

# PERSONAL STATEMENT ON MY RESEARCH: EVALUATING AND ENABLING HUMAN-AI COLLABORATION

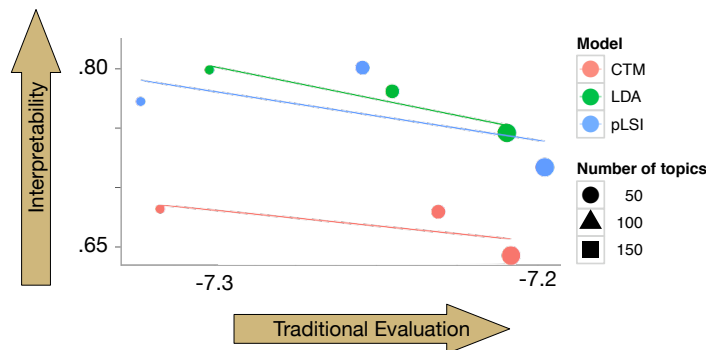
JORDAN BOYD-GRABER, UNIVERSITY OF MARYLAND

Artificial intelligence<sup>1</sup> (AI) is ubiquitous: detecting spam e-mails, flagging fraudulent purchases, and providing the next movie in a Netflix binge. But they do not exist in a vacuum: as Shneiderman [34] argues, AI must exist alongside humans. My goal is to create metrics to measure whether AI methods make sense to users, helping users craft examples to advance AI, and applying AI to illuminate complex social science applications.

## 1. EVALUATING INTERPRETABILITY

My journey with evaluating interpretability began over ten years ago with topic models. Topic models are sold as a tool for understanding large data collections: lawyers scouring Nordstream e-mails for a smoking gun, journalists making sense of Wikileaks, or humanists characterizing the oeuvre of Lope de Vega. But topic models' proponents never asked what those lawyers, journalists, or humanists needed. Instead, they optimized *held-out likelihood*.

When my colleagues and I developed the *interpretability* measure to assess whether topic models' users understood their outputs, interpretability and held-out likelihood were negatively correlated [6]! The topic modeling community (including me) had fetishized complexity at the expense of usability... and topic modeling is not alone.



Since this humbling discovery, I've built topic models that are a collaboration between humans and computers. The computer starts by proposing an organization of the data. The user separates confusing clusters or joins similar clusters together [18], an improvement over the "take it or leave it" philosophy of most machine learning algorithms.

Focusing on collaboration also requires algorithms that are low latency (not just high throughput). We extended the geometric interpretations of admixture models developed by Arora et al. [1] to multi-anchor topics [25] and multi-lingual topics [41]. These are much faster than traditional probabilistic topic models—they can handle millions of documents in seconds—but they are less well understood theoretically and less used in practice. Thus, we also developed better understanding of the projections of multilingual representations via graph theory [12] and the convergence of alternating projections [42].

After we proposed our "reading tea leaves" evaluation, it's heartening that Lau et al. [24] and their "machine reading tea leaves" (which correlate with our human measures) became a standard topic model evaluation: in a survey of forty recent topic modeling papers, **all but four** use a form of their coherence evaluation. However, as we argue in Hoyle et al. [17], you cannot just use this evaluation forever and forget about humans. In that same survey, **none** of those papers do a human evaluation. As topic models evolve (e.g., incorporating neural components), you need to validate that these automatic metrics still correlate with whether it is useful for a human-computer collaboration.

*Date:* Updated May 2023.

<sup>1</sup>I take a broad interpretation of AI; some of my examples might be better characterized machine learning. But rather than distracting boundary policing, I will embrace the general term but will be specific in describing particular tools/models.

## 2. TEAMING AS AN EVALUATION

Within the HCI community, we have argued for the foundations of what should go into human–computer collaborations: computers that incorporate users’ suggestions [23]; explanations with accountability [36]; and stable explanations [37].

In addition to these human-centered understanding of users’ needs and desires, we’ve developed machine learning approaches to measure how well users complete a task. For example, for a question answering task, we measured how much the accuracy of the human–computer *team* increases with different explanations and found that explanations help all users but that novices are easily overwhelmed [10]. In follow-on work, we learned how to explicitly optimize explanations for individual users [11].

## 3. CONNECTING WITH SOCIAL SCIENCE: PEDAGOGY, FRAMING, AND DECEPTION

The reverse of cooperation is human competition; it also has much to teach computers. I’ve increasingly looked at language-based games whose clear goals and intrinsic fun speed research progress. For example, in the board game *Diplomacy*, users chat with each other while marshaling armies for world conquest. Alliances are fluid: friends are betrayed and enemies embraced as the game develops. However, users’ conversations let us predict when friendships break.

Thus, we argued that *Diplomacy* would be an exciting testbed for natural language processing, and our 2015 paper is—to the best of our knowledge—the first NLP research on *Diplomacy*. Before a betrayal, betrayers write ostensibly friendly messages and become more polite, stop talking about the future, and change how *much* they write [31]. In follow-on work, we developed a dataset that predict both when users lie to each other and when recipients of lies detect deception [32]. *Diplomacy* may be a nerdy game, but it is a fruitful testbed to teach computers to understand messy, emotional human interactions. We are continuing to look into these questions with researchers from across the nation in a new DARPA program: SHADE, which focuses on *Diplomacy* as a testbed for understanding deception.

Recently, the use of NLP methods in the game of *Diplomacy* has been the subject of highly-publicized papers by DeepMind in Nature Communications [22] and Meta in Science [3]. The Meta paper, like our 2020 paper, used a classifier to detect deceptive statements. The DeepMind paper built a game theoretic understanding of when betrayal should happen, building on our descriptive investigation of deception in human games.

A game with higher stakes is politics. However, just like *Diplomacy*, the words that people use reveal their underlying goals; computational methods can help expose the “moves” political players can use. With collaborators in political science, we’ve built models that: show when politicians in debates strategically change the topic to influence others [27, 29]; frame topics to reflect political leanings [28]; use subtle linguistic phrasing to express their political leaning [20]; or create political subgroups with larger political movements [30].

Because political discourse is built on a common set of commonly accepted facts, we have focused on developing fact checking: datasets for general knowledge fact checking [8] and climate change fact checking [7]. However, because fact checking is part of an information arms race, we need to build these examples as part of a human-in-the-loop adversarial process, which I’m exploring in an ongoing collaboration with journalism professor Naeemul Hassan that extends my question answering work, which I talk about next.

## 4. HUMAN-IN-THE-LOOP ADVERSARIAL EXAMPLES

One of the most fun aspects of my research has been building trivia-playing robots [5, 19, 21]; beyond research papers, our system has faced off against former Jeopardy champions in front of hundreds high school students<sup>2</sup> and against researchers at NeurIPS 2015 (which won the best demonstration award). But after defeating some of the smartest trivia players, did I actually believe that our system was better at question answering? No!

Adversarial examples first came out of the vision community: add a small epsilon to an example and suddenly a object detector calls a turtle a gun [2].<sup>3</sup> While others have attempted to create adversarial examples for language using paraphrasing, it’s hard to know if the changes are perceptually negligible (“who wrote the invisible man”—a question with the answer H.G. Wells—is fundamentally different from “who wrote the man you can’t see”—an ill-formed questions—as is “who wrote the book invisible man”—a question with the answer Ralph Ellison) and it’s hard to “add epsilon” to a discrete word.

---

<sup>2</sup><https://www.youtube.com/watch?v=LqsUaprYMOw>

<sup>3</sup>Point of personal pride: I mentored Kevin on another research project [15], but I myself had nothing to do with this later adversarial work.

Consistent with the theme of my research, my NSF CAREER grant added a *human in the loop* to generate novel adversarial language examples that can provide new training examples to make AI more robust and to expose what AI cannot (yet) do. With Eric Wallace, an undergraduate student, we built a system that could help an expert trivia question writer to stump a computer: as the author writes the question, it shows the author what the system is thinking [38]. And it worked, even generalizing across models [39] (an example written with an IR model still stumps a neural model). After we introduced human-in-the-loop adversarial example generation, Meta/Facebook adopted this framework with gusto [4] in their Dynabench framework, the Dynamic Adversarial Data Collection workshop and call for proposals (which I’m grateful is funding our continuing research in this area).

## 5. BUT WAIT, THERE’S MORE!

Many of our best-cited papers are “traditional” papers that do better on some task:

- We developed deep averaging networks [21, DAN], a simple model still used in the transformer age [40].
- In question answering, we proposed new evaluation mechanisms for knowing if an answer is correct [35] and improved information retrieval to answer complicated questions [9, 33, 43].
- We also introduced reinforcement learning to *simultaneous machine interpretation* [14], a language-based task that requires significant human intuition, insight, and—for those who want to become interpreters—training.<sup>4</sup> We learned tricks from professional human interpreters—passivizing sentences and guessing the verb—to translate sentences sooner [16], letting speakers and algorithms cooperate together and enabling more natural cross-cultural communication. We also use reinforcement learning to learn machine translation feedback from noisy supervision such as star ratings on a webpage [26].

This work doesn’t *yet* fit nicely into the human–computer collaboration narrative, but these more complex tasks are part of my broader vision for where my research will go: state-of-the-art models built to support human decisions, not replace them. And that requires the low-latency models built to react to input “like a human” described above.

## 6. FUTURE WORK

I hope that I’ve convinced you that to get more effective AI, we need to measure how well it works (and plays) *with* humans. This is a problem that requires solutions in the form of new data, algorithms, and evaluations, and I hope to advance in all three areas.

**Data.** Existing datasets are not diverse and do not reflect the kinds of interactions people from diverse backgrounds have with AI systems. In question answering, Google’s Natural Questions, SQuAD, and other datasets contain entities that are overwhelmingly male and either American or British [13]. In newly funded NIH research, we’re working with answering questions specifically in complicated, code-switched environments with less educated users trying to navigate the healthcare landscape. These linguistically, acoustically, and socio-economically diverse samples will help provide more realistic data for question answering applications, replete with false presuppositions, ambiguity, and shifting information needs. We are also working with NIST to develop datasets that better reflect a specific *cultural* context: detecting when questions are answered differently in Ghana than in the us or topics that only people in Bhutan care about (and would be neglected by us-centric datasets).

**Algorithms.** In research recently funded by Meta, we’re capitalizing on the realization that the best way to get good datasets that reflect the strengths and weaknesses of current AI is to have users in the loop to generate adversarial data. While our previous adversarial datasets focused on a single system—find a question *this* neural system cannot answer—we are expanding our search to *learn* how to help users best engineer adversarial examples, a natural extension of our work using bandit algorithms to improve explanations. While our goal in that work was to get the human–computer team to answer a question correctly, we flip the script to create explanations that help the user craft examples the computer cannot answer.

However, those algorithms are for *supervised* tasks where there is a known answer. For *unsupervised* tasks like understanding the themes in a document collection, we need algorithms that are not just transparent but that are *fast*: the primary goal is exploration of the dataset, and the user needs to be able to quickly provide feedback to the algorithm. The algorithm then needs to update quickly given that feedback.

To develop more robust algorithms, we are building training approaches that—as a conversation develops—explicitly build representations that capture the theory of mind of conversation participants and the current “question under discussion”. We plan to apply this both in our Diplomacy game that captures strategic negotiations (detecting when a negotiation partner is angling for a particular resolution implicitly or when they are negotiation

---

<sup>4</sup>This framework—using reinforcement learning to capture human strategies—was featured in Liang Huang’s ACL keynote.

in bad faith) and in the information seeking question answering context (where a question can be posed imperfectly, revealing a false assumption or a lack of background information).

While eventually we will need to connect to expensive, difficult to compute neural models, many user updates can be approximated by fast spectral or probabilistic algorithms. We will design fast, browser-based Javascript approximations of these complicated neural models, to allow users to quickly interact with the models, get a result, and continue before reconciling the solution later.

**Evaluation.** Where computers have strengths that go beyond what a human can do, we need to know if it would be wise to build evaluations that can we will build interactive systems that help users come to a correct answer through a process they trust, measure that process, and optimize to encourage this cooperation. This requires basic engineering—ensuring all of the components of a system are efficient and low latency—and user modeling, as we cannot assume that every user will have the same knowledge and capabilities.

As these systems become more capable and usable, we can no longer assume that our model of the user should remain static: the user will learn and adapt to the system. This makes modeling more complicated, but it also allows for employing these models in educational settings through examples ordered in a curriculum: expanding the frontier of what the user knows, reinforcing weaker knowledge, and using strategies to both educate and explain information from the AI.

But interactions with individual people are not how AI will be a part of 21<sup>st</sup> society: it will be interactions with *populations*. Thus, we need to have models and systems that capture population-level interactions. Helping detect misinformation online, working collaboratively with authors to craft effective counter-measures, and to propagate that within a social network. This builds on our fake news systems and our deception detection work, but will also require deeper collaboration with social scientists and journalists to develop the interfaces and the models to build human-computer AI that informs and helps society as a whole.

## REFERENCES

- [1] Arora, S., Ge, R., Moitra, A.: Learning topic models—going beyond SVD. In: Proceedings of Foundations of Computer Science (2012)
- [2] Athalye, A., Engstrom, L., Ilyas, A., Kwok, K.: Synthesizing robust adversarial examples. In: Proceedings of the International Conference of Machine Learning (2018)
- [3] Bakhtin, A., Brown, N., Dinan, E., Farina, G., Flaherty, C., Fried, D., Goff, A., Gray, J., Hu, H., Jacob, A.P., Komeili, M., Konath, K., Kwon, M., Lerer, A., Lewis, M., Miller, A.H., Mitts, S., Renduchintala, A., Roller, S., Rowe, D., Shi, W., Spisak, J., Wei, A., Wu, D., Zhang, H., Zijlstra, M.: Human-level play in the game of Diplomacy by combining language models with strategic reasoning. *Science* 378(6624), 1067–1074 (2022)
- [4] Bartolo, M., Roberts, A., Welbl, J., Riedel, S., Stenetorp, P.: Beat the AI: Investigating adversarial human annotation for reading comprehension. *Transactions of the Association for Computational Linguistics* 8, 662–678 (2020)
- [5] Boyd-Graber, J., Satinoff, B., He, H., Daume III, H.: Besting the quiz master: Crowdsourcing incremental classification games. In: Proceedings of Empirical Methods in Natural Language Processing (2012)
- [6] Chang, J., Boyd-Graber, J., Wang, C., Gerrish, S., Blei, D.M.: Reading tea leaves: How humans interpret topic models. In: Proceedings of Advances in Neural Information Processing Systems (2009)
- [7] Diggelmann, T., Boyd-Graber, J., Bulian, J., Ciaramita, M., Leippold, M.: CLIMATE-FEVER: A dataset for verification of real-world climate claims. In: NeurIPS Workshop on Tackling Climate Change with Machine Learning (2020)
- [8] Eisenschlos, J.M., Dhingra, B., Bulian, J., Börschinger, B., Boyd-Graber, J.: Fool me twice: Entailment from wikipedia gamification. In: Conference of the North American Chapter of the Association for Computational Linguistics (2021)
- [9] Elgohary, A., Peskov, D., Boyd-Graber, J.: Can you unpack that? learning to rewrite questions-in-context. In: Proceedings of Empirical Methods in Natural Language Processing (2019)
- [10] Feng, S., Boyd-Graber, J.: What AI can do for me: Evaluating machine learning interpretations in cooperative play. In: International Conference on Intelligent User Interfaces (2019)
- [11] Feng, S., Boyd-Graber, J.: Learning to explain selectively: A case study on question answering. In: Proceedings of Empirical Methods in Natural Language Processing (2022)
- [12] Fujinuma, Y., Paul, M., Boyd-Graber, J.: A resource-free evaluation metric for cross-lingual word embeddings based on graph modularity. In: Proceedings of the Association for Computational Linguistics (2019)
- [13] Gor, M., Webster, K., Boyd-Graber, J.: Toward deconfounding the influence of subject’s demographic characteristics in question answering. In: Proceedings of Empirical Methods in Natural Language Processing. p. 6 (2021)
- [14] Grissom II, A., He, H., Boyd-Graber, J., Morgan, J.: Don’t until the final verb wait: Reinforcement learning for simultaneous machine translation. In: Proceedings of Empirical Methods in Natural Language Processing (2014)
- [15] He, H., Boyd-Graber, J., Kwok, K., Daumé III, H.: Opponent modeling in deep reinforcement learning. In: Proceedings of the International Conference of Machine Learning (2016)
- [16] He, H., Grissom II, A., Boyd-Graber, J., Daumé III, H.: Syntax-based rewriting for simultaneous machine translation. In: Proceedings of Empirical Methods in Natural Language Processing (2015)
- [17] Hoyle, A., Goel, P., Peskov, D., Hian-Cheong, A., Boyd-Graber, J., Resnik, P.: Is automated topic model evaluation broken?: The incoherence of coherence. In: Neural Information Processing Systems (2021)
- [18] Hu, Y., Boyd-Graber, J., Satinoff, B., Smith, A.: Interactive topic modeling. *Machine Learning* 95(3), 423–469 (Jun 2014), <http://dx.doi.org/10.1007/s10994-013-5413-0>
- [19] Iyyer, M., Boyd-Graber, J., Claudino, L., Socher, R., Daumé III, H.: A neural network for factoid question answering over paragraphs. In: Proceedings of Empirical Methods in Natural Language Processing (2014)
- [20] Iyyer, M., Enns, P., Boyd-Graber, J., Resnik, P.: Political ideology detection using recursive neural networks. In: Proceedings of the Association for Computational Linguistics (2014)
- [21] Iyyer, M., Manjunatha, V., Boyd-Graber, J., Daumé III, H.: Deep unordered composition rivals syntactic methods for text classification. In: Proceedings of the Association for Computational Linguistics (2015)
- [22] Kramár, J., Eccles, T., Gemp, I., Tacchetti, A., McKee, K.R., Malinowski, M., Graepel, T., Bachrach, Y.: Negotiation and honesty in artificial intelligence methods for the board game of Diplomacy. *Nature Communications* 13(1) (2022)
- [23] Kumar, V., Smith, A., Findlater, L., Seppi, K., Boyd-Graber, J.: Why didn’t you listen to me? comparing user control of human-in-the-loop topic models. In: Proceedings of the Association for Computational Linguistics (2019)
- [24] Lau, J.H., Newman, D., Baldwin, T.: Machine reading tea leaves: Automatically evaluating topic coherence and topic model quality. In: Proceedings of the European Chapter of the Association for Computational Linguistics (2014)
- [25] Lund, J., Cook, C., Seppi, K., Boyd-Graber, J.: Tandem anchoring: A multiword anchor approach for interactive topic modeling. In: Proceedings of the Association for Computational Linguistics (2017)
- [26] Nguyen, K., Boyd-Graber, J., Daumé III, H.: Reinforcement learning for bandit neural machine translation with simulated human feedback. In: Proceedings of Empirical Methods in Natural Language Processing (2017)
- [27] Nguyen, V.A., Boyd-Graber, J., Resnik, P.: SITS: A hierarchical nonparametric model using speaker identity for topic segmentation in multiparty conversations. In: Proceedings of the Association for Computational Linguistics (2012)
- [28] Nguyen, V.A., Boyd-Graber, J., Resnik, P.: Lexical and hierarchical topic regression. In: Proceedings of Advances in Neural Information Processing Systems (2013)
- [29] Nguyen, V.A., Boyd-Graber, J., Resnik, P., Cai, D., Midberry, J., Wang, Y.: Modeling topic control to detect influence in conversations using nonparametric topic models. *Machine Learning* 95, 381–421 (2014)
- [30] Nguyen, V.A., Boyd-Graber, J., Resnik, P., Miler, K.: Tea party in the house: A hierarchical ideal point topic model and its application to Republican legislators in the 112th Congress. In: Proceedings of the Association for Computational Linguistics (2015)

- [31] Niculae, V., Kumar, S., Boyd-Graber, J., Danescu-Niculescu-Mizil, C.: Linguistic harbingers of betrayal: A case study on an online strategy game. In: Proceedings of the Association for Computational Linguistics (2015)
  - [32] Peskov, D., Cheng, B., Elgohary, A., Barrow, J., Danescu-Niculescu-Mizil, C., Boyd-Graber, J.: It takes two to lie: One to lie and one to listen. In: Proceedings of the Association for Computational Linguistics (2020)
  - [33] Shi, T., Zhao, C., Boyd-Graber, J., Daumé III, H., Lee, L.: On the potential of lexico-logical alignments for semantic parsing to sql queries. In: Findings of EMNLP (2020)
  - [34] Shneiderman, B.: Human-Centered AI: A New Synthesis. Springer-Verlag, Berlin, Heidelberg (2021)
  - [35] Si, C., Zhao, C., Boyd-Graber, J.: What's in a name? answer equivalence for open-domain question answering. In: Proceedings of Empirical Methods in Natural Language Processing (2021)
  - [36] Smith, A., Boyd-Graber, J., Fan, R., Birchfield, M., Wu, T., Weld, D., Findlater, L.: No explainability without accountability: An empirical study of explanations and feedback in interactive ml. In: International Conference on Human Factors in Computing Systems (2020)
  - [37] Smith, A., Kumar, V., Boyd-Graber, J., Seppi, K., Findlater, L.: Digging into user control: Perceptions of adherence and instability in transparent models. In: International Conference on Intelligent User Interfaces (2020)
  - [38] Wallace, E., Boyd-Graber, J.: Trick me if you can: Adversarial writing of trivia challenge questions. In: ACL Student Research Workshop (2018)
  - [39] Wallace, E., Rodriguez, P., Feng, S., Yamada, I., Boyd-Graber, J.: Trick me if you can: Human-in-the-loop generation of adversarial question answering examples. Transactions of the Association for Computational Linguistics 10 (2019)
  - [40] Ye, Q., Khabsa, M., Lewis, M., Wang, S., Ren, X., Jaech, A.: Sparse distillation: Speeding up text classification by using bigger student models. In: Conference of the North American Chapter of the Association for Computational Linguistics (2022)
  - [41] Yuan, M., Van Durme, B., Boyd-Graber, J.: Multilingual anchoring: Interactive topic modeling and alignment across languages. In: Neural Information Processing Systems (2018)
  - [42] Zhang, M., Xu, K., Kawarabayashi, K.i., Jegelka, S., Boyd-Graber, J.: Are girls neko or shōjo? Cross-lingual alignment of non-isomorphic embeddings with iterative normalization. In: Proceedings of the Association for Computational Linguistics (2019)
  - [43] Zhao, C., Xiong, C., Boyd-Graber, J., Daumé III, H.: Distantly-supervised dense retrieval enables open-domain question answering without evidence annotation. In: Proceedings of Empirical Methods in Natural Language Processing (2021), [http://umiacs.umd.edu/~jbg/docs/2021\\_emnlp\\_weak\\_dpr.pdf](http://umiacs.umd.edu/~jbg/docs/2021_emnlp_weak_dpr.pdf)
-