

Shi Feng

shifeng@cs.umd.edu • [Google Scholar](#)

4108 Iribe Center • 8125 Paint Branch Drive • College Park, MD 20742

Education

| | |
|--|----------------------------------|
| University of Maryland PhD candidate, Computer Science Advised by Prof. Jordan Boyd-Graber, member of CLIP lab. | College Park, Maryland 2016 – |
| New York University Visiting Student, Center of Data Science Hosted by Prof. He He. | New York City Fall 2019 |
| Shanghai Jiao Tong University B.S. in Computer Science Member of the ACM Honor Class. | Shanghai, China 2012 – 2016 |

Research Interest

I'm interested in interpretability, AI alignment, and in general trustworthy machine learning for the people, with a focus on the language domain.

Publications

| | |
|--|------|
| Pathologies of Neural Models Make Interpretation Difficult Empirical Methods in Natural Language Processing (Long paper, oral presentation) Shi Feng, Eric Wallace, Alvin Grissom II, Mohit Iyyer, Pedro Rodriguez, Jordan Boyd-Graber | 2018 |
| What can AI do for me: Evaluating Machine Learning Interpretations in Cooperative Play ACM Intelligent User Interface (Long paper, oral presentation) Shi Feng, Jordan Boyd-Graber | 2019 |
| Universal Adversarial Triggers for Attacking and Analyzing NLP Empirical Methods in Natural Language Processing (Long paper, oral presentation) Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, Sameer Singh | 2019 |
| Customizing Triggers with Concealed Data Poisoning In submission Eric Wallace*, Tony Z. Zhao*, Shi Feng, Sameer Singh | 2020 |
| Misleading Failures of Partial-input Baselines The 57th annual meeting of the Association for Computational Linguistics (Short paper) Shi Feng, Eric Wallace, Jordan Boyd-Graber | 2019 |
| Understanding Impacts of High-Order Loss Approximations and Features in Deep Learning Interpretation International Conference on Machine Learning (Long paper) Sahil Singla, Eric Wallace, Shi Feng, Soheil Feizi | 2019 |
| Trick Me If You Can: Human-in-the-loop Generation of Adversarial Examples for Question Answering Transactions of the Association for Computational Linguistics (Presented at ACL) Eric Wallace, Pedro Rodriguez, Shi Feng, Jordan Boyd-Graber | 2019 |
| How Pre-trained Word Representations Capture Commonsense Physical Comparisons Commonsense Inference in NLP workshop Pranav Goel, Shi Feng, Jordan Boyd-Graber | 2019 |
| Quizbowl: The Case for Incremental Question Answering In submission Pedro Rodriguez, Shi Feng, Mohit Iyyer, He He, Jordan Boyd-Graber | 2019 |

| | |
|--|------|
| Interpreting Neural Networks with Nearest Neighbors BlackboxNLP Workshop at EMNLP Eric Wallace*, Shi Feng* , Jordan Boyd-Graber | 2018 |
| The UMD Neural Machine Translation Systems at WMT17 Bandit Learning Task The Second Conference on Machine Translation Amr Sharaf, Shi Feng , Khanh Nguyen, Kianté Brantley, Hal Daumé III | 2017 |
| Improving Attention Modeling with Implicit Distortion and Fertility for Machine Translation International Conference on Computational Linguistics (Long paper) Shi Feng , Shujie Liu, Nan Yang, Mu Li, Ming Zhou, Kenny Q. Zhu | 2016 |
| Knowledge-Based Semantic Embedding for Machine Translation The 54th annual meeting of the Association for Computational Linguistics (Long paper) Chen Shi, Shujie Liu, Shuo Ren, Shi Feng , Mu Li, Ming Zhou, Xu Sun, Huofeng Wang | 2016 |

On-going Projects

KAR³L: Spaced repetition meets representation learning karl.qanta.org

We want to see if representation learning can help us improve human memorization, and more generally, human learning. We implemented this flashcard app as a testbed for this idea. In traditional spaced repetition learning systems, all flashcards are treated as equal, so are all the users. This over-simplified model ignores useful signals that can help us infer the state of the user's memory: if the user correctly answers a question about Mozart, this should tell us something about the his/her knowledge about classical music, and in turn the probability of correctly answering a question about Beethoven. Our proposed algorithm uses representation learning to exploit connections like this, and is currently deployed on this interface. We have a paper in preparation for this project.

Play With QANTA: Human-computer Cooperative QA play.qanta.org

We want to see if post-hoc explanations improves human-AI cooperation, and built this online interface is a testbed for this idea. Each human player on the interface is assisted with a human-level question answering AI, where the AI communicated its predictions via several post-hoc explanations. Our **IUI'19** paper is based on experiments conducted using this interface. We have an on-going work that investigates whether we could adapt to each user and intelligently select which explanation to show in order to maximize the human-AI team performance.

QANTA: Human-level Quizbowl System github.com/pinafore/qb

At HSNCT'17 we beat *top* human players for the first time ([video](#)). I'm mainly responsible for the *buzzer* of QANTA, which controls when to buzz and when to wait. The buzzer was trained with reinforcement learning using game history collected from Protobowl. This RL buzzer was first introduced to the system for HSNCT'17 and turned out to be crucial to the victory against human.

Talks

| | |
|------------------------|-------------|
| NLP Highlights Podcast | Apr 25 2019 |
| Invited talk at UPenn | Mar 25 2019 |
| Invited talk at UCSD | Mar 19 2019 |
| Invited talk at UCI | Mar 18 2019 |

Awards and Service

Best reviewer award, EMNLP'18, NeurIPS'20

Reviewer: EMNLP'18'19'20, ACL'19'20, AAAI'20, CoNLL'20, NeurIPS'20, ICLR'21, NAACL'21

Work Experience

| | |
|--|-----------------|
| Salesforce Research , <i>Research Intern</i> • Advised by Bryan McCann | 2020.6 – 2020.8 |
|--|-----------------|

- Pretrain poisoning. We show that malicious unlabeled data during pretraining can lead to biases and backdoors in the downstream model based on the pretrained representations. In particular, we show that injecting a few thousand sentences to GPT-2's unlabeled training set exacerbates the gender bias of a sentiment classifier based on GPT-2 by 10% absolute difference. The security implications of this vulnerability is immense: unlabeled data for pretraining is almost always collected from the web with very minimal data cleaning. Poisoning against the unlabeled data is thus easy to carry out by an adversary without a lot of resources. We have a paper in preparation for this work.

Microsoft Research, *Research Intern*

2018.6 – 2018.8

- Health AI team
- Domain adaptation for machine translation. A boosting approach to safely select in-domain data to adapt a general translation system to the medical domain.

Microsoft Research Asia, *Research Intern*

2015.8 – 2016.2

- Natural Language Computing Group
- Built the first neural machine translation system with Theano for NLC group.
- Improved the attention mechanism, results published at COLING'16.
- Experimented sequence-to-sequence for many other tasks, including pos tagging, parsing, and Chinese couplet completion ([link](#)).