

Learning to Explain Selectively: A Case Study on Question Answering

Shi Feng

University of Chicago
shif@uchicago.edu

Jordan Boyd-Graber

University of Maryland
jbg@umiacs.umd.edu

Abstract

Explanations promise to bridge the gap between humans and AI, yet it remains difficult to achieve consistent improvement in AI-augmented human decision making. The usefulness of AI explanations depends on many factors, and always showing the same type of explanation in all cases is suboptimal—so is relying on heuristics to adapt explanations for each scenario. We propose learning to explain “selectively”: for each decision that the user makes, we use a model to choose the best explanation from a set of candidates and update this model with feedback to optimize human performance. We experiment on a question answering task, Quizbowl, and show that selective explanations improve human performance for both experts and crowdworkers.

1 Introduction

Recent advances in machine learning (ML) (Silver et al., 2017; Brown et al.; Jumper et al., 2021) sparked new interest in **intelligence augmentation**—the vision that computers are not mere number-crunching tools but also interactive systems that can *augment* humans at problem solving and decision making (Engelbart, 1962). The hope is to combine the complementary strengths of human and machine, and to fully harness these models with human intuitions and oversight (Dafoe et al., 2020; Amodei et al., 2016). But this agenda is obstructed by the many counterintuitive traits of neural networks (NNs) (Szegedy et al., 2014; Goodfellow et al., 2015; Zhang et al., 2017) and our lack of theoretical understanding (Belkin et al., 2019): these models are not interpretable to humans by default, and it is difficult to foresee when they will fail. This lack of interpretability also amplifies the risk of model bias (Angwin et al., 2016; Bolukbasi et al., 2016; Caliskan et al., 2017), making it difficult to use NN-powered AIs in real-world decision making (Poursabzi-Sangdeh et al., 2021).

To bridge the gap between human and machine, various methods attempt to explain model predictions in human-interpretable terms, e.g., by providing more context to the model’s uncertainty estimation (Gal et al., 2016; Bhatt et al., 2021), by highlighting the most important part of the input (Ribeiro et al., 2016; Lundberg and Lee, 2017; Ebrahimi et al., 2017), and by retrieving relevant training examples (Renkl, 2014; Koh and Liang, 2017). Grounded in psychology (Lombrozo, 2006; Kulesza et al., 2012), these explanations promise to augment human decision making (Pu and Chen, 2006; Rader et al., 2018). But in application-grounded evaluations—with real problems and real humans (Doshi-Velez and Kim, 2018), it proves difficult for any single explanation method to achieve consistent improvement in disparate contexts (Bansal et al., 2021; Buçinca et al., 2020).

A major contributor to this problem is the breadth of context that the explanation methods are applied in. Internally, the explanation method can falter when the model reacts badly to novel inputs (Goodfellow et al., 2015; Liu et al., 2021); externally, it faces human users with diverse levels of expertise (Feng and Boyd-Graber, 2019), engagement (Sidner et al., 2005), and general trust in AI (Dietvorst et al., 2015). Our current use of explanations demands a one-size-fits-all solution, but existing methods cannot provide that as they are largely oblivious to the above-mentioned variables.

Selective explanations. Each person is unique, and the right explanation might vary from one example to another, so we propose to show explanations selectively to maximize their utility as decision support. Concretely, we assume a given set of explanation methods, but instead of showing all of them for every decision that the human user makes, we use a selector to choose a subset of the explanations to display. We can think of the selector as controlling an on/off switch for each explanation

method. The selector is allowed, for example, to show three types of explanations for one example but withhold all of them for the next one.

Online optimization. For the policy to accurately estimate the utility of explanations in each context, its training data must provide reasonable coverage over the joint distribution of all types of explanations, human users, and examples. This means the dataset will have to include cases where the users receive suboptimal decision support, e.g., with excessive explanations causing information overload (Doshi-Velez and Kim, 2018). We focus on the online setting—new information is constantly hurtling toward the user with limited time to carefully select which explanation to pay attention to—which represents real-world scenarios where the opportunity cost of giving suboptimal support cannot be ignored. In this setting, a good policy must balance the trade-off between sticking to tried and true combinations of explanations and exploring new ones; we model this trade-off with a multi-armed bandit formulation (Robbins, 1952).

We evaluate selective explanations on Quizbowl—explained in more detail in Section 2—using the same platform as Feng and Boyd-Graber (2019). To mimic real-world decision making, we recruit twenty trivia enthusiasts and ran a multi-player, real-time Quizbowl tournament. We compare our method head-to-head against baselines such as showing all explanations for all examples. Selective explanations outperform all other strategies, including the best subset of explanations identified by Feng and Boyd-Graber (2019). We also evaluate our method with mechanical Turkers—amateurs whose performance without assistance is far worse than the AI. Explanations significantly boost their performance, but only selective explanations can help them reach performance comparable with the AI.

2 Explaining Selectively to Support Human Decisions

Explanations have many uses in human-AI cooperation; this paper focuses on using explanations as *decision support*—to improve the quality of human decisions under machine assistance. Not all problems benefit from machine assistance (Doshi-Velez and Kim, 2018)—in this section, we identify three criteria for decision support testbeds. We then introduce our setup based on Quizbowl (Rodriguez et al., 2019), a competitive trivia game.

2.1 Criteria for decision support testbeds

It is not uncommon to use low-stakes and synthetic tasks to evaluate machine assistance, but it’s important to find tasks where results can generalize. Building on existing work (e.g. Lee and See, 2004; Lim et al., 2009; Yin et al., 2019), we identify the three criteria for suitable tasks.

Clear objectives. The task must have well-defined metrics, and the standard for good decisions must be clear to all participants. With unreliable metrics, well-optimized decision support will still fail to improve decision quality (Amodei et al., 2016).

Diversity of context. A reliable testbed should be diverse in terms of both participants (e.g., their skill levels) and test examples (e.g., their difficulty level). As discussed in Section 1, the lack of diversity contributes to the inconsistent results.

Incentives to be engaged. The participants must be incentivized to pay attention to model outputs to establish proper reliance (Lim et al., 2009). As a corollary, the model should demonstrate complementary strengths and provide information that participants cannot extract by themselves. In terms of the setup, engagement can also be improved by imposing time limits (Doshi-Velez and Kim, 2018) and introducing competition (Bitrián et al., 2021).

We choose Quizbowl (Rodriguez et al., 2019)—a task that roughly satisfies all three criteria—as our testbed. Compared to previous work that uses Quizbowl to evaluate explanations (Feng and Boyd-Graber, 2019), we make several changes to the setup for evaluating online selective explanations. In the following, we first introduce the most basic setting with only human players and build up our system one component at a time.

2.2 Human-only Quizbowl

We start with the most basic (and traditional) setting: Quizbowl with only human players. Quizbowl is a trivia game popular in the English-speaking world where players compete to answer questions from all areas of academic knowledge, including history, literature, science, sports, and more.¹ Each Quizbowl question consists of four to

¹While these games often have collaboration on bonus questions, we consider only individual players on tossup (US) or starter (UK/INDIA) questions. Likewise, throughout this paper, we assume each human-AI team has a single human player. The extension to multiple humans is non-trivial and is thus left for future work.

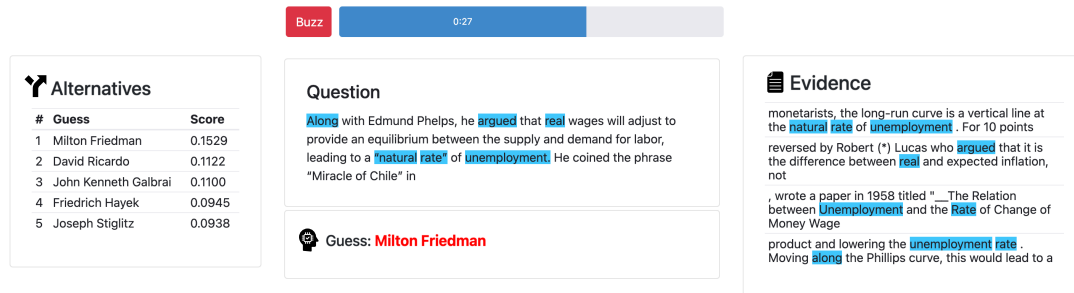


Figure 1: Our Quizbowl web interface when all four explanations are displayed. In the middle we show the question word-by-word; below, we show the current best model guess, which is colored red when the Autopilot is confident, otherwise gray; on the left, we show Alternatives, including confidence scores; on the right, we show snippets of relevant training examples as Evidence; finally, we show Highlights for the question and the evidence, respectively.

five clues. The clues are organized by their difficulty in each question: starting with the clue that's most difficult and obscure, and finishing with the one that's easiest and most telling. The clues are presented to all players *word-by-word* in real-time, verbally, or in text (e.g. web interfaces). And players compete to answer as early as possible.

To signal that they know the answer, players interrupt the question with a *buzz*, which takes its name from the signaling device's sound. Whoever buzzes needs to answer: ten points for a correct answer, and a five-point penalty for a wrong one. A player only gets one chance at each question.

To win Quizbowl, you need to answer quickly *and* correctly. This game requires not only trivia knowledge but also an accurate assessment of confidence and risk (He et al., 2016). We formally discuss the evaluation metric in Section 2.6.

2.3 Human + AI + Explanations

In our Quizbowl games, human players augmented with AI decision support compete against each other. In each human-AI team, the human player is still in charge of making decisions of when to buzz and what to answer, but now with the help of a machine learning *guesser* which predicts an answer given a question (we provide details about the guesser in Section 3). In addition to showing the guesser's current best guess, we show four types of explanations: Alternatives (Lai and Tan, 2019), salient word Highlights (Ribeiro et al., 2016), relevant training examples as Evidence (Wallace et al., 2018), and a new explanation, Autopilot. As the name suggests, Autopilot assumes the role of the human player and makes suggestions on *whether* to buzz

or to wait (details in Section 2.5). We build our interface (Figure 1) by extending the interface of Feng and Boyd-Graber (2019). We discuss these changes in detail next and in Section 3.

2.4 Human + AI + Selective Explanations

With selective explanations, decision support is customized for each player and question. For each new question, a selector policy (or *selector* for short) switches each explanation on and off. We refer to a combination of explanations as a *configuration*; for example, one configuration could be showing Highlights and Evidence but hiding Alternatives. A configuration is selected at the beginning of each question and kept constant throughout the question, but the content of each explanation is still updated dynamically. For example, Highlights will always show when they are turned on for a question, but the exact highlighted words change as more clues are revealed.

We make two important changes to the setting of Feng and Boyd-Graber (2019) to accurately estimate the effect of selectivity.

- **The guesser prediction is always available.** We make this design choice to better isolate the effect of the explanations.
- **Separate highlights for the question and the evidence.** Highlights can be applied to both the question and the evidence. In Feng and Boyd-Graber (2019), the two are treated as one explanation. However, their experiments confirm that question highlights alone are already effective. In this paper, we separate the two and the policy can control them individually. Table 1 lists the available configurations for Highlights and Evidence.

#	Evidence		Highlights	
	Question	Evidence	Question	Evidence
1				
2		✓		
3		✓	✓	
4		✓	✓	✓
5			✓	

Table 1: Each configuration is a set of visualizations shown to users, and our policy learns which configuration helps users the most. Most visualizations turn on or off independently, but some only make sense with others, e.g. we cannot highlight evidence without evidence. This table summarizes configurations for two visualizations: *Autopilot* and *Highlights*, which are dependent on each other. Combined with two independent explanations—*Alternatives* and *Autopilot*—there are twenty possible configurations.

2.5 A new explanation: 🤖 *Autopilot*

While most of our explanations appeared in previous work, we introduce a more assertive explanation, the *Autopilot*. At each step of the question, *Autopilot* gives the human player one bit of information: should you buzz or not. The suggestion is based on the binary prediction of whether the guesser’s current top answer is correct or wrong, just as how human players assess their confidence.

An autonomous AI could use *Autopilot* to decide when to buzz. But in a human-AI team, it’s just a suggestion, and the decision is still left to the human. If the human blindly follows the suggestion, the human-AI team reduces to an autonomous AI trying to win by itself, hence the name.

Both *Autopilot* and the selector try to maximize the chance of winning. While *Autopilot* optimizes for the AI only, the selector optimizes for the team. And this is no coincidence—*Autopilot* tests if selective explanation goes beyond implicit calibration: the hope is to decide to buzz better than both human-*Autopilot* teams and a fully-autonomous AI using *Autopilot*.

We use a simple, threshold-based model for *Autopilot* similar to Yamada et al. (2018): it looks at the normalized confidence scores of the top five guesses and recommends buzzing if the gap between the top two is larger than 0.05 (a threshold tuned on the dev set from Rodriguez et al. (2019)). Despite its simplicity, this model is very efficient at choosing the right time to buzz (Yamada et al., 2018; Rodriguez et al., 2019).

2.6 Evaluating accuracy and efficiency using one metric without an opponent

Winning in Quizbowl requires you to answer correctly before your opponent. In real Quizbowl games with two or more players, a high score is proof that a player is both accurate and efficient—in the sense that they require little information to get the answer right. In a perfect assessment of a Quizbowl player, we would control for factors such as question topics and have a head-to-head competition between every pair of players. In an ideal evaluation of decision support, we need to control for confounders such as player skill and have a head-to-head comparison between every possible pair of differently-augmented players, e.g., a strong player with no support vs. a weak player with selective explanations, and vice versa. However, this is infeasible due to the number of confounders.

Thus, we would like a single metric to evaluate both accuracy and efficiency without running head-to-head competitions. Accuracy is trivial to evaluate by itself, but efficiency is not as simple as counting the number of words that the player sees when they answer a question correctly because not all words have the same value: answering earlier by one word is much more difficult at the beginning of the question than at the end. The reward for answering earlier should be proportional to the increase in the chance of beating an opponent.

The expected wins (EW) metric implements this idea. Concretely, it assigns a weight to each correct answer depending on the percentage of the question revealed. The higher the percentage, the lower the assigned weight. For example, answering correctly halfway through the question counts as 0.3 points in EW, while a correct answer at the end only counts as 0.05 points. We use weights provided by Rodriguez et al. (2019) which are estimated using maximum likelihood from previous game data (Boyd-Graber et al., 2012).

2.7 Online optimization of the explanation selector to maximize cumulative EW

Our goal is to build effective human-AI teams whose cooperation requires the selector to select which explanations to show to the human. This section describes the machine learning model—learned from users’ preferences in behavioral data—which lets the selector pick user-specific explanations to show the user. Finally, to model the exploration-exploitation trade-off, we formulate

#	Description
1	Confidence of current top guesses.
2	Previous confidence of current top guesses.
3	Change in confidence of top guesses.
4	Gap in confidence between top guesses.
5	If top guesses maintained their rank.
6	If top guesses appear in previous step.
7	User’s accuracy.
8	User’s average relative buzzing position.
9	User’s average EW score.
10	Gap in EW compared to optimal buzzer.
11	Portion of words highlighted in question.
12	Portion of words highlighted in evidence.
13	Longest highlighted substring in question.
14	Longest highlighted substring in evidence.

Table 2: The user model uses the above features in addition to BERT representations of the questions. The three categories capture information about the guesser’s current prediction, the user, and the explanations. These features let the selector predict which explanations will be most useful for a human-AI team.

the online optimization of the selector as a multi-armed bandit problem and maximize the team’s cumulative EW score.

Given a human player, a question, and one of the available explanation configurations, the user model predicts the EW score received from this question. To model the human player as well as the properties of each specific question, the user model uses both manually crafted features and BERT representations. Table 2 shows the full list of features. The user model can also be viewed as a value function in reinforcement learning.

Our goal is to empower humans to complete the task at hand as accurately and as efficiently as possible. Given a new question, the selector should choose the best configuration based on its model of the user; however, to learn this model, the selector needs to test how well each configuration works for each type of question. This presents an exploration-exploitation trade-off, which we model with multi-armed bandits (Robbins, 1952). As the user plays, new information gathered about the user is incorporated into the user model via features (Table 2), and we train their personalized selector using LinUCB (Auer, 2002).

3 Experiments

We run two experiments with real human participants: a single-player experiment with crowdworkers and a multi-player real-time Quizbowl tournament with experts. This section first introduces the metric for evaluating Quizbowl competency, then provides details about the human players, the AI player, the explanation methods, and the baselines. We show that selective explanation provides personalized decision support and leads to the best augmented human performance.

3.1 Setup: Crowdworkers

We recruit twenty crowdworker players through Amazon Mechanical Turk. Each crowdworker player answers a set of sixty Quizbowl questions, and the questions are randomly permuted for each player. Each player is randomly assigned to either the experimental group with selective explanations or a control group with a baseline; more on these conditions later.

Before the user answers questions, we familiarize the user with the interface. During that period, the user can explore the interface without restriction (e.g., they can turn explanations on and off), and we switch to the assigned setting after the user clicks a button to indicate that they are ready.

3.2 Setup: Experts

We recruit twenty expert Quizbowl players from online trivia enthusiast forums. For these experts, we use a newly commissioned set of 144 questions no participant has seen before. The questions are randomly divided into six rounds.

Unlike the crowdworker experiment, the experts play a real multi-player Quizbowl game. To make sure that our game is fair and competitive, we divide players into three rooms. The initial assignment uses players’ self-reported skill levels. We subsequently adjust the assignment at the end of each round by promoting the top 20% players in each room and relegating the bottom 20%.






3.3 Setup: AI guesser and explanations

The human player is assisted by an AI guesser. Given a question, the guesser produces a multinomial distribution over the set of possible answers (Boyd-Graber et al., 2012); we update this prediction after every four question words. We use the BERT-based guesser from Rodriguez et al. (2019), and refer readers to that paper for model

Condition	Description
None	Display no explanation.
Everything	Display all explanations.
Autopilot	Display Autopilot suggestions only.
AI-only	Autopilot replaces human player.
Random	Choose a new random configuration for each question.
Selective	Selector chooses the configuration for each question.

Table 3: Conditions in the randomized controlled trial. Under `Random` and `Selective` conditions, the enabled configurations could change from one question to another; under all other conditions, one configuration is used for all questions. In all conditions, the human player has access to the guesser’s prediction. In the baseline `AI-only` condition, no human player is involved.

details and standard evaluation results. Next we discuss how we generate explanations for the guesser.

-  **Alternatives:** We show the guesser’s current top five predictions along with their confidence scores.
-  **Evidence:** We retrieve four training examples that are most similar to the current question. To measure similarity we use cosine distance between question representations by the guesser (Wallace et al., 2018).
-  **Highlights on question:** We use HotFlip (Ebrahimi et al., 2017) and show tokens with a normalized attribution score higher than 0.15.
-  **Highlights for evidence:** We search for the highlighted question tokens in the retrieved training examples, and highlight them.
-  **Autopilot:** We colorize the guesser’s prediction based on the Autopilot’s current decision: red if buzzing, gray if not. When Autopilot is disabled, the color is always black.

Hyperparameters of an explanation (e.g. the number of highlighted tokens) affect its effectiveness. Here we choose a fixed set of hyperparameters based on feedback from internal trial runs. However, the choice of hyperparameters can also be considered as part of the explanation configuration. In such a setup, we could use the selector with an expanded action space to, for example, also choose the number of tokens to highlight. We discuss this more in Section 4.3.

3.4 Setup: Conditions and baselines

Table 3 lists the conditions of our randomized controlled trial. The experimental condition is selective

explanations. The control conditions include baseline policies such as using a fixed explanation configuration for all questions. To control the number of conditions, we omit conditions with fixed configurations, e.g. `Alternatives+Evidence`. Instead, we include `Everything`, which Feng and Boyd-Graber (2019) show to be most effective at improving user accuracy.

The guesser’s accuracy is on par with the experts. So if the crowdworker players are willing and able to *blindly* and *precisely* follow `Autopilot`, they could achieve good scores. But we consider this as a degenerate solution to human-AI cooperation.

To account for this issue, we include two special settings. In `Autopilot`, we display `Autopilot` suggestions as the only explanation for the human player. In `AI-only`, we *replace* the human player with `Autopilot` to make decisions. Using these two settings, we can quantify to what degree the human player follows `Autopilot`.

In our forum post to recruit experts, we promise an “interface to augment human players with explanations of AI predictions”. To stay true to this promise and ensure a good experience for the experts (who participate in Quizbowl for fun), we omit the `None` baseline in our expert experiments. This omission should not affect our results since the baseline is compared to other conditions in Feng and Boyd-Graber (2019).

3.5 Evaluation: Does selector improve EW?

We use the mean cumulative EW score over the course of the game (144 questions for experts and 60 for crowdworkers) for our quantitative comparison. If the human-AI team with a tailored selector can improve in EW, this suggests explanations are helping the users more than other conditions.

Figure 2 shows how the mean EW score un-

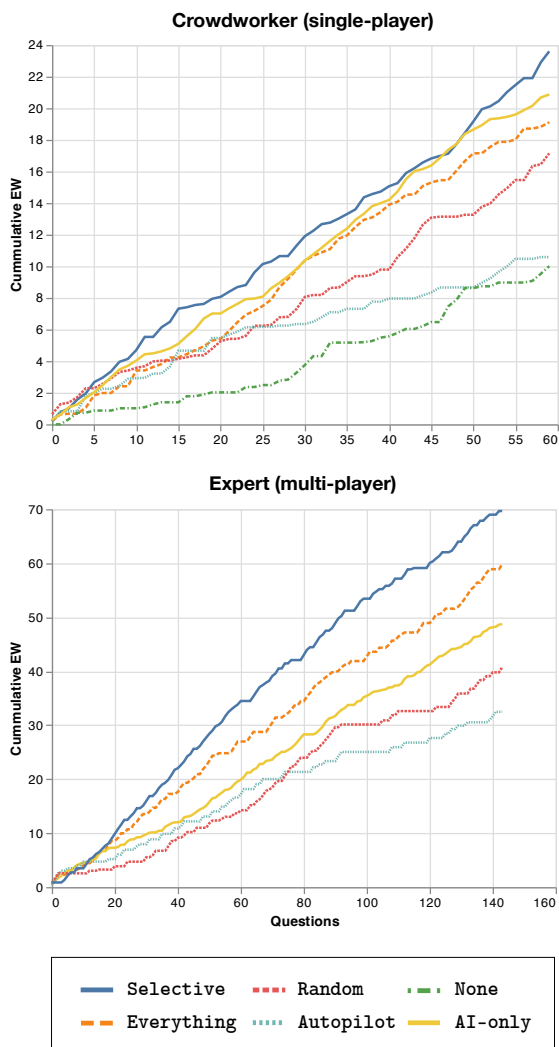


Figure 2: Mean cumulative EW score over number of questions under each condition by crowdworkers (top) and experts (bottom). The selective condition gets the highest score among all human-AI cooperative settings.

der each condition increases as the players answer more questions. Among all human-AI cooperative settings, the *Selective* condition is the best. Especially for experts, selective explanation by the selector is better than both showing all explanations and AI-only. Importantly, as our model acquires more data for each user with more questions (and as the user acclimatizes to their teammate), the gap between *Selective* and *Everything* grows.

Without explanations, crowdworkers are much worse than AI-only. With selective explanations, crowdworkers are comparable to AI-only and only slightly better than showing all explanations.

Under the *Autopilot* condition, if players blindly follow the AI’s suggestion—buzz when the *Autopilot* says so and provide the AI prediction

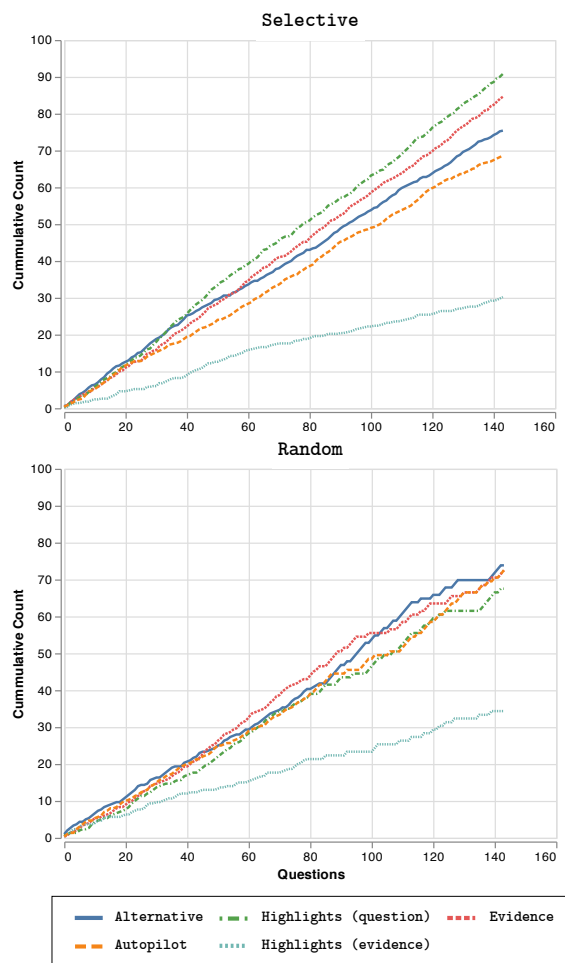


Figure 3: Mean cumulative count of explanations over time shown to **experts**. We compare the explanations shown by the selector (top) and by random (bottom). Based on the frequency, the selector learns a ranking of explanations consistent with the effectiveness reported in [Feng and Boyd-Graber \(2019\)](#): question highlights are most effective, then evidence, then alternatives.

as the answer—they should match the AI-only baseline. However, both experts and crowdworkers lose to the AI-only condition. This indicates that the other conditions evince a synergy: humans are not simply blindly following AI suggestions. Rather, the diverse and selective explanations allow the players to better decide when to follow and when to use their own knowledge.

3.6 Evaluation: What does the selector show?

We are interested in what the selector learns to be most effective and what it chooses to show to players. Figure 3 visualizes the evolving distribution of configurations selected by the bandit selector and the random selector.

First, the selector did not learn to show all expla-

nations for all questions—it learned, as the name suggests, to be selective. And compared to the random selector, the selector formed a clear preference among explanations. In fact, at the end of the game, the selector—learning purely from interaction—recovers the ranking of individual explanations reported by [Feng and Boyd-Graber \(2019\)](#): highlight > evidence > alternatives. Interestingly, the selector did not converge to this ranking until the players finished about 60 questions: initially the list of alternatives was the preferred explanation, possibly because it is easier for the players to interpret than the others. Eventually, as the players get more used to the other explanations and the selector continues to learn about the players, it converges.

4 Discussion and Related Work

4.1 Who should drive?

Clearly defining the shared obligations of the team is crucial to the success of the team. By design, we keep ultimate control of decision making with the human. However, this may not be optimal; a distracted, overloaded, or hesitant human might be better served by an AI “taking the wheel” if it is certain. The most relevant work to ours is [Gao et al. \(2021\)](#), which similarly uses bandit feedback to optimize team performance. While our policy chooses from the set of explanation configurations, their policy makes a binary decision: whether to delegate a decision to the human or leave it to the AI. Our `Autopilot` explanation can be seen as “soft” delegation. Future work should compare selective explanation with more methods for delegation and deferral ([Madras et al., 2018](#); [Lubars and Tan, 2019](#); [Kamath et al., 2020](#); [Lai et al., 2022](#)).

4.2 Alignment and learning to optimize human objectives

Typically, ML algorithms optimize automatic metrics: how well can a machine replace or emulate a human. However, this is inconsistent with how humans and machines interact in the real world; often models need to be personalized to users ([Zhou and Brunskill, 2016](#)). The research area that deals with the general problem of optimizing human objectives is alignment ([Amodei et al., 2016](#)). Specifically for human-AI teams, an unsettled question is how to optimize for that partnership; while we optimize for short-term accuracy, a reasonable alternative would be to optimize for longer-term learning ([Bragg and Brunskill, 2020](#)). An interesting

direction would be to take a real-world task and directly optimize the underlying model (not just the selector) to create tailored explanations, as [Lage et al. \(2018\)](#) did for synthetic tasks.

4.3 Expanding the selector’s action space

We present this work as another step towards learned explanations that are more aligned with human values ([Amodei et al., 2016](#)). Our method seeks to maximize a human objective, not heuristic proxies of that ([Doshi-Velez and Kim, 2018](#)), and not the objective of the solo machine. In this work, we focus on a simplified setting with a limited design and action space, but our experimental setting closely mimics how a human-AI team would operate in a real-world task; in particular, our testbed, Quizbowl, bears merits that are essential for a task to have to benefit from this idea.

We focus on this restricted selector to keep the sample complexity for online optimization under control. In principle, the selector could be more fine-grained if we allow it to dynamically change the configuration as the clues in the question are revealed. Despite challenges concerning sample complexity, we believe that this expansion of action space is a logical next step.

5 Conclusion: Explanations Tailored for Users

Users benefit from collaborating with an AI, and this collaboration can be improved by explaining the AI well. Moreover, the benefit is not universal, some users need or thrive with different explanations. However, finding the right combination is not easy; while our bandit approach can find useful explanations, it requires both the user to become acclimatized to human-AI teaming and the bandit to explore the space of explanations. As human-AI collaborations become more common, we must continue to search for better signals and methods to help the teaming minimize stress and acclimation but maximize fun and productivity.

Limitations

As we discussed in Section 1, a major contributor to the inconsistency of human-AI experimental results is the large number of factors that can influence cooperative effectiveness. One of those factors that’s relatively easy to model is the human’s skill level. In theory, selective explanation should

be able to model that: if we optimize selective explanation jointly for experts and crowdworkers, the selector should be able to learn and choose different explanations for the two different groups of players. Unfortunately, we couldn't have done that experiment because Quizbowl is too challenging for crowdworkers without any assistance, and when they compete head-to-head the game is made more difficult by the element of competition.

There are other factors of human-AI cooperation that have been identified by previous work but we couldn't model: the level of human agency (Lai and Tan, 2019; Bansal et al., 2021), the model's predictive accuracy (Bansal et al., 2020), the user's mental model of machine learning (Bansal et al., 2019), and the amount of interactivity (Smith-Renner et al., 2020a,b). Even within limited interactions, there is significant variation about the optimal modality of explanations (Gonzalez et al., 2020). Other factors, such as the distribution of test examples and model architecture, affect the quality of output from various post-hoc explanation methods (Ghorbani et al., 2019; Jones et al., 2020).

Another major limitation of our evaluation is that we only experimented with one question answering problem, Quizbowl. Our method is generally applicable to decision making problems. But finding another suitable task and adapting our infrastructure, experiment design, and incentive structures is highly non-trivial. We are actively looking for other problems to experiment on and hope to conduct more extensive experiments in the future.

Ethics Statement

All of our work is conducted under the supervision and approval of an institutional review board. Our annotators are fairly compensated for their work, and we have attempted to make the activity as intrinsically rewarding as possible.

Acknowledgments

We thank the anonymous reviewers and meta-reviewer for their suggestions and comments. Boyd-Graber is supported by NSF Grant IIS-1822494. Any opinions, findings, conclusions, or recommendations expressed here are those of the authors and do not necessarily reflect the view of the sponsors.

References

- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul F. Christiano, John Schulman, and Dan Mané. 2016. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine bias: Risk assessments in criminal sentencing.
- Peter Auer. 2002. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422.
- Gagan Bansal, Besmira Nushi, Ece Kamar, Eric Horvitz, and Daniel S Weld. 2020. Is the most accurate AI the best teammate? optimizing ai for teamwork. In *Association for the Advancement of Artificial Intelligence*.
- Gagan Bansal, Besmira Nushi, Ece Kamar, Walter S Lasecki, Daniel S Weld, and Eric Horvitz. 2019. Beyond accuracy: The role of mental models in human-ai team performance. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*.
- Gagan Bansal, Tongshuang Wu, Joyce Zhu, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel S Weld. 2021. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. In *International Conference on Human Factors in Computing Systems*.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. 2019. Reconciling modern machine-learning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences*.
- Umang Bhatt, Javier Antorán, Yunfeng Zhang, Q Vera Liao, Prasanna Sattigeri, Riccardo Fogliato, Gabrielle Melançon, Ranganath Krishnan, Jason Stanley, Omesh Tickoo, et al. 2021. Uncertainty as a form of transparency: Measuring, communicating, and using uncertainty. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 401–413.
- Paula Bitrián, Isabel Buil, and Sara Catalán. 2021. Enhancing user engagement: The role of gamification in mobile apps. *Journal of Business Research*, 132:170–185.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of Advances in Neural Information Processing Systems*.
- Jordan L. Boyd-Graber, Brianna Satinoff, He He, and Hal Daumé III. 2012. Besting the quiz master: Crowdsourcing incremental classification games. In

- Proceedings of Empirical Methods in Natural Language Processing.*
- Jonathan Bragg and Emma Brunskill. 2020. Fake it till you make it: Learning-compatible performance support. In *Proceedings of Uncertainty in Artificial Intelligence.*
- TB Brown, B Mann, N Ryder, M Subbiah, J Kaplan, P Dhariwal, A Neelakantan, P Shyam, G Sastry, A Askell, et al. Language models are few-shot learners. arxiv 2020. In *Proceedings of Advances in Neural Information Processing Systems.*
- Zana Bućinca, Phoebe Lin, Krzysztof Z Gajos, and Elena L Glassman. 2020. Proxy tasks and subjective measures can be misleading in evaluating explainable ai systems. In *International Conference on Intelligent User Interfaces.*
- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Allan Dafoe, Edward Hughes, Yoram Bachrach, Tatum Collins, Kevin R McKee, Joel Z Leibo, Kate Larson, and Thore Graepel. 2020. Open problems in cooperative ai. *arXiv preprint arXiv:2012.08630.*
- Berkeley J Dietvorst, Joseph P Simmons, and Cade Massey. 2015. Algorithm aversion: people erroneously avoid algorithms after seeing them err. *Journal of Experimental Psychology: General*, 144(1):114.
- Finale Doshi-Velez and Been Kim. 2018. Towards a rigorous science of interpretable machine learning. *Springer Series on Challenges in Machine Learning.*
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2017. HotFlip: White-box adversarial examples for text classification. In *Proceedings of the Association for Computational Linguistics.*
- Douglas C Engelbart. 1962. Augmenting human intellect: A conceptual framework. *Menlo Park, CA*, page 21.
- Shi Feng and Jordan Boyd-Graber. 2019. What can AI do for me: Evaluating machine learning interpretations in cooperative play. In *International Conference on Intelligent User Interfaces.*
- Yarin Gal, Yutian Chen, Roger Frigola, S. Gu, Alex Kendall, Yingzhen Li, Rowan McAllister, Carl Rasmussen, Ilya Sutskever, Gabriel Synnaeve, Nilesh Tripuraneni, Richard Turner, Oriol Vinyals, Adrian Weller, Mark van der Wilk, and Yan Wu. 2016. *Uncertainty in Deep Learning*. Ph.D. thesis, University of Oxford.
- Ruijiang Gao, Maytal Saar-Tsechansky, Maria De-Arteaga, Ligong Han, Min Kyung Lee, and Matthew Lease. 2021. Human-AI collaboration with bandit feedback. In *International Joint Conference on Artificial Intelligence.*
- Amirata Ghorbani, Abubakar Abid, and James Y. Zou. 2019. Interpretation of neural networks is fragile. In *Association for the Advancement of Artificial Intelligence.*
- Ana Valeria Gonzalez, Gagan Bansal, Angela Fan, Robin Jia, Yashar Mehdad, and Srinivasan Iyer. 2020. Human evaluation of spoken vs. visual explanations for open-domain qa. *arXiv preprint arXiv:2012.15075.*
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and harnessing adversarial examples. In *Proceedings of the International Conference on Learning Representations.*
- He He, Jordan L. Boyd-Graber, Kevin Kwok, and Hal Daumé III. 2016. Opponent modeling in deep reinforcement learning. In *Proceedings of the International Conference of Machine Learning.*
- Erik Jones, Shiori Sagawa, Pang Wei Koh, Ananya Kumar, and Percy Liang. 2020. Selective classification can magnify disparities across groups. *arXiv preprint arXiv:2010.14134.*
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. 2021. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589.
- Amita Kamath, Robin Jia, and Percy Liang. 2020. Selective question answering under domain shift. In *Proceedings of the Association for Computational Linguistics.*
- Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *Proceedings of the International Conference of Machine Learning.*
- Todd Kulesza, Simone Stumpf, Margaret Burnett, and Irwin Kwan. 2012. Tell me more? the effects of mental model soundness on personalizing an intelligent agent. In *International Conference on Human Factors in Computing Systems.*
- Isaac Lage, Andrew Slavin Ross, Been Kim, Samuel J Gershman, and Finale Doshi-Velez. 2018. Human-in-the-loop interpretability prior. *arXiv preprint arXiv:1805.11571.*
- Vivian Lai, Samuel Carton, Rajat Bhatnagar, Q Vera Liao, Yunfeng Zhang, and Chenhao Tan. 2022. Human-ai collaboration via conditional delegation: A case study of content moderation. In *International Conference on Human Factors in Computing Systems.*
- Vivian Lai and Chenhao Tan. 2019. On human predictions with explanations and predictions of machine learning models: A case study on deception detection. In *Proceedings of ACM FAT**.

- John D Lee and Katrina A See. 2004. Trust in automation: Designing for appropriate reliance. *Human factors*, 46(1):50–80.
- Brian Y Lim, Anind K Dey, and Daniel Avraami. 2009. Why and why not explanations improve the intelligibility of context-aware intelligent systems. In *International Conference on Human Factors in Computing Systems*.
- Han Liu, Vivian Lai, and Chenhao Tan. 2021. Understanding the effect of out-of-distribution examples and interactive explanations on human-ai decision making. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):1–45.
- Tania Lombrozo. 2006. The structure and function of explanations. *Trends in cognitive sciences*.
- Brian Lubars and Chenhao Tan. 2019. Ask not what ai can do, but what ai should do: Towards a framework of task delegability. In *Proceedings of Advances in Neural Information Processing Systems*.
- Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Proceedings of Advances in Neural Information Processing Systems*.
- David Madras, Toni Pitassi, and Richard Zemel. 2018. Predict responsibly: improving fairness and accuracy by learning to defer. In *Proceedings of Advances in Neural Information Processing Systems*.
- Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Wortman Vaughan, and Hanna Wallach. 2021. Manipulating and measuring model interpretability. In *International Conference on Human Factors in Computing Systems*.
- Pearl Pu and Li Chen. 2006. Trust building with explanation interfaces. In *International Conference on Intelligent User Interfaces*.
- Emilee Rader, Kelley Cotter, and Janghee Cho. 2018. Explanations as mechanisms for supporting algorithmic transparency. In *International Conference on Human Factors in Computing Systems*.
- Alexander Renkl. 2014. Toward an instructionally oriented theory of example-based learning. *Cognitive science*, 38(1):1–37.
- Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. “why should i trust you?”: Explaining the predictions of any classifier. In *Knowledge Discovery and Data Mining*.
- Herbert Robbins. 1952. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*.
- Pedro Rodriguez, Shi Feng, Mohit Iyyer, He He, and Jordan Boyd-Graber. 2019. Quizbowl: The case for incremental question answering. *arXiv preprint arXiv:1904.04792*.
- Candace L Sidner, Christopher Lee, Cory D Kidd, Neal Lesh, and Charles Rich. 2005. Explorations in engagement for humans and robots. *Artificial Intelligence*.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. 2017. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359.
- Alison Smith-Renner, Ron Fan, Melissa Birchfield, Tongshuang Wu, Jordan Boyd-Graber, Daniel S Weld, and Leah Findlater. 2020a. No explainability without accountability: An empirical study of explanations and feedback in interactive ml. In *International Conference on Human Factors in Computing Systems*.
- Alison Smith-Renner, Varun Kumar, Jordan Boyd-Graber, Kevin Seppi, and Leah Findlater. 2020b. Digging into user control: perceptions of adherence and instability in transparent models. In *International Conference on Intelligent User Interfaces*.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. 2014. Intriguing properties of neural networks. In *Proceedings of the International Conference on Learning Representations*.
- Eric Wallace, Shi Feng, and Jordan Boyd-Graber. 2018. Interpreting neural networks with nearest neighbors. In *EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*.
- Ikuya Yamada, Ryuji Tamaki, Hiroyuki Shindo, and Yoshiyasu Takefuji. 2018. Studio ousia’s quiz bowl question answering system. *arXiv preprint arXiv:1803.08652*.
- Ming Yin, Jennifer Wortman Vaughan, and Hanna Wallach. 2019. Understanding the effect of accuracy on trust in machine learning models. In *International Conference on Human Factors in Computing Systems*.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2017. Understanding deep learning requires rethinking generalization. In *Proceedings of the International Conference on Learning Representations*.
- Li Zhou and Emma Brunskill. 2016. Latent contextual bandits and their application to personalized recommendations for new users.