

# The Myth of Diagnosis as Classification: Examining the Effect of Explanation on Patient Satisfaction and Trust in AI Diagnostic Systems

Lamia Alam (✉ [lalam@mtu.edu](mailto:lalam@mtu.edu))

Michigan Technological University

Shane Mueller

Michigan Technological University

---

## Research Article

**Keywords:** Artificial Intelligence, healthcare, medical diagnosis

**Posted Date:** January 7th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-138628/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# **The Myth of Diagnosis as Classification: Examining the Effect of Explanation on Patient Satisfaction and Trust in AI Diagnostic Systems**

**Lamia Alam<sup>1</sup>, Shane Mueller<sup>2</sup>**

Michigan Technological University, Houghton, MI 49931, USA

lalam@mtu.edu<sup>1</sup> shanem@mtu.edu<sup>2</sup>

## Abstract

**Background:** Artificial Intelligence has the potential to revolutionize healthcare, and it is increasingly being deployed to support and assist medical diagnosis. One potential application of AI is as the first point of contact for patients, replacing initial diagnoses prior to sending a patient to a specialist, allowing health care professionals to focus on more challenging and critical aspects of treatment. But for AI systems to succeed in this role, it will not be enough for them to merely provide accurate diagnoses and predictions. In addition, it will need to provide explanations (both to physicians and patients) about why the diagnoses are made. Without this, accurate and correct diagnoses and treatments might otherwise be ignored or rejected.

**Method:** It is important to evaluate the effectiveness of these explanations and understand the relative effectiveness of different kinds of explanations. In this paper, we examine this problem across two simulation experiments. For the first experiment, we tested a re-diagnosis scenario to understand the effect of local and global explanations. In a second simulation experiment, we implemented different forms of explanation in a similar diagnosis scenario.

**Results:** Results show that explanation helps improve satisfaction measures during the critical re-diagnosis period but had little effect before re-diagnosis (when initial treatment was taking place) or after (when an alternate diagnosis resolved the case successfully). Furthermore, initial “global” explanations about the process had no impact on immediate satisfaction but improved later judgments of understanding about the AI. Results of the second experiment show that visual and example-based explanation integrated with rationales had a significantly better impact on patient satisfaction and trust than no explanations, or with text-based rationales alone. As in Experiment 1, these explanations had their effect primarily on immediate measures of satisfaction during the re-diagnosis crisis, with little advantage prior to re-diagnosis or once the diagnosis was successfully resolved.

**Conclusion:** These two studies help us to draw several conclusions about how patient-facing explanatory diagnostic systems may succeed or fail. Based on these studies and the review of the literature, we will provide some design recommendations for the explanations offered for AI systems in the healthcare domain.

## Introduction

### Background

AI systems are increasingly being fielded to support diagnoses and healthcare advice for patients [1]. Although these systems are still in their infancy, they have the potential to serve as a first point-of-contact for patients, and eventually may produce diagnoses and predictions about patient’s health, perform routine tasks, and provide non-emergency

medical advice. This has the potential to provide innovative solutions for improved healthcare outcomes at a reduced cost.

However, in order to replace or supplement human diagnosis from physicians and health care professionals, it may not be enough for the AI diagnosis system to just be accurate. An accurate diagnosis without justification or explanation might be ignored, even from a competent physician. This was perhaps first noted in the early days of medical diagnosis systems, Teach and Shortliffe [2] found that when considering AI diagnostic systems, the most important desire of both physicians and non-physicians was that it should be able to explain its diagnostic decisions. In contrast, avoiding incorrect diagnoses and erroneous treatments were rated among the least important properties. Recently, Holzinger et al. [3] argued that Explainable AI (XAI) may help to facilitate transparency and trust for the implementation of AI in the medical domain, so we expect that any successful patient-focused AI diagnoses system will also provide explanations and justifications of that diagnosis so that the patient can understand why a diagnosis is made or a treatment plan is recommended. It is even possible that an average diagnosis system with better explanation will lead to better healthcare outcomes than a perfect diagnosis system without explanation.

A variety of algorithms have been identified for providing explanations of AI diagnostic systems, both within and outside the field of healthcare. For example, early expert

systems provided rule-based logical explanations that were tightly coupled to the knowledge the systems used to make diagnoses [4]–[8]. More recently, researchers have focused on visualizing elements of the classification algorithms being used to make a diagnosis (e.g., heat-map image analysis), and visualizing decision trees or complex additive models [9]–[13]. Other researchers have explored using case-based explanations, providing examples, and compelling support for the systems’ conclusions [14]–[20]. Consequently, there are several algorithmic approaches to both diagnosis and explanation of diagnoses that have been explored in medical AI. However, it is not clear which methods are effective, whether a single method is sufficient, or whether the explanations need to be tailored to individual patients, situations, or different timepoints during diagnosis. Furthermore, the literature on XAI in healthcare focuses primarily on algorithms, rather than the impact the algorithms have on patients (such as their satisfaction, understanding, or willingness to use the system in the future).

One approach we have pursued to study explainable AI (XAI) in healthcare is to understand the types of explanations real physicians offer when they interact with patients. For example, Alam [21] conducted an interview study with physicians to document how they explained diagnoses to the patients. The results suggest that physicians use a variety of explanation methods, which are dependent on context, including time (i.e., early or later in diagnosis) and the patient or patient’s advocates identity (including cultural,

education, age, and other concerns). The explanations identified included the use of logical arguments, examples, test results, imagery, analogies, and emotional appeals. The results of this study also suggest that physicians tend to provide different types of explanations at different points of diagnosis. Although many of these explanations have been explored in the XAI literature [22], few systems have acknowledged the variety and contextual aspects of the different explanation types.

### **Methods for Providing Explanations**

In the present paper, we will report on two experiments we conducted that explore how different types of explanations may impact satisfaction and trust in a simulated AI diagnostic system. In these studies, we will examine several different types of explanations that have been proposed and explored in the XAI community.

One such type of explanation is whether the goal of the explanation is to inform about the diagnostic process, versus justify why a particular diagnosis was made. These explanation types are respectively referred to as “global” and “local” explanations [23]–[26]. In general, Alam [21] found that physicians report using both methods; sometimes they explain how a particular disease or diagnostic process works; other times they justify why a particular diagnosis is given based on evidence (symptoms, test results, history, etc.).

Another important distinction is the means by which an explanation is provided. Alam [21] also found that physicians’ explanations mapped onto many of the explanation

types studied in the XAI literature, including case-based information and examples [27], [28], analogies [29], [30], logical arguments [31], [32], visualizing imagery and highlighting important aspects [33], [34]. For imagery, AI healthcare systems may use graphs to show the relative probability of different outcomes or the relative importance of different symptoms for those outcomes, which is more akin to how the LIME algorithm [35] works for diagnostic features. Physicians may present visualization differently from how AI systems offer visual explanations, but even the use of x-rays and other test reports are generally accompanied by explanations highlighting the location of critical signs indicating a diagnosis—with a similar goal as gradient-based heatmaps [36]–[39] in XAI systems.

Next, we will report the results of two studies in which we tested a variety of explanation methods and approaches in a simulated diagnostic situation. Rather than testing a single explanation of an isolated case, we designed a garden path scenario in which symptoms initially pointed to one diagnosis, but later it became clear that another diagnosis was correct. This provided an emerging diagnosis, which we believe is particularly well-suited to understand how patients both trust and understand an AI diagnostic system.

## **Experiment 1**

Hoffman et al. [40] argued that elements of satisfaction and trust follow from an improved understanding of an AI system that might be gathered from different kinds of

explanations. Consequently, we hypothesize that explanations will induce greater satisfaction, trust, understanding, and perceptions of accuracy. To investigate this, we tested participants interacting with a simulated AI system that initially gives the most likely but incorrect diagnosis, but later it changes the diagnosis to the correct disease once further testing is complete. This provides an important case for understanding explanation, because, at all times, the AI can be judged to be behaving optimally given its information—even when its diagnosis is incorrect.

## **Method**

### **Participants**

Eighty undergraduate students at Michigan Technological University took part in the study in exchange for partial course credit.

### **Procedure**

We created a diagnosis scenario in which a simulated AI system gives a most likely but incorrect diagnosis but later changes the diagnosis to the correct disease. The scenario involved gastrointestinal disorders and symptoms, which are often difficult to diagnose in real-world situations. The participants played the role of patients in the scenario, instructed to say they were suffering from specific symptoms (abdominal pain, cramps, diarrhea, fatigue, and joint pain). A simulated AI system (called MediBot.ai) provided diagnostic information about the scenario, initially concluding that the patient was suffering from Irritable Bowel Syndrome (IBS), and advised patients to follow a specific diet

chart and come back for follow up next week. After one week, the participants were told that they had begun to feel better, but the symptoms started getting worse after that. When the patient did not feel good even after three consecutive weeks, MediBot determined that the patient might not be suffering from the “most likely” condition IBS and changed its diagnosis, ordered additional diagnostic tests, and determined the patient was suffering from Celiac disease, which occurs due to gluten allergy (the ‘ground truth’ of the scenario). Participants had to communicate with MediBot through six simulated weeks, but the study took around 20 minutes to complete. All participants experienced the same basic scenario with identical symptoms and diagnoses. To maintain certain intervals between the simulated weeks, they were given brief crosswords to solve during the intervals. After they solved one crossword, they were asked to start following up with MediBot and play their role as patients again.

Participants were divided into three groups: Control, global explanation, and local explanation. The control group received no explanation of why MediBot was making any decision in any week. The global explanation group received an initial tutorial describing how MediBot does diagnosis in general and focusing especially on how the AI follows a most-likely diagnostic approach, which means that it may make errors in particular cases. This included two examples: 1) A success case of the first diagnosis, and 2) A failure case for the first diagnosis, but eventually a

successful second diagnosis. The Local explanation group received local justifications about each decision and prediction of MediBot throughout the scenario. Local justifications explained why the MediBot made a particular decision for a particular case. For this group, MediBot showed a probability chart of the disease likelihood of the patient in each week (see Figure 1), representing the probability or likelihood of different outcomes visually, and including de-

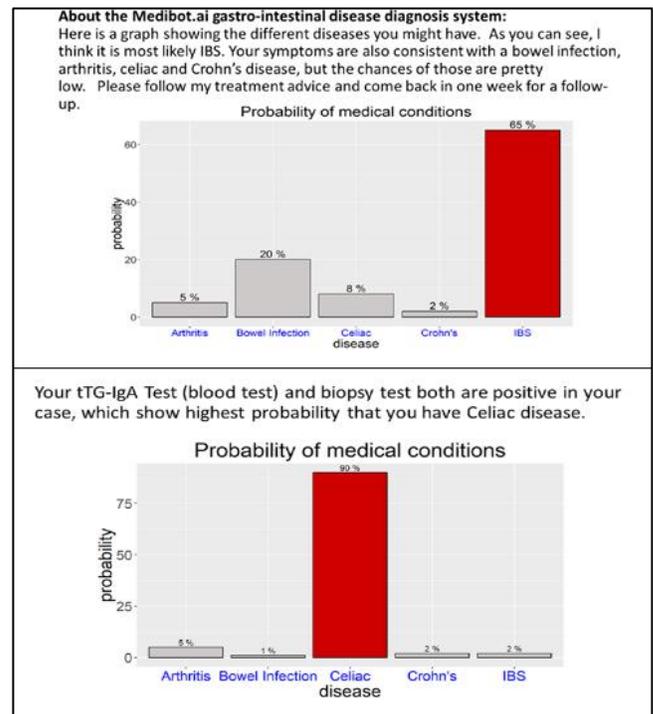


Figure 1: Week 4 (top panel) and 5 (bottom panel) Probability Chart and Explanation. Initial diagnosis of IBS changes as information emerges, and the explanations constitute the relative certainty of each disease in these bar charts and text description

scriptive text explanation about why it was making a particular decision.

The Appendix lists the entire scenario for a patient across six weeks of diagnosis. After each simulated week, participants were asked to rate their satisfaction, trust, perception of accuracy, sufficiency, usefulness, and completeness for

the explanations received from MediBot. These are some of the key attributes of explanations identified in the literature and are referred to as “Explanation Satisfaction Scale” attributes [40]. At the end of the study, participants also rated their agreement about their understanding of four 5-point Likert scale statements (see Table 2).

## Results

Both the control and the global explanation groups expressed less satisfaction, trust, perception of accuracy,

sufficiency, usefulness, and completeness than the local explanation group, as shown in Figure 2.

The control group and global explanation groups received the same scenario with no local explanations, and only differed in whether they saw an initial global explanation of the AI, and so the fact that they did not differ from one another on these ratings suggest that the satisfaction ratings focus

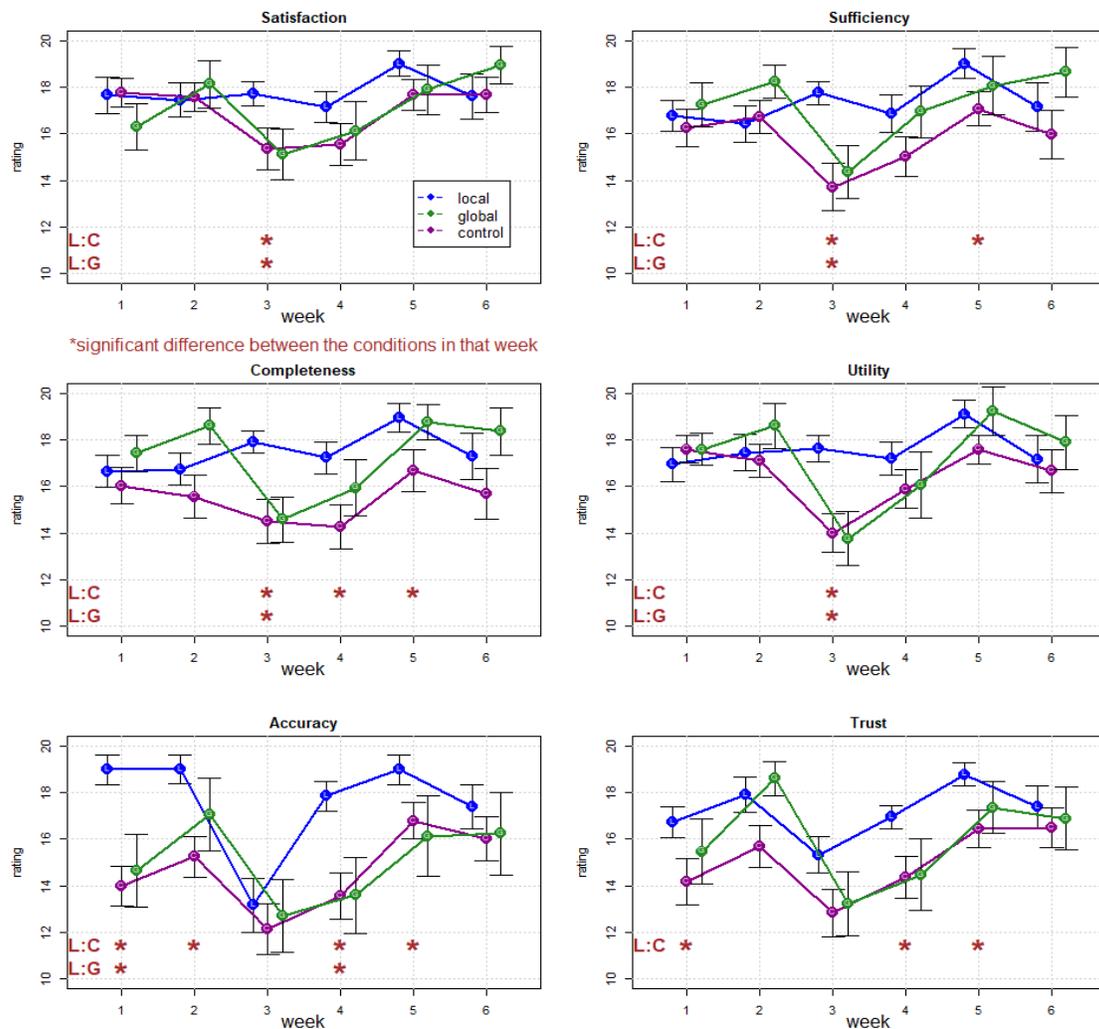


Figure 2: Results for explanation satisfaction scales

the user on the immediate situation and are not impacted by global understanding.

We examined the rating for each dimension of explanation satisfaction scales with a Type-III factorial ANOVA examining the main effects of time, explanation condition (local, global, and control), and their interaction using the R package ‘ez’ [41]. The Type-III ANOVA examines the main effects AFTER the interaction has been accounted for, allowing us to identify residual effects of explanation types across all time points. The results are shown in Table 1.

	Time	Explanation	Explanation: Time
Satisfaction	$F(5,385) = 8.20$ $p < 0.05$ $\eta_p^2 = 0.04$	$F(2,77) = 0.54$ $p = 0.58$ $\eta_p^2 = 0.01$	$F(10,385) = 2.28$ $p < 0.05$ $\eta_p^2 = 0.02$
Sufficiency	$F(5,385) = 7.52$ $p < 0.05$ $\eta_p^2 = 0.03$	$F(2,77) = 1.63$ $p = 0.20$ $\eta_p^2 = 0.03$	$F(10,385) = 3.14$ $p < 0.05$ $\eta_p^2 = 0.03$
Completeness	$F(5,385) = 6.20$ $p < 0.05$ $\eta_p^2 = 0.03$	$F(2,77) = 2.95$ $p = 0.06$ $\eta_p^2 = 0.04$	$F(10,385) = 2.28$ $p < 0.05$ $\eta_p^2 = 0.02$
Usefulness	$F(5,385) = 11.27$ $p < 0.05$ $\eta_p^2 = 0.06$	$F(2,77) = 0.95$ $p = 0.39$ $\eta_p^2 = 0.01$	$F(10,385) = 2.83$ $p < 0.05$ $\eta_p^2 = 0.03$
Accuracy	$F(5,385) = 17.40$ $p < 0.05$ $\eta_p^2 = 0.08$	$F(2,77) = 4.16$ $p < 0.05$ $\eta_p^2 = 0.06$	$F(10,385) = 2.18$ $p < 0.05$ $\eta_p^2 = 0.02$
Trust	$F(5,385) = 13.31$ $p < 0.05$ $\eta_p^2 = 0.07$	$F(2,77) = 3.03$ $p = 0.05$ $\eta_p^2 = 0.04$	$F(10,385) = 0.87$ $p = 0.50$ $\eta_p^2 = 0.01$

Table 1: Results from Type- III factorial ANOVA for explanation satisfaction scales

All interactions of Time and Explanation were significant at the  $p < .05$  level except for the trust response. The tests also showed the main effects of explanation were statistically significant for accuracy indicating that it was deemed better for explanation conditions across the entire duration of the experiment. Finally, the main effects of time were seen for

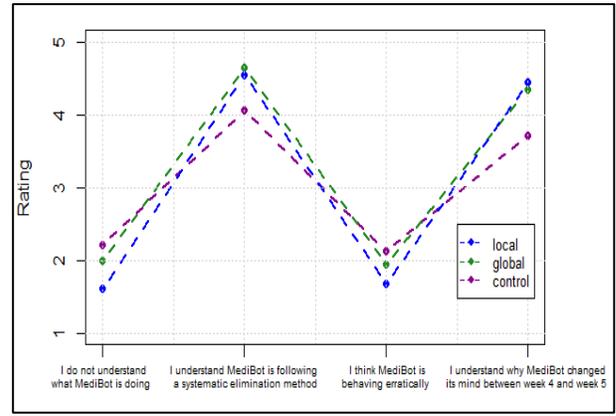


Figure 3: Results from statement ratings

all measures, indicating that the scenario was potent enough to manipulate subjective measures of trust as it moved through initial diagnosis to rediagnosis to resolution.

Welch t-test was conducted for comparing each pair of explanations at each week, the local explanation group with the control group, local explanation group with a global explanation group, and global explanation with the control group. The significant differences between each pair at each week are shown in Figure 2.

In contrast to the satisfaction ratings, the ratings of understanding elicited at the end of the scenario did in fact lead to differences between global explanation and the other conditions (see Figure 3). A one-way ANOVA showed that the three explanation conditions were significantly different ( $p < 0.05$ ) for the statements “I understand MediBot is following a systematic elimination method” ( $F(2,77) = 8.7, p < 0.05$ ) and “I understand why MediBot changed its mind between week 4 and week 5” ( $F(2,77) = 8.3, p < 0.05$ ), but they were not significantly different for the statements “I do not understand what MediBot is doing” ( $F(2,77) = 2.6, p =$

0.08) and “I think MediBot is behaving erratically” ( $F(2,77) = 2.3, p = 0.11$ ).

We used a posthoc Tukey test at a  $p < .05$  significance level on the three groups to examine pairwise differences (see Table 2). The global explanation condition produced ratings that were significantly better than the local explanation for statements 1 and 2, and both the local and global conditions were rated better than the control for statements 2 and 4. There were no differences between groups on statement 1 and 3. Thus, although the initial global explanation was not helpful for improving satisfaction during the scenario, it provided a better overall understanding of the general method of diagnosis by the AI system.

	Control-Local Ex	Control-Global Ex	Local-Global
1. I do not understand what MediBot is doing	$p = 0.06$	$p = 0.77$	$p = 0.45$
2. I understand MediBot is following a systematic elimination method	$p < 0.09$	$p < 0.05$	$p = 0.83$
3. I think MediBot is behaving erratically	$p = 0.09$	$p = 0.74$	$p = 0.55$
4. I understand why MediBot changed its mind between week 4 and week 5	$p < 0.05$	$p < 0.05$	$p = 0.9$

Table 2: Post-hoc Analysis for final understanding. Each pairwise comparison was performed with a pairwise Tukey HSD test

## Discussion and Summary

### Impact of Local Explanation/Justification

In this study, we examined how a re-diagnosis event impacted satisfaction and trust, and how different kinds of explanations impacted satisfaction, trust, and understanding of an AI system. Overall, the study showed that satisfaction and trust are harmed at the critical points during rediagnosis, even when the system is making the best diagnosis based on available information. Interestingly, the global explanation, which attempted to inoculate participants by teaching them that this very situation might occur, did little to reduce the impact of the rediagnosis on immediate measures of satisfaction and trust. Local justifications had effects throughout the scenario, but their greatest effect was at the point of re-diagnosis, in which they typically prevented a significant decline in subjective ratings of trust and satisfaction; and maintained this higher level of satisfaction until the end of the scenario when the diagnosis was resolved.

Thus, we found that local justifications were effective, but their effect is time-sensitive. During a critical situation or when AI was making errors, local justifications were very effective and powerful explanations for the patients.

### Impact of global explanations

In contrast, pre-test global explanations using example diagnoses do not show the same benefits. The global explanation did not help to raise satisfaction measures during the diagnosis in comparison to the control group that received no explanations. However, the global explanation brought

significant changes to the perception of the overall understanding of the AI system.

This study shows an initial demonstration of the time course of trust, satisfaction, and understanding during an unfolding diagnostic scenario. In the study, we used very simple visual explanations—bar charts describing the probabilities of different outcomes, with accompanying text. It is important to note that these explanations appeared effective, even though they are much simpler than many current explanatory algorithms that have been proposed for similar situations.

There are a number of alternative methods that have been explored for the explanation of classification and diagnosis. One approach attempts to focus attention on important causal factors in a classification decision or diagnosis [42]. Although like our study, the relative likelihood of different outcomes is typically shown, algorithms also often try to identify the importance of different features in making the diagnosis. For example, in the IBS/Celiac scenario, a symptom of joint pain supported celiac better than IBS, but not enough to override the higher base rate of IBS (especially because joint pain could arise from other sources and thus be attributed to something else). A single test for a gluten allergy would have been sufficient to change the diagnosis from the higher-base-rate IBS to the low-base-rate Celiac but could not impact the diagnosis if the test is not run. It might be important to let the patient understand how different signs and symptoms feed into the overall diagnosis.

A second method for explanation has been to rely on judiciously chosen examples. Examples and cases are known to be important methods for reasoning and persuasion [14], [15], [43] and have been extensively explored in the XAI literature [22].

To understand how more complex explanations might impact satisfaction and trust, we conducted a second study using a similar diagnosis scenario to investigate how different forms of local explanations affect patient satisfaction, trust, and perception of accuracy during diagnosis, which we will report next. In this study, we will examine and compare feature-highlighting approaches with case-based approaches.

## Experiment 2

Explanations in AI diagnostic systems may come in different forms such as text-based rationales, visualizations, examples, or contrasts. The goal of this study was to investigate whether different forms of explanation in an AI diagnostic system affect patient satisfaction, trust, and perception of accuracy. We implemented three forms of explanation: written rationales, visuals + rationales, and examples + rationales, in a diagnosis scenario similar to the one in Experiment 1. Again, a simulated AI system gave a most likely but incorrect diagnosis, but later it changed the diagnosis to the correct disease.

## Method

### Participants

One hundred and thirteen undergraduate students at Michigan Technological University took part in the study in exchange for partial course credit.

### Procedure

The study was conducted online, and it took 15-20 minutes to complete. Participants gave their consent online before taking part in the study. They played the role of a patient suffering from a gastrointestinal disorder interacting with the simulated AI system slightly modified from the Experiment 1 scenario.

This time, the patient suffered from abdominal pain, cramps, bloating, diarrhea, fatigue, and joint pain and had no family history of gastrointestinal diseases but had recently been exposed to a natural water source, making an initial diagnosis of Giardia likely. When tests for this came back negative, MediBot predicted that it might be IBS and asked to follow the IBS diet. The patient's condition was inconsistent for a few weeks following the diet, then eventually MediBot resolved the diagnosis as Celiac disease and confirmed it with tests.

Participants were randomly assigned into one of four groups, each receiving a different form of explanation: 1) text rationales as explanation; 2) visual + rationales explanation; 3) example-based rationales, and 4) a control group.

Rationales are the narrative justifications of how MediBot made decisions. Visual explanations include figures of the

likelihood of each suspected disease based on features MediBot used to make decisions as shown in the top panel of Figure 4. These visualizations were akin to the LIME algorithm [44], but were generated via a simple probabilistic Bayesian model (the symptom likelihood visualizations was given by the conditional probability of each disease given the symptom). We also showed the equivalent probability chart provided in Experiment 1, which showed the relative probability of each disease.

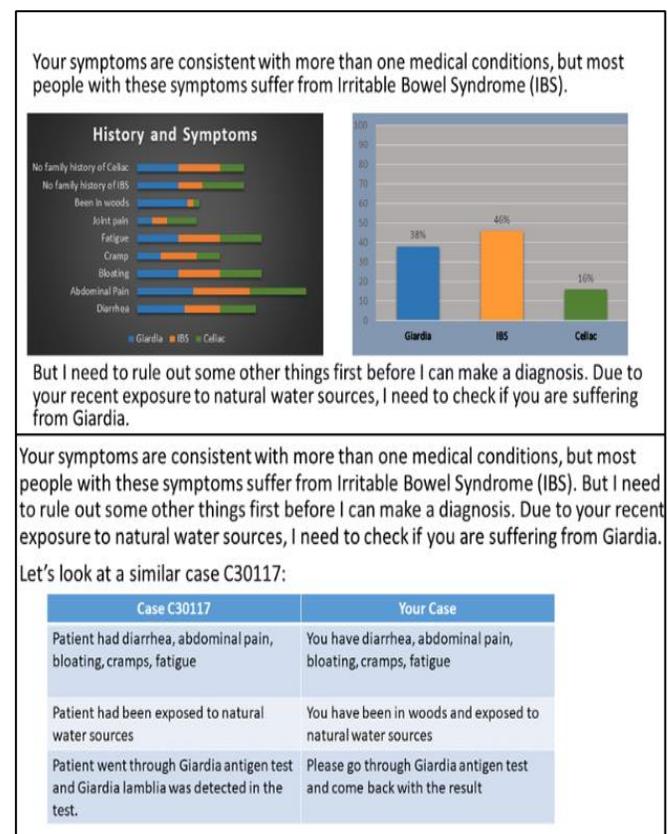


Figure 4: Sample explanations used in Experiment 2. Top panel shows visualizing feature weights and rationale; bottom panel shows example-based explanations

The example-based explanation included examples of similar cases diagnosed by MediBot in the past, as illustrated in the bottom panel of Figure 4. We used example-

based explanations where the system gave an example of a previous case and explained how it was diagnosed.

In week 5, instead of showing a positive example, it used an example that explains why it did not consider Celiac disease the most-likely condition at the beginning of the consultation. The rationales-only group saw all the justifications included in the visual and example-based explanation, only the figures and examples were removed from the explanation.

As in Experiment 1, participants interacted with MediBot for six simulated weeks and received an explanation about its prediction and diagnosis each week. After each simulated week, participants were asked to rate their satisfaction, trust, perception of accuracy, sufficiency, usefulness, and completeness for the explanations, as in Experiment 1.

## Results

In order to simplify the presentation of the results, we organized the ratings for all six weeks into three sets: Week 1

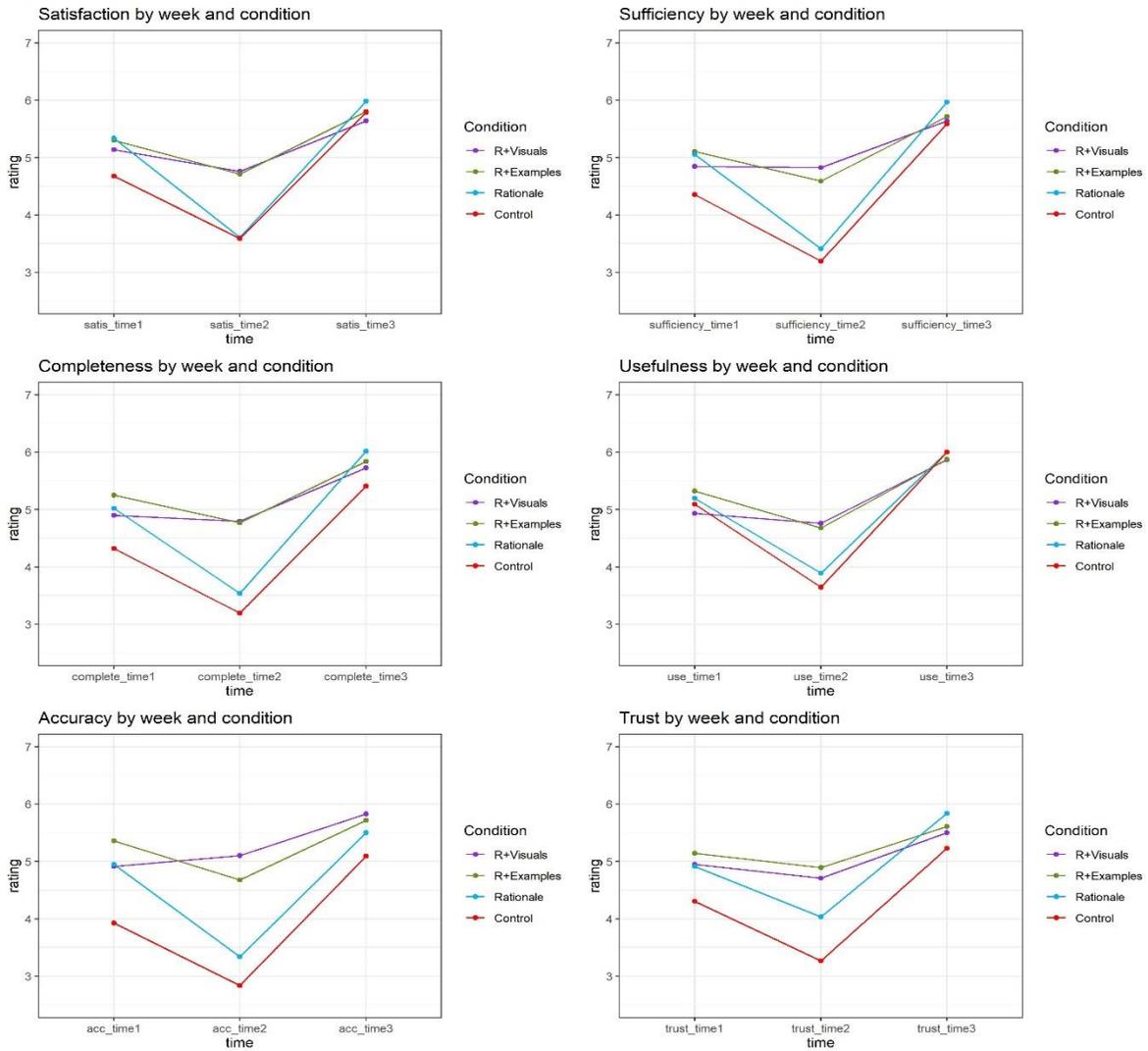


Figure 5: Rating for explanation satisfaction scales

and 2 averaged into Time 1 (initial diagnosis), Weeks 3 and 4 averaged into Time 2 (critical rediagnosis) and Weeks 5 and 6 averaged into Time 3 (resolution of diagnosis). The mean rating for all six attributes (satisfaction, trust, perception of accuracy, sufficiency, usefulness, and completeness) across conditions are shown in Figure 5.

	Time	Explanation	Explanation: Time
Satisfaction	$F(2,218) = 126.52$ $p < 0.05$ $\eta_p^2 = 0.25$	$F(3,109) = 1.97$ $p = 0.12$ $\eta_p^2 = 0.04$	$F(6,218) = 8.01$ $p < 0.05$ $\eta_p^2 = 0.06$
Sufficiency	$F(2,218) = 114.65$ $p < 0.05$ $\eta_p^2 = 0.25$	$F(3,109) = 3.38$ $p < 0.05$ $\eta_p^2 = 0.06$	$F(6,218) = 8.78$ $p < 0.05$ $\eta_p^2 = 0.07$
Completeness	$F(2,218) = 104.24$ $p < 0.05$ $\eta_p^2 = 0.24$	$F(3,109) = 4.85$ $p < 0.05$ $\eta_p^2 = 0.08$	$F(6,218) = 6.54$ $p < 0.05$ $\eta_p^2 = 0.06$
Usefulness	$F(2,218) = 110.36$ $p < 0.05$ $\eta_p^2 = 0.25$	$F(3,109) = 0.82$ $p = 0.49$ $\eta_p^2 = 0.02$	$F(6,218) = 5.06$ $p < 0.05$ $\eta_p^2 = 0.05$
Accuracy	$F(2,218) = 88.26$ $p < 0.05$ $\eta_p^2 = 0.20$	$F(3,109) = 9.95$ $p < 0.05$ $\eta_p^2 = 0.16$	$F(6,218) = 8.14$ $p < 0.05$ $\eta_p^2 = 0.07$
Trust	$F(2,218) = 64.71$ $p < 0.05$ $\eta_p^2 = 0.16$	$F(3,109) = 4.71$ $p < 0.05$ $\eta_p^2 = 0.08$	$F(6,109) = 4.10$ $p < 0.05$ $\eta_p^2 = 0.04$

Table 3: Results from Type- III factorial ANOVA for explanation satisfaction scales

We examined the rating for each dimension of explanation satisfaction scales with a Type-III factorial ANOVA examining the main effects of time, explanation condition, and their interaction using the R package ‘ez’ [41]. The Type-III ANOVA examines the main effects AFTER the interaction has been accounted for, allowing us to identify residual effects of explanation types across all time points. The results are shown in Table 3. The test of the interaction is the primary indicator of the effectiveness of an explanation

because different conditions began with little difference and converged by the end of the study.

All interactions of Time and Explanation were significant at the  $p < .05$  level. The tests also showed the main effects of explanation were statistically significant for sufficiency, completeness, accuracy, and trust indicating that these were deemed better for explanation conditions across the entire duration of the experiment. Finally, the main effects of time were seen for all measures, indicating that the scenario was potent enough to manipulate subjective measures of trust as it moved through initial diagnosis to rediagnosis to resolution.

	Time 1	Time 2	Time 3
<b>Satisfaction</b>	None	Visual; Examples > Rationale; Control	None
<b>Sufficiency</b>	None	Visual; Examples > Rationale; Control	None
<b>Completeness</b>	None	Visual; Examples > Rationale; Control	None
<b>Usefulness</b>	None	Visuals were better than Control	None
<b>Accuracy</b>	Example > Control	Visual; Examples > Rationale; Control	None
<b>Trust</b>	None	Visuals; examples > Control	None

Table 4: Significant differences between conditions at each Set according to the Tukey test, any pairing not mentioned was not significantly different for that Set.

To understand the differences between the Explanation Conditions at each time set, we conducted Tukey posthoc tests for each of the six scales using the R package *agricolae* [45]. The results are shown in Table 4.

For Time 1, there are no significant differences between any pair of explanation types overall six dimensions except accuracy. At Time 2, there are no significant differences between rationales + visuals and rationales + examples for

satisfaction, sufficiency, completeness, trust, and accuracy. But they both were better than control and rationales for satisfaction, sufficiency, completeness, and accuracy. At Time 3, there are no differences between any of the explanation types, indicating that the resolution of the scenario produced uniformly high satisfaction. Only during the rediagnosis crisis weeks, when the system was noticeably wrong, there were statistically significant differences between explanation conditions.

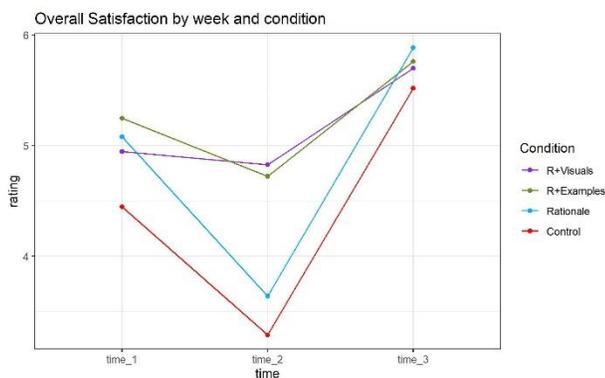


Figure 6: Mean rating for Overall Satisfaction.

Figure 6 summarizes these 6 measures with a single grand average that encapsulates the basic effect of explanation conditions in our scenario. A Type-III factorial ANOVA on the average score showed a statistically significant difference in overall satisfaction by time ( $F(2,218) = 144.68, p < 0.05$ ), condition ( $F(3,109) = 4.43, p < 0.05$ ) and the time by condition interaction ( $F(6,109) = 9.14, p < 0.05$ ).

A posthoc Tukey test showed there were no significant differences between any pair of explanation types at Time 1; there are no significant differences between rationales + visuals and rationales + examples but they both are better

than control and rationales at Time 2, and there were no differences between any of the explanation types at Time 3.

## Discussion and Summary

This study demonstrated several important results. First, like Experiment 1, explanations only appear to matter substantially during crisis weeks. It must be noted that this crisis was not due to a specific mistake or error on the part of the AI, but it was a consequence of making a most-likely diagnosis based primarily on the relative base rate of two diseases that have similar symptomology. Second, we found that richer explanations (visuals + rationales and examples + rationales) are the most effective at these critical points, but otherwise do not differ substantially from the control group. Next, for the majority of measures, rationales alone were no better than the control group. Additionally, although the visualization was substantially different from example-based explanations, we found no evidence that one method was more effective than the other. Finally, once the system came to a resolution the explanation no longer mattered and participants gave high satisfaction ratings.

Notably, this experiment did not test several conditions that might also be interesting. First, because of the lack of impact of global explanation on satisfaction measures, we did not compare global explanations in this study, either alone or accompanying the local explanations. We have no data on whether the global explanation would improve local justifications in this scenario but suspect that they would have little impact here as well. We also did not examine

whether together, examples, and visualizations would be better than either individually. The fact that subjective ratings improved at Time 3 versus Time 2 shows that there would certainly be room for improvement in the score, but given that neither were substantially impaired during Time 2 in comparison to the baseline Time 1 suggests that satisfaction may be as high as it can be under the circumstances of a disease that has not yet been cured. Finally, we examined only a single method of selecting examples. This method was sufficient enough to increase self-rated satisfaction of the system, but it may be the case that there is a variety of example types that could provide better or worse explanations. Our examples were chosen specifically to provide similar cases in the past that produced similar outcomes; another approach would be to use contrastive examples that highlight a critical aspect of the symptoms that led to the current diagnosis.

### **Discussion**

The two studies reported here allow us to draw several conclusions about how patient-facing explanatory diagnostic systems may succeed or fail. Overall, they show the importance of context on explanations. For example, justifying a decision is important to maintain satisfaction in the system; different kinds of explanations impact the patient differently, and the timing of explanations is also critical. We will examine the main lessons from these studies next.

### **Lesson: Explanations are time-sensitive**

These studies found that explanations are differentially effective at different timepoints. At the critical times when the AI is making errors, explanations can be very helpful for improved patient satisfaction, whereas they were often no different from control when things at non-critical points. This suggests that to manage patient attention and focus, developers may wish to avoid burdening patients with explanations when none are needed. Not only can this be distracting, but an explanation for something that is already understood may make the patient think they misunderstand something (why else would it need to be explained). Consequently, explanations should be used judiciously at appropriate times.

The impact of explanations at the critical Time 2 is important because this is the point at which real patients might start abandoning the system, seeking second opinions, or failing to adhere to recommendations. The type of error seen in this scenario is especially pernicious, because the diagnosis was in some sense optimal, even though it is wrong. The study shows that under the right circumstances, an explanation may mean the difference between seeing this and thinking that the diagnosis system is fundamentally unreliable or inaccurate.

### **Lesson: Significance of global explanation**

Global explanations were not as effective as local justifications for immediate measures of patient satisfaction and trust. Nevertheless, they showed significant improvement in

some post-scenario measures—ones related to the perception of global or overall understanding of the diagnosis by the AI system. And so, that should not be ignored if developers are really trying to build an XAI system for the patients. Thus, not only are different explanations effective at different times, but they also impact different aspects of their assessment of the system.

### **Lesson: Effectiveness of local justification**

These studies showed the power of local justification/explanation on immediate measures of satisfaction. When used at the right time, a local justification could be a powerful improvement for diagnostic systems. Our results suggest that system developers should concentrate on investing more effort into explanations in cases where the system may be wrong, and especially when a diagnosis is changing. Most XAI systems currently focus on a single time point explanation, but if a system can detect that its predictions are changing in a single case, this is an especially important point to use explanation.

### **Lesson: The format of explanation matters**

Across the two studies, we examined several different formats of explanation. We found that even a simple visualization showing the likelihood of different outcomes was effective (Exp. 1), as were more complex visualizations (Exp. 2) and examples (Exp. 2). However, a written logic-based narrative explanation alone (Exp. 2) did not improve subjective assessments of satisfaction. Not only that, but a detailed global explanation anticipating the type of mistakes

the AI would make had its greatest impact on post-scenario ratings of knowledge and not immediate measures of satisfaction. There are various forms of visual explanations and case-based or example-based explanations offered in XAI literature. Our studies suggest that instead of asking simple comparisons about “which kind of explanation is better”, researchers should start addressing questions about when and how different kinds of explanations are effective and helpful.

### **The myth of diagnosis as classification**

One final observation we make is that it is a mistake to think about AI diagnosis as merely a classification problem—determining what diseases or conditions the patients are suffering from considering their symptoms and signs. [21] identified several ways this is true for physicians diagnosing a variety of disorders, and this is also true for AI systems. This problem involves diagnosing, but also requires explaining why and how the AI is making the diagnosis. In many cases of diagnosis, an error is not necessarily an actual mistake, as it might follow the most likely outcome that happens to be wrong for an individual. In other cases, the course of treatment may not simply be following the most-likely option. Instead, a treatment (e.g., antibiotics) may be pursued even if it is not the most likely if it has little risk, but the consequences of not treating are large. Still, other cases involve several possibilities that each could be the source of the symptoms.

## Conclusion

To improve patient satisfaction and trust at such points, building AI systems with higher accuracy might not be enough, and may not even be possible. In critical situations, AI systems may offer an erroneous diagnosis in the process of determining the most-likely disease or condition, but patients would not understand the reason behind this if they do not get exposed to the explanations and justifications. Incorporating appropriate explanations with the AI systems may help a patient understand the diagnosis better in these situations and make them satisfied with the diagnosis as well.

## Declarations

### Ethics approval and consent to participate

This study was approved by the MTU Institutional Review Board Human Subjects Committee (M1966) and all methods in the study were carried out following the relevant guidelines and regulations. Participants were prompted with an informed consent form at the beginning of the study also outlining that by participating they acknowledge that the data obtained from this study will be used for publication. We sought explicit consent. For the online version, participants had to click a specific button and for the in-person version they had to sign the consent form to indicate they consented after reading the aforementioned form.

### Consent for publication

Not applicable.

## Availability of data and materials

The datasets used and/or analyzed during the current study are available from the corresponding author on reasonable request.

## Competing interests

Not applicable.

## Funding

This research was developed with funding from the Defense Advanced Research Projects Agency (DARPA). This research was developed with funding from the Defense Advanced Research Projects Agency (DARPA).

## Author's contribution

LA conceived and contributed to the design of the study, contributed to the statistical analysis and writing of the manuscript text. STM supervised the study, contributed to the interpretation of the data, and made critical contributions to the manuscript including revisions. Both the authors agreed on the final version of the manuscript.

## Acknowledgments

We would like to express our gratitude to every member of our lab for providing valuable suggestions and taking part in the pilot studies.

## References

- [1] I. Team, "Forbes Insights: AI And Healthcare: A Giant Opportunity," *Forbes*. <https://www.forbes.com/sites/insights-intelai/2019/02/11/ai-and-healthcare-a-giant-opportunity/> (accessed Nov. 23, 2020).
- [2] R. L. Teach and E. H. Shortliffe, "An analysis of physician attitudes regarding computer-based clinical consultation

- systems,” *Comput. Biomed. Res.*, vol. 14, no. 6, pp. 542–558, 1981.
- [3] A. Holzinger, C. Biemann, C. S. Pattichis, and D. B. Kell, “What do we need to build explainable AI systems for the medical domain?,” *ArXiv Prepr. ArXiv171209923*, 2017.
- [4] K.-P. Adlassng and M. Akhavan-Heidari, “Cadiag-2/gall: An experimental expert system for the diagnosis of gallbladder and biliary tract diseases,” *Artif. Intell. Med.*, vol. 1, no. 2, pp. 71–77, 1989.
- [5] W. J. Clancey, “The epistemology of a rule-based expert system—a framework for explanation,” *Artif. Intell.*, vol. 20, no. 3, pp. 215–251, 1983.
- [6] R. A. Miller, H. E. Pople Jr, and J. D. Myers, “Internist-I, an experimental computer-based diagnostic consultant for general internal medicine,” *N. Engl. J. Med.*, vol. 307, no. 8, pp. 468–476, 1982.
- [7] E. H. Shortliffe, “MYCIN: a rule-based computer program for advising physicians regarding antimicrobial therapy selection,” STANFORD UNIV CALIF DEPT OF COMPUTER SCIENCE, 1974.
- [8] D. W. Hasling, W. J. Clancey, and G. Rennels, “Strategic explanations for a diagnostic consultation system,” *Int. J. Man-Mach. Stud.*, vol. 20, no. 1, pp. 3–19, 1984.
- [9] Z. Che, S. Purushotham, R. Khemani, and Y. Liu, “Interpretable deep models for ICU outcome prediction,” in *AMIA Annual Symposium Proceedings*, 2016, vol. 2016, p. 371.
- [10] S. Kundu, S. Kolouri, K. I. Erickson, A. F. Kramer, E. McAuley, and G. K. Rohde, “Discovery and visualization of structural biomarkers from MRI using transport-based morphometry,” *ArXiv170504919 Cs*, May 2017, Accessed: Jun. 04, 2020. [Online]. Available: <http://arxiv.org/abs/1705.04919>.
- [11] S. Nemat, A. Holder, F. Razmi, M. D. Stanley, G. D. Clifford, and T. G. Buchman, “An Interpretable Machine Learning Model for Accurate Prediction of Sepsis in the ICU,” *Crit. Care Med.*, vol. 46, no. 4, pp. 547–553, 2018.
- [12] C.-S. Yu *et al.*, “Clustering Heatmap for Visualizing and Exploring Complex and High-dimensional Data Related to Chronic Kidney Disease,” *J. Clin. Med.*, vol. 9, no. 2, p. 403, 2020.
- [13] Z. Zhang, Y. Xie, F. Xing, M. McGough, and L. Yang, “Mdnnet: A semantically and visually interpretable medical image diagnosis network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6428–6436.
- [14] P. Cunningham, D. Doyle, and J. Loughrey, “An evaluation of the usefulness of case-based explanation,” in *International Conference on Case-Based Reasoning*, 2003, pp. 122–130.
- [15] D. Doyle, A. Tsymbal, and P. Cunningham, “A review of explanation and explanation in case-based reasoning,” Trinity College Dublin, Department of Computer Science, 2003.
- [16] W. Goodridge, H. Peter, and A. Abayomi, “The Case-Based Neural Network Model and its use in medical expert systems,” in *Joint European Conference on Artificial Intelligence in Medicine and Medical Decision Making*, 1999, pp. 232–236.
- [17] C.-C. Hsu and C.-S. Ho, “A new hybrid case-based architecture for medical diagnosis,” *Inf. Sci.*, vol. 166, no. 1–4, pp. 231–247, 2004.
- [18] P. Koton, “A medical reasoning program that improves with experience,” *Comput. Methods Programs Biomed.*, vol. 30, no. 2–3, pp. 177–184, 1989.
- [19] M. Kwiatkowska and M. S. Atkins, “Case representation and retrieval in the diagnosis and treatment of obstructive sleep apnea: a semio-fuzzy approach,” in *Proceedings of the 7th European Conference on Case-Based Reasoning*, 2004, pp. 25–35.
- [20] L. S. Ong *et al.*, “The colorectal cancer recurrence support (CARES) system,” *Artif. Intell. Med.*, vol. 11, no. 3, pp. 175–188, 1997.
- [21] L. Alam, “Investigating the Impact of Explanation on Repairing Trust in AI Diagnostic Systems for Re-Diagnosis,” (Publication No. 28088930) [Master’s Thesis, Michigan Technological University]. ProQuest Dissertations Publishing, 2020.
- [22] S. T. Mueller, R. R. Hoffman, W. Clancey, A. Emrey, and G. Klein, “Explanation in human-AI systems: A literature meta-review, synopsis of key ideas and publications, and bibliography for explainable AI,” *ArXiv Prepr. ArXiv190201876*, 2019.
- [23] D. C. Berry and D. E. Broadbent, “Explanation and Verbalization in a Computer-Assisted Search Task,” *Q. J. Exp. Psychol. Sect. A*, vol. 39, no. 4, pp. 585–609, Nov. 1987, doi: 10.1080/14640748708401804.
- [24] F. Doshi-Velez and B. Kim, “A Roadmap for a Rigorous Science of Interpretability,” *ArXiv*, vol. abs/1702.08608, 2017.
- [25] P. Lipton, “Contrastive explanation,” *R. Inst. Philos. Suppl.*, vol. 27, pp. 247–266, 1990.
- [26] M. R. Wick and W. B. Thompson, “Reconstructive expert system explanation,” *Artif. Intell.*, vol. 54, no. 1–2, pp. 33–70, 1992.
- [27] Z. C. Lipton, “The mythos of model interpretability,” *ArXiv Prepr. ArXiv160603490*, 2016.
- [28] P. Shafto and N. Goodman, “Teaching games: Statistical sampling assumptions for learning in pedagogical situations,” in *Proceedings of the 30th annual conference of the Cognitive Science Society*, 2008, pp. 1632–1637.
- [29] R. R. Hoffman, “AI models of verbal/conceptual analogy,” *J. Exp. Theor. Artif. Intell.*, vol. 10, no. 2, pp. 259–286, 1998.
- [30] R. J. Spiro, P. J. Feltovich, R. L. Coulson, and Daniel K Anderson, “Multiple Analogies for Complex Concepts: Antidotes for Analogy-Induced Misconception in Advanced Knowledge Acquisition. Technical Report No. 439,” 1988, Accessed: Sep. 18, 2017. [Online]. Available: <https://eric.ed.gov/?id=ED301873>.
- [31] F. H. George, “Logical constructs and psychological theory,” *Psychol. Rev.*, vol. 60, no. 1, pp. 1–6, 1953, doi: 10.1037/h0057812.
- [32] S. Wachter, B. Mittelstadt, and C. Russell, “Counterfactual explanations without opening the black box: Automated decisions and the GDPR,” 2017.
- [33] Y. Goyal, A. Mohapatra, D. Parikh, and D. Batra, “Interpreting visual question answering models,” in *ICML Workshop on Visualization for Deep Learning*, 2016, vol. 2, Accessed: Aug. 19, 2017. [Online]. Available: <https://pdfs.semanticscholar.org/72ce/bd7d046080899703ed3cd96e3019a9f60f13.pdf>.

- [34] L. A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell, "Generating visual explanations," in *European Conference on Computer Vision*, 2016, pp. 3–19, Accessed: May 18, 2017. [Online]. Available: [http://link.springer.com/chapter/10.1007/978-3-319-46493-0\\_1](http://link.springer.com/chapter/10.1007/978-3-319-46493-0_1).
- [35] M. T. Ribeiro, S. Singh, and C. Guestrin, "Model-agnostic interpretability of machine learning," *ArXiv Prepr. ArXiv160605386*, 2016.
- [36] W. Nie, Y. Zhang, and A. Patel, "A theoretical explanation for perplexing behaviors of backpropagation-based visualizations," *ArXiv Prepr. ArXiv180507039*, 2018.
- [37] N. F. Rajani and R. J. Mooney, "Using explanations to improve ensembling of visual question answering systems," *Training*, vol. 82, pp. 248–349, 2015.
- [38] V. Ramanishka, A. Das, J. Zhang, and K. Saenko, "Top-down visual saliency guided by captions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7206–7215.
- [39] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [40] R. R. Hoffman, S. T. Mueller, G. Klein, and J. Litman, "Metrics for explainable AI: Challenges and prospects," *ArXiv Prepr. ArXiv181204608*, 2018.
- [41] M. A. Lawrence and M. M. A. Lawrence, "Package 'ez,'" *R Package Version*, vol. 4, no. 0, 2016.
- [42] G. A. Klein, L. Rasmussen, M.-H. Lin, R. R. Hoffman, and J. Case, "Influencing preferences for different types of causal explanation of complex events," *Hum. Factors*, vol. 56, no. 8, pp. 1380–1400, 2014.
- [43] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *ArXiv Prepr. ArXiv14126572*, 2014.
- [44] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why Should I Trust You?: Explaining the Predictions of Any Classifier," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, pp. 1135–1144, Accessed: May 18, 2017. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2939778>.
- [45] F. de Mendiburu and M. F. de Mendiburu, "Package 'agricolae,'" *R Package Version*, pp. 1–2, 2019.

# Figures

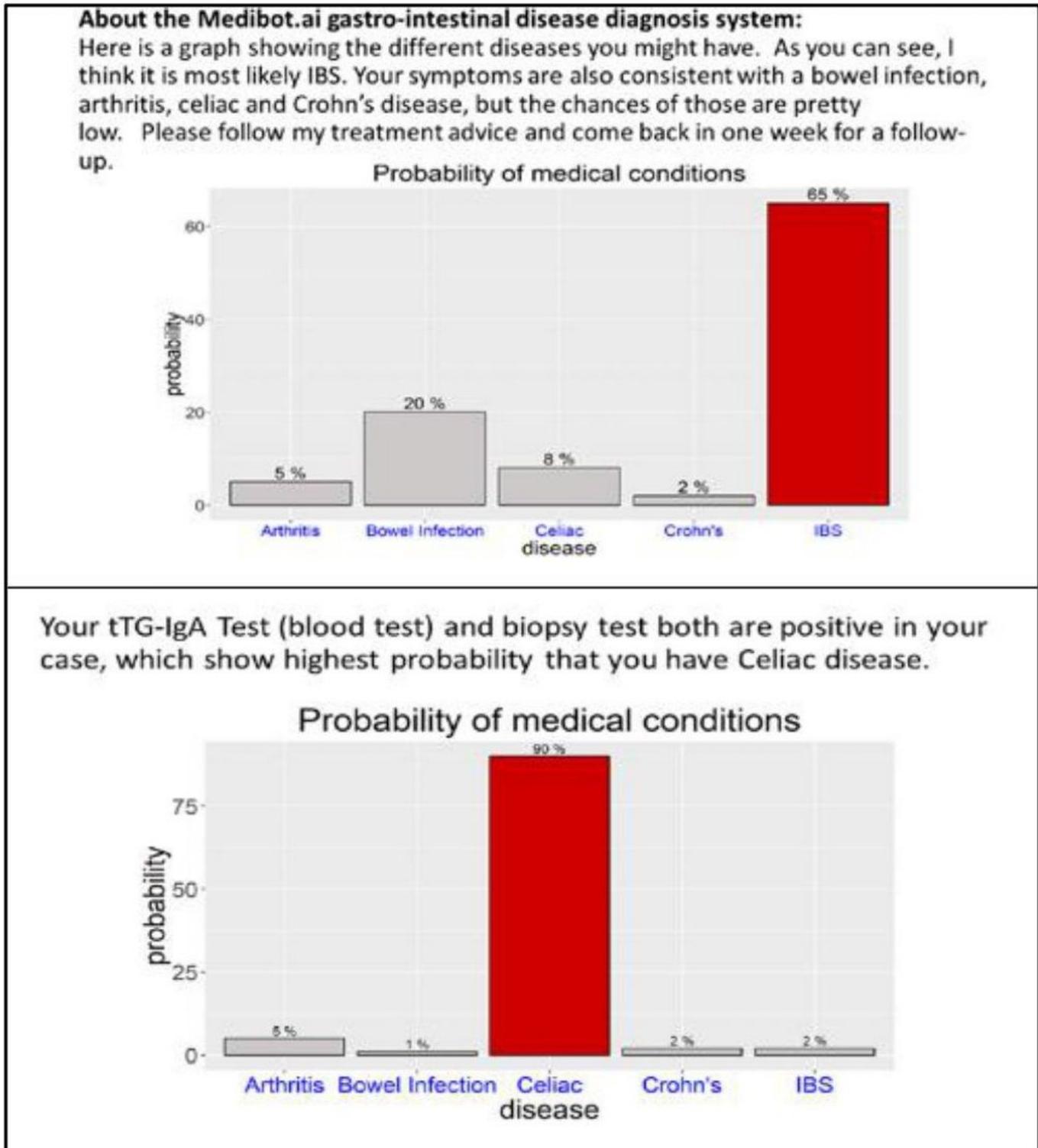
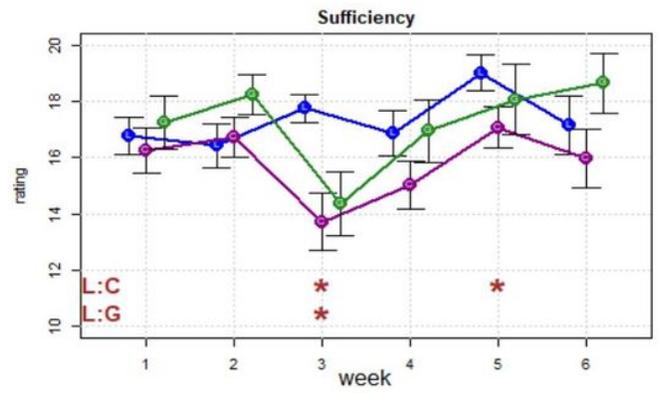
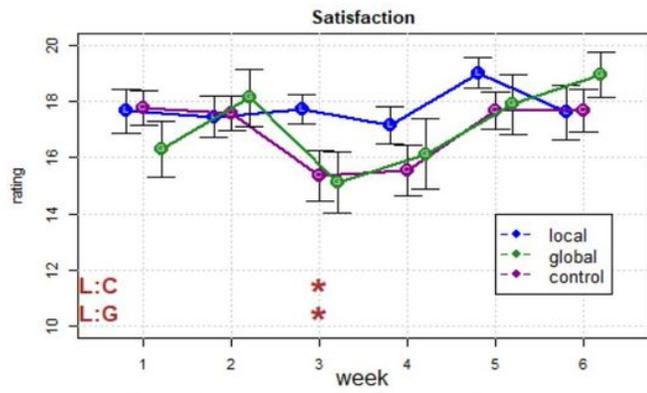


Figure 1

Week 4 (top panel) and 5 (bottom panel) Probability Chart and Explanation. Initial diagnosis of IBS changes as information emerges, and the explanations constitute the relative certainty of each disease in these bar charts and text description



\*significant difference between the conditions in that week

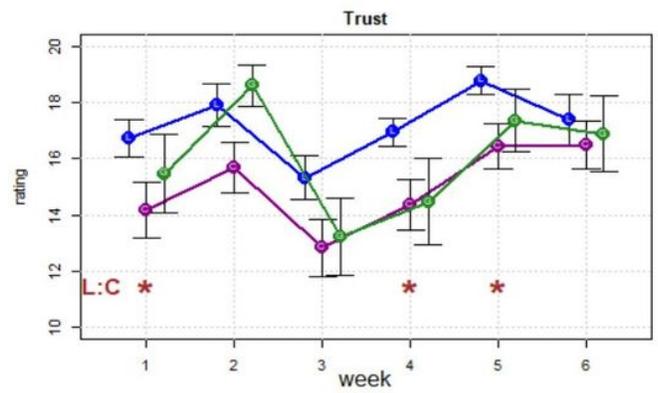
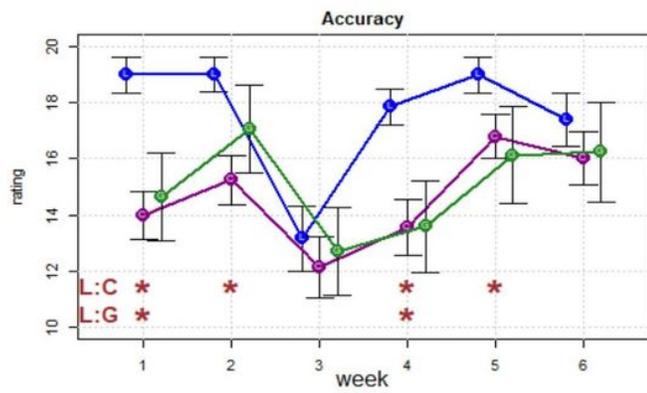
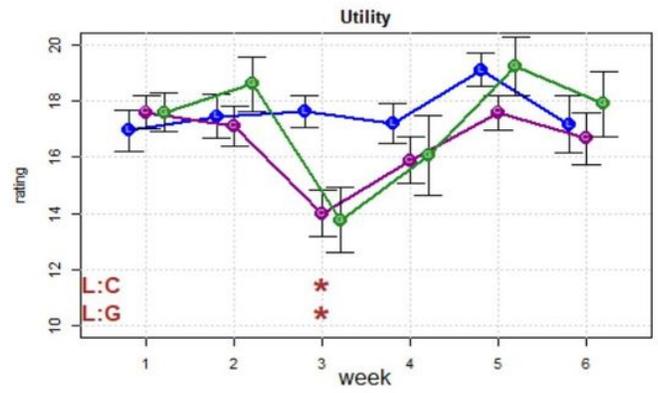
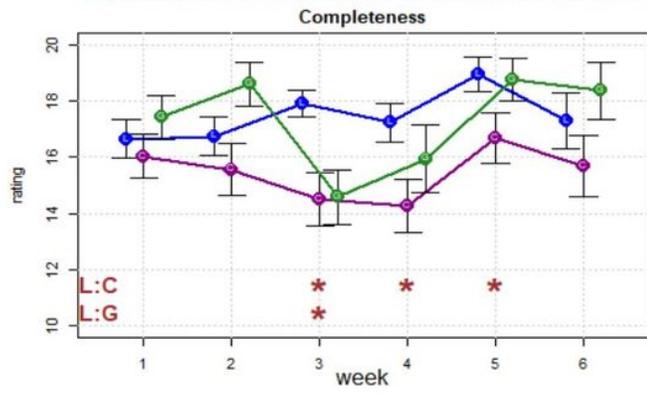


Figure 2

Results for explanation satisfaction scales

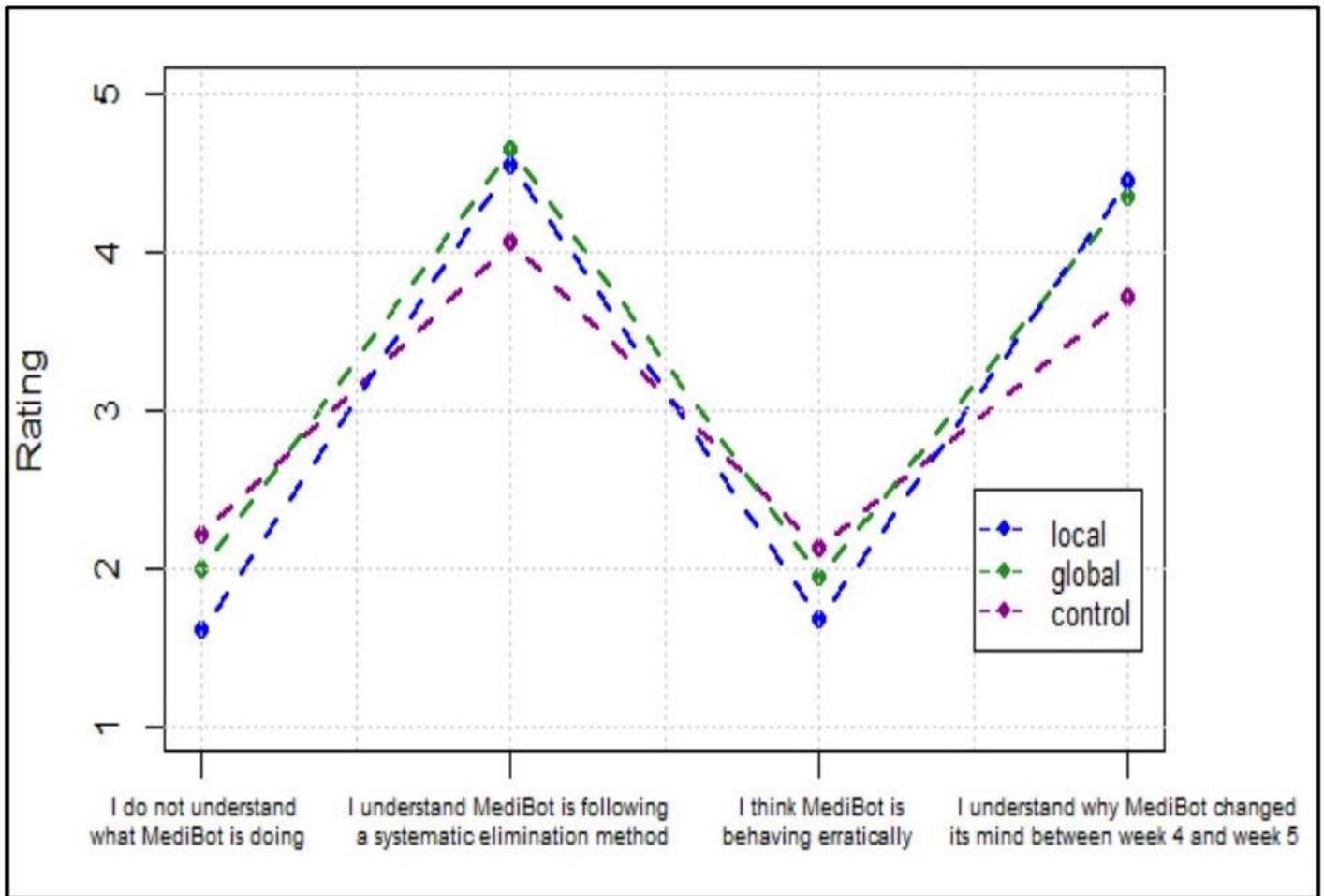
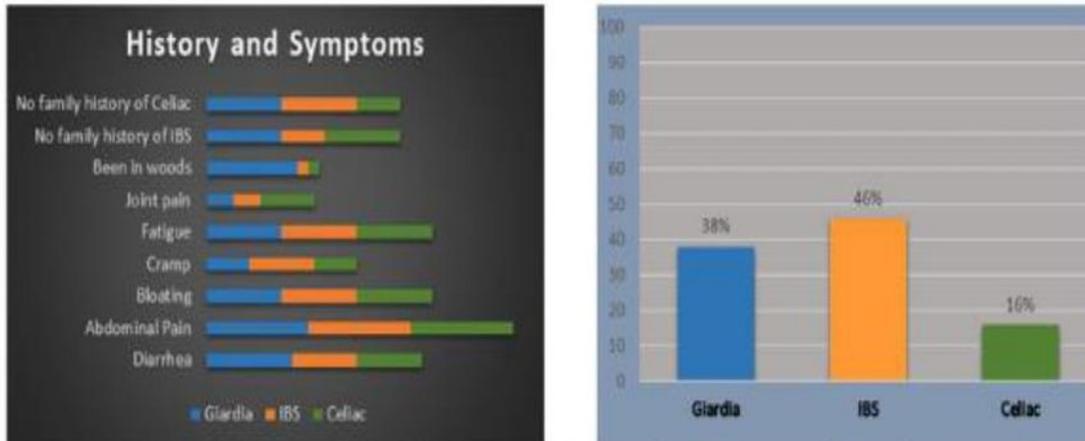


Figure 3

Results from statement ratings

Your symptoms are consistent with more than one medical conditions, but most people with these symptoms suffer from Irritable Bowel Syndrome (IBS).



But I need to rule out some other things first before I can make a diagnosis. Due to your recent exposure to natural water sources, I need to check if you are suffering from Giardia.

Your symptoms are consistent with more than one medical conditions, but most people with these symptoms suffer from Irritable Bowel Syndrome (IBS). But I need to rule out some other things first before I can make a diagnosis. Due to your recent exposure to natural water sources, I need to check if you are suffering from Giardia.

Let's look at a similar case C30117:

Case C30117	Your Case
Patient had diarrhea, abdominal pain, bloating, cramps, fatigue	You have diarrhea, abdominal pain, bloating, cramps, fatigue
Patient had been exposed to natural water sources	You have been in woods and exposed to natural water sources
Patient went through Giardia antigen test and Giardia lamblia was detected in the test.	Please go through Giardia antigen test and come back with the result

Figure 4

Sample explanations used in Experiment 2. Top panel shows visualizing feature weights and rationale; bottom panel shows example-based explanations

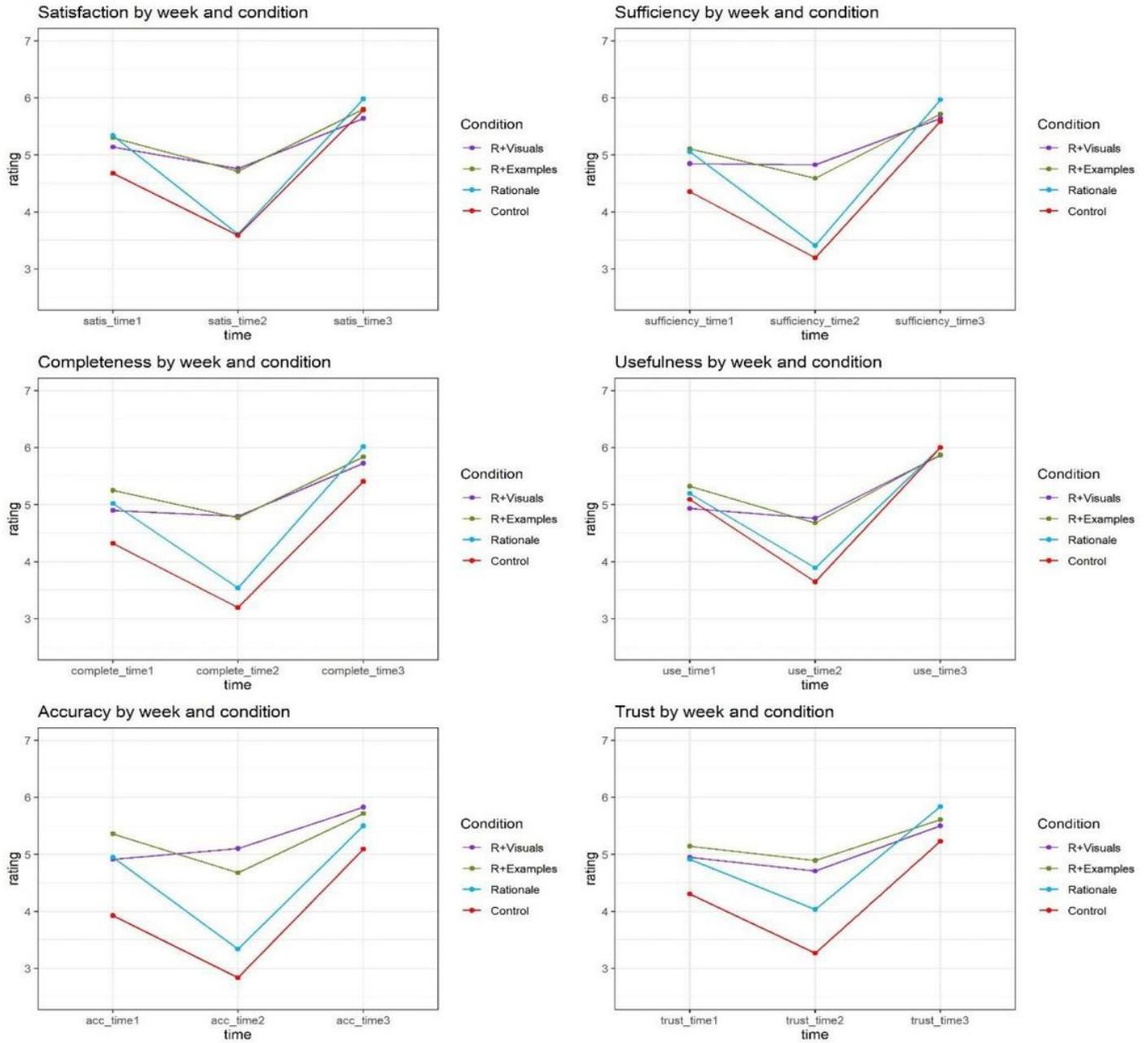
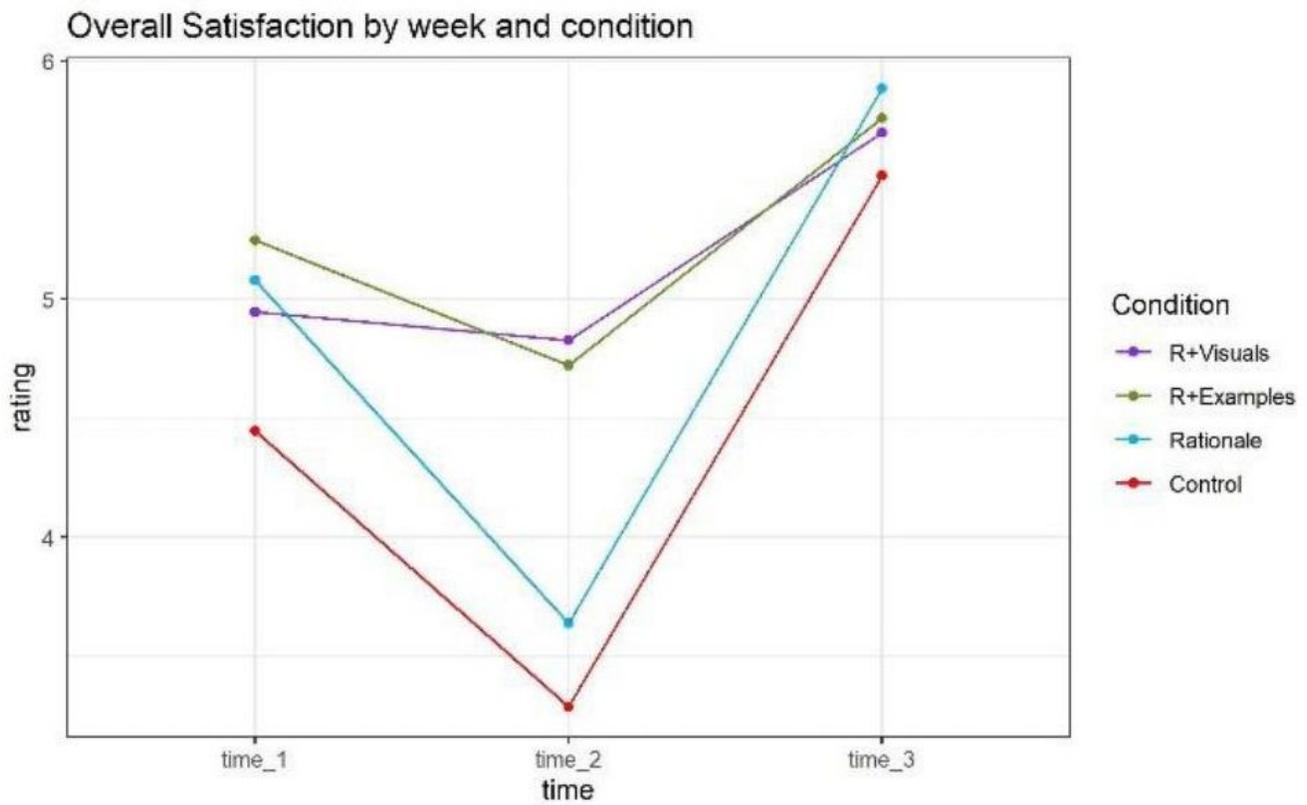


Figure 5

Rating for explanation satisfaction scales



**Figure 6**

Mean rating for Overall Satisfaction.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryAppendix.pdf](#)