

---

# Are Visual Explanations Useful?

## A Case Study in Model-in-the-Loop Prediction

---

**Eric Chu**  
MIT Media Lab  
echu@media.mit.edu

**Deb Roy**  
MIT Media Lab  
dkroy@media.mit.edu

**Jacob Andreas**  
MIT CSAIL  
jda@mit.edu

### Abstract

We present a randomized controlled trial for a model-in-the-loop regression task, with the goal of measuring the extent to which (1) good explanations of model predictions increase human accuracy, and (2) faulty explanations decrease human trust in the model. We study explanations based on visual saliency in an image-based age prediction task for which humans and learned models are individually capable but not highly proficient and frequently disagree. Our experimental design separates model quality from explanation quality, and makes it possible to compare treatments involving a variety of explanations of varying levels of quality. We find that presenting model *predictions* improves human accuracy. However, visual *explanations* of various kinds fail to significantly alter human accuracy or trust in the model—regardless of whether explanations characterize an accurate model, an inaccurate one, or are generated randomly and independently of the input image. These findings suggest the need for greater evaluation of explanations in downstream decision making tasks, better design-based tools for presenting explanations to users, and better approaches for generating explanations.

## 1 Introduction

While significant research effort has been devoted to automatically explaining decisions from machine learning models, it remains an open question to what extent these explanations are useful for humans in downstream applications. One fundamental assumption underlying much interpretable machine learning research is that more faithful and accurate explanations help people use models more effectively—explanations indicating that models have identified features relevant to the target prediction should increase human confidence in predictions, and explanations indicating that models have latched onto noise or irrelevant features should decrease trust [35, 40, 29, 13]. However, most evaluation of explanations has focused on their intrinsic relation to model properties [8, 30, 35, 40] rather than their effect on human decision-making. Here we investigate (1) whether explanation quality actually impacts model-in-the-loop human performance, and (2) whether explanation quality impacts human trust in the model.

We present a randomized controlled trial (RCT) involving model-in-the-loop decision making in a nontrivial perception problem with a modern neural prediction architecture and interpretability method. Our work follows recent RCTs studying related questions [7, 16, 17, 25, 26, 34, 41, 49], but critically, our setup allows us to isolate the effect of explanations of varying quality in a scenario with more complex models and inputs. We include in our experiments both faithful explanations of a high-quality model, via integrated gradients [40], and a variety of “faulty explanations”—saliency maps from a model trained on data with spurious correlations and completely uninformative random explanations, as shown in Figure 1. We find that neither faulty explanations nor accurate ones significantly affect task accuracy, trust in the model predictions, or understanding of the model. Counter-intuitively, even the obviously unreliable explanations shown in Figure 1c and 1d fail to

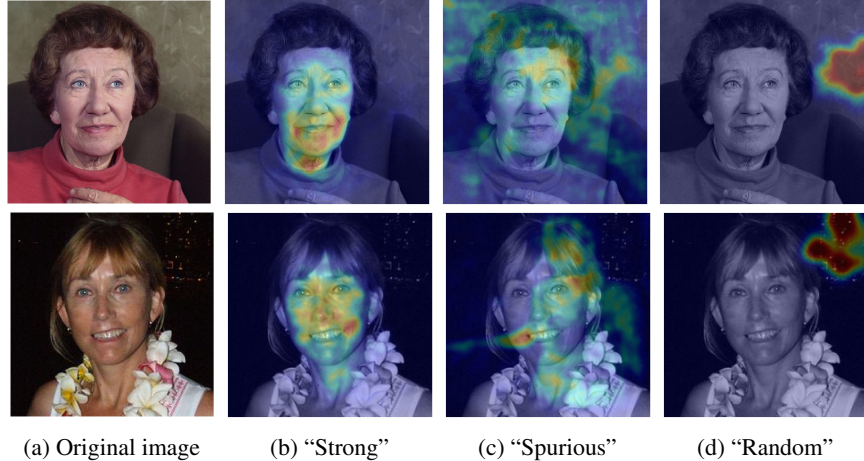


Figure 1: Images and varying saliency maps used in explanation-based treatments in our age prediction task. Participants are shown an image, a prediction from our strong model, and either a “strong”, “spurious”, or “random” explanation. The strong explanations are derived from the strong model and focus on details of the face. The spurious explanations are derived from a model trained on data with spurious correlations and tend to focus on both the face and the background. The random explanations are input-agnostic and focus on the background of the image.

significantly decrease trust in the model. All of these explanation-based treatments are comparable to each other, and to other treatments such as personifying the model with a name and an image that reveal nothing about its behavior. Ultimately, our studies point to the limited effectiveness of pixel-level visual explanations in this model-in-the-loop setting, and motivate research on better explanations, communication of the meaning of explanations, and training of users to interpret them.

## 2 Experimental Design

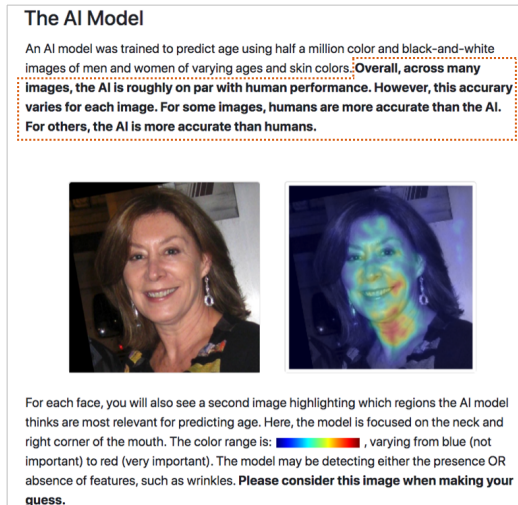
### 2.1 Task

We study a model-in-the-loop scenario, where participants are shown an input and a model’s prediction, and then asked to make their own guess. Our study examines the age prediction task, where users guess a person’s age given an image of their face. We chose this task because (a) both models and humans can perform the task with some proficiency, but with far from perfect accuracy, and (b) it is representative of high-stakes real-world scenarios that use similar models in passport recognition and law enforcement [37, 1]. A discussion of the ethical concerns of tasks that involve predicting demographic features can be found in the Broader Impact section.

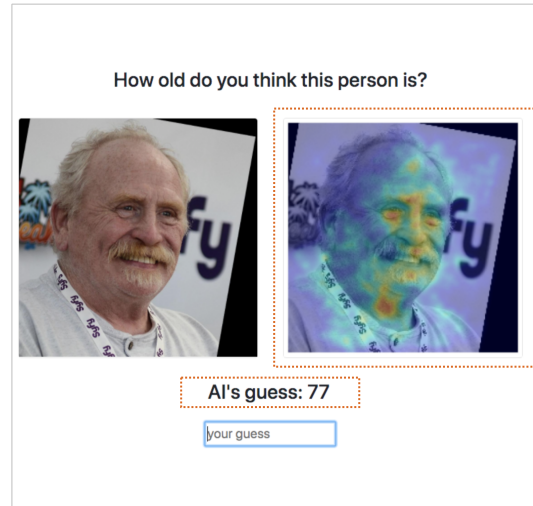
Users are shown images from the validation and test of the APPA-REAL dataset [14], which contains 7,591 images with real age labels. We sample images uniformly across ages, and many images are seen by multiple users. We also balance the dataset shown to users by combining equal proportions of images for which the model is more accurate than previously collected human guesses (available in the APPA-REAL data), and vice versa. This allows us to have greater coverage over the more interesting decision space in which the human and model disagree. The model in our model-in-the-loop task is a Wide Resnet [48] pre-trained on the IMDB-WIKI dataset [36] and fine-tuned on the APPA-REAL dataset. Trained to predict the real ages, this model gets near human-performance with a 5.24 mean absolute error (MAE) on the test set of APPA-REAL, and is more accurate than the previously collected guesses 46.9% of the time. We call this the “strong” model.

### 2.2 Treatments

Users are randomly placed into different experimental conditions, or treatment arms, allowing us to measure the effect of specific interventions. Different treatments vary aspects of the human-AI system, such as whether the model prediction is shown and what kind of explanation accompanies



(a) Description of the model and guidelines for interpreting and using the explanations.



(b) Users are asked to guess a person's age.

Figure 2: Webapp used to run the human study. Elements boxed in orange are shown or not shown depending on the treatment arm. The model prediction is shown in all cases except the control, and in additional treatments discussed in the Appendix. The saliency map in (b) is shown in the Explanation-based treatments. The description of the model performance in (a) is shown in all treatments involving the model, except for two additional treatments discussed in the Appendix.

it. Elements of the experiment are shown in Figure 2, including how model predictions is presented to the participants. The full set of treatments are listed in Table S2. Our main experiments focus on the effect of showing explanations of varying quality alongside model predictions. Additional experiments further exploring showing explanations alone, and the effect of the description of the model performance, are described in the Appendix.

**Baseline treatments** Our two main baselines are (a) the *Control* treatment, in which the user guesses without the help of the model, and (b) the *Prediction* treatment, in which the user guesses with the model prediction shown but without an explanation shown.

**Explanation-based treatments** These treatments show explanations in addition to the model prediction. Our explanations are pixel-wise importances, or saliency maps, calculated using integrated gradients [40]. For each pixel, we sum the absolute value of the channel-wise attributions, then normalize these pixel attribution scores by the 98th percentile of scores across that image. Users are given a guide for how to interpret the saliency maps (Figure 2a) and also explicitly asked to consider the explanation when making their own guess.

Explanations of varying quality are shown in Figure 1. In addition to the strong explanations, we also tested two versions of lower quality saliency maps. Importantly, these varying explanations are all shown with predictions from the same strong model, which isolates the effect of the explanation.

- *Explain-strong*. These are saliency maps from the same model whose predictions are shown to users. These explanations tend to focus on, in decreasing order, (a) the areas around the eyes, (b) lines anywhere on the face, and (c) regions around the mouth and nose.
- *Explain-spurious*. We train a “spurious” model by modifying the APPA-REAL dataset to contain spurious correlations between the background and the label. The area outside the face bounding box is modified by scaling the pixel values by  $\alpha$ , which is a linear mapping  $f(\text{age})$  from the  $[0,100]$  age range to a value in  $[0.25, 5.0]$ . Saliency maps from the spurious model often focus on both the face and the background. As with all explanation-based treatments, we show the spurious model’s saliency map with the predictions from the strong model.

- *Explain-random*. We also test completely uninformative, input-agnostic saliency maps that do not focus on the face. To generate these attributions, we first sample an anchor point around the border of the image. We then sample 50 points in a window around the anchor point, which are used as the centers for 2D Gaussians that we then combine. These are similarly normalized and mapped to the same colorscale as the integrated gradients attributions.

“Algorithmic aversion” refers to human loss of trust in a model after seeing it make a mistake [12, 47]. Our question is whether faulty explanations could act as a similar deterrent. A model may be accurate overall but still have undesired behavior. Explanations could expose those deficiencies and serve as a nudge to not trust the model.

**Design-based treatments** We also compare against existing and novel design-based treatments, which vary the description of the model and the way the model’s predictions are shown. We are interested in whether these simple framing approaches can be as effective at increasing accuracy or trust as explanation-based treatments. The *Delayed Prediction* treatment tests the effect of anchoring bias [22] and was previously shown to work well for improving accuracy in [17]. We record the initial guess, show the model prediction, then ask for a final guess. The *Empathetic* treatment personifies the model as an AI named Pat, shown in Figure S1 in the appendix. When a human perceives a computer to be more similar to themselves, they may be more cooperative and find information from the computer to be of higher quality [31]. We state “Like every teammate, Pat isn’t perfect” and show Pat next to every prediction. The *Show Top-3 Range* treatment tests a form of uncertainty by showing the range of the models’s top 3 predicted ages. The user is told “The AI’s guess will be shown as a range, e.g. 27-29. Sometimes the range may be wider, e.g. 27-35. When the range is wider, this means the AI is more uncertain about its guess.”

### 2.3 Metrics

We measure and analyze four quantities: (1) the **error** of the user’s guess (the absolute value of the difference between the guess and the *ground-truth* age label), (2) **trust** (quantified as the absolute value of the difference between the guess and the *model’s* prediction), (3) the time spent making the guess, and (4) answers to post-survey questions on how the model prediction was used in their decision-making process and how reasonable the explanations seemed. Our definition of trust follows previous work operationalizing trust as *agreement* with model predictions [47, 49, 26].

We use a mixed-effects regression model to estimate error and trust as defined above. The model includes fixed-effect terms  $\beta_{\text{image\_age}}$ , the age of the person in the image (which is correlated with the absolute error,  $\rho = 0.21$ ,  $p < 2.2e^{-16}$ ), and  $\beta_{\text{treatment}}$ , for each of the treatments. We also include random-effect intercept terms  $z_{\text{user}}$  and  $z_{\text{image}}$  to capture effects specific to each image and user. The model is defined as follows, where  $\langle \text{target} \rangle$  is the error or trust defined above.

$$y_{\langle \text{target} \rangle} = \beta_0 + \beta_{\text{treatment}} \cdot x_{\text{treatment}} + \beta_{\text{image\_age}} \cdot x_{\text{image\_age}} + z_{\text{user}} \cdot x_{\text{user}} + z_{\text{image}} \cdot x_{\text{image}} + \epsilon \quad (1)$$

### 2.4 Experiment Details

We ran experiments on Amazon Mechanical Turk, with 1,058 participants. Participants were incentivized to perform well on the task by paying top performers a 100% bonus. Prior to data collection, we conducted a two-tailed power analysis using the mean and standard deviation of previously collected guesses in the APPA-REAL dataset. In order to detect 1 year differences between treatments, we needed to collect 546 guesses per treatment. We ultimately collected 827.5 guesses (82.75 participants) per treatment, which would allow us to detect a 1 year difference of means at the  $p < 0.05$  level with probability 93%.

## 3 Analysis and Results

The overall mean absolute errors per treatment are shown in Table 1. For more detailed analysis, we control for image- and user-specific effects using the regression model in Equation 1, shown in Figures 3 and 4. This allows us to more precisely quantify the additional effect of explanation-based treatments on error relative to a human-only control (Figure 3) and the effect of explanations on trust relative to a prediction-only control (Figure 4). In addition to a regression on the overall data, we also



analyzed subsets of the data, defined by the mean of all human and model errors (8.65). For example, the “good human, bad model” subset, denoted **Human+Model-** and **H+M-** for short, is the set of datapoints where the human guess was more accurate and the model prediction was less accurate than the average human and model guess. There are 1802 and 1721 guesses in the **H+M-** and **H-M+** settings, respectively. Using these experiments, we aim to understand the effect of explanations on human use and perception of model predictions.

### 3.1 How do model predictions and explanation quality affect model-in-the-loop accuracy?

**Participants in the Empathetic, Show Top-3 Range, and Explain-strong treatments performed best and outperformed humans without AI.** Participants shown the model prediction are more accurate by 2 years on average overall, as seen in Figure 3. The predictions generally help whenever the human is inaccurate (**H-M+**, **H-M-**), but can hurt when the human is accurate and the model is inaccurate (**H+M-**).

The top-performing treatments have similar effects overall. However, their effects vary in different settings. For example, Show Top-3 Range is potentially more helpful in the **H-M+** setting, with a 4.55 improvement in accuracy, and also the only top treatment to not have a statistically significant harmful effect in the **H+M-** setting. However, we note that Tukey HSD testing indicates that the pairwise differences between all of these top treatments is not statistically significant.

The results are also a reminder of the importance of the *design* of the human-AI system, with the Empathetic and Show Top-3 Range treatments equally or more effective as our explanation-based treatments. Understanding these approaches may also help design better explanation-based approaches, as there may be underlying mechanisms that affect both. For example, the Empathetic and Explain-strong treatments both increase trust, and it could be that the former does so through an emotional approach while the latter does so through more logical means.

Improved accuracy is not attributable to the amount of time spent making guesses: participants in the Control took 5.7 seconds per image, compared to Empathetic (6.4), Explain-strong (5.9), Show Top-3 Range (5.4), and Prediction (5.2).

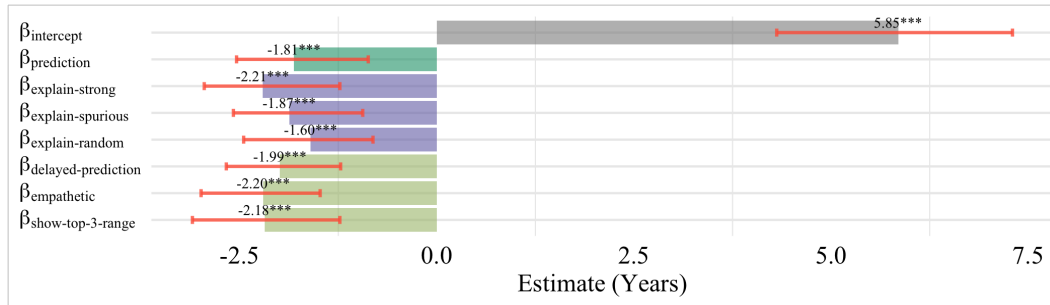
**The addition of explanations did not improve accuracy.** As seen in Figure 3a, the Explain-strong treatment has a similar effect size to the Prediction treatment. This, along with additional explanation-without-prediction treatments detailed in Appendix Sections A.2 and B, which resulted in small and non statistically significant effects, indicate the limited utility of these explanations for such a visual classification task. Survey responses indicate that participants did indeed examine the highlighted areas and focus on important features such as location-specific wrinkles, but could not extract information that would boost their performance.

It is possible that our results would change if users were extensively trained to interpret and use the saliency maps, instead of only being presented with our short guide. We do note, however, that prior work on training users to interpret saliency map explanations in a different task did not increase performance when model predictions were also shown [26, 25]. We believe nevertheless that one broader takeaway remains the same — designers of human-AI systems should question the utility of pixel-level saliency explanations when designing ML-powered tools.

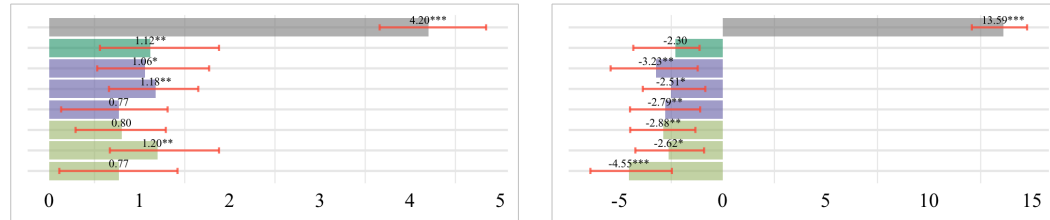
**The quality of explanations had little effect on accuracy.** Though directionally what one might expect overall (explain-strong < explain-spurious < explain-random), the differences are small and

Table 1: Mean absolute error of guesses. Bootstrapped 95% confidence intervals in parentheses. Results for all treatments are in Appendix Section B. We note that the MAE is higher for the “Model Alone” case than the MAE stated in Section 2 because we sample uniformly across the entire age range, and errors are often much larger when the person is older.

Treatment Arm	MAE
Control (Human Alone)	10.0 (9.4 - 10.5)
Model Alone	8.5 (8.3 - 8.7)
Prediction	8.4 (7.8 - 9.0)
Explain-strong	8.0 (7.5 - 8.5)
Explain-spurious	8.5 (8.0 - 9.1)
Explain-random	8.7 (8.1 - 9.2)
Delayed Prediction	8.5 (8.0 - 9.0)
Empathetic	8.0 (7.6 - 8.5)
Show Top-3 Range	8.0 (7.4 - 8.5)



(a) Overall



(b) Human+Model-

(c) Human-Model+

Figure 3: Estimates for  $y_{\text{error}}$  regression. **Intercept represents control treatment (human alone); other estimates are relative to the intercept. Lower values are better, indicating reduced error.** Bootstrapped 95% confidence intervals are shown. Starred items are statistically significant, calculated with Bonferroni correction: \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ . Additional treatments and settings are in Appendix Section B. While showing the model prediction increased accuracy, the explanation-based treatments did not differ significantly from showing the prediction alone. Pairwise differences between the top treatments were also not statistically significant, in all three settings shown.

not statistically significant in any setting. This is likely related to how explanation quality had little impact on the *trust* participants placed in the model predictions, discussed in the following section.

### 3.2 How does explanation quality affect human trust and understanding?

**Faulty explanations did not significantly decrease trust in model predictions.** We use the same regression model as in Equation 1 but with “trust”, the absolute value of the difference between the model’s prediction and the user’s guess, as the outcome variable. Smaller differences would indicate that explanations can increase the degree to which humans believe model predictions are accurate. The results are in Figure 4. Strong explanations could *increase* trust up to 1.08 years (CI lower bound) relative to the Prediction treatment (no explanations), while the random explanations could *decrease* trust up to 1.48 years (CI upper bound). These effect sizes could be sizable (30–40% relative to the intercept mean), but no treatment was statistically significant. Moreover, the difference between the spurious and random saliency maps is small, and none of the pairwise differences are statistically significant. These findings hold even when the model prediction is inaccurate (Model-), indicating that users are not learning to trust the model based on its predictions alone.

These findings, coupled with the large differences in accuracy between the H+M- and H-M+ settings and the decrease in accuracy in the H+M- setting (Figure 3), raise the question to what degree humans can identify and ignore erroneous model predictions. For example, a model that is accurate 98% of the time but makes large errors in the remaining 2% can be very dangerous in a high-stakes scenario if humans default to trusting all model predictions. Importantly, these subsets are typically unknown on test sets. Two possible approaches to alleviate this problem are (a) surfacing uncertainty measures as touched upon in the Show Top-3 Range treatment, and (b) training users by showing common model errors on a validation set. We also briefly expand upon the complementarity of human guesses and model decisions in Appendix Section C.2.

**Most participants claimed that explanations appeared reasonable, even when they were obviously not focused on faces.** Responses out of 7 to the post-survey question, “How easy to under-

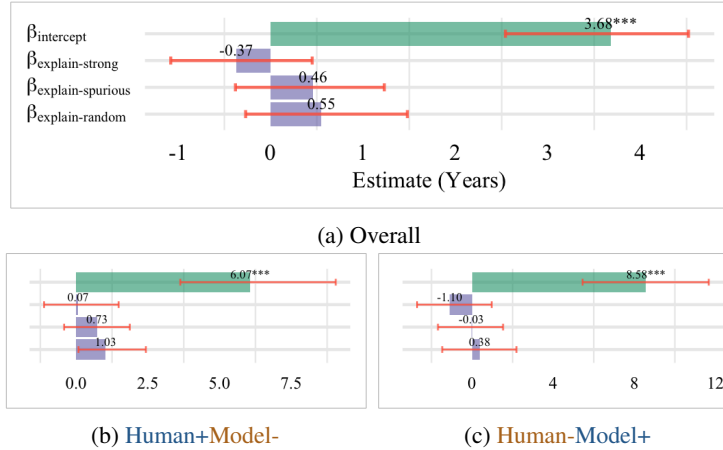


Figure 4: Estimates for  $y_{trust}$  regression. **Intercept represents the Prediction treatment; other estimates are relative to the intercept. Smaller values indicate more trust**, i.e. smaller difference between model prediction and human guess. Bootstrapped 95% confidence intervals are shown. Starred items are statistically significant, calculated with Bonferroni correction: \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ . Additional stats are in Appendix Section B. Pairwise differences between explanation-based treatments were not significant, and faulty explanations did not significantly decrease trust.

stand were the AI’s explanations? Did the explanations seem reasonable?”, are shown in Figure 5. Despite the clear flaws in the saliency maps shown in Figure 1, there were only small differences between the treatments. Participants shown a strong explanation rated the explanations 5.37 / 7, versus 5.05 and 4.71 for the spurious and random explanations. In the next section, we provide qualitative examples to shed some light into these counterintuitive findings.

### 3.3 How did humans incorporate model predictions and explanations into their decision making process?

Responses to the post-survey question “How did you use the model’s prediction in your decision-making process?” provide some clues to the workings of our human-AI system. Participants in the Explain-random treatment did highlight the randomness of the saliency maps, with one saying “I took it slightly into consideration but didn’t weight it heavily because it looked like the model was picking inaccurate locations...”. However, many others focused simply on the accuracy of the prediction, with one stating “Well it did a poor job recognizing face or features but the ages sound mostly correct so i sort of went with it”. We again note, however, that faulty explanations did not significantly decrease trust even in the presence of *inaccurate* model predictions (Figure 3b). Participants in the Explain-spurious treatment were similar, sometimes noting that the explanations were “totally out of whack, like on the wall”, but with only a few, explicit statements of these explanations mediating judgment, such as “If it was close to an area that might actually matter (neck, under eyes, edges of mouth etc) I took it into consideration, but if not, I dismissed it.”

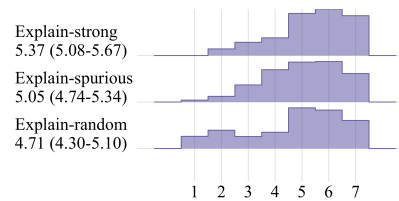


Figure 5: Distribution, mean, and 95% confidence intervals on the intelligibility of explanations. Response on a 1-7 Likert scale. Participants rated explanations similarly, regardless of the quality of the explanations.

We also examined responses for the top 75 guessers in terms of mean error, (collective MAE of 5.19). Answers to how the model prediction was used were bucketed into 6 categories: (1) 10.3% ignored it, (2) 17.9% considered it, (3) 5.1% used it if they were unsure, (4) 28.2% used the model prediction as a starting point, (5) 21.8% had their own initial guess and adjusted based on the model prediction, (6) 9.0% used the model prediction if it seemed reasonable; otherwise gave their own guess.

## 4 Related Work

**Interpretable machine learning** There has been a wide range of methods, including instance-level explanations [35] vs. global explanations based on feature representations across the dataset [4], glass-box methods with access to model gradients [40, 39] vs. black box methods [35, 32], and input feature attributions [35, 40, 39, 32, 4] vs. natural language explanations [20, 27] vs. counterfactuals [43]. Input feature attributions has been a common approach, which we use for our experiment.

Evaluating explanations has included human assessments of correctness, attempts to derive insights from explanations, and seeing if explanations help humans predict model behavior [40, 35, 11, 30]. Recent work has also suggested that popular interpretability methods exhibit undesirable properties, such as being unreliably and unreasonably sensitive to minor changes in the input [21, 23, 16].

**Model-in-the-Loop experiments** Previous experiments have mixed results on the effect of explanations. In a deceptive review detection task, explanations alone helped human performance slightly. However, the best setting was simply showing the prediction alone while informing users of the model’s accuracy [26]. Other work has shown feature-wise importances for low-dimensional inputs to both slightly increase [17] and decrease accuracy [49]. Explanations have also been found to increase trust and understanding [26, 28], but can hurt if they too overtly convey the model’s limitations [46]. Model complexity was also examined in [34], which found that less complex, transparent, linear models did not help users in an apartment price prediction task. In fact, transparency hindered people from detecting model errors, echoing other work on the risk of cognitive overload [2, 24].

We expand upon these prior works, as their limitations include (a) linear, relatively simple, or non-neural models [34, 16, 7, 24, 26], (b) small input feature spaces, making explanations simpler [17, 34, 49], (c) imbalance of task performance between humans and AI, e.g. human performance is 51% on a binary classification task (near random), vs. 87% for the model in [26, 25], (d) no investigation into using explanations for certification (identifying faulty or biased models) [7, 16, 17, 25, 26, 34, 41, 49].

**Certification with explanations** In an ad recommendation setting, explanations allowed users to detect that models were highly limited [15]. The authors of [35] found users able to identify the better of two models using their interpretable method with 90% accuracy. [3] introduces “model parameter randomization” and “data randomization” tests to analyze whether explanations are specific to the model and input. However, there have not been extensive human studies of certification.

**Design of the human-AI system** Varying the stated accuracy of the model can greatly increase trust in the model, though this trust will decrease if the observed accuracy is low [47]. There are also potential benefits for displaying uncertainty estimates for each prediction [42], though [49] found no improvement when showing confidence scores in textual form. “Design” covers an even broader category of elements, such as incentives to use a ML-driven tool, level of human agency and interactivity, and the “unremarkability” of the system [45]. These can often determine success in real-life settings [9, 38, 10, 5], but are out of the scope of our study.

## 5 Conclusion

Ideally, contextualizing model predictions with explanations would help improve people’s decision-making process in model-in-the-loop settings. Randomized control trials on this age prediction task, however, found that additionally showing explanations with model predictions did not have a significant effect on accuracy, trust, or understanding. Moreover, the quality of the explanations were unimportant, and even faulty explanations did not significantly decrease human trust in the model.

Existing interpretable ML methods have largely been used as a debugging tool for researchers and industry engineers, rather than a mechanism for communicating with end users [6]. Our findings are a reminder of this point and suggest input feature attributions, while helpful for machine learning researchers, may be less useful in downstream decision-making problems. Other interpretability methods, such as counterfactuals or natural language explanations, may be more effective but should be properly motivated and evaluated by their downstream use if model-in-the-loop usage is a potential goal. This echoes broader trends in human-computer interaction [18, 19] and grounded ML research, which argue for greater situating of ML models and methods in real-world scenarios.

## Broader Impact

We chose our image classification task precisely because there are analogous models used in high-stakes situations. There are risks to examining such a problem, as these findings could also be used to improve ethically-dubious systems involving prediction of demographic attributes such as facial recognition for surveillance systems or criminal identification. Complementary human and model biases may also be amplified in a human-AI system, further harming already disproportionately marginalized communities.

However, we believe the machine learning community and designers of ML-powered tools may immediately benefit from the findings of our study, as it motivates more useful explanations and ways to design human-AI systems. Ultimately, we hope this will improve the legibility of automated decision systems for non-technical users and improve outcomes for those affected by ML-powered tools.

## Acknowledgements

Thank you to Nabeel Gillani, Martin Saveski, Doug Beeferman, Sneha Priscilla Makini, and Nazmus Saquib for helpful discussion about the design and analysis of the RCT, as well as Jesse Mu for feedback on the paper.

## References

- [1] Passport facial recognition checks fail to work with dark skin, 2019. <https://www.bbc.com/news/technology-49993647>, Last accessed on 2020-06-01.
- [2] Ashraf Abdul, Christian von der Weth, Mohan Kankanhalli, and Brian Y Lim. Cogam: Measuring and moderating cognitive load in machine learning model explanations. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2020.
- [3] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515, 2018.
- [4] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 3319–3327. IEEE, 2017.
- [5] Emma Beede, Elizabeth Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Paisan Ruamvi-boonsuk, and Laura M Vardoulakis. A human-centered evaluation of a deep learning system deployed in clinics for the detection of diabetic retinopathy. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2020.
- [6] Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Yunhan Jia, Joydeep Ghosh, Ruchir Puri, José MF Moura, and Peter Eckersley. Explainable machine learning in deployment. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 648–657, 2020.
- [7] Reuben Binns, Max Van Kleek, Michael Veale, Ulrik Lyngs, Jun Zhao, and Nigel Shadbolt. 'it's reducing a human being to a percentage' perceptions of justice in algorithmic decisions. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2018.
- [8] Carrie J Cai, Jonas Jongejan, and Jess Holbrook. The effects of example-based explanations in a machine learning interface. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, pages 258–262, 2019.
- [9] Carrie J Cai, Emily Reif, Narayan Hegde, Jason Hipp, Been Kim, Daniel Smilkov, Martin Wattenberg, Fernanda Viegas, Greg S Corrado, Martin C Stumpe, et al. Human-centered tools for coping with imperfect algorithms during medical decision-making. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2019.
- [10] Carrie J Cai, Samantha Winter, David Steiner, Lauren Wilcox, and Michael Terry. "hello ai": Uncovering the onboarding needs of medical practitioners for human-ai collaborative

- decision-making. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–24, 2019.
- [11] Arjun Chandrasekaran, Deshraj Yadav, Prithvijit Chattopadhyay, Viraj Prabhu, and Devi Parikh. It takes two to tango: Towards theory of ai’s mind. *arXiv preprint arXiv:1704.00717*, 2017.
- [12] Berkeley J Dietvorst, Joseph P Simmons, and Cade Massey. Algorithm aversion: People erroneously avoid algorithms after seeing them err. *Journal of Experimental Psychology: General*, 144(1):114, 2015.
- [13] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. 2017.
- [14] S Escalera X Baro I Guyon R Rothe. E Agustsson, R Timofte. Apparent and real age estimation in still images with deep residual regressors on appa-real database. In *12th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG), 2017*. IEEE, 2017.
- [15] Motahhare Eslami, Sneha R Krishna Kumaran, Christian Sandvig, and Karrie Karahalios. Communicating algorithmic process in online behavioral advertising. In *Proceedings of the 2018 CHI conference on human factors in computing systems*, pages 1–13, 2018.
- [16] Shi Feng and Jordan Boyd-Graber. What can ai do for me? evaluating machine learning interpretations in cooperative play. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, pages 229–239, 2019.
- [17] Ben Green and Yiling Chen. The principles and limits of algorithm-in-the-loop decision making. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–24, 2019.
- [18] Steve Harrison, Phoebe Sengers, and Deborah Tatar. Making epistemological trouble: Third-paradigm hci as successor science. *Interacting with Computers*, 23(5):385–392, 2011.
- [19] Steve Harrison, Deborah Tatar, and Phoebe Sengers. The three paradigms of hci. In *Alt. Chi. Session at the SIGCHI Conference on human factors in computing systems San Jose, California, USA*, pages 1–18, 2007.
- [20] Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. Generating visual explanations. In *European Conference on Computer Vision*, pages 3–19. Springer, 2016.
- [21] Sarthak Jain and Byron C Wallace. Attention is not explanation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3543–3556, 2019.
- [22] Daniel Kahneman, Stewart Paul Slovic, Paul Slovic, and Amos Tversky. *Judgment under uncertainty: Heuristics and biases*. Cambridge university press, 1982.
- [23] Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. The (un) reliability of saliency methods. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 267–280. Springer, 2019.
- [24] Isaac Lage, Emily Chen, Jeffrey He, Menaka Narayanan, Been Kim, Samuel J Gershman, and Finale Doshi-Velez. Human evaluation of models built for interpretability. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 7, pages 59–67, 2019.
- [25] Vivian Lai, Han Liu, and Chenhao Tan. Why is “chicago” deceptive? towards building model-driven tutorials for humans. *arXiv preprint arXiv:2001.05871*, 2020.
- [26] Vivian Lai and Chenhao Tan. On human predictions with explanations and predictions of machine learning models: A case study on deception detection. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 29–38, 2019.
- [27] Tao Lei, Regina Barzilay, and Tommi Jaakkola. Rationalizing neural predictions. *arXiv preprint arXiv:1606.04155*, 2016.
- [28] Brian Y Lim, Anind K Dey, and Daniel Avrahami. Why and why not explanations improve the intelligibility of context-aware intelligent systems. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 2119–2128, 2009.
- [29] Zachary C Lipton. The mythos of model interpretability. *arXiv preprint arXiv:1606.03490*, 2016.



- [30] Menaka Narayanan, Emily Chen, Jeffrey He, Been Kim, Sam Gershman, and Finale Doshi-Velez. How do humans understand explanations from machine learning systems? an evaluation of the human-interpretability of explanation. *arXiv preprint arXiv:1802.00682*, 2018.
- [31] Clifford Nass, BJ Fogg, and Youngme Moon. Can computers be teammates? *International Journal of Human-Computer Studies*, 45(6):669–678, 1996.
- [32] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*, 2018.
- [33] P Jonathon Phillips, Amy N Yates, Ying Hu, Carina A Hahn, Eilidh Noyes, Kelsey Jackson, Jacqueline G Cavazos, Géraldine Jeckeln, Rajeev Ranjan, Swami Sankaranarayanan, et al. Face recognition accuracy of forensic examiners, superrecognizers, and face recognition algorithms. *Proceedings of the National Academy of Sciences*, 115(24):6171–6176, 2018.
- [34] Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Vaughan, and Hanna Wallach. Manipulating and measuring model interpretability. *arXiv preprint arXiv:1802.07810*, 2018.
- [35] Marco Tulio Ribeiro, Sameer Singh, and Carlos 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*, pages 1135–1144. ACM, 2016.
- [36] Rasmus Rothe, Radu Timofte, and Luc Van Gool. Dex: Deep expectation of apparent age from a single image. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 10–15, 2015.
- [37] Tom Schuba. Cpd using controversial facial recognition program that scans billions of photos from facebook, other sites, 2020. <https://chicago.suntimes.com/crime/2020/1/29/21080729/clearview-ai-facial-recognition-chicago-police-cpd>, Last accessed on 2020-06-01.
- [38] Mark Sendak, Madeleine Clare Elish, Michael Gao, Joseph Futoma, William Ratliff, Marshall Nichols, Armando Bedoya, Suresh Balu, and Cara O’Brien. "the human body is a black box" supporting clinical decision-making with deep learning. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 99–109, 2020.
- [39] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- [40] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3319–3328. JMLR. org, 2017.
- [41] Sarah Tan, Julius Adebayo, Kori Inkpen, and Ece Kamar. Investigating human+ machine complementarity for recidivism predictions. *arXiv preprint arXiv:1808.09123*, 2018.
- [42] Anne Marthe van der Bles, Sander van der Linden, Alexandra LJ Freeman, James Mitchell, Ana B Galvao, Lisa Zaval, and David J Spiegelhalter. Communicating uncertainty about facts, numbers and science. *Royal Society open science*, 6(5):181870, 2019.
- [43] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.
- [44] Dayong Wang, Aditya Khosla, Rishab Gargeya, Humayun Irshad, and Andrew H Beck. Deep learning for identifying metastatic breast cancer. *arXiv preprint arXiv:1606.05718*, 2016.
- [45] Qian Yang, Aaron Steinfeld, and John Zimmerman. 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*, pages 1–11, 2019.
- [46] Michael Yeomans, Anuj Shah, Sendhil Mullainathan, and Jon Kleinberg. Making sense of recommendations. *Journal of Behavioral Decision Making*, 32(4):403–414, 2019.
- [47] Ming Yin, Jennifer Wortman Vaughan, and Hanna Wallach. Understanding the effect of accuracy on trust in machine learning models. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2019.
- [48] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.

- [49] Yunfeng Zhang, Q Vera Liao, and Rachel KE Bellamy. Effect of confidence and explanation on accuracy and trust calibration in ai-assisted decision making. *arXiv preprint arXiv:2001.02114*, 2020.

## A Experimental Setup

### A.1 Datasets, Models, and Training Details

The IMDB-WIKI dataset used to pretrain our strong model consists of 524,230 images and age labels, and was partitioned into 90% train, 5% validation, and 5% test splits. The train, validation, and test splits for the APPA-REAL dataset were included in the original dataset<sup>1</sup>. We also trained a “weak” model for use in one additional treatment, which was not pre-trained on the IMDB-WIKI dataset and not trained until convergence on the APPA-REAL dataset. Saliency maps from the weak model tended to focus on similar features, but were often more diffuse, with fewer red spots. The strong and weak models used pretrained ImageNet weights, while the spurious model was trained from scratch in order to be more tuned to the spurious correlations. The model MAEs are listed in Table S1.

Table S1: Age prediction model performance in terms of mean absolute error (MAE). The MAE for the Spurious model is on the modified APPA-REAL dataset and has slightly lower MAE than the Strong model precisely because it is tuned into additional spurious correlations.

Model	Valid MAE	Test MAE
Strong	4.62	5.24
Weak	5.84	6.86
Spurious	3.33	4.58

Each model was trained on one GeForce GTX 1080 Ti using Pytorch. We used the Adam optimizer with learning rate 0.001 (after sweeping over [0.01, 0.005, 0.001, 0.0005]). Images were scaled to  $224 \times 224$ . We performed data augmentation during training by applying additive Gaussian noise 12.5% of the time, Gaussian blur 12.5% of the time, 20 degree rotations, 5% scaling and translation operations, horizontal flipping, and hue and saturation adjustments.

### A.2 Treatment Arms

The avatar used for the *Empathetic* treatment is shown in Figure S1. The icon was made by Freepik from [www.flaticon.com](http://www.flaticon.com) and can be found at [https://www.flaticon.com/free-icon/robot\\_1587565](https://www.flaticon.com/free-icon/robot_1587565).



Figure S1: Pat the AI.

**Additional treatments** We also tested the effect of explanations alone, hypothesizing that they may have a small, but slight effect. The *Explain-strong, No Pred* and *Explain-weak, No Pred* treatments show the saliency maps from the strong and weak model, respectively, *without the prediction*.

We also tested a global, model-level explanation in the *Explain-guide, No Pred* treatment. Before the task begins, users are shown a grid of saliency maps and told that important regions are: “(1) the areas around the eyes, and (2) lines anywhere on the face. The next two most important regions are around the mouth and nose.” The researchers manually went through 200 saliency maps and tallied regions in red in order to determine these features. Users are reminded of these guidelines at every image. This approach is similar in spirit to [25].

For the faulty saliency maps, we additionally tested not stating the model’s performance in the *No Acc* treatments. We hypothesized that allowing users to come to their own conclusion about the model’s ability would result in the faulty explanations having a larger effect.

<sup>1</sup>Download link: <http://chalearnlap.cvc.uab.es/dataset/26/data/45/description/>

Table S2: Full list of treatment arms. All model predictions are from the same “strong” model.

	<b>Treatment Arm</b>	<b>Shorthand Description</b>
	Control	Ask for age without help of model
	Prediction	User is shown model’s prediction
<i>Design</i>	Delayed Prediction	User guesses before and after seeing model’s prediction
	Empathetic	Model is described as a fallible “Pat the AI”
	Show Top-3 Range	Prediction shown as range of top-3 values, e.g. 28-32
<i>Explanations</i>	Explain-strong	Show strong model’s saliency map
	Explain-spurious	Show spurious model’s saliency map
	Explain-random	Show random saliency map
	Explain-strong, No Pred	Show strong model’s saliency map, hide prediction
	Explain-weak, No Pred	Show weak model’s saliency map, hide prediction
	Explain-guide, No Pred	Show summary of feature importances, hide prediction
	Explain-spurious, No Acc	Show spurious model’s saliency map, hide model’s accuracy
	Explain-random, No Acc	Show random saliency map, hide model’s accuracy

### A.3 Participant Demographics

40.3% of the participants were female, with a mean age of 36.5 (standard deviation 10.8).

## B Results

The full regression results for MAE, accuracy, and trust are shown in Tables S3, S4, and S5. We briefly cover findings from the additional treatments as follows:

- The instance-level explanation-only treatments (*Explain-strong, No Pred* and *Explain-weak, No Pred*) had small, non-statistically significant effects on accuracy.
- The model-level explanation-only treatment, *Explain-guide, No Pred*, was helpful overall (1.2 years increase in accuracy, compared to 1.8 years increase in the *Prediction* treatment). It was also as helpful as the top performing treatments when both humans and models were inaccurate, i.e. in the **Human-Model-** setting.
- Contrary to our hypothesis, hiding the model performance did not significantly increase or decrease the effect of the faulty explanations. Directionally, however, it appeared to *increase trust* in the **Human-Model-** setting. It may be that the statement on model performance actually emphasized the fallability of the model.

Table S3: Mean absolute error of guesses for all treatments. Bootstrapped 95% confidence intervals in parentheses.

Treatment Arm	MAE
Control (Human Alone)	10.0 (9.4 - 10.5)
Model Alone	8.5 (8.3 - 8.7)
Prediction	8.4 (7.8 - 9.0)
Explain-strong	8.0 (7.5 - 8.5)
Explain-spurious	8.5 (8.0 - 9.1)
Explain-random	8.7 (8.1 - 9.2)
Delayed Prediction	8.5 (8.0 - 9.0)
Empathetic	8.0 (7.6 - 8.5)
Show Top-3 Range	8.0 (7.4 - 8.5)
Explain-strong, No Pred	9.7 (9.2 - 10.3)
Explain-weak, No Pred	10.2 (9.5 - 10.9)
Explain-guide, No Pred	9.4 (8.9 - 10.0)
Explain-spurious, No Acc	8.2 (7.6 - 8.8)
Explain-random, No Acc	8.0 (7.6 - 8.5)

Table S4: Estimates for  $y_{error}$  regression. **Intercept represents control treatment (human alone); other estimates are relative to the intercept. Lower values are better, indicating reduced error.** Bootstrapped 95% confidence intervals are shown in parentheses; starred items are statistically significant, calculated with Bonferroni correction: \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ .

Overall	
$\beta_{intercept}$	5.9 (4.6,7.2)***
$\beta_{prediction}$	-1.8 (-2.6,-1.1)***
$\beta_{delay-model-pred}$	-2.0 (-2.8,-1.2)***
$\beta_{empathetic}$	-2.2 (-3.0,-1.2)***
$\beta_{show-top-3-range}$	-2.2 (-2.9,-1.3)***
$\beta_{explain-strong}$	-2.2 (-3.2,-1.3)***
$\beta_{explain-spurious}$	-1.9 (-2.7,-1.2)***
$\beta_{explain-random}$	-1.6 (-2.5,-0.9)***
$\beta_{explain-strong\_no-pred}$	-0.5 (-1.5,0.3)
$\beta_{explain-weak\_no-pred}$	0.2 (-0.7,0.9)
$\beta_{explain-guide\_no-pred}$	-1.2 (-2.0,-0.4)*
$\beta_{explain-spurious\_no-acc}$	-1.9 (-2.7,-1.0)***
$\beta_{explain-random\_no-acc}$	-2.0 (-2.6,-1.2)***
$\beta_{image-age}$	0.1 (0.0,0.1)***

(a) Overall

	Human+Model+	Human+Model-	Human-Model+	Human-Model-
$N$	4940	1802	1721	2812
$\beta_{intercept}$	3.3 (2.9,3.7)***	4.2 (3.6,4.8)***	13.6 (12.1,14.8)***	12.7 (10.3,15.9)***
$\beta_{prediction}$	-0.3 (-0.6,0.0)	1.1 (0.5,1.8)**	-2.3 (-4.3,-1.1).	-2.0 (-3.4,-0.7)*
$\beta_{delay-model-pred}$	-0.1 (-0.4,0.3)	0.8 (0.3,1.3).	-2.9 (-4.5,-1.3)**	-2.3 (-3.5,-1.0)**
$\beta_{empathetic}$	-0.4 (-0.7,-0.0)	1.2 (0.6,2.0)**	-2.6 (-4.2,-0.9)*	-2.7 (-3.7,-1.6)***
$\beta_{show-top-3-range}$	-0.4 (-0.8,-0.0)	0.8 (0.3,1.4)	-4.5 (-6.4,-2.5)***	-1.7 (-2.8,-0.8)*
$\beta_{explain-strong}$	-0.2 (-0.6,0.1)	1.1 (0.5,1.6)*	-3.2 (-5.4,-1.2)**	-2.4 (-3.4,-1.2)**
$\beta_{explain-spurious}$	-0.2 (-0.5,0.1)	1.2 (0.7,1.7)**	-2.5 (-3.9,-0.9)*	-1.5 (-2.8,-0.3)
$\beta_{explain-random}$	-0.2 (-0.5,0.2)	0.8 (0.2,1.3)	-2.8 (-4.5,-1.1)**	-1.7 (-3.0,-0.5).
$\beta_{explain-strong\_no-pred}$	0.4 (0.0,0.9)	0.2 (-0.3,0.8)	-1.1 (-2.8,0.4)	-1.4 (-2.7,0.0)
$\beta_{explain-weak\_no-pred}$	0.1 (-0.3,0.7)	-0.1 (-0.8,0.8)	1.0 (-0.9,2.4)	-0.6 (-1.8,0.9)
$\beta_{explain-guide\_no-pred}$	0.0 (-0.2,0.4)	-0.2 (-0.8,0.3)	-1.1 (-2.6,0.8)	-2.4 (-3.9,-1.1)**
$\beta_{explain-spurious\_no-acc}$	-0.3 (-0.7,0.1)	1.1 (0.5,1.9)**	-2.3 (-3.7,-0.8)	-1.2 (-2.4,0.1)
$\beta_{explain-random\_no-acc}$	-0.2 (-0.5,0.3)	1.1 (0.6,1.8)**	-3.0 (-4.7,-1.6)**	-2.2 (-3.2,-1.1)**
$\beta_{image-age}$	0.0 (0.0,0.0)*	-0.0 (-0.0,0.0)	0.0 (0.0,0.1)	0.1 (0.0,0.1)*

(b) Splits



Table S5: Estimates for  $y_{trust}$  regression. **Intercept represents Prediction treatment; other estimates are relative to the intercept. Lower values indicate greater trust.** Bootstrapped 95% confidence intervals are shown in parentheses; starred items are statistically significant, calculated with Bonferroni correction: \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ .

Overall	
$\beta_{intercept}$	3.7 (2.5,4.5)***
$\beta_{explain-strong}$	-0.4 (-1.1,0.4)
$\beta_{explain-spurious}$	0.5 (-0.4,1.2)
$\beta_{explain-random}$	0.5 (-0.3,1.5)
$\beta_{explain-spurious\_no-acc}$	0.3 (-0.5,1.0)
$\beta_{explain-random\_no-acc}$	0.4 (-0.5,1.0)
$\beta_{image-age}$	0.0 (0.0,0.0)*

(a) Overall

	Human+Model+	Human+Model-	Human-Model+	Human-Model-
$N$	4940	1802	1721	2812
$\beta_{explain-strong}$	-0.1 (-0.5,0.5)	0.1 (-1.1,1.5)	-1.1 (-2.7,1.0)	0.1 (-1.5,1.6)
$\beta_{explain-spurious}$	0.1 (-0.3,0.7)	0.7 (-0.4,1.9)	-0.0 (-1.7,1.5)	1.0 (-0.3,2.1)
$\beta_{explain-random}$	0.0 (-0.4,0.4)	1.0 (0.1,2.4)	0.4 (-1.5,2.2)	1.4 (-0.2,2.6)
$\beta_{explain-spurious\_no-acc}$	0.3 (-0.1,0.7)	0.7 (-0.2,1.9)	1.2 (-1.2,3.1)	-0.2 (-1.8,1.1)
$\beta_{explain-random\_no-acc}$	0.3 (-0.1,0.8)	0.7 (-0.3,1.9)	-0.3 (-2.0,1.5)	0.2 (-1.2,1.2)
$\beta_{image-age}$	0.0 (-0.0,0.0)	0.0 (-0.0,0.1)	0.0 (-0.0,0.1)	0.0 (-0.0,0.1)

(b) Splits

## C Additional Analyses

### C.1 The *economic cost* of Model-in-the-Loop predictions

We also consider the *economic cost* of a prediction, calculated simply as error  $\times$  time. Under this metric, the treatment arms are similar or worse than simply showing the model prediction: Control (56.7), Empathetic (50.6), Explain-strong (47.2), Show Model Pred (43.5), Show Top-3 Range (42.5). We use this simply as an illustrative example, as we believe time and cognitive load are important considerations when designing human-AI systems, especially with possibly complex explanations of high-dimensional inputs.

### C.2 Combining human guesses and model predictions

We investigate the possible gain of combining the two predictions in simple hybrid models, whose input features are simply the model’s prediction and the human’s guess. Such models have been found to outperform either human or machine alone [44, 33], though this may not always be the case [41]. We perform cross-fold validation and hyperparameter search, with the test MAE results for the top-performing treatments shown in Table S6.

Follow-up questions to our work include: (a) how the differences between treatments can be related to the complementarity of human and model predictions, and (b) how to best design human-AI systems where the AI complements existing human capabilities. One experiment could be to derive coarse strategies from these hybrid models, such as important decision tree rules, and test whether these strategies could help further improve accuracy (i.e. “Trust the model if you think the person is above 60, and the model’s prediction is significantly greater than yours.”)

Table S6: MAE on test split of hybrid models that combine human guesses and model predictions

	Prediction	Explain-strong	Empathetic	Show Top-3 Range
Human guess (with AI assistance)	8.4	8.0	8.0	8.0
Model prediction	8.4	8.4	8.3	8.5
Logistic Regression	8.1	8.3	7.4	9.4
Decision Tree	7.1	6.3	8.4	6.0