

Finding AI's Faults with AAR/AI: An Empirical Study

ROLI KHANNA, JONATHAN DODGE, ANDREW ANDERSON, RUPIKA DIKKALA, JED IRVINE, ZEYAD SHUREIH, KIN-HO LAM, CALEB R. MATTHEWS, ZHENGXIAN LIN, MINSUK KAHNG, ALAN FERN, and MARGARET BURNETT, Oregon State University

Would you allow an AI agent to make decisions on your behalf? If the answer is “not always,” the next question becomes “in what circumstances”? Answering this question requires human users to be able to assess an AI agent—and not just with overall pass/fail assessments or statistics. Here users need to be able to *localize* an agent’s bugs so that they can determine when they are willing to rely on the agent and when they are not. After-Action Review for AI (AAR/AI), a new AI assessment process for integration with Explainable AI systems, aims to support human users in this endeavor, and in this article we empirically investigate AAR/AI’s effectiveness with domain-knowledgeable users. Our results show that AAR/AI participants not only located significantly *more* bugs than non-AAR/AI participants did (i.e., showed greater recall) but also located them more *precisely* (i.e., with greater precision). In fact, AAR/AI participants outperformed non-AAR/AI participants on every bug and were, on average, almost six times as likely as non-AAR/AI participants to find any particular bug. Finally, evidence suggests that incorporating labeling into the AAR/AI process may encourage domain-knowledgeable users to abstract above individual instances of bugs; we hypothesize that doing so may have contributed further to AAR/AI participants’ effectiveness.

CCS Concepts: • **Human-centered computing** → **Empirical studies in HCI**;

Additional Key Words and Phrases: AAR/AI, after-action review for AI, explainable AI (XAI)

ACM Reference format:

Roli Khanna, Jonathan Dodge, Andrew Anderson, Rupika Dikkala, Jed Irvine, Zeyad Shureih, Kin-Ho Lam, Caleb R. Matthews, Zhengxian Lin, Minsuk Kahng, Alan Fern, and Margaret Burnett. 2022. Finding AI’s Faults with AAR/AI: An Empirical Study. *ACM Trans. Interact. Intell. Syst.* 12, 1, Article 1 (February 2022), 33 pages. <https://doi.org/10.1145/3487065>

1 INTRODUCTION

Explainable AI (XAI) has recently begun to expand its scope. Besides simply explaining AI to its users, some XAI researchers are focusing on explanation-based systems to help users *assess* an AI system’s decisions (e.g., [12, 23, 35, 38, 48, 60, 61, 70]).

The reviewing of this article was managed by associate editor Paternò, Fabio.

This material is based upon work supported by the AI Research Institutes program supported by NSF and USDA-NIFA under the AI Institute’s Agricultural AI for Transforming Workforce and Decision Support (AgAID) award #2021-67021-3534, and by DARPA #N66001-17-2-4030. Any opinions, findings, and conclusions or recommendations expressed are the authors’ and do not necessarily reflect the views of DARPA, NSF, the USDA, the Army Research Office, or the U.S. government.

Authors’ address: R. Khanna, J. Dodge, A. Anderson, R. Dikkala, J. Irvine, Z. Shureih, K.-H. Lam, C. R. Matthews, Z. Lin, M. Kahng, A. Fern, and M. Burnett, School of EECS, Kelley Engineering Center, Oregon State University, Corvallis, OR, 97331 USA; emails: {khannaro, dodgej, anderson2, dikkalar}@oregonstate.edu, irvine@eeecs.oregonstate.edu, {shureihz, lamki, matthea, linzhe, minsuk.kahng, Alan.Fern}@oregonstate.edu, burnett@eeecs.oregonstate.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

2160-6455/2022/02-ART1 \$15.00

<https://doi.org/10.1145/3487065>

Imagine “Pat,” a user knowledgeable in some domain who is trying to make an educated decision about whether or when to rely upon a particular intelligent agent in a *particular* situation important to them. In the domain of AI-powered aviation, Pat might be a test pilot, deciding on the circumstances in which a human pilot should override AI-generated decisions in a new AI-powered airplane. In the domain of health care, Pat might be a remote aging-in-place caregiver, deciding when to follow an AI system’s recommendations on actions their grandparent may need from them right away. These particular situations matter to Pat and to the person or people Pat affects. No matter how thoroughly trained an AI system is, for Pat the dilemma is not about the AI system’s overall correctness statistics—it is about their responsibility for making the most appropriate decision for this particular case. The European Commission’s European Group on Ethics in Science and New Technologies put it this way: “[Autonomous systems] must not impair [the] freedom of human beings to set their own standards and norms and be able to live according to them” [21, 50].

In assessing domains like these, no objective “ground truth” is available, because (1) only Pat knows *this* airplane’s or *their* grandparent’s behaviors in the context in which they happened, and (2) there may be multiple response paths Pat could take that would produce positive outcomes.¹

To enable domain experts such as Pat to assess an AI agent and find its faults in circumstances like these, we devised **After-Action Review for AI (AAR/AI)** [19, 46], which we detail in Section 1.3. AAR/AI aims to help domain experts assess AI agents in sequential decision-making environments. This work evaluates how well AAR/AI achieves this aim.

1.1 The Domain

Real-Time Strategy (RTS) games are a popular sequential decision-making domain for AI research. In RTS games, players attempt to strategically maneuver through a plethora of choices to win the game. Sometimes in AI research, the RTS player is an AI agent maneuvering on behalf of a human,² which is the case in this work.

For our study, we used a model-based **Reinforcement Learning (RL)** agent that played an RTS game. The game was the same StarCraft 2 “Tug-of-War” custom game as was used in earlier AAR/AI publications [19, 46]. Tug-of-War games entail two evenly matched players, a *Friendly AI* and *Enemy AI* pursuing the same goal. In our game, tugs of war occur in the top and bottom “lanes” of a game (Figure 1), over the course of a maximum of 40 **Decision Points (DPs)** or rounds. Players perform actions in either lane at each DP, depending on affordability (e.g., how much they spent and earned prior to the current DP). Actions include purchasing troop production buildings and/or purchasing pylons to increase income. Figure 1 shows a screenshot of the game replay as it appeared to our study’s participants.

In the game, troops battle to win, with troop types in a rock-paper-scissors relationship: marines are effective against immortals, immortals against banelings, and banelings against marines. Different troops cost different amounts, depending on these capabilities. Troops spawn behind the nexus (the player’s base, represented as gold star-shaped objects on the gameboard in Figure 1). Once spawned, they march down the lane and attack enemies in pre-programmed fashion, as with the marines shown in Figure 1. Players can win in two ways: (1) destroy one of its opponent’s nexuses, or (2) if all the nexuses remain after 40 DPs, the player whose nexus has the lowest health loses.

1.2 An AI RTS Player’s Failures and Faults

What if an AI agent, such as the one playing this game, makes a flawed decision? An AI agent’s flawed decisions are analogous to the software engineering concept of “failures.” Ammann and

¹Groce et al. [23] also pointed out that no objective ground truth exists in many human/AI decision domains—there is only “good enough for my purposes.”

²For example, as a proxy for dangerous strategy-centered situations such as military operations.

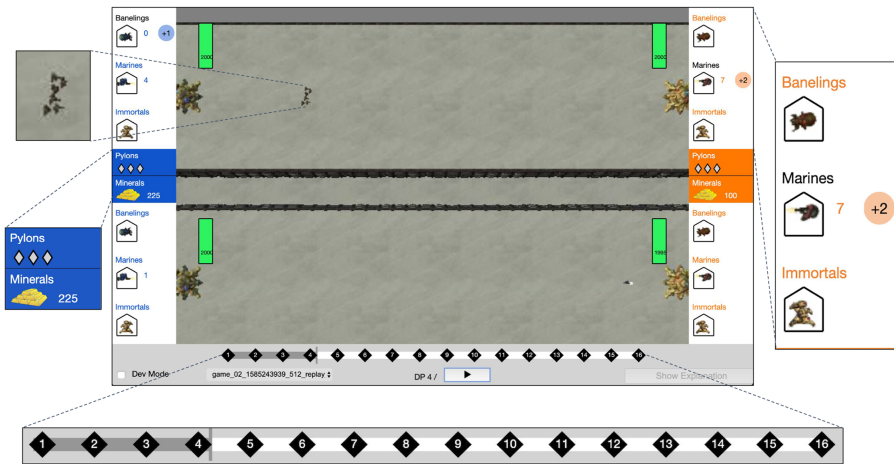


Fig. 1. The participants’ replay view of the game just past DP 4 (fourth diamond; see callout at bottom). The game board has two lanes where action takes place: a top lane and a bottom lane. The (blue) Friendly AI agent’s “home” is the left side, and the (orange) Enemy is the right side. Each lane’s troop inventories are shown in a side panel for that lane (e.g., the callout at right blows up the side panel for the Enemy’s top lane). Both players’ side panels also summarize resources (e.g., the blue callout (middle left) shows the Friendly AI’s resources). On the game board, a group of Friendly AI marines in the top lane are currently moving toward the Enemy’s nexus (top left callout).

Offutt [3] define a “failure” as “external, incorrect behavior with respect to the requirements.” In the RTS domain, our analogous requirement is the AI agent deciding upon good “enough” actions (according to a human knowledgeable in the domain), so we define failures as user-visible decisions the AI agent makes that are not adequate³ according to that particular user’s standards.

Still, a failure is only a symptom of something going wrong under the hood. Ideally, an interactive XAI system could not only help Pat spot such symptoms but also locate the root causes of those symptoms. Only in this way can Pat know which decisions the AI is making for acceptable reasons, so as to avoid “lucky guesses,” ward off ethical concerns, or defend against potential legal challenges (e.g., a malpractice suit) [21, 37].

In medicine, the causes of symptoms are diseases; in software engineering literature, the causes of symptoms (failures) are termed *faults*. Avizienis et al. [6] define a fault to be the underlying cause or condition that may lead to a failure, and “fault localization” to be the act of identifying the *locations* of faults. Building upon these definitions, in an RTS game with XAI support, we define a fault to be erroneous reasoning by the AI agent—ideally revealed to the users in the explanations—and fault localization to be finding the component of the explanation that reveals the erroneous reasoning.⁴ In this article, we also use the term *bug* synonymously with *fault*.

1.3 Human Users Assessing an AI RTS Player’s Failures and Bugs with/without AAR/AI

Finding an AI’s failures and bugs involves users understanding the AI’s behaviors. However, numerous empirical studies have reported difficulties users face in building this understanding, even in the presence of explanations (e.g., [4, 14, 40, 60]). One possibility is that presenting explanations

³As with Ammann and Offutt’s definition, an AI agent’s failure is not always a “show-stopper.” In other words, a bad decision is a failure even if the AI agent later makes good enough decisions to overcome the initial bad decision.

⁴In medicine, identifying the disease is a necessary step toward a cure but still may not be sufficient to produce an effective cure. Similarly, in software engineering, localizing a fault is necessary but still may not be sufficient to produce a fix.

to users is not enough—because explanations unscaffolded by a *process* require the user to not only build their understanding but also to build their own process for doing so.

In this work, we empirically investigate whether a process known as AAR/AI [19, 46] can improve XAI users’ ability to understand an AI agent well enough to localize its bugs. AAR/AI is a new member of the After-Action Reviews (AARs) family. AARs were originally devised by the U.S. Army [64] for assessing human decisions. The AAR is a facilitated, team-based debriefing method. In the military, it has been used to assess soldier training sessions. Sawyer and Deering [59] characterize the AAR process using the acronym “DEBRIEF”: Define rules, Explain objectives, Benchmark performance, Review what was supposed to happen, Identify what happened, Examine why, and Formalize learning [59]. More generally, AAR has been used for decades to assess human decisions in the military (e.g. [24]) and has also been adapted to manned-unmanned teams [9]. Civilians have also used AAR processes in transportation [45], medical treatment [54, 59], and emergency response [18, 30, 42]. A recent meta-analysis of 61 studies by Keiser and Arthur [33] found that using AAR produced beneficial and practical effects.

AAR/AI is the first use of AAR in AI. The original AAR/AI publication [46] describes AAR/AI in seven steps, adapted from the preceding DEBRIEF sequence. These steps are conducted with a facilitator and one or more assessors as follows:

- (1) The facilitator defines the domain.
- (2) The facilitator explains the agent’s objective.
Next begins an “inner loop” for each decision to be assessed:
- (3) The facilitator reviews what was supposed to happen.
- (4) An assessor identifies what happened.
- (5) An assessor describes why it happened.
- (6) An assessor formalizes learning from this decision.
Finally, at the end of these iterations:
- (7) An assessor formalizes learning holistically from every decision they analyzed.

The AAR/AI process allows flexibility in the details within each step, to allow customization to the assessors’ purpose in their domain.

Although there is some qualitative evidence revealing some of AAR/AI’s strengths [19, 46], AAR/AI has not been empirically compared with *not* using AAR/AI. Only a comparison of with-AAR/AI vs. without-AAR/AI can measure causality—whether using AAR/AI leads human users to significantly greater effectiveness at assessing an AI system than they would achieve without AAR/AI. To address this need, we prototyped an AAR/AI-supported XAI system (which will be illustrated in Section 3.2) and used it to conduct a controlled lab experiment comparing the effectiveness of domain-knowledgeable users localizing faults (bugs) using AAR/AI vs. without AAR/AI. In both treatments, participants could use information in the game itself and a full explanation of the AI system’s reasoning. Our study investigated the following research questions:

- RQ1: Does the AAR/AI process help domain-knowledgeable users to localize faults in an XAI-based system?
- RQ2: Does the type of fault interact with RQ1?
- RQ3: When supported by AAR/AI, do users somehow abstract beyond individual instances of faults? If so, how?

2 BACKGROUND AND RELATED WORK

A substantial body of research has investigated human users *understanding*, *finding* (e.g., through testing), and/or *debugging/improving* AI systems, all of which relate to humans localizing an AI agent’s faults.

Fault localization in an AI agent requires the human doing the localizing to have at least a partial understanding of how the AI agent reasons. XAI aims at exactly this goal. One of its aims is to improve people's mental models [4, 5, 38, 39]—representations people construct in their heads about how something works from whatever they have experienced with it [49].

However, affecting someone's mental model is not always straightforward. Although people have mental models about most things, their mental models are not always accurate and sometimes are not be very malleable. For example, Tullio et al. [63] found explanations helped clarify some misconceptions, but overall mental model structure went largely unchanged. This exposes a central challenge XAI faces when trying to help people understand an AI system, which Yang et al. [69] describe as “[humans’] uncertainty surrounding AI’s capabilities [and] AI’s output complexity” [69].

To address this problem, some XAI researchers have drawn from social science the strategy of helping humans generate “self-explanations,” a process that has been shown to support knowledge acquisition [56]. A user might generate self-explanations as a result of being prompted to do so or might do so on their own accord [26].

As an example of XAI work involving self-explanations, Chi et al. [13] showed that learners relying more heavily on examples had worse outcomes, which they credited to inability to engage in self-explanation. Another example occurred in our early qualitative AAR/AI results [46], in which participants’ uses of self-explanation, in combination with other factors, produced high levels in Bloom’s learning taxonomy [8]. Still another example of encouraging self-explanations is the use of counterfactuals; Byrne [10] offers evidence that counterfactuals enable people to explain how events relate to one another such as identifying cause-effect or reason-action relationships.

XAI consumers correspond to human learners, in that the humans consuming XAI are doing so to learn how the AI reasoning went. This correspondence opens the possibility of drawing from **Cognitive Load Theory (CLT)** for insights. CLT models tasks as having three kinds of load: intrinsic (“nature of the material”), extraneous (“manner in which the material is presented”), and germane (“reflects the effort that contributes to the construction of schemas”) [62]. The recommendation of much of CLT work is to increase germane load and decrease extraneous load where possible.

Where can XAI researchers and developers turn to find concrete XAI-pertinent guidance to fulfill recommendations like these? For non-AI systems, when faced with the challenge of helping people form more accurate mental models, UI designers can draw upon substantial work distilling research results into usability fundamentals and practical guidelines. Unfortunately, however, few such works yet exist for XAI. Although Hoffman et al. [27] recently conducted a large-scale literature survey on guidance for empiricists *measuring* XAI’s effects, explanation design was not covered. In a very large-scale literature survey by Abdul et al. [1] on XAI with an HCI (human-computer interaction) perspective, none of the usability papers reported were tailored for XAI. Soon after the paper of Abdul et al. [1], Amershi et al. [2] created 18 usability-centric guidelines for human-AI interaction. A few of these guidelines are applicable to XAI, but as one of the first works in the direction of usability guidelines in AI, the guidelines mainly contribute a set of design goals to achieve for usable AI, not how to achieve them. For example, Guideline 11 is “Make clear why the system did what it did,” which is an important design goal but is not guidance on how to do so. Complementing the work of Amershi et al. [2], Wang et al. [65] presented a theory-centric framework connecting social science fundamentals on human reasoning and human biases to XAI techniques. This work produced six XAI lessons learned from the social science research. These six, like Amershi’s 18, are at the design goal level (e.g., “support hypothesis generation”), but unlike Amershi’s 18, these six also drill down one more level of theory. For example, one recommendation is to support hypothesis generation via contrastive reasoning,

hypothetico-deductive reasoning, and abductive reasoning. However, for XAI researchers not adept with social science concepts on how these kinds of reasoning work in humans, more concrete operationalizations may still be needed.

Given the paucity of XAI-specific usability fundamentals, many researchers have turned to advancing community knowledge through empirical studies. Taxonomies and related sets of principles are ways to build upon these researchers' individual empirical results—they abstract above individual experiments, thus providing intellectual tools for understanding the dimensions of XAI that researchers have been investigating.

For example, Kulesza et al. [38, 40], taxonomize XAI research via two proposed principles—soundness and completeness—illustrated in the phrase, “the whole truth (completeness) and nothing but the truth (soundness).” Understanding explanations' attributes according to these principles have implications for XAI's consumers. For example, when completeness is too low, explanation consumers may perceive it as “sneaking,” a UI dark pattern adapted to XAI [15]. However, if completeness is too high, users' searches for “the right” information can so become onerous that finding failures or localizing faults in an explanation may be reminiscent of finding the proverbial needle in a haystack.

In their “intelligibility types” taxonomy, Lim [43] and Lim and Dey [44] categorized explanations according to the kinds of questions they answer (e.g., What, Why) and its relationship to the system (e.g. Inputs, Model, Outputs). The relative importance of each intelligibility types can vary by domain. For example, Lim and Dey [44] found that users wanted Why Not information when they perceived flaws, whereas other researchers found a heavy emphasis on What information in domains like smart homes and RTS games [11, 52]. Research has reported that supporting the “right” intelligibility types for a particular situation or domain improved users' confidence in the system [17].

Although most current XAI research focuses on helping people interpret models' inner workings (e.g. [28, 66]), some tools in the closely related area of interactive **Machine Learning (ML)** are intended for failure detection and/or fault localization. Examples include using scalable query-based approaches for NLP [68], clustering around user-selected example-based “anchors” [12], or “covering” different input/output combinations [23]. Others support fault localization (“visual debugging” [61]) by revealing system internals—in this case, latent vectors for sequence-to-sequence models for translation. The explanations in our study also include revealing system internals but in a model-based agent's search tree.

For this work, the domain is a complex, sequential decision-making environment based on RTS games. Ontañón et al. [51] have pointed to the gap in research about human needs for understanding AI for RTS, and researchers have been working to fill this gap. As a few examples, Metoyer et al. [47] contributed formative work via human explanations of RTS games used within expert-novice pairs, Dodge et al. [20] and Penney et al. [52] investigated how expert broadcasters explain RTS, Kim et al. [34, 35] investigated human responses to Human vs. AI battles, and Penney et al. [52, 53] investigated pairs of AI players making sense of “simulated AI” behavior. Although some of these RTS investigations included participants *noticing AI failures* (symptoms of faults), none except the AAR/AI work [46] offer insights into humans attempting to *localize AI faults* in this domain. To help fill this gap, in this article we present the first quantitative evaluation of humans' effectiveness with the AAR/AI process.

3 METHODOLOGY

To investigate the effectiveness of the AAR/AI process for localizing AI's faults/bugs, we conducted an empirical study with domain-knowledgeable participants using AAR/AI vs. without AAR/AI. Due to COVID-19, we conducted sessions over teleconference (Zoom) and a browser-based custom

Table 1. Participant Demographics as Per Their Questionnaire Responses to Their Gender Identification (Including a Free-Form Response), Student/Non-Student Status, and Age

	AAR/AI	Non-AAR/AI	Total
Man	25	24	49
Woman	7	8	15
Transgender	1	0	1
Undergrad	13	14	27
Grad	10	7	17
Non-Student	10	11	21
Total	33	32	65

Median age was 24 years (minimum: 17; maximum: 48), with about half below and half above. (One 17-year-old claimed to meet the ≥ 18 inclusion criterion before the study, then gave their actual age on the questionnaire.) AAR/AI vs. non-AAR/AI participant demographics were similar for all categories.

combination of the platform (game and explanation system, including AAR/AI features for the AAR/AI treatment) and questionnaires. The participants were experienced with RTS games but had no AI or ML background.

3.1 Participants and Procedure

We required participants to be at least 18 years of age, and to have 10+ hours of prior experience with RTS games to ensure they would understand our domain. In addition, we excluded respondents who had taken any AI or ML class before. (We later disqualified 1 participant who became persistently inattentive during the study session.) Of the final 65 participants, 49 self-identified as men, 15 as women, and 1 as transgender (Table 1). Participants were randomly assigned by flipping a coin to one of two treatments: AAR/AI and non-AAR/AI. Each zoom session had 1 to 7 participants. Upon completing the study, they received a \$20 Amazon gift card as compensation.

All participants observed an AI agent playing the web-based RTS game described in Section 1.1. Participants' task was to localize the AI agent's bugs. Participants in both the treatments saw the same explanation (which we describe in Section 3.2)—the only difference between the treatments was the presence/absence of the AAR/AI supports. Data collected were participants' responses to a pre-task demographic questionnaire; their in-task answers to the AI agent's reasoning, and where in the explanation they saw these problems (Figure 2); a click log of their interactions; and their responses to a post-task NASA/**Task Load index (TLX)** questionnaire [25]. Additional details about the study procedures can be found in Appendix B, and all questionnaires are included in the supplemental documents accompanying this article.

The study proceeded as follows. The participants agreed to an IRB-approved informed consent form, then filled out the pre-task demographic questionnaire, then performed the following steps 1 through 3 (also illustrated in Figure 3), and finally filled out the NASA/TLX questionnaire and were compensated.

Step 1: Tutorial. The researcher began the tutorial by informing participants that (1) they would observe a game between Friendly and Enemy AI players, (2) the Friendly AI would lose, and (3) their main task was to find "problems" in the Friendly AI's actions. ("Problems" was the vocabulary we used with participants to encourage them to find any/all of the Friendly AI agent's failures and faults/bugs as defined in Section 1.)

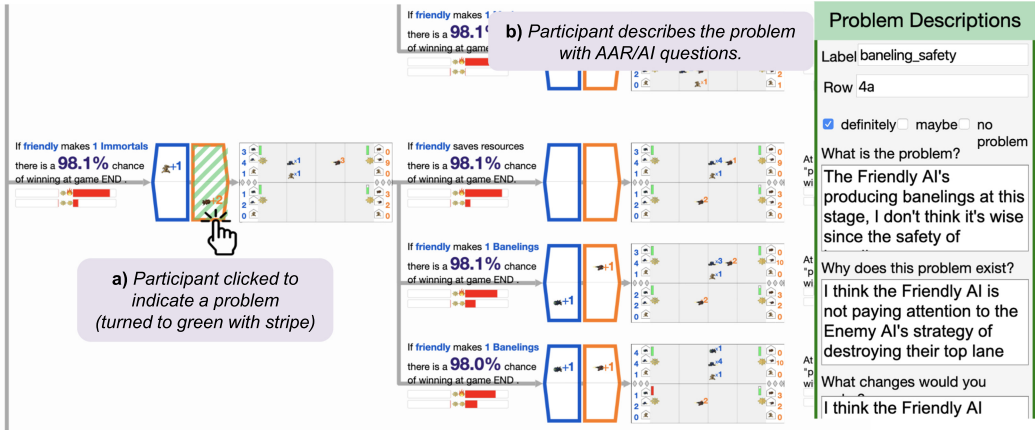


Fig. 2. How a (hypothetical) participant could mark up the Explanation UI for the AAR/AI treatment. (a) The participant selects what they think is a problem on the diagram. (b) The participant describes the problem by including a label, location, level of certainty, and responses to the What, Why, and What changes questions.

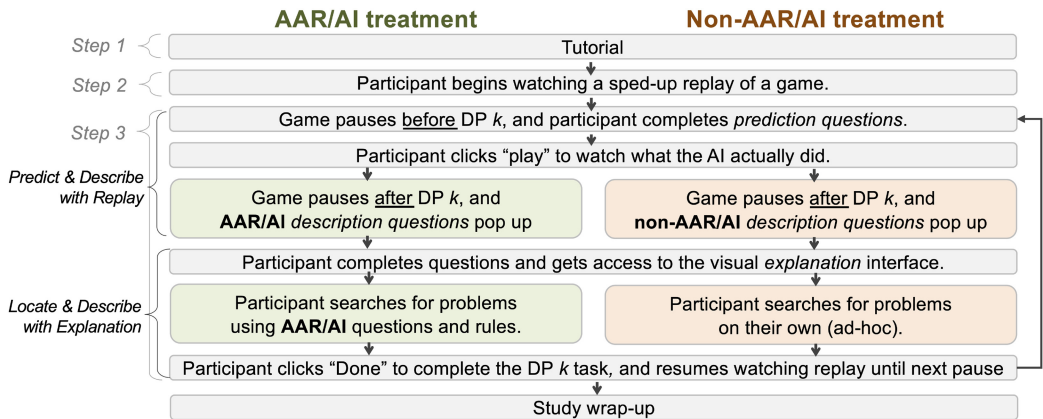


Fig. 3. Summary of study procedure.

The researcher guided the participants through working with the interface to familiarize them with the game interface and explanations for 30 to 40 minutes. The tutorial included example problems, such as “violation of game rules for buying troops in both the top and bottom lanes,” but ultimately they were told, “If you think it’s a problem, it’s a problem.” This setup and tutorial handled the first two steps of AAR/AI: (1) defining the rules and (2) explaining the agent’s objective.

Step 2: Entering the game interface. After the hands-on tutorial, the participants got access to the main task’s interface (Figure 4(A)), and they began watching a sped-up replay of a game with the Friendly AI competing against the Enemy AI. Participants could only observe the game; they did not play the game.

From here onward, the researcher had real-time access to the actions taken by all the participants through a dashboard. This ensured that the researcher could track signs of inattention or inappropriate actions taken by the participants, and deal directly with the participant about them. (One participant’s inattention could not be resolved, and we ultimately discarded their data; this is in addition to the 65 participants reported in this article.)

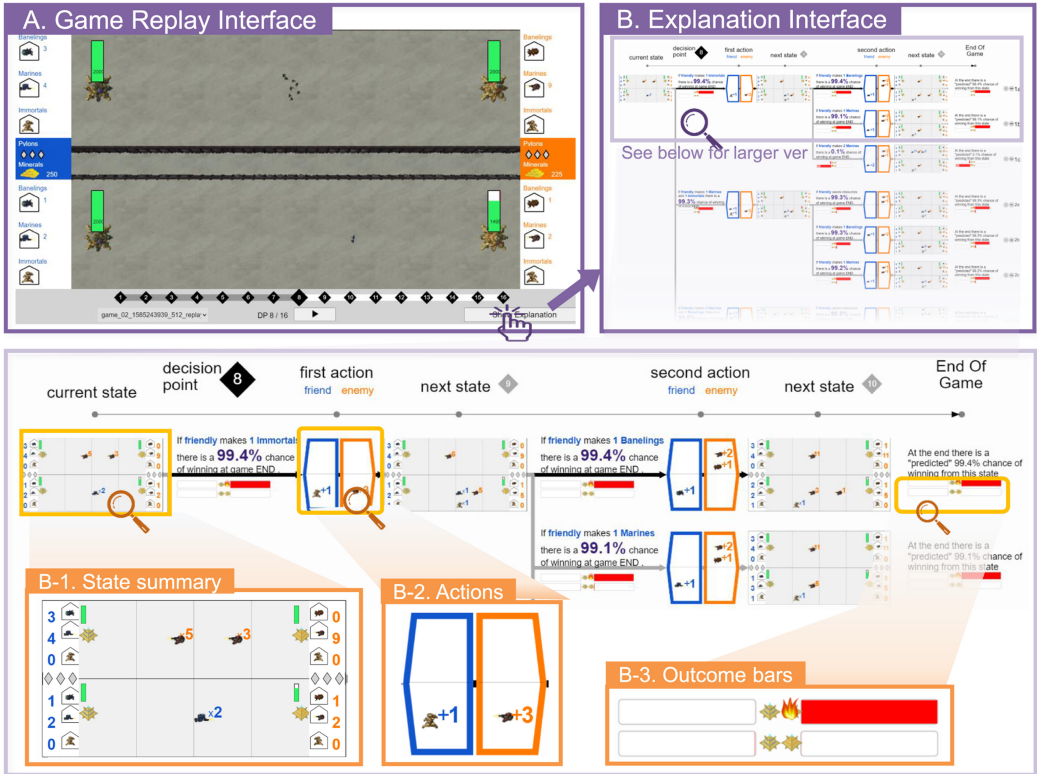


Fig. 4. (A) The game interface that participants used to watch the game in action. (B) This interface visually explains Friendly AI’s explanation for its actions. The screenshot shows top two next actions. The top row represents the best next action among multiple actions. (B-1) Current state is graphically represented. (B-2) AI’s predicted action pair (Friendly in blue and Enemy in orange). In addition, its child state, the AI’s next predicted action pair, and grandchild state are shown. (B-3) At right, the outcome bars are shown to represent how the probability is calculated.

Step 3: Main Task Loop (Predict and Describe with Replay; Locate and Describe with Explanation).

Predict and Describe with Replay: The game automatically paused at a DP before the one they would analyze, and participants provided written answers to what they thought the Friendly AI would do by the next DP. Specifically, participants said which lane it would build in, and whether it would make any marines, any banelings, any immortals, and/or a pylon. The purpose of these questions was to get the participant active in trying to figure out the AI’s reasoning.

The participants then watched the AI’s decision and answered a set of questions. The participants in the non-AAR/AI group answered a question asking what the Friendly AI had just done; the AAR/AI participants answered three questions as part of the AAR/AI process: *what* had the Friendly AI just done, *why* they thought it made those decisions, and *what changes* they would make in the Friendly AI’s decisions.

Locate and Describe with Explanation: After they had answered the initial questions, participants were able to see the explanations (Figure 4(B)). Participants were then told to locate problems in the AI’s reasoning using the explanations by clicking parts of the explanation they found problematic

and answering a set of questions. AAR/AI participants also had the option of labeling the problems they found via a free-form textbox.⁵ We captured all their interactions in a click log.

AAR/AI participants could start at any row they wanted, but once they had started a row, they had to finish locating bugs in that row and describing them via the AAR/AI questions as in Figure 2. Once they said they were finished with the row (by clicking on “Done with this row”), they moved on to whatever next row they wanted. (They could later go back to review any previous row, but they could not change it after they had said they had finished it.) In contrast, non-AAR/AI participants could move freely among rows, tackling the task of locating and describing them however they pleased.

Once participants completed finding and describing problems in one DP, they then could click on the “Done” button to indicate the completion of the task. They then could resume watching the game. The game would pause again at another DP, and participants repeated the same process in this new DP. The two DPs participants worked with were DP 8 and DP 15, selected for the bugs they exhibited. Participants could spend a minimum of 10 minutes and a maximum of 40 minutes per DP.

3.2 Explanations

Figure 4(B) shows a visual explanation of the agent for a given DP. It visualizes the internal search tree the agent made to find the best actions. The leftmost node (also shown in Figure 4(B-1)) graphically represents a current state of the game (i.e., root of the search tree). Note that the state in Figure 4(B-1) is an approximate thumbnail of the gameboard (Figure 4(A)). The tree expansion to the right of the current state shows different combinations of actions and states the agent predicts could happen next. Figure 4(B-2) shows the Friendly AI’s action in the blue box and the Enemy AI’s action in orange. Next is the predicted “next state” (child), followed by another pair of actions, and the predicted grandchild state. The explanation interface shows 5 of the 20 searched actions: the top 2, median, and bottom 2 (Figure 4(B)).

Outcome predictions, shown in Figure 4(B-3), appear in two ways: a sentence describing the win probability associated with that action and a visualization decomposing that win probability into four stacked bars, one for each nexus. Each bar shows two probabilities: one for the nexus being destroyed (shown in red) and the other for that nexus having the lowest health at the end of a game (shown in pink). The sum of all eight probabilities is 100% (the game has to end in one of eight ways); thus, a single player’s win probability is that sum minus the four probabilities that they lose (encoded by the total size of the red and pink bars on the right side), shown by the large bold number.

3.3 The RL Agent

In our study, the Friendly AI “player” is powered by a model-based RL agent,⁶ which determines the Friendly AI’s next action by predicting future states. This section summarizes the agent, and Appendix C explains its architecture in more detail.

The agent makes its predictions using the following neural network driven functions:

- (1) **Action Ranking Function (ARF)** with type signature $(\text{State}, \text{Action}) \rightarrow \text{float}$: answers “How good is taking this Action in this State?”

⁵In formative investigations and pilots, we observed that one benefit of AAR/AI seemed to be its encouragement of consistent search practices [7]. We added labeling to AAR/AI to further encourage these practices.

⁶This agent was also used in the work of Mai et al [46].

- (2) **Transition Function (TF)** with type signature (State, Action1, Action2) \rightarrow State: answers “What State will arise if I (Friendly AI) take Action1 and the opponent (Enemy AI) takes Action2?”
- (3) **Leaf Evaluation Function (LEF)** with type signature (State) \rightarrow (outcome-probabilities): answers “How good is this State?”

Using these components, the agent internally builds a search tree to select the best action for the Friendly AI player. It first enumerates all actions available in the current state, applying the ARF to each and pruning all but the top 20 actions (shown in blue in Figure 4(B)). It then applies the ARF again from the opponent's perspective, pruning all but the top 10 actions (the top one shown in orange next to the blue one in Figure 4(B)). Then for each combination of promising moves, the agent applies the TF, predicting the resultant child state. By this point, one level of the game tree has been built. It then builds one more level in the same fashion, starting from the child state, with smaller numbers of actions. Although we could repeat this process indefinitely, we stop the prediction here, because in many domains searching enough to reach a terminal state would be intractable. Thus, after the agent applies the LEF to each grandchild state, it propagates resulting values back up the tree via minimax search.⁷

3.4 The Bugs

The participants' task was to identify the AI agent's bugs. We based our bugs on the agent's naturally occurring bugs. For example, one of the bugs we found was that the AI predicted that a health value of the nexus increases over time, which cannot occur in a real game.

To harvest these bugs, we wrote scripts to find similar cases that were objectively wrong (e.g., ones that violate game rules, like the preceding nexus health example), then hand-validated the results. After rigorous analysis [41] of these naturally occurring bugs and iteratively trying them in our preliminary pilots, we harvested 10 of these bug instances, selected DPs that contained those bugs, and exaggerated those whose effects were too small to notice easily. Appendix B's Table 6 enumerates the complete list of bugs and any exaggerations we made. All participants' explanations contained all of the bugs enumerated. Five bug instances were in each of two DPs. Half of the bugs were in the agent's TF, and half were in the agent's LEF.

For example, one of the LEF bugs is Bug ID #1 at DP 8; this bug is shown in detail later in Figure 7. At one point (which will be called out in row 1C in the figure), the agent predicts that the Friendly AI will *lose* with its *bottom* nexus being destroyed, but it consistently predicts that the Friendly AI will *win* by destroying the Enemy AI's *top* lane nexus in other rows. In addition, although the actions in 1B and 1C are very similar, their outcomes are radically different, which is not possible. This is a bug with the win probabilities flipped for both the Friendly and Enemy AI's top and bottom lanes.

An example of a TF bug is Bug ID #4 at DP 8, which is shown later in Figure 8. The agent predicts that the Enemy AI will have two immortals in the next state; however, this is not possible because there exists only one building for producing immortals and each building can produce only one in a single round.

4 RESULTS

4.1 RQ1 Results: Does AAR/AI Help Localize Faults?

RQ1 asks whether the AAR/AI process helps domain-knowledgeable users localize faults. To answer this question, we measured participants' ability to find and describe the 10 bugs enumerated

⁷A full discussion on minimax and game tree search can be found in Chapter 5 in the work of Russell and Norvig [58].

Table 2. Coding Rules We Used to Code Participants' Problem Reports

	No Credit (0)	Partial Credit (+0.5)	Full Credit (+1)
Location: Participant... (A)... marked location correctly	Not (A)	N/A	(A)
Description: Participant... (A)... completely described bug correctly (B)... partially described bug correctly	Neither (A) nor (B)	(B)	(A)

For Location, if a participant's location markings were combined in a way that introduced ambiguity, we disambiguated by looking for location information in their free-form descriptions.

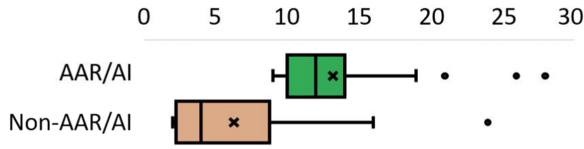


Fig. 5. Problem report count per participant. AAR/AI: Mean = 13.182, SD = 4.565; non-AAR/AI: Mean = 6.281, SD = 5.050. The AAR/AI participants submitted significantly more problem reports than their non-AAR/AI counterparts.

in Section 3.4. To code their efforts, two researchers independently coded 20% of the data corpus and achieved an inter-rater reliability (IRR) of 80.6% (Jaccard index [31]). Given this level of agreement, they then split up the remaining coding. The code set followed a scoring system. Participants could earn up to 2 points for each bug they reported: if they correctly *located* a bug, they could earn up to 1 point, and if they correctly *described* the bug, they could earn another point. Table 2 details the coding rules for no credit, partial, or full credit for locating and describing bugs.

First, we compare sheer volume of AAR/AI vs. non-AAR/AI participants' problem reports. Figure 5 shows the distributions of how many problems participants reported. AAR/AI participants reported significantly more problems than non-AAR/AI participants (t -test, $t(63) = 5.7829$, $p < .0001$).⁸ In fact, even the AAR/AI participant with the fewest problem reports (9) still submitted more than 75% of the participants in the non-AAR/AI treatment did ($Q3 = 8.25$).

To evaluate the correctness of their problem reports, we use the term *bug* to refer to the bugs in Table 6 and the term *problem* to denote whatever participants reported as problematic. We use these concepts to compute two metrics commonly used in ML: recall and precision.⁹ Recall measures the proportion of the system's 10 bugs the participants reported, and precision measures the proportion of participants' problem reports that were actually bugs. An "ideal" participant whose problem reports would show perfect recall and precision would include all of the bugs in Table 6 (perfect recall) and nothing else (perfect precision).

Using these measures, the AAR/AI participants had both significantly greater average recall (Welch's t -test, $t(55.666) = 4.5479$, $p < .0001$) and precision (t -test, $t(63) = 2.0358$, $p = .04598$) than

⁸Levene's test for equal variance determined when to use a standard t -test vs. Welch's t -test; we point out Welch's whenever we use it.

⁹See Equations (5) and (6) in Appendix A for details of how we computed recall and precision for each participant. Because this data labeling would be considered multi-class and multi-label, we could not use the basic formulas. Further, although Zhang and Zhou [71] offer Equations (3) and (4) for such labelings, they do not incorporate other issues present in our data corpus. For example, few negative examples are present in the corpus because most participants only reported bugs they thought to be present.



Fig. 6. Participants’ recall (left) and precision (right). Recall AAR/AI: Mean = 0.233, SD = 0.160; non-AAR/AI: Mean = 0.080, SD = 0.105. Precision AAR/AI: Mean = 0.179, SD = 0.129; non-AAR/AI: Mean = 0.116, SD = 0.121 over all 10 bugs. AAR/AI participants performed significantly better than non-AAR/AI participants with both measures.



Fig. 7. Example of an LEF bug (Bug ID #1), present at DP 8. The game outcome for row 1C is suspicious. Since the Friendly AI’s action (i.e., two marines) is similar to that for 1A and 1B (especially 1B: one marine), we can expect that the Friendly AI would win (>99% chance of winning) by destroying the top Enemy nexus (as in row 1B); however, the agent predicts that Friendly AI will lose (0.1% chance of winning), which implies that the win probabilities for row 1C have likely been flipped (i.e., LEF bug).

the non-AAR/AI participants (Figure 6). Cohen’s d showed a large effect size ($d = 1.121$) for the recall difference and a medium effect size ($d = .505$) for the precision difference.¹⁰

Together, these three results suggest that the AAR/AI process not only encouraged participants to report significantly more problems (Figure 5) but also encouraged them to report problems that were indeed bugs, as measured by their significantly higher recall and precision (Figure 6). These results are especially encouraging given that none of the participants had backgrounds in AI/ML.

4.2 RQ2 Results: Does the Type of Fault Matter?

RQ2 raises the question of whether AAR/AI vs. non-AAR/AI participants’ success differences depended on which particular bugs or types of bugs they were pursuing. We begin by considering the types of bugs: LEF bugs vs. TF bugs.

LEF bugs occur when the neural network provides an inaccurate game outcome for an input state, such as in Figure 7. TF bugs occur when the neural network predicts an inaccurate future state, given a current state and actions, such as in Figure 8. The experiment’s 10 bugs were evenly split between 5 LEF and 5 TF bugs. RQ2 asks whether AAR/AI vs. non-AAR/AI played out differently for these two bug types.

¹⁰We consider Cohen’s $d \in [0, 0.2)$ to be no effect, $d \in [0.2, 0.5)$ to be small, $d \in [0.5, 0.8)$ to be medium, and $d \in [0.8, 1.4)$ to be large, by convention [16].



Fig. 8. Example of a TF bug (Bug ID #10), present at DP 8. The bug (in the highlighted box) is that the agent predicts there will be two immortals in the bottom lane (solid arrows), even though there is only one immortal production building (dashed arrow).

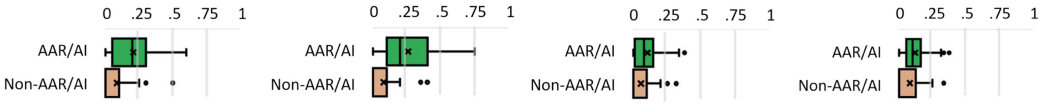


Fig. 9. AAR/AI participants’ vs. non-AAR/AI participants’ (left to right): Recall for LEF bugs only (AAR/AI: Mean = 0.208, SD = 0.185; non-AAR/AI: Mean = 0.075, SD = 0.105), recall for TF bugs only (AAR/AI: Mean = 0.258, SD = 0.188; non-AAR/AI: Mean = 0.086, SD = 0.135), precision for LEF bugs only (AAR/AI: Mean = 0.100, SD = 0.105; non-AAR/AI: Mean = 0.068, SD = 0.109), and precision for TF bugs only (AAR/AI: Mean = 0.122, SD = 0.101; non-AAR/AI: Mean = 0.078, SD = 0.122).

To answer this question, we analyzed recall and precision separately for LEF bugs and TF bugs. (Note that the number of target bugs is now split into two for analysis, which affects the distributions.) Figure 9 shows the results, with recall in the left two pairs (LEF and TF bugs, respectively) and precision in the right two pairs. The AAR/AI vs. non-AAR/AI recall differences for both bug types were significant. Specifically, AAR/AI participants found significantly greater proportions of both LEF bugs (t -test, $t(63) = 3.0358, p = .0035$) and TF bugs (Welch’s t -test, $t(51.341) = 4.7479, p < .001$). Average precision differences in AAR/AI participants vs. non-AAR/AI participants were suggestive but did not reach significance for either LEF bugs (t -test, $t(63) = 1.1891, p = .2389$) or TF bugs (t -test, $t(63) = 1.5878, p = .1173$).

As these LEF vs. TF recall and precision results show, bug type did not determine when AAR/AI participants were more effective than non-AAR/AI participants—AAR/AI participants performed at least as well as non-AAR/AI participants on both bug types. In fact, as Figure 10 shows, AAR/AI participants outperformed non-AAR/AI participants on every bug. On average, AAR/AI participants were about six times as likely as non-AAR/AI participants to find each bug (odds ratio for two sample proportions [55], $\hat{\phi} = 5.889$).

4.3 RQ3 Results: Labeling and Abstractions

This study’s third research question explored whether adding labeling to the AAR/AI process would facilitate participants’ ability to spot patterns of bugs and potentially develop abstractions capturing these patterns. Toward that end, our interface enabled AAR/AI participants to label the bugs they localized as they went along. They had free rein to label any way they chose or not to label at all. We then analyzed the kinds of labels participants used and whether their use of different kinds of labels related to their successes at localizing bugs.

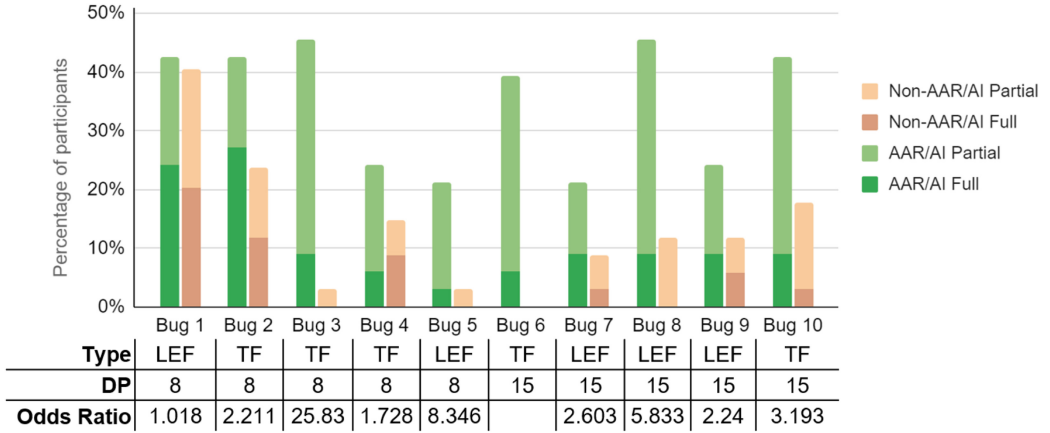


Fig. 10. Percentage of participants who found each bug. More AAR/AI participants than non-AAR/AI participants found every one of the bugs. Odds ratios show AAR/AI participants' increased likelihood of finding each bug. (Bug 6's odds ratio is undefined because no non-AAR/AI participants found that bug.)

4.3.1 What Kinds of Labels Did They Devise? Participants used a wide variety of labels to characterize the bugs they found. Some seemed to use labels simply as a way to group similar instances (e.g., “1,” “problem2”), some used location information in their labeling schemes (e.g., “first row”), and some used labels for lighter purposes (e.g., “bored”). However, some participants' labels abstracted above individual instances into concepts, either from a “naive-AI” perspective (e.g., “bad-prediction”) or from a domain perspective (e.g., “marines to be made”).

We coded the labels into categories. Two researchers generated the code set using a process similar to the summative content analysis of Hsieh and Shannon [29], where keywords are identified before and during data analysis to form the code set. This generated the categories of labels shown in Table 3. When a participant's label was applicable to more than one category, we coded it in all the applicable categories. For example, “battlefield counting issues” was coded as an instance of both Domain Concepts and Counting/Math. The researchers independently coded 20% of the data with 85% agreement (Jaccard index [31]), coding the labels directly or in the context of the participants' reports when necessary to disambiguate labels. Given these researchers' high level of code consistency, only one researcher was needed to complete the rest of the coding.

As Figure 11 shows, participants' use of labels most frequently tended toward abstractions relating to concepts of the game (Domain Concepts) or concepts of AI and/or AI Explanations (AI/XAI Concepts), with more than 120 instances of each. The next most frequent was using labels as simply grouping mechanisms (Identified Groups) (e.g., “ai3,” “problem2”); there were more than 80 instances of these. Participants also sometimes used “Not a Problem After All” labels as a way to document “all clear” diagnoses after perusing the Explanation UI; there were 70 instances of this category. The Bug Location category (e.g., “first row”) and “Un-category” category (labels not even attempting to categorize; e.g., “abc,” “can't understand”) were also somewhat common, with more than 30 instances of each. Lowest in frequency was the No Time category, in which participants used labels to document where they ran out of time, with 6 instances.

4.3.2 Which of Their Labels Correlated with Success? Of these categories, the three codes showing positive correlations with participants' success (scores) were AI/XAI Concepts, Counting/Math, and Domain Concepts (Figure 12). Each of these had a correlation coefficient r between [0.45, 0.54],

Table 3. Code Set Used to Categorize AAR/AI Participants' Labels on the Faults They Localized

Code: Description	Examples of Participants' Labels
AI/XAI Concepts: Concepts/terminology related to XAI/AI	health prediction, incorrect decision, success outcome change, overestimation of ability...
Bug location: Location on the Explanation UI	next state 17, outcome 2, second action, error row 5a, 2b not possible...
Count/Math: Related to counting, math, or calculation	battlefield counting issues, nexus health calculation, too many enemies... marines...
Domain Concepts: Concepts/terminology related to the game domain (troops, lanes, nexus)	lane 1 nexus, unit disappear, nexus randomly dies, suddenly immortals, banelings in the bottom lane...
Identified Groups: Evidence of an "ID" of some kind used repeatedly to group similar instances	ai, ai2,... 1issue, 2 issue,... problem, problem2,...
No Time: Did not have time to complete search for problems	no time, out of time
Not a problem after all: Location was marked as a (potential) problem, but added a label indicating no problem after all	no problem, n/a, no, ignore this, no issues...
"Un-category": Did not attempt to categorize; labels ranged from gibberish to messages to the researcher	error, bored, alex, cant understand, abc...

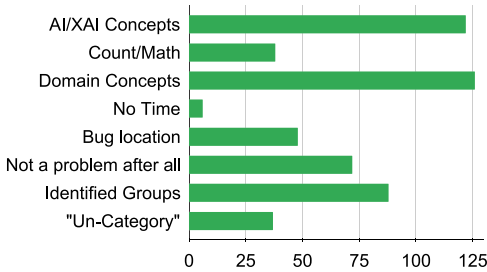


Fig. 11. Frequencies of AAR/AI participants' ways of labeling faults into these categories. (Non-AAR/AI participants did not have a labeling feature.) Domain Concepts and AI/XAI Concepts were the most common.

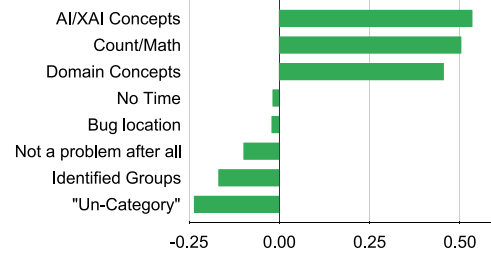


Fig. 12. Correlation between the scores in each categorization scheme and the participant's total score.

which is generally considered to be a moderate correlation [22]. In contrast, the remaining categories had small *negative* correlations with participants' success.

AI/XAI Concept labels had the highest correlation with the participants' success scores ($r = 0.535$). Among the participants' labels hinting at AI/XAI-concept abstractions were "bad prediction" (P158, score 7; P116, score 5), "winning percentage" (P130, score 14), and "overconfidence top lane" (P106, score 5.5). Participants' label usage in this category, when matched up with their click histories, suggested that they may have located bugs by walking through the explanation tree showing the AI's reasoning, in a "reasoning walkthrough" somewhat analogous to a code walkthrough. Figure 13 illustrates one such walkthrough. Another example excerpted from a participant's report is:

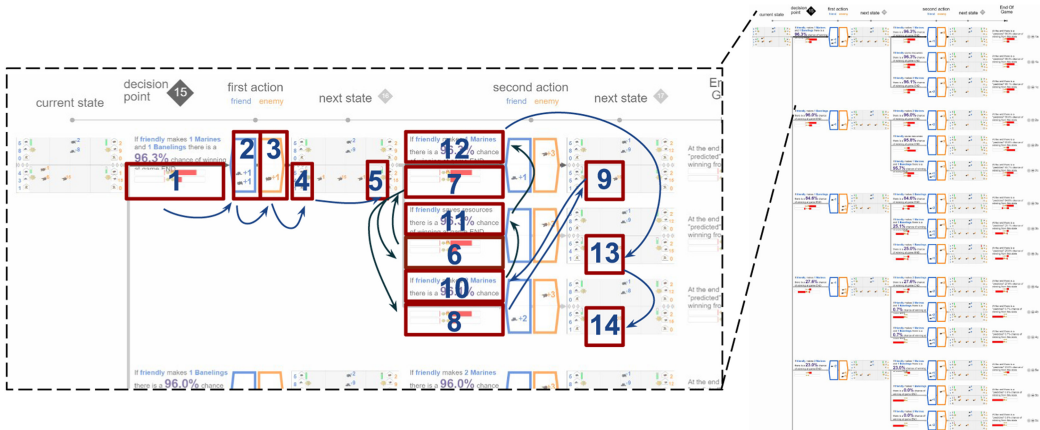


Fig. 13. P119’s search path leading to a problem report they labeled “bad decisions.” To report this problem, P119 walked down to the second level, then perused each node at the second level, drilling down further if the node seemed potentially problematic (e.g., node 8), then returned to their progression through the second level until the need arose to drill down again (e.g., node 12).

P120 (Labeled “Action-prediction incompatibility”) What: “AI adds a friendly baneling to the top row while assuming the opponent will save money. Despite this addition to units, the prediction for success does not increase from 98%. I would think it would increase.”

Counting/Math labels also correlated with participants’ successful outcomes ($r = .504$). Participants’ labels in the Counting/Math category pointed out where some number (in the game UI) or calculation (made by the AI player and shown in the explanation tree) was incorrect. Some of these referred to only to math (P102, score 3: “number off”), but some also brought in AI/XAI Concepts (P123, score 11: “chances <probabilities> are wrong”) or Domain Concepts (P111, score 12: “battlefield counting issues”). Participants’ reports with Counting/Math labels seemed to suggest that participants were localizing bugs by “auditing” the counting/math via the explanation tree:

P128 (Labeled “number”): “the numbers are inconsistent. the total marines are 9 but it only shows 6 in the game area:”

P104 (Labeled “troop_calculation”): “during action, friendly will have 4 marines, 1 each of immortal and baneling on map, enemy has 5 marines, 1 baneling. 3 friendly marines and 1 baneling was destroyed while only 1 enemy baneling was destroyed. does not seem to add up if it was a fair game.”

Domain Concept labels also were correlated with participants’ successful outcomes ($r = .455$). Participants’ use of Domain Concept labels often showed that they were localizing bugs by mapping what they saw in the explanations to game concepts, then using them to spot portions of the explanation running counter to the logic of the game. For example, P111 and P137 spotted bugs via game-logic contradictions, which they both labeled using Domain Concepts:

P111 (Labeled “evaluated battlefield error”) What: “there are 3 more marines that are alive the next time than the previous. Though the situations have not changed, so neither should the speculated battlefield.”

Why: “I’m not actually sure why, other than a misprediction.”

What changes: *“Change the 10 state bottom lane to have 3 enemy marines to match the previous prediction.”*

P137 (Labeled “Nexus health”) What: *“The health bar for Enemy bottom Nexus isn’t consistent. It’s low in the first state and then becomes full.”*

Although some examples in the Domain Concept category, were solely Domain Concept labels, about half the instances in this category also related to the AI/XAI category and/or the Counting/Math category. For example, the “evaluated battlefield error” entry above by P111 (score 12) was also in the AI/XAI Concept category. Another example was P120’s “improper unit count,” which was in both the Domain Concept category and the Counting/Math category (P120, score 5). These co-occurrences mark instances in which participants may have been reasoning about a single bug in multiple ways.

Taken together, the correlations between success and the AI/XAI Concepts, Counting/Math audits, and Domain Concepts suggest that adding labeling to the AAR/AI process may be a potentially powerful aid. Perhaps the participants’ labeling effort facilitated forms of self-explanation, in which participants were able to make sense of the individual instances of bugs by (self-)explaining them via other patterns of reasoning, such as concepts of math, RTS gameplay, or how they assumed AI agents work.

5 DISCUSSION

Did using AAR/AI impose an extra cognitive load on the AAR/AI participants beyond the load experienced by non-AAR/AI participants? We expected it to, because in our past XAI research [4], participants paid a statistically significant cognitive load cost for their successes with the most efficacious explanations. Thus, we analyzed participants’ perceptions of cognitive load via the NASA TXL, to see if adding the AAR/AI process on top of the explanations imposed a cognitive load “tax.”

To our surprise, none of the five NASA TLX dimensions gathered¹¹ showed statistically significant differences between the AAR/AI and non-AAR/AI participants; all p -values fell within [.133, .878]. That said, the temporal demand dimension was particularly interesting.

As with the other NASA/TLX questions, AAR/AI vs. non-AAR/AI participants did not differ statistically in the amount of temporal demand that they reported (TLX Temporal dimension: t -test, $t(49) = 1.5276$, $p = .1331$). However, somewhat contradicting their self-reports, AAR/AI participants actually *spent* significantly more time than non-AAR/AI participants did, for both DP8 (Welch’s t -test, $t(50) = 2.3095$, $p = .0251$) and DP15 (Welch’s t -test, $t(49) = 2.5874$, $p = .0127$). Given this contradiction between participants’ self-reported perceptions of time and actual time taken, and in light of the remaining non-significant differences in the TLX responses, further studies are required to answer whether adding AAR/AI to XAI adds cognitive load to domain-knowledgeable users’ efforts to assess their AI agents—and if so, whether the added cognitive load is desirable (e.g., participants better scrutinizing the AI) or undesirable (e.g., participants wasting time due to AAR/AI features).

As with all empirical studies, our study has limitations and threats to validity [36, 67]. One challenge with controlled experiments involving human participants in XAI is controlling *which* portions of the explanations that participants see. If participants are allowed to explore freely through a huge virtual explanation space, no two participants see the same explanations. In a quantitative experiment, such uncontrolled variations in participants’ experiences would introduce too

¹¹The “physical” dimension was not gathered because the study was online rather than in the lab. Fourteen participants (21.6%) did not complete the questionnaire: 8 out of 33 AAR/AI participants and 6 out of 32 of the non-AAR/AI participants.

much experimental noise for inferential statistics to be useful. For example, in a qualitative study in the RTS domain, Penney et al. [52, 53] showed that different participants focused on different things.

To ensure that participants in both treatments could start the experiment seeing exactly the *same* subset of the virtual explanation space, we pruned the XAI explanation tree they could view.¹² Constraining their view had the benefit of keeping their attention in parts of the explanation space where we had previously confirmed bugs. Participants could then pan/zoom/collapse the amount of information on the screen but could not expand beyond the original subset. We selected the explanation subset such that it included information pertinent to every bug. This involved pruning away large portions of the virtual explanation space, which raises an ecological validity threat. This is an example of a classic trade-off for most controlled experiments, in which some ecological validity must be traded off to achieve the controls necessary to isolate an independent variable [36] (in our study, to remove other factors so as to isolate the AAR/AI vs. non-AAR/AI independent variable).

Another threat to ecological validity is the bugs that participants needed to locate. Our bugs were naturally occurring bugs—they turned up without our help in the AI-learned model. As Ko et al. [36] point out, naturally occurring bugs increase ecological validity. However, our pilot participants revealed that some of these natural bugs were so subtle that participants rarely could locate them. This could have led to “floor effects,” in which the task is so difficult that no participant can complete them in the allotted time no matter which tool they use [57], and no effects can be revealed by statistics. To address this issue, we exaggerated the size¹³ of some of these naturally occurring bugs, as enumerated in Table 6. These exaggerations likely affected participants’ success rates; however, this threat was equally present in both treatments.

Another potential threat is in how we calculated the measures of participants’ success rates with recall and precision. Calculating recall and precision with bug identification is normally straightforward, because a single area of code either is or is not faulty, and a single bug report usually describes a single issue. However, in our experiment, a single problem report could (and sometimes did) point at multiple bugs. In addition, when two of the bugs were co-located in a single node of the explanation (Bug ID #1 and #2), attribution of participants’ bug reports to one of those two bugs was even more difficult. These complexities break the one-to-one correspondence between report and bug, yielding a multi-class, multi-label problem. Thus, we had to derive our own calculations for recall and precision, as detailed in Appendix A. The uniqueness of our recall and precision calculations affect the ability to precisely compare them against simpler precision and recall calculations used in other fault localization literature.

Ideally, participants would have been given enough time to find all the bugs (high recall), and to do so carefully enough to achieve high precision (with few false positives). In our study, neither of these measures’ averages reached 25%. A likely contributor to this was our need to control the amount of time a participant could spend in the study. Given the practicality of people’s schedules, we settled on 2-hour session times, which does not seem to have been enough time to complete the task. However, this constraint applied equally to both treatments, so it does not threaten the validity of our comparisons between treatments.

¹²Some of the virtual explanation space was not available even to us, because like many AI agents, the AI system proactively pruned away unpromising portions of its potential solution space—and, as a side effect, the explanation space—to reduce calculation time.

¹³As an additional safeguard, we included multiple exaggeration amounts for the same bug type in our experiment (e.g., some bugs being a small, medium, or large exaggeration of another bug type, as detailed in Table 6). Participants’ results did not reveal any patterns as to whether the size of the exaggeration mattered.

Another threat to participants' ability to localize the bugs is that we avoided defining what a fault/bug/problem is, because we did not want to influence participants' assessment efforts. Even when participants asked if the problem they found was really a problem, the researcher did not answer for ecological validity reasons—in software debugging, there is no “oracle” monitoring someone's debugging efforts to tell them whether or not a line of code they are puzzling over is problematic. However, this design choice introduced a threat: how could participants find a bug without knowing what constitutes one? We attempted to head off this threat by telling them that if they thought a problem was a problem, they should report it. We also provided “definitely,” “maybe,” and “never mind” bug reporting options (see Figure 2) to encourage participants to report everything they thought was even potentially problematic. If we had defined to participants what did and did not count as a bug, participants' success rates would probably have been different.

Limitations like these can be addressed only by additional studies across a spectrum of empirical methods, to isolate different independent variables of study and to establish generality of findings over different explanation styles, different bugs, different measures, different AI algorithms, different domains, and different populations attempting to find an AI agent's problematic behaviors.

6 CONCLUSION

In this article, we presented the results of an empirical study comparing domain-knowledgeable users' attempts to find an AI agent's bugs using AAR/AI vs. not using AAR/AI. The results showed the following:

- AAR/AI participants' recall rate on the bugs they reported was significantly higher than that of non-AAR/AI participants, with a large effect size. This indicates that AAR/AI participants found more of the actual bugs—one of the key goals of assessment in XAI domains.
- AAR/AI participants also reported a significantly larger number of problems. This result would be worrying if it showed that they achieved high recall simply by reporting that everything was wrong. However, the results showed that this was not the case, as explained in the next entry.
- AAR/AI participants also showed significantly higher precision than non-AAR/AI participants, with a medium effect size. Typically, there is a tradeoff in precision vs. recall, so increasing both is a very strong improvement.
- When considering the bugs one by one, we saw no evidence of AAR/AI's advantages being particular to specific bugs or types of bugs—rather, AAR/AI participants outperformed non-AAR/AI participants on *every* bug. On average, AAR/AI participants were almost six times as likely as non-AAR/AI participants to find any particular bug.
- The AAR/AI participants' labeling behaviors suggest that incorporating labeling into the AAR/AI process may bring important benefits. In this study, some AAR/AI participants used labels to abstract above individual instances of bugs, using concepts from the domain or from (X)AI. Others used labels in ways suggestive of auditing. Use of these types of labels correlated with higher recall rates in finding the bugs.

Finally, recall that the only difference in treatments was AAR/AI vs. no AAR/AI—all participants consumed the *same* explanations. These results suggest the importance of integrating explanations with an assessment process such as AAR/AI, to enable domain-knowledgeable users to make informed decisions about when to follow an AI agent's recommendations and when *not* to.

APPENDICES

A ANALYSIS MATH

Because our participants could select multiple nodes in the diagram, our analysis was faced with a multi-class, multi-label problem. In classification problems with two classes, recall and precision can be computed via familiar expressions like the following:

$$Recall_{Basic}(TP, FP, TN, FN) = \frac{TP}{TP + FN} \quad (1)$$

$$Precision_{Basic}(TP, FP, TN, FN) = \frac{TP}{TP + FP} \quad (2)$$

To find a multi-class, multi-label analog, we consulted a review by Zhang and Zhou [71], which describes two flavors of multi-label analysis. One equally weights each data *instance*, and another equally weights each *label*. Either one is well defined for us, so we picked the former approach, using the recall/precision equations from their Section 2.2.2, as Equations (3) and (4) provided verbatim here, barring a single notational change, where

- $h(\cdot)$ is the classification function, which returns labels given the i th data point as a feature vector \vec{x}_i ;
- d is the number of data instances (Zhang and Zhou used p , but we will use that later; this notation swap is the only change from their equation); and
- Y_i is the set of ground truth labels associated with the i th data point.

$$Recall_{Zhang}(h) = \frac{1}{d} \sum_{i=1}^d \frac{|Y_i \cap h(\vec{x}_i)|}{|Y_i|} \quad (3)$$

$$Precision_{Zhang}(h) = \frac{1}{d} \sum_{i=1}^d \frac{|Y_i \cap h(\vec{x}_i)|}{|h(\vec{x}_i)|} \quad (4)$$

We started from these expressions, and cast the summands into our situation via the following steps. First, $|Y_i|$ is just the number of bugs present (in this case, 10). Second, since our bugs were all known to be present, the \vec{Y} for a particular DP is a 1-vector, so the intersection becomes a sum of the bugs a participant found in *all* their reports. Third, since $|h(\vec{x}_i)|$ is intended to model the “number of shots fired at targets,” we use the number of problem reports that participant submitted as the denominator in precision (considering that one could consider 0 reports to be 0 precision because bugs were known to be present, participants provided 2 reports at minimum). Importantly, in this framing, the denominator has no dependence on the summation, and so it can be pulled outside the summation.

Next, we need to manipulate the summation part to handle reports not being independent. In our case, we do not get negative data instances¹⁴ because we did not request any certifications that (regions of) the explanation were free of bugs, although a few participants saw fit to submit such reports. This means we need to interpret silence on a bug as a *FN*, but we can *only* know if the participant was silent on a bug after they have finished with that subtask. Further, we will never get a *TN* because bugs were known to be present. The other kind of non-independence we need to handle is best illustrated by Table 4: many participants reported the same problem multiple times. To handle this without awarding excess credit, we create Table 5 so that it is sliced per *participant*

¹⁴If confused by this terminology, consult the confusion matrix at <https://developers.google.com/machine-learning/crash-course/classification/true-false-positive-negative>.

Table 4. Mockup of Data Post Labeling, Presented Per *Problem Report*

ReportID	PID	Bug A		Bug B		Bug C	
		Describe	Select	Describe	Select	Describe	Select
1	Alice	0	0	1	1	0	1
2	Alice	0	0	0	0	1	1
3	Alice	0	0	0	0	0	0
4	Bob	0	0	0	0	0	0
5	Bob	1	0	0	0	0	0
6	Cindy	0	0	0	1	0	0
7	Cindy	0	1	1	1	0	0
8	Cindy	0	0	1	1	0	0
9	Cindy	0	0	0	1	0	0

In our labeling, each bug could be selected and/or described properly, so 2 points available per bug. Notably, as illustrated here, participants often found the same issues repeatedly, so we devised Equations (5) and (6) to handle not awarding additional credit for these repeat finds. In the example provided, we used binary indicator variables for simplicity of presentation, although our formulation naturally handles partial credit with no modification.

Table 5. Mockup of Data Post Labeling, Per Participant (i.e., After Taking the Max Across Reports from That Person)

PID	#Reports	Bug A		Bug B		Bug C	
		Describe	Select	Describe	Select	Describe	Select
Alice	3	0	0	1	1	1	1
Bob	2	1	0	0	0	0	0
Cindy	4	0	1	1	1	0	0

To compute recall/precision, one can take the sum across rows of this table, normalize by 2 (because there are 2 points per bug), then divide by the number of bugs or problem reports, respectively.

by aggregating each participant's part of columns of Table 4 by taking a max over the reports from each participant.

Putting the pieces together yields Equations (5) and (6), for participant p 's recall/precision, where

- B is the set of bugs (in our case, there are 10);
- M is the maximum score available per bug (in our case, 2);
- j indexes particular labels (of which there are BM in each report, because each bug could be described and/or labeled correctly); and
- $R(\cdot)$ is the report function, which returns the set of problem reports from a given participant.

$$Recall_{participant}(p) = \frac{1}{|B|M} \sum_j^{|B|M} \max_{r \in R(p)} r_j \quad (5)$$

$$Precision_{participant}(p) = \frac{1}{|R(p)|M} \sum_j^{|B|M} \max_{r \in R(p)} r_j \quad (6)$$

Calculating recall per *bug* is very similar, with the summation going down columns instead of across rows, but precision becomes a bit more complicated. Because some reports could be attributed to bugs while others could not, determining the correct number of reports to divide

by requires checking to see if a report was attributed to a *different* bug. Mathematically, we use Equations (7) and (8) to compute the recall of a label b (e.g., “Bug A Select”), where

- P is the total number of participants;
- p indexes participants; and
- $R(\cdot, \cdot)$ is an overloaded function returning the set of problem reports for a participant that could possibly be attributed to a particular bug (either by actually being attributed to that bug *or* by being unable to be attributed to any bug).

$$Recall_{Label}(b) = \frac{1}{P} \sum_p \max_{r \in R(p)} r_b \quad (7)$$

$$Precision_{Label}(b) = \frac{1}{\sum_p |R(p, b)|} \sum_p \max_{r \in R(p)} r_b \quad (8)$$

Using these as the basis, we can analyze a bug (or kind of bug, e.g., LEF bugs) by running this on multiple columns and normalizing results appropriately.

To illustrate usage of Equations (5) and (6), one can compute the *recall* of a participant using Table 5; we can sum across the row and divide by the number of bugs (the max was already incorporated moving from Table 4 to Table 5). Similarly, for *precision*, we use the same sum across the row but divide by the number of problem reports the participant provided. In the concrete example, $|B| = 3$ and $M = 2$, so Alice’s recall is a whopping $\frac{4}{6} \approx 67\%$ with the same precision (Alice submitted three bug reports). Meanwhile, Bob’s recall was only $\frac{1}{6} \approx 17\%$ with precision a bit higher at $\frac{1}{2 * M} = 25\%$ because he did not provide many reports. Cindy, however, had recall at $\frac{3}{6} = 50\%$ but with lower precision ($\frac{3}{4 * M} \approx 38\%$) because of the higher report volume.

Next, to illustrate usage of Equation (7), one of the three participants (Bob) found “Bug A Describe,” so the recall for that bug would be 33%. Similarly, one of the three participants selected Bug A, leading to a 33% recall overall for Bug A.

Finally, to illustrate usage of Equation (8), we compute the precision for “Bug B Describe.” Alice went 1-for-2 (report 1 hit, report 2 was attributable to Bug C, and report 3 was unable to be attributed and so possibly intended for that target). Meanwhile, Bob went 0-for-1 (report 4 was possibly intended for that target, whereas report 5 was attributable to Bug A). And last, Cindy went 1-for-4 (report 6 missed but was intended for Bug B, report 7 hit, report 8 hit but was duplicate, and report 9 missed but was intended for Bug B). Combining these results yields $\frac{2}{7} \approx 29\%$.

B STUDY DESIGN ADDITIONAL DETAILS

B.1 Details of the StarCraft II Game

In our StarCraft II setup, both players were RL-powered agents. We refer to them as the “Friendly AI” (referring back to Figure 1’s left) and the “Enemy AI” (right). Participants were shown the game from the Friendly AI’s perspective.

As Figure 1 shows, the game has a top and a bottom lane, each of which is a separate battlefield between the AI players. Each AI player has two nexuses, which is the AI’s base in this game. A nexus is represented by the golden/yellow star structure at the corner of each lane. Next to each nexus is a health bar with that nexus’s corresponding **health points (HPs)**.

The two ways an AI player can win are by (1) destroying one of the opposing nexuses before 40 rounds, by bringing the nexus health to 0, or (2) having the lowest nexus health if all nexuses are standing at the end of 40 rounds. Thus, in trying to win, throughout the game the AIs generate troops behind their respective nexus to cause damage to the opposing nexus in their lane.

Table 6. The 10 Bugs: 5 in DP 8 (Above the Line) and 5 in DP 15 (Below the Line)

ID-Type	Description (<i>and how we engineered them</i>)
1-LEF	<i>Why a bug:</i> Because the actions preceding state at DP 10 (i.e., predicted future states from DP 8) differ by 1 marine from its sibling actions, yet the expected outcome is radically different (flipped) than all other sibling actions. A correct state would have similar outcome expectations to the sibling states. <i>Exaggerated by changing:</i> Friendly AI's win by destroying top Enemy nexus to Friendly AI's loss by Enemy AI destroying Friendly AI's bottom nexus.
2-TF	<i>Why a bug:</i> Because the Friendly agent has 4 marine-producing buildings in the top lane, but there are 21 total Friendly marines expected to be in the top lane, State at DP 10, row 1C. A correct state would have 4 or fewer marines in the top lane. <i>Exaggerated by changing:</i> Friendly marines top grid #1 to 19. Friendly marines top grid #2 to 2.
3-TF	<i>Why a bug:</i> Because base (nexus) health cannot heal. In expected state DP 9 row 2A, the Enemy bottom base is expected to incur a significant amount of damage as shown by the red bar. However, in predicted child state DP 10 row 2C, that damage is not reflected as the base's health in DP 10 row 2C is greater than its health in DP 9 row 2A. A correct state would not have greater Friendly Base HP in DP 10 row 2C. <i>Exaggerated by changing:</i> Friendly Base HP D9 row 2A to 20; Friendly Base HP D10 row 2B to 100.
4-TF	<i>Why a bug:</i> Because the Enemy has built one immortal-producing building in the preceding action and can therefore have at most one immortal troop on the field. This state depicts two immortal troops in adjacent game grid. A correct state would be one or fewer Enemy immortals in the bottom lane. <i>Exaggerated by changing:</i> Enemy immortal bottom grid #4 to 1; Enemy immortal grid #3 to 1.
5-LEF	<i>Why a bug:</i> Because the Friendly base is expected to have 0 HP meaning, it has been destroyed; therefore, per game rules, the Friendly agent has lost the game at DP 10. However, the expected outcome shows that the Friendly agent expects to destroy the Enemy's top and bottom base. A correct state would have been for the Enemy win 100% by destroying a Friendly bottom base. <i>Exaggerated by changing:</i> Friendly bottom base HP to 0.
6-TF	<i>Why a bug:</i> Recall the game rule that the top and bottom lanes are independent; troops built in one lane will not affect the outcome of what happens in the other. From DP 16 row 1A to all child actions and states, the Enemy builds three marine-producing buildings in the top lane and the Friendly does not build anything in the top lane for all three actions from DP 16 to DP 17 (rows 1A, 1B, and 1C). Since the action for the top lane is the same for all three sibling predictions, the expected state in the top lane in DP 17 in all three sibling states should all be the same. However, state DP 17 row 1C differs from its siblings, and one fewer marine is expected in DP 17 row 1C. <i>No exaggeration needed.</i>
7-LEF	<i>Why a bug:</i> Because all sibling child states reflect a likely win in the top lane with a less likely win by purchasing troops in the bottom lane, but the bars in row 2B depict a high expectation to win in the bottom lane with no purchases in the bottom lane. <i>Exaggerated by changing:</i> Friendly win by destroying Enemy bottom to 76%, Friendly win by destroying Enemy top to 20%.
8-LEF	<i>Why a bug:</i> Because the Enemy agent has zero immortal-producing buildings in state DP 17, yet it is expected to have one immortal troop in the bottom lane. <i>Exaggerated by changing:</i> Enemy immortals bottom grid #2 to 1.
9-LEF	<i>Why a bug:</i> Because in this predicted state the game is guaranteed to end with the Enemy top base getting destroyed, but outcome bars show the Friendly expecting to lose. <i>Exaggerated by changing:</i> Enemy top base HP to 0 in row 3C.
10-TF	<i>Why a bug:</i> Because the game is guaranteed to end with the Friendly bottom base getting destroyed sibling states in rows 5B and 5C, both have the correct outcome expectations (0% win, lose by Friendly bottom base destroy); however, in row 5A the Friendly expects to win by destroying Enemy top. <i>Exaggerated by changing:</i> Friendly win by destroying Enemy top to 23%. State D10, rows 5A, 5B, 5C set Friendly bottom base HP to 0.

LEF, in LEF outputs; TF, in TF outputs (Section 3.3). More information about the bugs and their locations can be found in the supplemental documents. All bugs occurred naturally, but we exaggerated some as marked.

The bar with diamonds at the bottom of the screen in Figure 1 is the game timeline. StarCraft II is played in rounds: each new round starts at one of the black diamonds (marked as D1, D2, etc). These are called *decision points* (DPs), each marking a point where the AI decides on an action before the next round begins. At a DP, the AI decides (1) what troops to buy, if any, and (2) which lane to place them in (top or bottom).

To engage in battle, the AIs need minerals, because minerals are akin to money: the AI can use them to buy troops or invest in pylons. At the start of the game, each AI is given 150 minerals and then receives 100 minerals in every subsequent round. A pylon generates an additional 75 minerals per round, and each AI can buy up to 3 pylons in the game. The pylons for the Friendly and Enemy AIs are represented by the three diamonds in the center of the column on each side. The AIs spend their minerals to buy troop production buildings, which produce troops. Once an AI buys a troop production building, one unit of that type is generated in each subsequent round. After being generated, these troops will run toward the opposing nexus and fight any units in their way. The AI cannot control what units attack the other, or troop formations; it can only buy the troop production buildings to generate troops.

The three types of troop production buildings in this game are marines, banelings, and immortals, all of which are shown in Figure 1. Marines are small, low-cost units. They have the lowest health of the three troops and attack with small, quick shots. Banelings are explosive bug units with a moderate cost. They have medium health, and they explode upon contact with an enemy unit. Immortals are large and are also the most expensive unit. They have the highest health, since they have a shield. They attack with slow shots that are effective against nexuses. There is a rock-papers-scissors relationship between the marines, banelings, and immortals. Marines are effective against immortals, immortals against banelings, and banelings against marines.

B.2 Experiment Session Walkthrough

Participants were given a tutorial detailing the game rules (specified earlier). Next, they were given access to the web interface with a unique ID and password. After entering their credentials, they were prompted by the facilitator to click on “Play” and start watching the game.

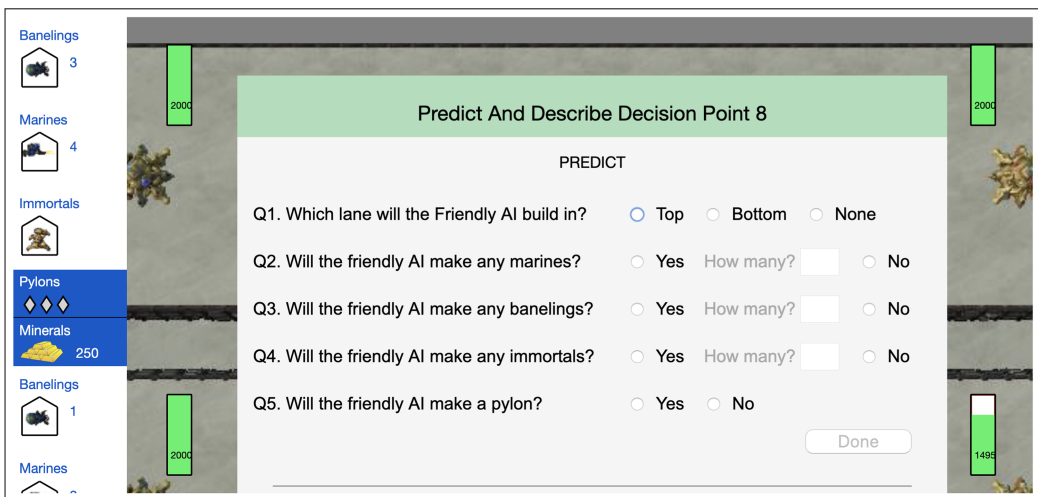


Fig. 14. Participants watch the game unfold until round 8, at which point the game pauses and the prediction questions pop up. The participant has to predict what they think the Friendly AI will do in round 9.

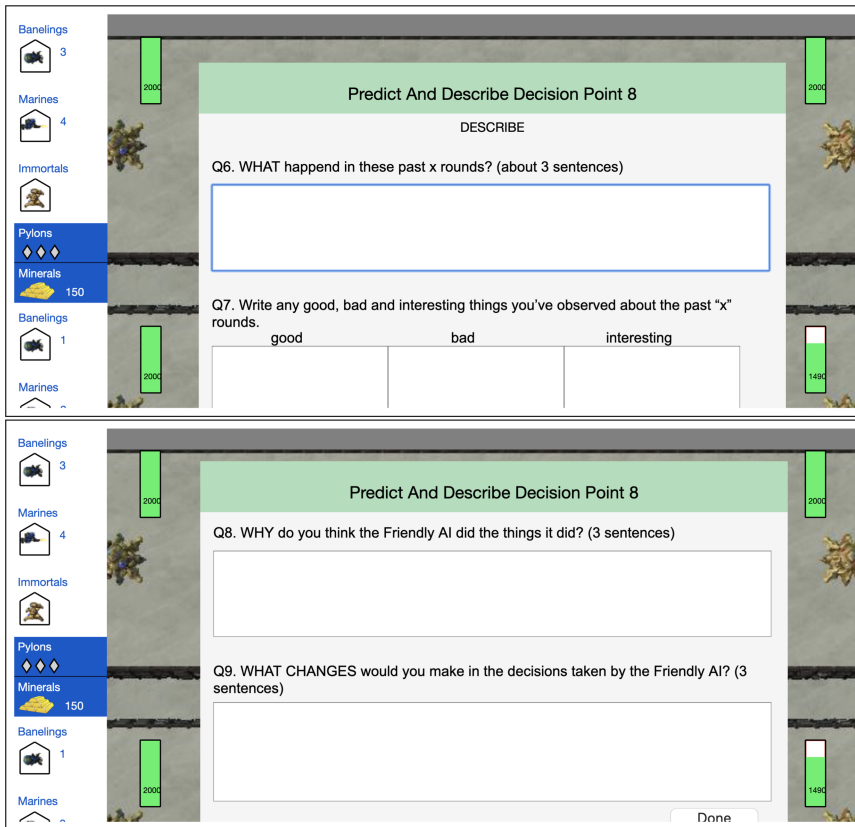


Fig. 15. Having predicted what the agent would do (Figure 14), participants then see what it really did. Here, AAR/AI participants watch the game until round 9, at which point the game pauses again and the description questions pop up. The participant first describes what happened in round 9 (top) and then why (bottom).

Participants in both treatments saw identical AI agents, game replays, and explanations. Their agents also exhibited identical bugs, which are detailed in Table 6.

The game automatically paused at DP 8, and participants in both treatments were shown identical “prediction” questions (Figure 14). After answering these questions, participants watched the game round. After the game round ended, participants answered “description” questions, where they described the Friendly AI’s actions in the round they had just watched. The AAR/AI group was given a guided process (Figure 15), whereas the non-AAR/AI group was given observational questions (Figure 16).

After answering the prediction and description questions, participants in both treatments then saw the Friendly AI’s explanation for its actions. Figure 17 shows the AAR/AI interface, and Figure 18 shows the non-AAR/AI interface. All participants saw the same explanations, but as these figures show, the different treatments asked different questions about the bugs participants found. In addition, AAR/AI participants had to finish the row they had selected to work on before moving onto another row as part of the AAR/AI process, whereas non-AAR/AI participants were allowed to navigate freely among rows.

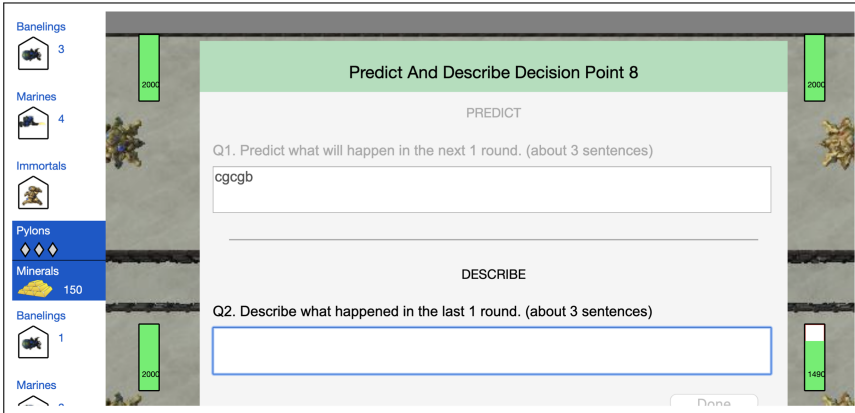


Fig. 16. Having predicted what the agent would do (Figure 14), participants then see what it really did. Here, non-AAR/AI participants watch the game until round 9, at which point the game paused again and the description questions pop up. The participant has to describe what happened in round 9.

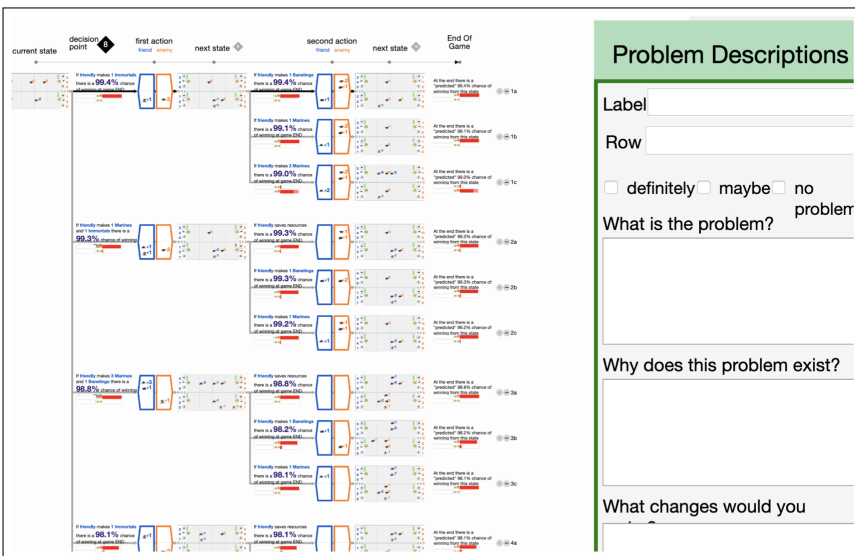


Fig. 17. AAR/AI participants answered the AAR/AI questions of “What-Why-What changes” for each problem they found, in addition to giving their problem descriptions labels.

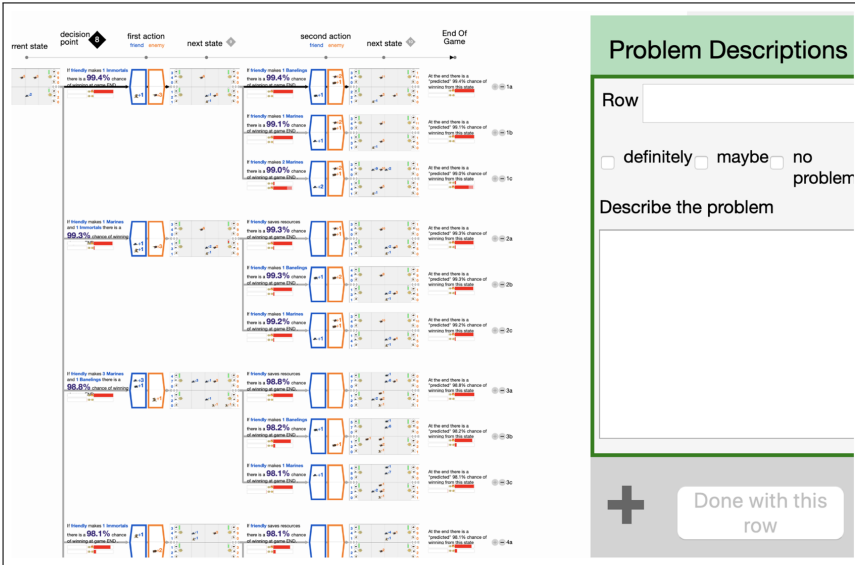


Fig. 18. Non-AAR/AI participants had less structure than their AAR/AI counterparts, simply indicating the location and severity, then offering a description.

C MODEL-BASED AGENT ARCHITECTURE

This section describes details about the model-based agent we used for the studies. Figure 19 illustrates the overall architecture of the agent. We constructed a minimax search tree by combining a *decomposed reward deep Q-network (drDQN)* [32], used for *action ranking* and *leaf evaluation*, and a transition model [41]. We describe the details of both in the next two sections.

Model training details. Both the drDQN and transition models are neural networks that were pre-trained. They were trained separately and frozen while the agent plays. The training process continued until the agent achieved a high win percentage against a pool of agents or until resources were expended. This took around 3 days on a consumer desktop machine. In other words, the

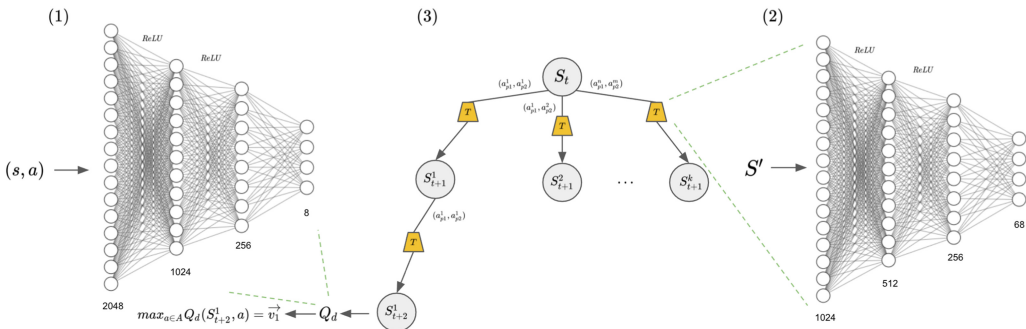


Fig. 19. The model-based agent consists of three parts. (1, left) The drDQN model, which takes a state-action pair (s, a) as input and outputs a decomposed Q-value vector. (2, right) The transition model, which outputs the estimated next game state by taking the after-state S' as input. (3, middle) The tree structure that utilizes (1) and (2) together. Here, the minimax algorithm assigns the Q-value vectors computed at the leaf all the way to the root.

model-based agent in its entirety does *not* have a training process. Both networks have three fully connected layers, and each hidden layer uses ReLU as its activation function. They have different numbers of neurons and output sizes, as indicated at the bottom of each layer in Figure 19. We used *mean squared errors* as the loss functions for both models. The learning rates of the drDQN and transition model were 10^{-4} and 5^{-4} , respectively.

C.1 Decomposed Reward Deep Q-Network (for Action Ranking and Leaf Evaluation)

The purpose of using a drDQN agent instead of a standard deep Q-network agent is that rather than only a single Q-value, it provides a more explanatory *vector* at the leaf nodes. In our case, we decomposed the Q-value (which is a scalar win probability) into an eight-element vector composed of the probability of each nexus being destroyed and the probability of each nexus having the lowest HP if the game reaches the tie-breaker. Therefore, we can compute the win probability for a single player by taking the sum of winning by destroying each opponent nexus (two elements) and by tie-breaking each opponent nexus (two elements). Further, the sum of all eight elements should be 1.0, since it represents a probability distribution.

The drDQN model was pre-trained via pool-based self-play learning to achieve a reasonable high win probability, which provides a meaningful decomposed Q-value vector for the leaf nodes of the minimax tree. Because the size of the minimax tree grows exponentially in its depth, we cannot expand it to the end of the game. To cope with this, we use the drDQN to prune the tree, declining to expand actions that do not look promising (this is the ARF referred to in Section 3.3). Therefore, evaluating leaf nodes with a neural network is important because it predicts the value of the future based on the leaf states—*without* expanding the tree further (this is the LEF referred to in Section 3.3). The decomposed Q-value function provides a discounted accumulation value vector predicting the future based on a state-action pair. The leaf node value vector $v = \max_{a \in A} Q(s, a)$ is back-propagated back to the root node, where A is the action space and (s, a) is the state-action pair. Thus, we use the *same* drDQN twice in same agent, for both ranking actions and evaluating leaf nodes.

C.2 Transition Model

The transition model was also pre-trained by supervised learning based on the dataset from running the drDQN agent playing against an opponent pool that includes several types of agent. It takes an after-state S' as input that combines the current state S , Player 1 (Friendly AI)'s action a_{p1} , and Player 2 (Enemy AI)'s action a_{p2} . Since actions in Tug-of-War correspond to an integer vector corresponding to the buildings the player is going to create, the action is deterministic and can be simply *added* up with the current state to produce an after-state corresponding to a tuple (s, a_{p1}, a_{p2}) . It outputs an estimated state that describes the game at the next DP. The estimated state has exactly elements as the input state has, which includes mineral resources, number of buildings, number of troops in different regions for both lanes, Nexus HP, and the wave number.

ACKNOWLEDGMENTS

We thank Tom Dietterich for words of wisdom during the research that produced this article.

REFERENCES

- [1] Ashraf Abdul, Jo Vermeulen, Danding Wang, Brian Y. Lim, and Mohan Kankanhalli. 2018. Trends and trajectories for explainable, accountable and intelligible systems: An HCIresearch agenda. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (CHI'18)*. ACM, New York, NY, Article 582, 18 pages.
- [2] Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, et al. 2019. Guidelines for human-AI interaction. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI'19)*. ACM, New York, NY, Article 3, 13 pages.

- [3] Paul Ammann and Jeff Offutt. 2016. *Introduction to Software Testing*. Cambridge University Press.
- [4] Andrew Anderson, Jonathan Dodge, Amrita Sadarangani, Zoe Juozapaitis, Evan Newman, Jed Irvine, Souti Chattopadhyay, Matthew Olson, Alan Fern, and Margaret Burnett. 2020. Mental models of mere mortals with explanations of reinforcement learning. *ACM Transactions on Interactive Intelligent Systems* 10, 2 (May 2020), Article 15, 37 pages. <https://doi.org/10.1145/3366485>
- [5] Andrew Anderson, Jonathan Dodge, Amrita Sadarangani, Zoe Juozapaitis, Evan Newman, Jed Irvine, Souti Chattopadhyay, Alan Fern, and Margaret Burnett. 2019. Explaining reinforcement learning to mere mortals: An empirical study. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI'19)*.
- [6] Algirdas Avizienis, J.-C. Laprie, Brian Randell, and Carl Landwehr. 2004. Basic concepts and taxonomy of dependable and secure computing. *IEEE Transactions on Dependable and Secure Computing* 1, 1 (2004), 11–33.
- [7] Adam T. Biggs and Stephen R. Mitroff. 2015. Improving the efficacy of security screening tasks: A review of visual search challenges and ways to mitigate their adverse effects. *Applied Cognitive Psychology* 29, 1 (2015), 142–148. <https://doi.org/10.1002/acp.3083>
- [8] Benjamin S. Bloom, Max D. Engelhart, Edward J. Furst, Walker H. Hill, and David R. Krathwohl. 1956. *Taxonomy of Educational Objectives*. Longmans, Green & Co. Ltd.
- [9] Ralph Brewer, Anthony Walker, E. Ray Pursel, Eduardo Cerame, Anthony Baker, and Kristin Schaefer. 2019. Assessment of manned-unmanned team performance: Comprehensive after-action review technology development. In *Proceedings of the 2019 International Conference on Human Factors in Robots and Unmanned Systems (AHFE'19)*. 119–130.
- [10] Ruth M. J. Byrne. 2019. Counterfactuals in explainable artificial intelligence (XAI): Evidence from human reasoning. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI'19)*. 6276–6282. <https://doi.org/10.24963/ijcai.2019/876>
- [11] Nico Castelli, Corinna Ogonowski, Timo Jakobi, Martin Stein, Gunnar Stevens, and Volker Wulf. 2017. What happened in my home? An end-user development approach for smart home data visualization. In *Proceedings of the ACM Conference on Human Factors in Computing Systems*. ACM, New York, NY, 853–866.
- [12] Nan-Chen Chen, Jina Suh, Johan Verwey, Gonzalo Ramos, Steven Drucker, and Patrice Simard. 2018. AnchorViz: Facilitating classifier error discovery through interactive semantic data exploration. In *Proceedings of the 23rd International Conference on Intelligent User Interfaces (IUI'18)*. ACM, New York, NY, 269–280.
- [13] Michelene T. H. Chi, Miriam Bassok, Matthew W. Lewis, Peter Reimann, and Robert Glaser. 1989. Self-explanations: How students study and use examples in learning to solve problems. *Cognitive Science* 13, 2 (4 1989), 145–182. https://doi.org/10.1207/s15516709cog1302_1
- [14] Michael Chromik, Malin Eiband, Felicitas Buchner, Adrian Krüger, and Andreas Butz. 2021. I think I get your point, AI! The illusion of explanatory depth in explainable AI. In *Proceedings of the 26th International Conference on Intelligent User Interfaces (IUI'21)*. ACM, New York, NY, 307–317. <https://doi.org/10.1145/3397481.3450644>
- [15] Michael Chromik, Malin Eiband, Sarah Theres Völkel, and Daniel Buschek. 2019. Dark patterns of explainability, transparency, and user control for intelligent systems. In *Proceedings of IUI Workshops*.
- [16] Jacob Cohen. 2013. *Statistical Power Analysis for the Behavioral Sciences*. Academic Press.
- [17] Kelley Cotter, Janghee Cho, and Emilee Rader. 2017. Explaining the news feed algorithm: An analysis of the “news feed FYI” blog. In *ACM CHI Conference Extended Abstracts on Human Factors in Computing Systems*. ACM, New York, NY, 1553–1560.
- [18] Robert Davies, Elly Vaughan, Graham Fraser, Robert Cook, Massimo Ciotti, and Jonathan E. Suk. 2019. Enhancing reporting of after action reviews of public health emergencies to strengthen preparedness: A literature review and methodology appraisal. *Disaster Medicine and Public Health Preparedness* 13, 3 (June 2019), 618–625. <https://doi.org/10.1017/dmp.2018.82>
- [19] Jonathan Dodge, Roli Khanna, Jed Irvine, Kin-Ho Lam, Theresa Mai, Zhengxian Lin, Nicholas Kiddle, et al. 2021. After-action review for AI (AAR/AI). *ACM Transactions on Interactive Intelligent Systems* 11, 3-4 (2021), Article 29, 35 pages.
- [20] Jonathan Dodge, Sean Penney, Claudia Hilderbrand, Andrew Anderson, and Margaret Burnett. 2018. How the experts do it: Assessing and explaining agent behaviors in real-time strategy games. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (CHI'18)*. ACM, New York, NY, Article 562, 12 pages.
- [21] Luciano Floridi, Josh Cows, Monica Beltrametti, Raja Chatila, Patrice Chazerand, Virginia Dignum, Christoph Luetge, et al. 2018. AI4People—An ethical framework for a good AI society: Opportunities, risks, principles, and recommendations. *Minds and Machines* 28, 4 (2018), 689–707.
- [22] David Freedman, Robert Pisani, and Roger Purves. 2007. *Statistics* (4th international student ed.). W. W. Norton, New York, NY.
- [23] A. Groce, T. Kulesza, C. Zhang, S. Shamasunder, M. Burnett, W. Wong, S. Stumpf, et al. 2014. You are the only possible oracle: Effective test selection for end users of interactive machine learning systems. *IEEE Transactions on Software Engineering* 40, 03 (March 2014), 307–323. <https://doi.org/10.1109/TSE.2013.59>

- [24] Samer Hanoun and Saeid Nahavandi. 2018. Current and future methodologies of after action review in simulation-based training. In *Proceedings of the 2018 Annual IEEE International Systems Conference (SysCon'18)*. IEEE, Los Alamitos, CA, 1–6.
- [25] Sandra G. Hart and Lowell E. Staveland. 1988. Development of NASA-TLX (task load index): Results of empirical and theoretical research. In *Human Mental Workload*, Peter A. Hancock and Najmedin Meshkati (Eds.). Advances in Psychology, Vol. 52. North-Holland, 139–183. [https://doi.org/10.1016/S0166-4115\(08\)62386-9](https://doi.org/10.1016/S0166-4115(08)62386-9)
- [26] Robert Hoffman, Gary Klein, and Shane Mueller. 2018. Explaining explanation for “Explainable AI.” *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 62 (Sept. 2018), 197–201. <https://doi.org/10.1177/1541931218621047>
- [27] Robert R. Hoffman, Shane T. Mueller, Gary Klein, and Jordan Litman. 2018. Metrics for explainable AI: Challenges and prospects. *CoRR* abs/1812.04608 (2018). arXiv:1812.04608 <http://arxiv.org/abs/1812.04608>.
- [28] Fred Hohman, Minsuk Kahng, Robert Pienta, and Duen Horng Chau. 2019. Visual analytics in deep learning: An interrogative survey for the next frontiers. *IEEE Transactions on Visualization and Computer Graphics* 25, 8 (2019), 2674–2693. <https://doi.org/10.1109/TVCG.2018.2843369>
- [29] Hsiu-Fang Hsieh and Sarah E. Shannon. 2005. Three approaches to qualitative content analysis. *Qualitative Health Research* 15, 9 (2005), 1277–1288.
- [30] Andrew Ishak and Elizabeth Williams. 2017. Slides in the tray: How fire crews enable members to borrow experiences. *Small Group Research* 48, 3 (March 2017), 336–364. <https://doi.org/10.1177/1046496417697148>
- [31] Paul Jaccard. 1908. Nouvelles recherches sur la distribution florale. *Bulletin de la Societe Vaudoise des Sciences Naturelles* 44, 163 (1908), 223–270.
- [32] Zoe Juozapaitis, Anurag Koul, Alan Fern, Martin Erwig, and Finale Doshi-Velez. 2019. Explainable reinforcement learning via reward decomposition. In *Proceedings of the IJCAI 2019 Workshop on Explainable Artificial Intelligence*. 47–53.
- [33] Nathanael Keiser and Winfred Arthur Jr. 2020. A meta-analysis of the effectiveness of the after-action review (or debrief) and factors that influence its effectiveness. *Journal of Applied Psychology* 106, 7 (Aug. 2020), 1007–1032. <https://doi.org/10.1037/apl0000821>
- [34] Man-Je Kim, Kyung-Joong Kim, SeungJun Kim, and Anind Dey. 2016. Evaluation of StarCraft artificial intelligence competition bots by experienced human players. In *Proceedings of the 2016 CHI Conference Extended Abstracts on Human Factors in Computing Systems (CHI EA'16)*. ACM, New York, NY, 1915–1921.
- [35] Man-Je Kim, Kyung-Joong Kim, SeungJun Kim, and Anind K. Dey. 2016. Evaluation of StarCraft artificial intelligence competition bots by experienced human players. In *ACM CHI Conference Extended Abstracts*. ACM, New York, NY, 1915–1921.
- [36] A. J. Ko, T. D. Latoza, and M. M. Burnett. 2015. A practical guide to controlled experiments of software engineering tools with human participants. *Empirical Software Engineering* 20, 1 (2015), 110–141.
- [37] Cliff Kuang. 2017. Can AI be taught to explain itself? *New York Times*. Retrieved December 26, 2017 from <https://www.nytimes.com/2017/11/21/magazine/can-ai-be-taught-to-explain-itself.html>.
- [38] T. Kulesza, M. Burnett, W. Wong, and S. Stumpf. 2015. Principles of explanatory debugging to personalize interactive machine learning. In *Proceedings of the ACM International Conference on Intelligent User Interfaces*. ACM, New York, NY, 126–137.
- [39] Todd Kulesza, Simone Stumpf, Margaret Burnett, Weng-Keen Wong, Yann Riche, Travis Moore, Ian Oberst, Amber Shinsel, and Kevin McIntosh. 2010. Explanatory debugging: Supporting end-user debugging of machine-learned programs. In *Proceedings of the 2010 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC'10)*. IEEE, Los Alamitos, CA, 41–48.
- [40] T. Kulesza, S. Stumpf, M. Burnett, S. Yang, I. Kwan, and W. K. Wong. 2013. Too much, too little, or just right? Ways explanations impact end users’ mental models. In *Proceedings of the 2013 IEEE Symposium on Visual Languages and Human Centric Computing (VL/HCC'13)*. 3–10. <https://doi.org/10.1109/VLHCC.2013.6645235>
- [41] Kin-Ho Lam, Zhengxian Lin, Jed Irvine, Jonathan Dodge, Zeyad T. Shureih, Roli Khanna, Minsuk Kahng, and Alan Fern. 2020. Identifying reasoning flaws in planning-based RL using tree explanations. In *Proceedings of the IJCAI-PRICAI 2020 Workshop on XAI*. <https://drive.google.com/file/d/1ihT39-S2SFpCsarJfzDsnMOU85YjCXK/view?usp=sharing>.
- [42] Adam Lareau and Brice Long. 2018. The art of the after-action review. *Fire Engineering* 171, 5 (May 2018), 61–64. <http://search.proquest.com/docview/2157468757/>.
- [43] Brian Y. Lim. 2012. *Improving Understanding and Trust with Intelligibility in Context-Aware Applications*. Ph.D. Dissertation. Carnegie Mellon University.
- [44] Brian Y. Lim and Anind K. Dey. 2009. Assessing demand for intelligibility in context-aware applications. In *Proceedings of the ACM International Conference on Ubiquitous Computing*. ACM, New York, NY, 195–204.

- [45] Sandra Deacon Lloyd Baird, Phil Holland. 1999. Learning from action: Imbedding more learning into the performance fast enough to make a difference. *Organizational Dynamics* 27 (1999), 19–32. [https://doi.org/10.1016/S0090-2616\(99\)90027-X](https://doi.org/10.1016/S0090-2616(99)90027-X)
- [46] Theresa Mai, Roli Khanna, Jonathan Dodge, Jed Irvine, Kin-Ho Lam, Zhengxian Lin, Nicholas Kiddle, et al. 2020. Keeping it “organized and logical”: After-action review for AI (AAR/AI). In *Proceedings of the 25th International Conference on Intelligent User Interfaces (IUI’20)*. ACM, New York, NY. <https://doi.org/10.1145/3377325.3377525>
- [47] Ronald Metoyer, Simone Stumpf, Christoph Neumann, Jonathan Dodge, Jill Cao, and Aaron Schnabel. 2010. Explaining how to play real-time strategy games. *Knowledge-Based Systems* 23, 4 (2010), 295–301.
- [48] Nicole Mirmig, Gerald Stollnberger, Markus Miksch, Susanne Stadler, Manuel Giuliani, and Manfred Tscheligi. 2017. To err is robot: How humans assess and act toward an erroneous social robot. *Frontiers in Robotics and AI* 4 (2017), 21.
- [49] Donald A. Norman. 1983. Some observations on mental models. *Mental Models* 7, 112 (1983), 7–14.
- [50] European Group on Ethics in Science and New Technologies. 2018. Statement on artificial intelligence, robotics and “Autonomous” Systems. Retrieved Novembewr 17, 2021 from https://ec.europa.eu/info/news/ethics-artificial-intelligence-statement-ege-released-2018-apr-24_en.
- [51] S. Ontañón, G. Synnaeve, A. Uriarte, F. Richoux, D. Churchill, and M. Preuss. 2013. A survey of real-time strategy game AI research and competition in StarCraft. *IEEE Transactions on Computational Intelligence and AI in Games* 5, 4 (Dec. 2013), 293–311. <https://doi.org/10.1109/TCIAIG.2013.2286295>
- [52] Sean Penney, Jonathan Dodge, Andrew Anderson, Claudia Hilderbrand, Logan Simpson, and Margaret Burnett. 2021. The shoutcasters, the game enthusiasts, and the AI: Foraging for explanations of real-time strategy players. *ACM Transactions on Interactive Intelligent Systems* 11, 1 (2021), Article 2, 46 pages. <https://doi.org/10.1145/3396047>
- [53] Sean Penney, Jonathan Dodge, Claudia Hilderbrand, Andrew Anderson, Logan Simpson, and Margaret Burnett. 2018. Toward foraging for understanding of StarCraft agents: An empirical study. In *Proceedings of the 23rd International Conference on Intelligent User Interfaces (IUI’18)*. ACM, New York, NY, 225–237. <https://doi.org/10.1145/3172944.3172946>
- [54] John Quarles, Samsun Lampotang, Ira Fischler, Paul Fishwick, and Benjamin Lok. 2013. Experiences in mixed reality-based collocated after action review. *Virtual Reality* 17, 3 (Sept. 2013), 239–252. <https://doi.org/10.1007/s10055-013-0229-6>
- [55] Fred Ramsey and Daniel Schafer. 2012. *The Statistical Sleuth: A Course in Methods of Data Analysis*. Cengage Learning.
- [56] Alexander Renkl, Robin Stark, Hans Gruber, and Heinz Mandl. 1998. Learning from worked-out examples: The effects of example variability and elicited self-explanations. *Contemporary Educational Psychology* 23, 1 (Jan. 1998), 90–108. <https://doi.org/10.1006/ceps.1997.0959>
- [57] Robert Rosenthal and Donald B. Rubin. 1978. Interpersonal expectancy effects: The first 345 studies. *Behavioral and Brain Sciences* 1, 3 (1978), 377–386.
- [58] Stuart J. Russell and Peter Norvig. 2016. *Artificial Intelligence: A Modern Approach*. Pearson Education Ltd., Malaysia.
- [59] Taylor Lee Sawyer and Shad Deering. 2013. Adaptation of the US army’s after-action review for simulation debriefing in healthcare. *Simulation in Healthcare* 8, 6 (Dec. 2013), 388–397. <https://doi.org/10.1097/SIH.0b013e31829ac85c>
- [60] James Schaffer, John O’Donovan, James Michaelis, Adrienne Raglin, and Tobias Höllerer. 2019. I can do better than your AI: Expertise and explanations. In *Proceedings of the 24th International Conference on Intelligent User Interfaces (IUI’19)*. ACM, New York, NY, 240–251. <https://doi.org/10.1145/3301275.3302308>
- [61] Hendrik Strobelt, Sebastian Gehrmann, Michael Behrisch, Adam Perer, Hanspeter Pfister, and Alexander Rush. 2019. Seq2Seq-Vis: A visual debugging tool for sequence-to-sequence models. *IEEE Transactions on Visualization and Computer Graphics* 25, 1 (2019), 353–363. <https://doi.org/10.1109/TVCG.2018.2865044>
- [62] John Sweller, Jeroen J. G. Van Merriënboer, and Fred Paas. 1998. Cognitive architecture and instructional design. *Educational Psychology Review* 10 (Sept. 1998), 251–296. <https://doi.org/10.1023/a:1022193728205>
- [63] J. Tullio, A. Dey, J. Chalecki, and J. Fogarty. 2007. How it works: A field study of non-technical users interacting with an intelligent system. In *Proceedings of the ACM Conference on Human Factors in Computing Systems*. ACM, New York, NY, 31–40.
- [64] U.S. Army. 1993. *Training Circular 25-20: A Leader’s Guide to After-Action Reviews*. Technical Report. Department of the Army, Washington, DC.
- [65] Danding Wang, Qian Yang, Ashraf Abdul, and Brian Y. Lim. 2019. Designing theory-driven user-centric explainable AI. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI’19)*.
- [66] J. Wang, L. Gou, H. Shen, and H. Yang. 2019. DQNViz: A visual analytics approach to understand deep Q-networks. *IEEE Transactions on Visualization and Computer Graphics* 25, 1 (2019), 288–298.
- [67] Claes Wohlin, Per Runeson, Martin Höst, Magnus C. Ohlsson, Björn Regnell, and Anders Wesslén. 2012. *Experimentation in Software Engineering*. Springer Science & Business Media.

- [68] Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2019. Errudite: Scalable, reproducible, and testable error analysis. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL'19)*. 747–763. <https://doi.org/10.18653/v1/P19-1073>
- [69] Qian Yang, Aaron Steinfeld, Carolyn Rosé, and John Zimmerman. 2020. Re-examining whether, why, and how human-AI interaction is uniquely difficult to design. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (CHI'20)*. ACM, New York NY, 1–13. <https://doi.org/10.1145/3313831.3376301>
- [70] Qian Yang, Aaron Steinfeld, and John Zimmerman. 2019. Unremarkable AI: Fitting intelligent decision support into critical, clinical decision-making processes. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI'19)*. 1–11.
- [71] M. Zhang and Z. Zhou. 2014. A review on multi-label learning algorithms. *IEEE Transactions on Knowledge and Data Engineering* 26, 8 (2014), 1819–1837. <https://doi.org/10.1109/TKDE.2013.39>

Received January 2021; revised May 2021; accepted August 2021