S. Swayamdipta**, E. Schlinger, Y. Tsvetkov

## WORKING PAPERS

Are We Automating the Joy Out of Work? Designing AI to Augment Work, Not Meaning

*Under Review at CHI 2026*

- J. Ranjit, K. Zhou, **S. Swayamdipta**, D. Quercia

Believing without Seeing: Quality Scores for Contextualizing Vision-Language Model Explanations

*Under Review at ARR 2026*

- K. He, T. Srinivasan, B. Joshi, X. Ren, J. Thomason, **S. Swayamdipta**

How Reliable is Language Model Micro-Benchmarking?

*Under Review at ICLR 2026*

- G. Yauney, S. S. Warrach, **S. Swayamdipta**

Every Language Model Has a Forgery-Resistant Signature

*Under Review at ICLR 2026*

- M. Finlayson, X. Ren, **S. Swayamdipta**

Uncovering and Mitigating Covert Dialect Bias in LLMs

*Under Review at FaccT 2026*

- K. Kondapally, C. Smerdon, P. Patel, O. Akoni, J. Torres, J. Ranjit, M. Finlayson, **S. Swayamdipta**

BenchBrowser: Retrieving Evidence of Intent Representation from Evaluation Benchmarks

*Under prep for ICML 2026*

- H. Diddee, G. Yauney, **S. Swayamdipta**, D. Ippolito

Mark My Words: Toward Unforgeable Language Model Signatures via Token Rankings

*Under prep for ICML 2026*

- M. Finlayson, A. Grivas, X. Ren, **S. Swayamdipta**

Representation Collapse in Language Models

*Under prep for ICML 2026*

- A. Kulkarni, A. Balasubramonian, J. Springer, **S. Swayamdipta**

On the Trustworthiness of Generative Foundation Models: Guideline, Assessment, and Perspective

*Under Review at TMLR 2026*

- Y. Huang, C. Gao, S. Wu, H. Wang, X. Wang, Y. Zhou, Y. Wang, J. Ye, J. Shi, and 57 more authors

Political-LLM: Large Language Models in Political Science

*Under Review at TMLR 2026*

- L. Li, J. Li, C. Chen, F. Gui, H. Yang, C. Yu, Z. Wang, J. Cai, J. A. Zhou, and 38 more authors

Shallow Syntax in Deep Water

*arXiv:1908.11047*

- **S. Swayamdipta**, M. Peters, B. Roof, C. Dyer and N. A. Smith

Frame-Semantic Parsing with Softmax-Margin Segmental RNNs and a Syntactic Scaffold

[arXiv:1706.09528](https://arxiv.org/abs/1706.09528)

- **S. Swayamdipta**, S. Thomson, C. Dyer and N. A. Smith

DyNet: The Dynamic Neural Network Toolkit

[arXiv:1701.03980](https://arxiv.org/abs/1701.03980)

- G. Neubig, C. Dyer, Y. Goldberg, A. Matthews, W. Ammar, A. Anastasopoulos, M. Ballesteros, D. Chiang, D. Clothiaux, T. Cohn, K. Duh, M. Faruqui, C. Gan, D. Garrette, Y. Ji, L. Kong, A. Kuncoro, G. Kumar, C. Malaviya, P. Michel, Y. Oda, M. Richardson, N. Saphra, **S. Swayamdipta**, P. Yin

## TUTORIAL AND ORGANIZATION PAPERS

Proceedings of the 3rd Workshop on Deep Learning Approaches for Low-Resource NLP  
(DeepLo 2022)

[NAACL 2022](#)

- C. Cherry, A. Fan, G. Foster, G. Haffari, S. Khadivi, N. Peng, X. Ren, E. Shareghi, **S. Swayamdipta**

Transfer Learning in Natural Language Processing

[NAACL 2019](#)

- S. Ruder, M. E. Peters, **S. Swayamdipta**, T. Wolf

Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP  
(DeepLo 2019)

[EMNLP 2019](#)

- C. Cherry, G. Durrett, G. Foster, R. Haffari, S. Khadivi, N. Peng, X. Ren, **S. Swayamdipta**

Frame Semantics across Languages: Towards a Multilingual FrameNet

[CoLing 2018](#)

- C. F. Baker, M. Ellsworth, M. R. L. Petrucci, **S. Swayamdipta**

## Mentorship

---

### PH.D. STUDENTS

Fall 2024-Present	<b>Xinyue Cui</b> , Ph.D.	
Fall 2024-Present	<b>Muru Zhang</b> , Ph.D.	<i>Co-Advisor: Robin Jia</i>
Fall 2024-Present	<b>Atharva Kulkarni</b> , Ph.D.	
Spring 2024-Present	<b>Brihi Joshi</b> , Ph.D., passed Proposal	<i>Co-Advisor: Xiang Ren</i>
Fall 2023-Present	<b>Matthew Finlayson</b> , Ph.D., passed Quals	<i>Co-Advisor: Xiang Ren</i>
Fall 2022-Present	<b>Jaspreet Ranjit</b> , Ph.D., passed Quals	

### MASTERS STUDENTS

Summer 2025-Present	<b>Harshavardhan Alimi</b> , Masters
Spring 2025-Present	<b>Dhruv Tarsadiya</b> , Masters
Spring 2024-Present	<b>Shahzaib Saqib Warrach</b> , Masters
Spring 2024-Present	<b>Xingjian Dong</b> , Masters

### UNDERGRADS

Fall 2023-Present	<b>Risha Surana</b> , BS / PDP
Fall 2025-Present	<b>Naysa Bhargava</b> , BS
Fall 2025-Present	<b>Thor Christoffersen Hochman</b> , BS
Fall 2025-Present	<b>Brennen Ho</b> , BS
Fall 2025-Present	<b>Mike Gee</b> , BS

### ALUMNI

Fall 2024-Summer 2025	<b>Gregory Yauney</b> , PostDoc	
Spring 2024-2025	<b>Ryan Wang</b> , BS	<i>Now at UC Berkeley</i>
Spring 2024-2025	<b>Joseph Liu</b> , BS	<i>Now at CMU</i>
Spring 2024-2025	<b>Keyu He</b> , BS	<i>Now at CMU</i>
Summer 2023-2025	<b>Aryan Gulati</b> , BS	<i>Now at Bloomberg</i>
Summer 2023-2024	<b>Catherine He</b> , BS	
Spring 2023-Spring 2024	<b>Yoonsoo Nam</b> , Masters	<i>Now at RenderWolf</i>
Summer 2023	<b>Ruyuan Zuo</b> , Masters	<i>Now at Google</i>
Spring 2022	<b>Hanjie Chen / Junlin Wang / Abdallah Bashir</b> , Ph.D., UVA / Masters, UCI / Masters, Saarland Universit	<i>AI2 Intern</i>
Fall 2021	<b>Jillian Fisher / Liwei Jiang</b> , Ph.D., UW	
Summer 2021	<b>Kawin Ethayarajh</b> , Ph.D., Stanford University	<i>AI2 Intern</i>
Summer 2021	<b>Sarah Wiegreffe</b> , Ph.D., George Institute of Technology	<i>AI2 Intern</i>
Summer 2021	<b>Ximing Liu</b> , BS, UW CSE	<i>AI2 Intern</i>
Fall 2020-2023	<b>Alisa Liu</b> , Ph.D., UW CSE	
Fall 2020	<b>Alon Jacovi</b> , Ph.D., Bar Ilan University	<i>AI2 Intern</i>
Summer 2020-1	<b>Jenny Liang</b> , BS, UW CSE	<i>AI2 Intern</i>
Spring 2020-2023	<b>Xuhui Zhou</b> , CLMS, UW Linguistics	
Fall 2019 - Spring 2020	<b>Yiben Yang</b> , Ph.D., Northwestern University	
Fall 2019 - Spring 2020	<b>Chaitanya Malaviya</b> , PYI, Allen Institute for AI	

## Professional Service

---

### EXTERNAL SERVICE

#### Senior Area Chair

- EMNLP 2025: Generalizability and Transfer
- EMNLP 2024: Machine Learning in NLP
- ACL 2024: Sentence-level Semantics

#### Area Chair

- ICLR 2026, COLM 2025, ICLR 2025, COLM 2024, ICLR 2024, EMNLP 2023: Machine Learning for NLP, ACL 2023: Interpretability, EMNLP 2022: Language Models, EMNLP 2021: Machine Learning for NLP, NAACL 2021: Sentence-Level Semantics, EACL 2021: Sentence-Level Semantics, ACL 2020: Semantics (Long)

#### Reviewer

- **Conferences:** ACL (ARR): 2015-2022, 2025; NAACL : 2015-2022; EMNLP : 2015-2022; NeurIPS : 2018-2022; ICML : 2015, 2019, 2020, 2023; EACL : 2017; AAAI : 2017-2020; CoNLL : 2017-2018, 2020
- **Journals:** Transactions of ACL : 2020 - 2022; Computational Linguistics : 2019 - 2022; Journal of AI Research : 2019

#### Co-organizer

- |  |      |
|--|------|
| • EMNLP 2025: Workshop on Uncertainty in NLP 2   | 2025 |
| • ICML 2024: 2nd Workshop on High-dimensional Learning Dynamics (HiLD): The Emergence of Structure and Reasoning | 2024 |
| • EACL 2024: Workshop on Uncertainty in NLP  | 2024 |
| • NAACL 2022: Workshop on Deep Learning in Low Resource NLP (DeepLo)   | 2022 |
| • EMNLP 2019: Workshop on Deep Learning in Low Resource NLP (DeepLo)   | 2019 |
| • West Coast NLP Workshop  | 2018 |

## Teaching

---

### CLASSES

Fall 2025	<b>CSCI 444: NLP</b> , Students: 14	<i>USC</i>
Fall 2024	<b>CSCI 544: Applied NLP</b> , Students: 260	<i>USC</i>
Spring 2024	<b>CSCI 499: Language Models in NLP</b> , Students: 28	<i>USC</i>
Fall 2023	<b>CSCI 499: Language Models in NLP</b> , Students: 30	<i>USC</i>
Fall 2022	<b>CSCI 699: Data-Centric NLP</b> , Students: 29	<i>USC</i>

## Invited Panels

---

Oct 18, 2024 **Alfred P. Sloan Foundation Science Seminar**,  
Jul 13, 2023 **Limitations of Large Language Models**, Rep4NLP Workshop  
Jun 29, 2023 **Generative AI and Education**, Roundtable  
Apr 5, 2023 **Alimaginings**, Polymathic Pizza  
Jul 14, 2022 **Adversarial Data Augmentations**, DADC Workshop  
Jun 17, 2022 **Role of LLMs**, Responsible AI Symposium

*USC School of Cinematic Arts*  
*ACL Toronto*  
*ASEE 2023*  
*USC Sidney Harman Academy*  
*for Polymathic Studies*  
*NAACL Seattle*  
*AILA*

## Invited Talks (Since 2020) ---

Oct 16, 2025	<b>Cornell Tech Learning Machine Seminar Series (LMSS)</b> , Simple Learning Recipes for Safe and Accountable Language Models	<b>Cornell Tech</b>
Sep 16, 2025	<b>USC School of Social Work</b> , The heaviest data: Can AI alleviate the weight of suicide mortality research?	<b>USC</b>
Apr 10, 2025	<b>Bloomberg CTO Data Science Speaker Series</b> , Rethinking Language Model Evaluation	<b>Bloomberg</b>
Apr 4, 2025	<b>Simons Institute: The Future of LMs and Transformers</b> , The Future of LMs: A Perspective on Evaluation	<b>UC Berkeley</b>
Feb 14, 2025	<b>Simons Research Fellow Introductions</b> , Rethinking the Evaluation of Large Language Models	<b>UC Berkeley</b>
Feb 12, 2025	<b>UC Berkeley AI and Society Mixer</b> , Building LLM Assistants for Social Workers	<b>UC Berkeley</b>
Nov 07, 2024	<b>NYU NLP/Text-as-Data Talk Series</b> , Ensuring Safety and Accountability in LLMs, Pre- and Post-Training	<b>New York University</b>
Sep 23, 2024	<b>Amazon AWS Bedrock AI Seminar</b> , Ensuring Safety and Accountability in LLMs, Pre- and Post-Training	<b>Amazon AWS, Virtual</b>
May 3, 2024	<b>ISI AI Seminar</b> , Understanding LLMs through their Generative Behavior, Successes and Shortcomings	<b>USC ISI</b>
Mar 26, 2024	<b>NSF-Open Source Generative AI Workshop</b> , Towards (Closed-Source) LLM Accountability via Logit Signatures	<b>Cornell Tech</b>
Mar 13, 2024	<b>Utah Data Science Seminar</b> , Understanding LLMs through their Generative Behavior, Successes and Shortcomings	<b>University of Utah</b>
Feb 13, 2024	<b>Intel Rising Stars Talk</b> , LLM Inference for Scaling Data Creation	<b>Intel Labs</b>
Feb 8, 2024	<b>Cambridge LTL Seminar</b> , Understanding LLMs via their Generative Successes and Shortcomings	<b>University of Cambridge</b>
Dec 15, 2023	<b>ATTRIB Workshop Keynote</b> , Understanding LLMs via their Generative Successes and Shortcomings	<b>NeurIPS 2023</b>
Nov 16, 2023	<b>UCLA Communications Symposium</b> , Understanding Online Discourse through Social Context and Structured Pragmatics	<b>UCLA</b>
Oct 3, 2023	<b>DataComp Workshop</b> , Understanding Data with $\mathcal{V}$ -Information	<b>ICCV 2023</b>
Sep 22, 2023	<b>LA County of Education CS Speaker Series</b> , The Role of Language Models in NLP	<b>Los Angeles CoE</b>
Aug 24, 2023	<b>CMU LTI Student Research Symposium</b> , Understanding Datasets and Explanations through $\mathcal{V}$ -Information	<b>Alumni Keynote: CMU, Pittsburgh</b>
Jul 13, 2023	<b>Rep4NLP Workshop</b> , Contextualizing Representations in Varied Annotator Perspectives	<b>ACL Toronto</b>
Jun 14, 2023	<b>Google Responsible Machine Learning</b> , Contextualizing Data in Varied Annotator Perspectives	<b>Google Brain</b>
Mar 14, 2023	<b>Cohere for AI Reading Group</b> , Understanding Dataset Difficulty with V-Usable Information	<b>Cohere for AI</b>
Feb 28, 2023	<b>Spotify Research Seminar</b> , Designing Controls and Filters for Dataset Generation	<b>Spotify Research Labs</b>
Feb 03, 2023	<b>Amazon Data-Centric AI Seminar</b> , What's in a Dataset? Interpreting datasets to enable better data creation	<b>Amazon</b>
Feb 01, 2023	<b>USC CAIS++ Seminar</b> , Contextualizing Bias in Hate Speech Detection	<b>USC</b>
Nov 18, 2022	<b>SoCal NLP Symposium</b> , Generating Datasets for Robust Generalization	<b>UC Santa Barbara</b>
Nov 09, 2022	<b>USC Center for AI in Society (CAIS) Seminar</b> , Contextualizing Bias in Hate Speech Detection through Annotator Perspectives	<b>USC</b>
May 23, 2022	<b>ACL Spotlight Talks for Young Rising Stars</b> , The Devil's in the Data: Mapping and Generating Datasets for Robust Generalization	<b>Dublin, Ireland</b>
May 13, 2022	<b>UC Irvine CS Seminar</b> , Mapping and Generating Datasets for Robust Generalization	<b>UC Irvine</b>
Mar 16, 2022	<b>UC Santa Cruz IFDS Ethics Seminar</b> , Rethinking Dataset Construction: : The Role of Generative Modeling and Annotator Perspectives	<b>UC Santa Cruz</b>
Oct 20, 2021	<b>Oracle ML Talks</b> , What's in your Data? Mapping Datasets and Exploring Data Usability	<b>Oracle</b>
Feb 24, 2020	<b>Georgetown NERT Seminar</b> , Addressing Biases for Robust, Generalizable AI	<b>Georgetown University</b>
Feb 12, 2020	<b>Georgia Tech NLP Seminar</b> , Addressing Biases for Robust, Generalizable AI	<b>Georgia Tech</b>
Nov 02, 2020	<b>Microsoft E+D Product Leads</b> , Responsible AI: Addressing Biases in Datasets and Models	<b>Microsoft</b>