

# DAVID HEINEMAN

[davidheineman.com](https://davidheineman.com)  
[davidh@allenai.org](mailto:davidh@allenai.org)  
[linkedin.com/in/david-heineman](https://linkedin.com/in/david-heineman)

## EDUCATION

### Georgia Institute of Technology

2020 - 2024

B.S. in Computer Science. GPA: 3.9, Highest Honors.

**Coursework:** Low-level Computer Architecture, Information Security, Deep Learning, Database Implementation

## INDUSTRY EXPERIENCE

### Allen Institute for AI | Predoctoral Young Investigator, AllenNLP

2024 - Present

- Ongoing work with the AllenNLP team on LLM training and evaluation.

### Amazon Web Services | Software Engineering Intern, EC2

Summer 2023

- Built prototype LLM application to query AWS documentation, account knowledge to solve root causes of CloudWatch alarms
- Used parameter efficient training (i.e., QLoRA) to fine-tune Falcon, LLaMA to align with developer-friendly responses
- Created first internal deployment of fine-tuned open-source 40B LLM on a Sagemaker endpoint, contributed to `bitsandbytes` to fix a GPU compatibility error when performing quantization in the G4 class of EC2 instances.
- Proposed new method of chaining T5 queries in parallel to explore account details in a graph-like structure which improved generation pre-processing time from 30s to 5s
- Built an internal manual inspection scheme to evaluate performance and automatic identification system for harmful responses

### Amazon Web Services | Software Engineering Intern, CloudWatch

Summer 2022

- Developed a new feature for CloudWatch Application Insights to monitor processes running on EC2 instances
- Worked with EC2 Windows experts to identify breakpoints for customers' Microsoft SQL and .NET workloads, predict these common errors using CloudWatch metrics and provide insights to reduce problem resolution time
- On-boarded existing AWS customers like Koch Industries, Moody's to monitor processes on critical Windows applications

### Patientco | Software Engineering Intern, ML

Summer 2021

- Created a large-scale data ingestion pipeline for training ML models to predict the likelihood of bill payments
- Used AWS Sagemaker to productionalize models to re-train and label 500,000 unique bills daily (5% of U.S. healthcare bills)
- Developed a novel sequence-based approach for detecting anomalies in insurance charges, used by providers like Piedmont to identify and fix errors in adjudication for patient bills

## RESEARCH EXPERIENCE

### Georgia Tech NLPx Lab | Undergraduate Research Assistant, Advisor: Prof. Wei Xu

2021 - Present

- **Controllable Text Diffusion:** Managed large scale training of text diffusion models on GPU clusters to replicate controllable text diffusion experiments. Developed a prototype controllable diffusion method for sequence-level control at decoding time by using an automatic metric as part of the control signal
- **Minimum Bayes Risk Prompting:** Created flexible decoding library for swapping out datasets, metrics and managing experiments, including custom multi-GPU distributed inference with open-source LLMs
- **Fine-grained Text Generation Analysis:** Performed the first large scale fine-grained text evaluation to exhaustively consider every linguistic transformation performed by the LLM. Built an open-source library for analysis, modeling and agreement calculations for fine-grained text generation, including a novel span-based evaluation metric. Contributed to the `lens-metric` library to load automatic simplification metrics directly from Huggingface
- **Rank & Rate Evaluation:** Built a multi-stage highlighting, ranking and rating interface for text generation, orchestrated large-scale in-house human evaluation across 26 text simplification systems used to develop state-of-the-art simplification metrics
- **Real-time Text Simplification:** Implemented efficient lookup algorithms for large n-gram datasets using KenLM to allow for real-time model inference.
- Led discussions of bi-weekly undergraduate reading group (Spring 2021, Fall 2021)

**Service:** Reviewer at TSAR 2022 (co-located at EMNLP 2022), panel speaker at TSAR 2024 (co-located at EMNLP 2024)

## PUBLICATIONS

### Signal and Noise: A Framework for Reducing Uncertainty in Language Model Evaluation

David Heineman, Valentin Hofmann, Ian Magnusson, Yuling Gu, Noah A. Smith, Hannaneh Hajishirzi, Kyle Lo, Jesse Dodge.

*under review at NeurIPS, 2025*

### Fluid Language Model Benchmarking

Valentin Hofmann, David Heineman, Ian Magnusson, Kyle Lo, Jesse Dodge, Maarten Sap, Pang Wei Koh, Chun Wang, Hannaneh Hajishirzi, Noah A. Smith.

*under review at COLM, 2025*

## Establishing Task Scaling Laws via Compute-Efficient Model Ladders

Akshita Bhagia\*, Jiacheng Liu\*, Alexander Wettig, **David Heineman**, Oyvind Tafjord, Ananya Harsh Jha, Luca Soldaini, Noah A. Smith, Dirk Groeneveld, Pang Wei Koh, Jesse Dodge, Hannaneh Hajishirzi.

under review at COLM, 2025

## Evaluating LLMs on Chinese Idiom Translation

Cai Yang, Yao Dou, **David Heineman**, Xiaofeng Wu, Wei Xu.

under review at COLM, 2025

## 2 OLMo 2 Furious

Evan Pete Walsh\*, Luca Soldaini\*, Dirk Groeneveld\*, Kyle Lo\*, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, ..., **David Heineman**, ..., Ali Farhadi, Noah A. Smith, Hannaneh Hajishirzi.

under review at COLM, 2025

## DataDecide: How to Predict Best Pretraining Data with Small Experiments

Ian Magnusson\*, Nguyen Tai\*, Ben Bogin\*, **David Heineman**, Jena D. Hwang, Luca Soldaini, Akshita Bhagia, Jiacheng Liu, Dirk Groeneveld, Oyvind Tafjord, Noah A. Smith, Pang Wei Koh, Jesse Dodge.

ICML, 2025

## Improving Minimum Bayes Risk Decoding with Multi-Prompt

**David Heineman**, Yao Dou, Wei Xu.

EMNLP, 2024

## Towards a Path Dependent Account of Category Fluency

**David Heineman**, Reba Koenen, Sashank Varma.

CogSci, 2024

## Thresh: A Unified, Customizable and Deployable Platform for Fine-Grained Text Evaluation

**David Heineman**, Yao Dou, Wei Xu.

EMNLP: System Demonstrations, 2023

## Dancing Between Success and Failure: Edit-level Simplification Evaluation using SALSA

**David Heineman**, Yao Dou, Mounica Maddela, Wei Xu.

EMNLP, 2023

## LENS: A Learnable Evaluation Metric for Text Simplification

Mounica Maddela\*, Yao Dou\*, **David Heineman**, Wei Xu.

ACL, 2023

## TEACHING

---

(CS 4650) **Natural Language Processing**

Fall 2023

Teaching Assistant (TA) | Instructor: Alan Ritter

(CS 3510) **Design & Analysis of Algorithms**

Fall 2022

Teaching Assistant (TA) | Instructor: Gerandy Brito & Dana Randall

(CS 3510) **Design & Analysis of Algorithms**

Fall 2021

Teaching Assistant (TA) | Instructor: Jacob Abernathy

## OPEN SOURCE CONTRIBUTIONS

---

**Thresh** | [thresh.tools](https://thresh.tools) | [github.com/davidheineman/thresh](https://github.com/davidheineman/thresh)

2023

- Created a highly customizable in-browser span alignment and annotation builder, with questions rendered in a recursive tree structure. Thresh currently supports 13 languages, paragraph annotation and a variety of question and span-selection options
- Implemented 12 well known fine-grained annotation interfaces into standardized YML templates and easy-to-load JSON datasets
- Built an integrated Python library to load arbitrary data collected with Thresh and manage large-scale annotation jobs

**Hashtag Segmentation API** | [github.com/davidheineman/hashtag-master-web-demo](https://github.com/davidheineman/hashtag-master-web-demo)

2022

- Deployed real-time inference for neural text segmentation, allowing projects to query an API for integration in future work
- Created interactive Vue app to demonstrate the hashtag segmentation task and new methodology

**Individual Contributions:** bitsandbytes, vLLM, OLMo-core

## AWARDS

---

- Georgia Tech College of Computing (CoC) Outstanding Undergraduate Research Award (1 of 3000+ CS undergraduates)
- 2nd place in AGI House Open Source AI Hackathon (\$1000 award from Prime Intellect)
- President's Undergraduate Research Award (PURA Research Grant)
- Top UserStyles Contributor (20,000 Installs)