# Explainable AI improves task performance in human-AI collaboration

Julian Senoner[1*], Simon Schallmoser[2,3*], Bernhard Kratzwald[1],

Stefan Feuerriegel[2,3], Torbjørn Netland[1†]

[1]ETH Zurich, Zurich, Switzerland

[2]LMU Munich, Munich, Germany

[3]Munich Center for Machine Learning (MCML), Munich, Germany

[*]Contributed equally

[†]Correspondence: Torbjørn Netland (tnetland@ethz.ch)

1

## Abstract

Artificial intelligence (AI) provides considerable opportunities to assist human work. However, one crucial challenge of human-AI collaboration is that many AI algorithms operate in a black-box manner where the way how the AI makes predictions remains opaque. This makes it difficult for humans to validate a prediction made by AI against their own domain knowledge. For this reason, we hypothesize that augmenting humans with explainable AI as a decision aid improves task performance in human-AI collaboration. To test this hypothesis, we analyze the effect of augmenting domain experts with explainable AI in the form of visual heatmaps. We then compare participants that were either supported by (a) black-box AI or (b) explainable AI, where the latter supports them to follow AI predictions when the AI is accurate or overrule the AI when the AI predictions are wrong. We conducted two preregistered experiments with representative, real-world visual inspection tasks from manufacturing and medicine. The first experiment was conducted with factory workers from an electronics factory, who performed $N = 9,600$ assessments of whether electronic products have defects. The second experiment was conducted with radiologists, who performed $N = 5,650$ assessments of chest X-ray images to identify lung lesions. The results of our experiments with domain experts performing real-world tasks show that task performance improves when participants are supported by explainable AI instead of black-box AI. For example, in the manufacturing setting, we find that augmenting participants with explainable AI (as opposed to black-box AI) leads to a five-fold decrease in the median error rate of human decisions, which gives a significant improvement in task performance.

**Keywords**: Explainable AI, task performance, decision-making, human-centered AI, human-AI collaboration

# Introduction

Artificial intelligence (AI) provides considerable opportunities to assist human work in various domains [1, 2]. For example, in manufacturing, AI is widely used to support humans when inspecting the quality of produced products to identify defects [3]. Similarly, in medicine, disease diagnosis now makes increasing use of AI systems. For instance, a recent survey found that AI is used by about 15% of radiologists at least weekly [4]. More broadly, an analysis showed that about 30% of all jobs in the United States are at high exposure to be assisted by AI [5]. Hence, the importance of human-AI collaboration is expected to grow in the near future.

However, many questions regarding the effective design of human-AI collaborations remain open. One particular challenge in the use of AI for human work is that state-of-the-art AI algorithms, which frequently involve millions of trainable parameters [6, 7], operate as "black-box" algorithms. The term "black-box" refers to the opacity of these systems, meaning that the internal workings and decision-making processes of these algorithms are not transparent or easily understandable by humans [8]. This can have crucial implications in practice as the lack of transparency makes it difficult – or even impossible – for humans to validate the predictions made by an AI against their domain knowledge. Hence, without being able to assess whether a prediction generated by an AI is accurate, humans will not be able to correct predictions of the AI, because of which the unique expertise of workers is essentially lost, which will make the collaboration between domain experts and AI largely ineffective.

Increasing efforts have been made to overcome the black-box nature of AI by developing methods that generate *explanations* for how AI algorithms reach their decisions [9, 10, 11, 12, 13, 14, 15]. Explainable AI refers to a set of methods that support humans in understanding how AI algorithms map certain inputs (e.g., lung X-rays, patient characteristics) to certain outputs (e.g., probability estimates for pneumonia) [16, 17]. Explainable AI can be broadly categorized into inherently interpretable models and post-hoc explanation techniques (see Supplement A for an extended literature review on explainable AI). For inherently interpretable algorithms, the decision-making of the algorithm can be inspected by humans, e.g., by inspecting the coefficients in linear regression or the splitting rules in decision trees [18]. Post-hoc explanation techniques are required when the inner workings of an AI algorithm become too complex to be understood by humans such as in deep neural networks. For example, one approach is to approximate the

behavior of a black-box AI with a simpler model (e.g., a linear model) that can be interpreted [19]. Other methods rely on game theory to estimate the contribution of each model input to the model output while considering possible interaction effects [20]. Common methods for explaining AI algorithms in computer vision include the use of heatmaps. These heatmaps visually highlight the areas that are most relevant to the predictions made by the AI [21, 22]. Such explanation techniques are commonly used by AI engineers in the development of AI algorithms. Hence, this literature stream is orthogonal to the use of explainable AI in our work, where we use post-hoc explanation techniques to improve decisions by domain experts in real-world job tasks.

Several works have studied behavioral dimensions of human-AI collaboration. For example, it has been examined whether humans are willing to delegate work to AI [23, 24, 25]. Another common dimension is algorithm aversion, where humans are averse to following decisions by algorithms and instead rely on their own (mis)judgment [26, 27, 28, 29, 30]. An antecedent to algorithm aversion is *trust in AI*, critically influencing whether humans adopt or reject AI recommendations [31, 32, 33]. Oppositely to algorithm aversion, overreliance is also a problem negatively impacting the effectiveness of human-AI collaboration [34, 35, 36]. That is, humans risk following AI predictions blindly without attentively performing the task. While all these dimensions are interesting from a behavioral perspective, the main outcome of interest for business and healthcare organizations is *task performance*. However, the impact of explainable AI on task performance in human-AI collaboration in real-world job tasks remains unclear.

We hypothesize that augmenting domain experts with explainable AI, as opposed to black-box AI, improves task performance in human-AI collaboration. Specifically, we treat explainable AI as a form of decision aid that supports domain experts in better understanding algorithmic decisions. Experts can then compare the explanations to their domain knowledge, thereby validating whether the AI is correct or overwriting the AI if is not correct. Here, explainable AI does not provide more information from an AI perspective (i.e., identical predictive performance). However, for domain experts, it gives rich additional information by making the AI predictions more accessible. Thus, we expect that domain experts supported by explainable AI will outperform those supported by black-box AI in two ways: (1) they are more likely to follow AI predictions when they are accurate, and (2) they are more likely to overrule AI predictions when they were wrong.

Previous research has studied the effect of explainable AI on task performance in human-AI

collaboration (see Supplement A for a detailed overview), yet with key limitations. In particular, existing works are typically restricted by either (i) recruiting laypeople or (ii) overly simplified tasks that are not representative of real job tasks [37, 38, 39, 40]. However, a realistic estimate of the effect of explainable AI on task performance requires a real-world task performed by domain experts. Such works that actually study real-world tasks with domain experts are on the other hand restricted by (i) comparing explainable AI vs humans alone [41, 42], (ii) using no real explainable AI [43], or (iii) research designs that do not isolate the effect of explainable AI on task performance [44, 45, 46]. In contrast, the strength of our work is that we study the effect of explainable AI on task performance relative to black-box AI in human-AI collaborations with *real-world tasks* and actual *domain experts*.

In this paper, we analyze the effect of augmenting domain experts with explainable AI on task performance in human-AI collaboration. For this, we conducted two preregistered experiments in which domain experts were asked to solve real-world visual inspection tasks in manufacturing (Study 1) and medicine (Study 2). We followed a between-subject design where we randomly assigned participants to two treatments: (a) black-box AI (i.e., where AI predictions are opaque) and (b) explainable AI (i.e., where AI predictions are explained). The latter thus offers not only the prediction from the AI but further shows explanations in the form of a visual heatmap as a decision aid. Heatmaps are frequently used and are considered state-of-the-art with respect to their localization performance across various settings [47, 48, 49, 50]. Study 1 was conducted in a manufacturing setting, where participants had to identify quality defects in electronic products. For this, we specifically recruited actual factory workers performing $N = 9,600$ assessments of electronic products at *Siemens*. Study 2 was conducted in a medical setting, where participants had to identify lung lesions on chest X-ray images. To that end, medical professionals, i.e., radiologists, were recruited and performed $N = 5,650$ assessments of chest X-ray images. In both studies, participants performed better when being supported by explainable AI as a decision aid.

The tasks of both experiments are representative of many real-world human-AI collaborations. The manufacturing task is an identical, one-to-one copy of a real-world job task at *Siemens* and, hence, highly representative of visual inspection tasks in manufacturing [51, 52]. Visual inspection tasks are standard in the manufacturing industry. Regardless of how much manufacturers have sought to build quality into products and processes, labor-intensive inspec-

tion tasks still abound [53]. In healthcare, visual inspection tasks are common across many different subdisciplines such as dermatology, radiology, pathology, ophthalmology, and dentistry, among many others. As concrete examples, physicians have to inspect, for instance, skin lesions in dermatology, tissues in pathology, and lung lesions in radiology [54, 55, 50]. Hence, establishing whether physicians benefit from explainable AI in visual inspection tasks is highly relevant for setting correct disease diagnosis and subsequent treatment.

## Results

To analyze the effect of explainable AI on task performance in human-AI collaboration, we conducted two randomized experiments across two different settings, i.e., in manufacturing (Study 1) and medicine (Study 2). In both experiments, participants had to perform a visual inspection task. In the manufacturing experiment, factory workers were asked to inspect electronic products and to identify defective products. In the medical experiment, radiologists were asked to decide whether lung lesions are visible in chest X-ray images. Participants were randomly assigned to one of two different treatments aiding them in the task: (a) black-box AI or (b) explainable AI (Figure 1). Participants with black-box AI received an opaque AI score as a decision aid. Participants with explainable AI received the same score and an additional decision aid: the explanation of the score in the form of a heatmap. The heatmap does not provide more information from an AI perspective (the score is identical) but allows users to verify the prediction made by the AI. However, heatmaps provide a clear and intuitive way of highlighting quality defects/lung lesions [56]. We hypothesized that explainable AI as a decision aid improves task performance of domain experts in human-AI collaboration. Details on both experiments are provided in the Methods section.

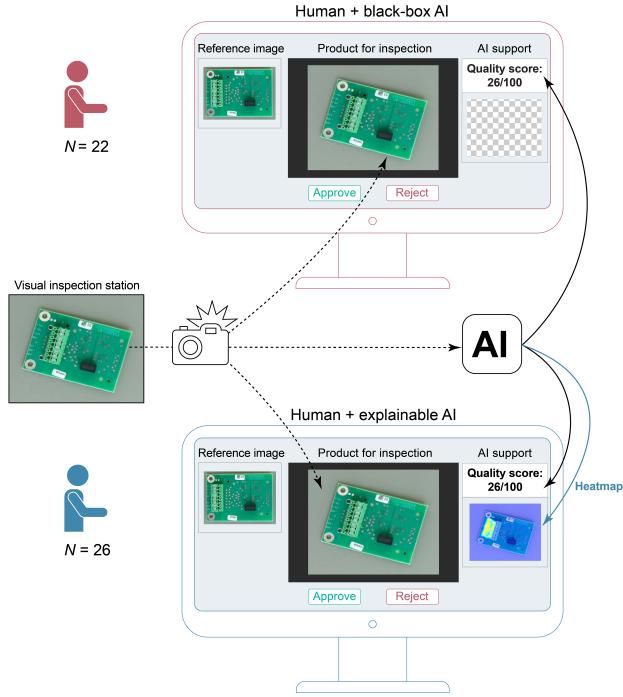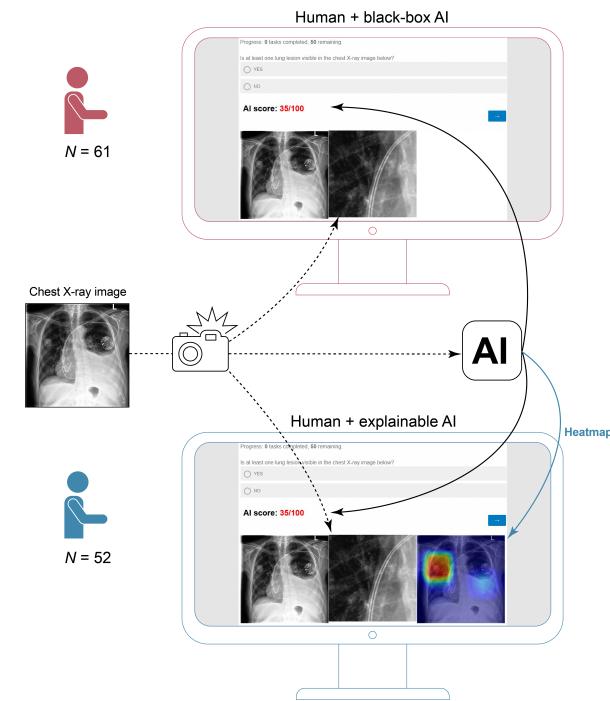**Figure 1: Overview of the experiments for assessing the effect of explainable AI on task performance.** (**A**) Experimental design of the manufacturing experiment where factory workers were asked to "approve" images of faultless products and to "reject" images of defective products through a computer interface. (**B**) Experimental design of the medical experiment where radiologists were asked to decide whether lung lesions are visible in the chest X-ray image. In both experiments, participants were randomly assigned to one of the two treatments: (a) black-box AI or (b) explainable AI.

## Study 1: Manufacturing experiment

The manufacturing experiment was conducted at a factory of *Siemens*, a global industrial conglomerate particularly known for its consumer and industrial electronics products. The objective of the task was to identify quality defects in electronic products (e.g., missing components, wrong components, and faulty components) with high accuracy. The task is representative of a real-world inspection task at *Siemens* and is analogous to other inspection tasks in manufacturing [51, 52]. Factory workers from *Siemens* were asked to visually inspect 200 images using a computer interface and label them as faultless or defective. The inspection task had to be completed within 35 minutes, which corresponds to realistic field conditions.

Workers were randomly assigned to two treatments where they were either supported by black-box AI or explainable AI (Figure 1**A**). In both treatment arms, workers received a reference image of a faultless product (note that this is common practice at *Siemens*). In addition, they were provided with AI predictions of the product quality given by a numerical "quality score" between 0 (most certainly defect) and 100 (most certainly faultless). Workers with explainable AI additionally received a tailored decision aid in the form of visual heatmaps that indicated the assumed location of potential quality defects. The quality scores in both treatment arms were identical, so that differences in task performance could only be attributed to the explanations (i.e., the heatmaps).

The objective of the field experiment was to obtain accurate estimates of the treatment effect under real-world conditions. Hence, we ran the experiment with actual domain experts rather than laypeople. The results are based on the entire available workforce of one shift of factory workers from *Siemens*. The factory workers performed in total $N = 9,600$ assessments of electronic products. We then analyzed the effect of explainable AI on task performance using the balanced accuracy and defect detection rate (i.e., proportion of correctly identified defective products among all actual defective products) based on the quality assessments in the visual inspection task.

We found that workers supported by explainable AI achieved a better task performance than workers supported by black-box AI. Workers with black-box AI achieved a balanced accuracy with a mean of only 88.6%, whereas workers with explainable AI treatment achieved a balanced accuracy with a mean of 96.3% (Figure 2**A**). We then estimated the treatment effect of explain-

able AI by regressing the balanced accuracy on the treatment (black-box AI = 0, explainable AI = 1). The regression results show that the treatment effect of explainable AI is statistically significant and large ($\beta = 7.653$, $SE = 2.178$, $P = 0.001$); that is, an improvement of 7.7 percentage points. Compared to the black-box AI, the explainable AI leads to a five-fold decrease in the median error rate.

Workers with explainable AI outperformed workers with black-box AI also with respect to the defect detection rate with a mean of 93.0% versus a mean of 82.0% (Figure 2**B**). The regression results again confirm that the treatment effect of explainable AI is statistically significant and large ($\beta = 11.014$, $SE = 3.680$, $P = 0.004$). All regression results remain statistically significant when including relevant control variables (demographics, tenure, self-reported IT skills, and decision speed) in the regression model (see Supplement H.1).

A detailed analysis of the workers' assessments revealed that workers with explainable AI followed accurate predictions more often than workers with black-box AI (mean = 93.5% for black-box AI, mean = 98.6% for explainable AI). In particular, workers supported by black-box AI were 3.6 times more likely to erroneously overrule an AI prediction, despite the prediction being accurate ($t = 2.437$, $P = 0.011$). Interestingly, 73.1% of the workers with explainable AI performed even better than the standalone AI algorithm. This suggests that the explanations (i.e., the heatmaps) not only improve adherence to accurate AI predictions, but also help humans make correct assessments when the AI predictions are wrong. We found that workers with explainable AI were, on average, able to identify and overrule 96.9% of the wrong AI predictions. For comparison, workers supported by black-box AI only overruled 86.4% of the wrong AI predictions. These results are highly relevant since – regardless of an AI's performance – wrong AI predictions can always occur due to external factors such as dust or different light conditions. The difference between both treatments is again statistically significant ($t = 2.631$, $P = 0.007$). These findings underscore the effectiveness of augmenting humans with explainable AI.
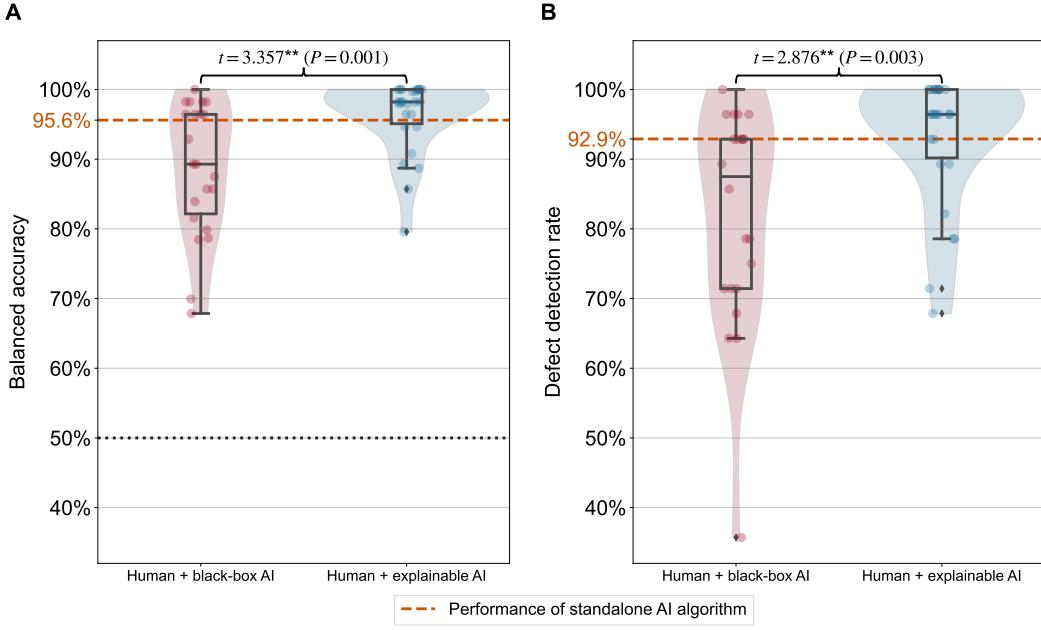
**Figure 2: Results of manufacturing experiment.** The boxplots compare the task performance between the two treatments: black-box AI and explainable AI. The task performance is measured by the balanced accuracy (**A**) and the defect detection rate (**B**) based on the quality assessment of workers and the ground-truth labels of the product images. A balanced accuracy of 50% provides a naïve baseline corresponding to a random guess (black dotted line). The standalone AI algorithm attains a balanced accuracy of 95.6% and a defect detection rate of 92.9% (orange dashed lines). Statistical significance is based on a one-sided Welch's $t$-test ($^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$). In the boxplots, the center line denotes the median; box limits are upper and lower quartiles; whiskers are defined as the 1.5x interquartile range.

Finally, we assessed whether workers with explainable AI spent more time on making their quality assessments. For this, we analyzed whether the workers' median decision speeds across the 200 product images differed. No statistically significant difference ($t = 0.308$, $P = 0.380$) was observed between both treatments (mean $= 5.01\,$s for black-box AI, mean $= 4.88\,$s for explainable AI). Therefore, explainable AI improved task performance without affecting the productivity of the workers.

## Study 2: Medical experiment

In the medical experiment, radiologists were asked to visually inspect 50 chest X-ray images and decide whether at least one lung lesion was visible (Figure 1**B**). Visual inspection tasks like ours are common in medicine across various subdisciplines [55, 54]. Analogous to the manufacturing task, radiologists had 35 minutes to complete the task and were randomly assigned to be either

supported by black-box AI or by explainable AI (Figure 1**B**). Both types of AI provided a score between 0 (most certainly a lung lesion visible) and 100 (most unlikely a lung lesion visible), which was identical in both treatment arms. In addition to that, radiologists with explainable AI were provided with a heatmap that highlights regions in the chest X-ray image that the AI finds most relevant for predicting lung lesions.

The results are based on a sample of $N = 5,650$ assessments of chest X-ray images performed by 113 radiologists from the United States. Again, task performance was analyzed using the balanced accuracy and the disease detection rate (i.e., the true negative rate, where chest X-ray images containing lung lesions were considered a negative sample) based on the assessments made in the visual inspection task.

Radiologists augmented with explainable AI outperformed peers with black-box AI. Radiologists with black-box AI achieved a balanced accuracy with a mean of only 79.1%, whereas radiologists with explainable AI achieved a balanced accuracy with a mean of 83.8% (Figure 3**A**). We again estimated the treatment effect of explainable AI by regressing the balanced accuracy on the treatment (black-box AI = 0, explainable AI = 1). The regression results show that the treatment effect of explainable AI is statistically significant and large ($\beta = 4.693$, $SE = 1.800$, $P = 0.01$); that is, an improvement of 4.7 percentage points. All results remain statistically significant when including relevant control variables (tenure, self-reported IT skills, and decision speed) in the regression model (see Supplement H.2). In contrast to the manufacturing experiment, no difference in task performance with respect to the disease detection rate was observed; radiologists in both treatment arms achieved a disease detection rate with a mean of 90.4% (Figure 3**B**). This was also observed when regressing the disease detection rate on the treatment ($\beta = -0.014$, $SE = 2.244$, $P = 0.995$). This can be expected since missing a lung lesion has more serious consequences than erroneously believing a lung lesion is visible; thus, leading to conservative decision-making of radiologists. Therefore, we additionally inspected precision as a task performance metric. We find that radiologists augmented with explainable AI were significantly more precise (improvement of 6.4 percentage points, $P = 0.014$) in identifying lung lesions compared to radiologists with black-box AI (see Supplement G).

As in Study 1, we found that radiologists with explainable AI followed accurate AI predictions more often than radiologists with black-box AI treatment (mean = 72.4% for black-box AI, mean = 82.1% for explainable AI). In particular, radiologists supported by black-box AI were

54.2% times more likely to erroneously overrule an AI prediction, although it was correct ($t = 3.084$, $P = 0.001$). We observed that radiologists with explainable AI only overruled 50.8% of the wrong AI predictions compared to 57.7% for radiologists with black-box AI treatment.
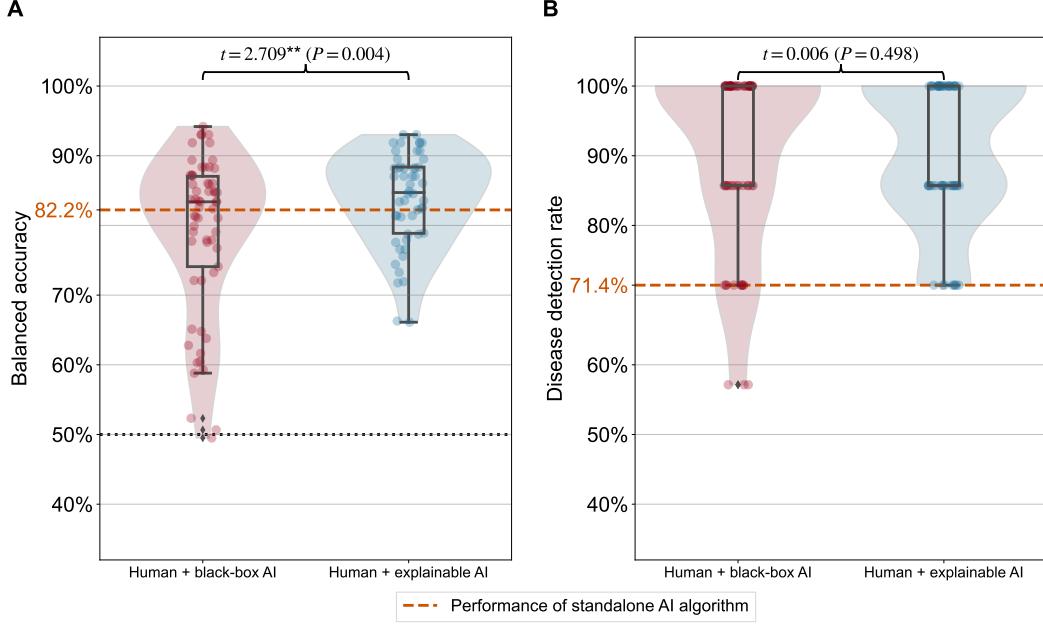


**Figure 3: Results of medical experiment.** The boxplots compare the task performance between the two treatments: black-box AI and explainable AI. The task performance is measured by the balanced accuracy (**A**) and the disease detection rate (**B**) based on the quality assessment of radiologists and the ground-truth labels of the chest X-ray images. A balanced accuracy of 50% provides a naïve baseline corresponding to a random guess (black dotted line). The standalone AI algorithm attains a balanced accuracy of 82.2% and a disease detection rate of 71.4% (orange dashed lines). Statistical significance is based on a one-sided Welch's $t$-test ($^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$). In the boxplots, the center line denotes the median; box limits are upper and lower quartiles; whiskers are defined as the 1.5x interquartile range.

Again, we assessed the decision speed of radiologists in both treatment arms. We found no significant difference ($t = 0.392$, $P = 0.348$) between both treatments (mean $= 10.71\,$s for black-box AI, mean $= 10.29\,$s for explainable AI). Thus, task performance was improved by explainable AI without reducing the productivity of the radiologists.

# Discussion

The "age of AI" redefines the way humans and machines collaborate, thus raising questions about how human-AI collaborations can be effectively designed. As we show, the effectiveness of

human-AI collaboration largely depends on the extent to which humans incorporate correct AI predictions and overrule wrong ones. However, many state-of-the-art AI algorithms operate as black-box, thus making it difficult for humans to compare the reasoning of the AI to their own domain knowledge. In this paper, we contribute a unique perspective by studying the impact of AI explainability on task performance of domain experts in human-AI collaboration, presenting empirical evidence from different domains with robust and generalizable results.

We conducted two preregistered experiments to estimate the effect of explainability in human-AI collaboration in real-world visual inspection tasks. Our results demonstrate that domain experts make subpar decisions when they are supported by a black-box AI algorithm with opaque predictions. In contrast, we find that explanations from an explainable AI are a powerful decision aid. Explanations were provided in the form of heatmaps, which provide a clear and intuitive way of highlighting areas that are determinants of AI predictions. The explanations do not provide more information from an AI perspective (i.e., the prediction performance is identical), but rather make the information more accessible to domain experts. Specifically, compared to black-box AI, augmenting domain experts with explainable AI improved the task performance by 7.7 percentage points in a manufacturing experiment and by 4.7 percentage points in a medical experiment. In the manufacturing experiment, 73.1% of the domain experts even outperformed the standalone AI algorithm when they were augmented with explainable AI. The prime reason was that domain experts supported by explainable AI were more likely to follow AI predictions when they were accurate and more likely to overrule them when they were wrong.

Improving the task performance of domain experts has practical implications in many fields such as manufacturing and medicine. For example, factory workers at *Siemens* augmented with explainable AI were able to identify 13% more defects than peers augmented with black-box AI. Thus, explainable AI could help to reduce downstream costs for manufacturing companies by filtering out defective products at the earliest possible stage. Similarly in medicine, task performance of physicians is crucial, especially for important tasks such as identifying possibly cancerogenous lung lesions. By showing that the results are consistent across different settings, we demonstrate that our insights are generalizable.

Our work contributes experimental evidence to the literature on human-AI collaboration [57, 58, 59, 60, 26, 27, 28, 61, 62, 63, 64, 65]. Algorithm aversion provides a barrier to the wider adoption of human-AI collaboration. Prior literature has presented several remedies, such

as describing the functional logic of an algorithm [60], giving users permission to modify an algorithm [27], or letting users integrate their own forecasts into an algorithm [62]. This paper presents evidence of an effective alternative; that is, explaining individual predictions from an otherwise opaque AI algorithm. Such explanations allow domain experts to validate how an AI arrives at a certain prediction. Interestingly, while many studies advocate for complete automation, we also show the importance of human-AI collaboration: domain knowledge can help to identify errors in the AI and lead to a better performance than an AI-only system.

A strength of our study is that we gather empirical evidence of improved task performance by explainable AI compared to black-box AI. In particular, by performing experiments in two different settings, we demonstrate that these results are generalizable. Unlike previous works on studying task performance in human-AI collaboration [37, 38, 39, 40], we (i) conducted experiments of two real-world job tasks in manufacturing and medicine and (ii) recruited domain experts for those tasks, i.e., factory workers and radiologists. When experiments use simplified decision tasks (e.g., object recognition) that are not representative of actual human work in the field, real-world validity is reduced. In contrast, our study has high external validity.

Our work is orthogonal to the literature on explainable AI in computer science, where the main goal is to develop and evaluate new methods for explaining black-box AI algorithms. Contrary, we are interested in a behavioral outcome, namely task performance in human-AI collaboration. A previous study on the effect of explainable AI on task performance found an improvement of 1.5 percentage points in accuracy relative to black-box AI [46]. However, their experiment was designed such that participants were not only shown the real explainable AI but also systemically biased explainable AI. This could have decreased the trust of the participants in the AI and, thus, explain the smaller treatment effect in comparison to our experiments. Prior work has also made use of expert annotations as a proxy for explainable AI [43]. However, this prevents any conclusion on whether real explainable AI improves task performance.

One limitation of our research is that we, in both experiments, studied one specific human-AI work setting (a visual inspection task) with one specific form of explainability (a heatmap indicating the location of potential quality defects or lung lesions). However, this experimental task is representative of many real-world human-AI work settings and heatmaps are standard in explaining AI predictions of images. We also show that other heatmap algorithms lead to similar results (see Supplement F). Still, we invite future research replicating our findings in other

work settings using other methods for explainability. We also acknowledge reservations against using explainable AI in general. Explainable AI can be fooled by adversarial attacks [66] or may itself generate explanations that are unreliable and thus lead to misleading conclusions [8]. Nevertheless, it is likely that performance improvements from explainable AI can be achieved in other settings, where explanations serve as a decision aid. Further, it is important to note that, in both experiments, 14% of the images contained quality defects/lung lesions. Thus, quality defects/lung lesions were more prevalent than what domain experts would typically encounter in their respective job. This discrepancy might have influenced their prior expectations while performing the task. However, as participants in both treatment arms were shown exactly the same images, this factor likely had minimal impact on our findings.

Policy initiatives in many countries aim to promote transparency in AI algorithms (e.g., the United States [67] and the European Union [68]). These efforts are usually motivated from the perspective of ethics, regulation, and safety [9, 69, 70, 71]. Our research suggests that the benefits of algorithmic transparency are more profound: augmenting domain experts with explainable AI can enable better decisions with benefits for individuals, organizations, and society.

# Methods

This work analyzes the effect of augmenting domain experts with explainable AI (as opposed to black-box AI) in human-AI collaboration. We preregistered our hypotheses (i.e., Study 1: https://osf.io/7djxb and Study 2: https://osf.io/69yqt; see also Supplement K), which were tested in two randomized experiments. We comply with all ethical regulations and the research design was approved by the Ethics Commission of ETH Zurich (EK 2021-N-34). All participants provided informed consent.

## Tasks

In the following, details of both visual inspection tasks are provided.

### Study 1: Manufacturing experiment

For the manufacturing task, we designed a representative, real-world visual inspection task in collaboration with *Siemens Smart Infrastructure* in Zug, Switzerland. The experimental task is representative of various domains in which workers have to make decisions under a limited time budget. During the experiment, workers were shown images of electronic products and were asked to label them as faultless or defective. Reassuringly, we emphasize that our experiment involved a real work scenario: we conducted it with real workers familiar with quality management practice, a real user interface for state-of-the-art quality management, realistic incentives, and real product images. All steps in the experiments were carried out via a computer interface that was designed analogously to the real-world quality inspection setup at *Siemens* (see Supplement D for details). In the experiment, we made sure that all workers have the exact same conditions (exact same product images, same computer setup, same time limits, etc.). Thereby, we can rule out confounding variables that would arise naturally during the usual work routines and thus ensure that the experiment is scientifically sound.

We obtained 200 images of four different types of electronic products (printed circuit boards) from *Siemens*. Example images are provided in Supplement B. All four different product types are of equal importance to *Siemens*. Each product type comprised 43 images with faultless products and 7 images with defective products (e.g., missing components, wrong components, and faulty components). The different defects are all considered equally bad by the partner

16

company, i.e., the products are considered to be either functional or non-functional. Hence, we considered defective products as scrap as best practice in quality management [51, 52].

We implemented an AI algorithm that computed an individual quality score for each image. The quality score gives a numerical value between 0 (most certainly defect) and 100 (most certainly faultless). Workers were instructed that a quality score below 90 suggests an increased likelihood of a quality defect and that the AI algorithm can make mistakes. As humans cannot understand how the AI algorithm arrives at the prediction, the quality score is regarded as opaque ("black-box AI"). When evaluating the quality score with a cutoff of 90 for mapping the numerical value onto a binary faultless/defect label, the standalone AI algorithm achieves a balanced accuracy of 95.6% and a defect detection rate of 92.9%. The prediction performance of the standalone AI algorithm was not communicated to the workers. The AI algorithm was trained on an additional set of product images that was not included in the experiment.

We used anomaly heatmaps [47] to explain the opaque quality from the AI algorithm. The heatmaps were computed with standard computer vision methods and highlighted image regions with suspected quality defects (i.e., deviations from a faultless product). We chose heatmaps as the explanation technique for our AI algorithm as they provide a clear and intuitive way of highlighting areas with quality defects. This is especially important since the recruited domain experts are typically not familiar with explanation techniques for AI algorithms. Further, heatmaps are frequently used for images and are considered state-of-the-art with respect to their localization performance across various settings [47, 48, 49, 50]. Details on the implementation of the AI algorithm and heatmaps are provided in Supplement C.

In the experiment, workers were randomly assigned to one of the two treatments: (a) black-box AI or (b) explainable AI. Workers in the black-box AI treatment arm were only supported by the opaque quality score. Workers in the explainable AI treatment arm had access to the same quality score but additionally received the heatmap that explained the otherwise opaque quality score. Of note, the explainable AI had the same accuracy as the black-box AI and did not carry more information from an AI perspective (i.e., the heatmaps were of the same predictive power).

The procedure of the experiment was as follows. Before starting the experiment, workers had to give written consent to participate and then pass a tutorial on how to use the interface. After that workers were randomly assigned to one of the two treatments, i.e., either black-box

AI or explainable AI. During the experiment, the 200 product images were consecutively shown in random order. For each image, the workers had to assess the quality; that is, to "approve" or "reject" the shown product. We tracked the decision speed and the quality assessment (i.e., labeled as faultless or defective) made by the worker. To match real-world conditions, the workers were given a maximum of 35 minutes to finish the inspection task of 200 product images (around 10 seconds per image). In total, $N = 9,600$ assessments of product images were performed by the workers. Finally, workers completed a post-experimental questionnaire (Supplement L).

*Study 2: Medical experiment*

For the medical task, radiologists had to identify lung lesions in real chest X-ray images. Lung lesions are common findings in chest X-ray images [72] and can be easily overlooked due to their frequently small size [73]. The radiologists were asked whether at least one lung lesion was visible in the X-ray image. The experiment was conducted via Qualtrics. To ensure a realistic experimental setup that resembles the same task in daily, medical practice, we implemented a zoom function, which allowed the radiologists to investigate an enlarged view of the image by moving their computer mouse over the image. Analogous to Study 1, we emphasize that our medical experiment involved a realistic work scenario: we conducted it with actual medical professionals, who were asked to investigate real chest X-ray images.

We used 50 chest X-ray images from the CheXpert dataset [74]. The dataset comprised 7 images with at least one lung lesion and 43 images without lung lesions. Example images are provided in Supplement B.

We implemented an AI algorithm that outputs the probability of whether a lung lesion is visible in the chest X-ray image. We transformed these probability outputs for lung lesions such that the AI score gives a numerical value between 0 (most certainly contains a lung lesion) and 100 (most certainly does not contain a lung lesion) to mirror the quality score from the manufacturing setting. Hence, the AI output can be interpreted as a risk score, which are widely used in medical practice. Radiologists were instructed that an AI score below 90 indicates that the AI algorithm suspects at least one lung lesion is visible and that the AI algorithm can make mistakes. When evaluating the AI score with a cutoff of 90 for mapping the numerical value onto a binary label (lung lesion visible yes/no), the standalone AI algorithm achieves a balanced accuracy of 82.2% and a disease detection rate of 71.4%. Analogously to the manufacturing

task, the prediction performance of the standalone AI algorithm was not communicated to the participants.

As in the manufacturing task, the black-box AI algorithm was converted into an explainable AI by explaining the AI score via a heatmap, which is a state-of-the-art explanation technique for chest X-ray images in medicine [50]. Further details about the implementation of the AI algorithm and the heatmap are provided in Supplement C.

The procedure was analogous to the manufacturing experiment. Before starting the experiment, radiologists had to confirm their area of specialization and give written consent to participate. Subsequently, the task was explained and the radiologists had to pass a tutorial on how to use the interface. After that radiologists were randomly assigned to one of the two treatments, i.e., either black-box AI or explainable AI. During the experiment, 50 chest X-ray images were randomly shown either in forward or reverse order. For each chest X-ray image, the radiologists had to answer whether at least one lung lesion is visible. The corresponding answers as well as the decision speed were tracked. Radiologists were given a maximum of 35 minutes to finish the inspection task of 50 chest X-ray images. Radiologists were given more time per image compared to factory workers in the manufacturing experiment to reflect the differences in manufacturing and clinical practice. In total, $N = 5,650$ assessments of chest X-ray images were performed by the radiologists. Finally, all radiologists completed a post-experimental questionnaire (Supplement L).

### Study populations

The inclusion were as follows. Participants had to be at least 18 years old. For the manufacturing task, participants additionally had to have no self-reported visual impairment. The exclusion criteria were preregistered and were as follows. In both studies, we excluded participants that failed the tutorial or did not finish the inspection task on time. Participants with obvious misbehavior were also excluded from our analyses. In the manufacturing task, this was the case for workers that approved all products (i.e., labeled no images as defective). In the medical task, this was the case for radiologists that assigned the same label for all 50 chest X-ray images. In both studies, participants whose performance with respect to balanced accuracy was more than three standard deviations worse than the mean of their respective treatment arm were excluded.

We performed randomization checks to confirm that all treatment arms were demographically unbiased (Supplement E).

*Study 1: Manufacturing experiment*

The manufacturing experiment was carried out from June 29 to July 8, 2021 on-site at a *Siemens* factory in Zug, Switzerland. The objective of the field experiment was to get a real-world estimate of the treatment effect based on a representative sample of actual factory workers. Therefore, we only considered factory workers who were experienced in quality control practices. The factory workers were well familiar with the shown products and the visual inspection task. Overall, 56 factory workers (consisting of manufacturing employees, quality engineers, and team leaders) participated in our study. Out of them, all workers passed the tutorial; 6 did not finish on time; and 2 were excluded due to obvious misbehavior. The final sample consisted of 48 factory workers with an average working experience of 13.8 years. A larger sample size in our manufacturing experiment was not possible because the entire available workforce in one shift did not exceed 56 workers. Still, the experiment is well-powered as the treatment effect in the field experiment is considerably large. No additional financial incentive was given beyond the base salary to be representative of many real-world tasks from domain experts (e.g., as in manufacturing at *Siemens*).

*Study 2: Medical experiment*

The medical experiment was carried out via an online interface from February 27 to March 31, 2024. Actual radiologists based in the United States were recruited via MSI-ACI (https://site.msi-aci.com/). MSI-ACI paid a financial compensation to radiologists regardless of performance and adheres to the federal minimum wage in the United States. Overall, 122 radiologists started the study. Out of them, all passed the tutorial; 4 did not complete the study; 2 did not finish on time; and 3 were excluded due to obvious misbehavior. Hence, the final sample consisted of 113 radiologists with an average tenure as radiologist of 13.5 years.

## Statistical analysis

In manufacturing, it is common that all defects (e.g., missing components, wrong components, and faulty components) are considered equally bad [51, 52], so that the product to inspect

could be either functional or non-functional. Analogously, in the medical setting, either lung lesions were present or not. Hence, in both settings, the outcomes were binary. Therefore, task performance in the visual inspection tasks between the participants' assessments and the ground-truth labels was computed via (1) balanced accuracy (i.e., average sensitivity across faultless and defective products) and (2) defect/disease detection rate. For (1), the balanced accuracy is calculated via $0.5 \times [TP/P + TN/N]$ with true positives $TP$, positives $P$, true negatives $TN$, and negatives $N$. Here, we used balanced accuracy since it accounts for imbalanced distributions of labels by equally weighing the performance on each label, thus following best practice [75]. In contrast, the standard accuracy score would not account for the imbalanced distribution of positive and negative labels encountered in both settings (i.e., 172 faultless products and 28 defective products in the manufacturing setting; 7 chest X-ray images with and 43 without lung lesions in the medical setting). For (2), the defect/disease detection rate is defined as $TN/N$, where defective products and chest X-ray images with lung lesions were defined as negatives. In our manufacturing setting, missing a defective product has more severe implications than labeling a faultless product as defective. In medicine, missing a lung lesion on a chest X-ray image has more severe implications than additionally performing a CT scan for a healthy patient. Hence, it is crucial to find the negative samples.

All statistical tests in the results are based on one-sided Welch's $t$-tests. We further used ordinary least square (OLS) regression models to estimate the treatment effect of explainable AI on task performance. The OLS models are estimated via

$$Y_i = \beta_0 + \beta_1 \ Treatment_i + \varepsilon_i, \tag{1}$$

where $Y_i$ is the observed task performance (i.e., balanced accuracy or defect/disease detection rate), $Treatment_i$ is a binary variable which equals 0 if participant $i$ received the black-box AI treatment and 1 if participant $i$ received the explainable AI treatment. A significance level of $\alpha = 0.05$ was preregistered.

## Robustness checks

We conducted the following robustness checks. First, we repeated our analyses using precision as an additional task performance metric (Supplement G). Second, we repeated the OLS re-

gression models with additional participant-specific controls to estimate the treatment effect of explainable AI (Supplement H). Third, we estimated the treatment effect with quasi-binomial regression (Supplement H). Fourth, we estimated the regression models including participants that were previously excluded due to obvious misbehavior or because they did not finish the inspection task in time (Supplement I). All robustness checks yielded conclusive findings.

To demonstrate that the heatmaps in our medical setting are robust with respect to the choice of algorithm, we used two additional, different algorithms to generate heatmaps. We find that different algorithms lead to similar heatmaps (Supplement F).

## Comparison to non-experts

Additionally, we repeated the manufacturing task with non-experts recruited from Amazon MTurk as a robustness check (Supplement J). Typically, non-experts can not leverage explanations in the same way as domain experts due to missing domain knowledge. Hence, we were interested whether the large treatment effect of explainable AI on task performance we observed with domain experts transfers also to non-experts.

We found that non-experts supported by explainable AI also achieved a higher task performance than non-experts supported by black-box AI. Task performance of non-experts augmented with explainable AI was improved by 6.3 percentage points with respect to balanced accuracy. However, the treatment effect of explainable AI is slightly smaller as compared to the experiment with domain experts (where the balanced accuracy increased by 7.7 percentage points). Furthermore, non-experts with explainable AI achieved a higher defect detection rate with an improvement of 11.3 percentage points. This is of a similar effect size as in the real-world experiment (11.0 percentage points). The treatment effects in both metrics were again statistically significant (balanced accuracy: $\beta = 6.252$, $SE = 1.733$, $P < 0.001$; defect detection rate: $\beta = 11.271$, $SE = 3.276$, $P = 0.001$). For more details, see Supplement J.

## References

[1] Brynjolfsson, E. & Mitchell, T. What can machine learning do? Workforce implications. *Science* **358**, 1530–1534 (2017).

[2] Perrault, R. & Clark, J. Artificial intelligence index report 2024. Human-Centered Artificial Intelligence. United States of America. Retrieved from https://policycommons.net/artifacts/12089781/hai_ai-index-report-2024/12983534/ on 26 Apr 2024. CID: 20.500.12592/h70s46h (2024).

[3] Bertolini, M., Mezzogori, D., Neroni, M. & Zammori, F. Machine learning for industrial applications: A comprehensive literature review. *Expert Systems with Applications* **175**, 114820 (2021).

[4] Scheetz, J., Rothschild, P., McGuinness, M., Hadoux, X., Soyer, H. P., Janda, M., Condon, J. J. J., Oakden-Rayner, L., Palmer, L. J., Keel, S. & van Wijngaarden, P. A survey of clinicians on the use of artificial intelligence in ophthalmology, dermatology, radiology and radiation oncology. *Scientific Reports* **11**, 5193 (2021).

[5] Cazzaniga, M., Jaumotte, F., Li, L., Melina, G., Panton, A. J., Pizzinelli, C., Rockall, E. J. & Tavares, M. M. Gen-AI: Artificial intelligence and the future of work. International Monetary Fund. Staff Discussion Notes 2024/001 https://www.imf.org/en/Publications/Staff-Discussion-Notes/Issues/2024/01/14/Gen-AI-Artificial-Intelligence-and-the-Future-of-Work-542379 (2024).

[6] Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations* (2015).

[7] He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 770–778 (2016).

[8] Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence* **1**, 206–215 (2019).

[9] Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F. & Pedreschi, D. A survey of methods for explaining black box models. *ACM Computing Surveys* **51**, 1–42 (2018).

[10] Murdoch, W. J., Singh, C., Kumbier, K., Abbasi-Asl, R. & Yu, B. Definitions, methods, and applications in interpretable machine learning. *Proceedings of the National Academy of Sciences* **116**, 22071–22080 (2019).

[11] Gunning, D., Stefik, M., Choi, J., Miller, T., Stumpf, S. & Yang, G.-Z. XAI—Explainable artificial intelligence. *Science Robotics* **4**, eaay7120 (2019).

[12] Cheng, H.-F., Wang, R., Zhang, Z., O'Connell, F., Gray, T., Harper, F. M. & Zhu, H. Explaining decision-making algorithms through UI. In *CHI Conference on Human Factors in Computing Systems*, 1–12 (2019).

[13] Liao, Q. V., Gruen, D. & Miller, S. Questioning the AI: Informing design practices for explainable AI user experiences. In *CHI Conference on Human Factors in Computing Systems* (2020).

[14] Bodria, F., Giannotti, F., Guidotti, R., Naretto, F., Pedreschi, D. & Rinzivillo, S. Benchmarking and survey of explanation methods for black box models. *Data Mining and Knowledge Discovery* **37**, 1719–1778 (2023).

[15] Dwivedi, R., Dave, D., Naik, H., Singhal, S., Omer, R., Patel, P., Qian, B., Wen, Z., Shah, T., Morgan, G. & Ranjan, R. Explainable AI (XAI): Core ideas, techniques, and solutions. *ACM Computing Surveys* **55**, 1–33 (2023).

[16] Samek, W., Montavon, G., Vedaldi, A., Hansen, L. K. & Müller, K.-R. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning* (Springer Nature, 2019).

[17] Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J. & Muller, K.-R. Explaining deep neural networks and beyond: A review of methods and applications. *Proceedings of the IEEE* **109**, 247–278 (2021).

[18] Molnar, C. *Interpretable machine learning: A guide for making Black Box Models interpretable* (Lulu. com, 2019).

[19] Ribeiro, M. T., Singh, S. & Guestrin, C. Why should I trust you? Explaining the predictions of any classifier. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016).

[20] Lundberg, S. & Lee, S.-I. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems* (2017).

[21] Simonyan, K., Vedaldi, A. & Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *Workshop at International Conference on Learning Representations* (2014).

[22] Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D. & Batra, D. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *IEEE International Conference on Computer Vision*, 618–626 (2017).

[23] Fügener, A., Grahl, J., Gupta, A. & Ketter, W. Cognitive challenges in human–artificial intelligence collaboration: Investigating the path toward productive delegation. *Information Systems Research* **33**, 678–696 (2022).

[24] Andreas Fügener, Alok Gupta, Jörn Grahl, Wolfgang Ketter & Anna Taudien. Exploring user heterogeneity in human delegation behavior towards AI. In *International Conference on Information Systems* (2021).

[25] Bauer, K., von Zahn, M. & Hinz, O. Please take over: XAI, delegation of authority, and domain knowledge. Preprint at *SSRN* https://doi.org/10.2139/ssrn.4512594 (2023).

[26] Dietvorst, B. J., Simmons, J. P. & Massey, C. Algorithm aversion: People erroneously avoid algorithms after seeing them err. *Journal of Experimental Psychology: General* **144**, 114–126 (2015).

[27] Dietvorst, B. J., Simmons, J. P. & Massey, C. Overcoming algorithm aversion: People will use imperfect algorithms if they can (even slightly) modify them. *Management Science* **64**, 1155–1170 (2018).

[28] Dietvorst, B. J. & Bharti, S. People reject algorithms in uncertain decision domains because they have diminishing sensitivity to forecasting error. *Psychological Science* **31**, 1302–1314 (2020).

[29] Burton, J. W., Stein, M.-K. & Jensen, T. B. A systematic review of algorithm aversion in augmented decision making. *Journal of Behavioral Decision Making* **33**, 220–239 (2020).

[30] Ben David, D., Resheff, Y. S. & Tron, T. Explainable AI and adoption of financial algorithmic advisors. In *AAAI/ACM Conference on AI, Ethics, and Society*, 390–400 (2021).

[31] Choung, H., David, P. & Ross, A. Trust in AI and its role in the acceptance of AI technologies. *International Journal of Human–Computer Interaction* **39**, 1727–1739 (2023).

[32] Nourani, M., Kabir, S., Mohseni, S. & Ragan, E. D. The effects of meaningful and meaningless explanations on trust and perceived system accuracy in intelligent systems. *AAAI Conference on Human Computation and Crowdsourcing* **7** (2019).

[33] Panigutti, C., Beretta, A., Giannotti, F. & Pedreschi, D. Understanding the impact of explanations on advice-taking: a user study for AI-based clinical decision support systems. In *CHI Conference on Human Factors in Computing Systems* (2022).

[34] Bansal, G., Wu, T., Zhou, J., Fok, R., Nushi, B., Kamar, E., Ribeiro, M. T. & Weld, D. Does the whole exceed its parts? The effect of AI explanations on complementary team performance. In *CHI Conference on Human Factors in Computing Systems* (2021).

[35] Vasconcelos, H., Jörke, M., Grunde-McLaughlin, M., Gerstenberg, T., Bernstein, M. S. & Krishna, R. Explanations can reduce overreliance on AI systems during decision-making. *ACM on Human-Computer Interaction* **7**, 129 (2023).

[36] Chen, V., Liao, Q. V., Wortman Vaughan, J. & Bansal, G. Understanding the role of human intuition on reliance in human-AI decision-making with explanations. *ACM on Human-Computer Interaction* **7**, 370 (2023).

[37] Buçinca, Z., Lin, P., Gajos, K. Z. & Glassman, E. L. Proxy tasks and subjective measures can be misleading in evaluating explainable AI systems. In *International Conference on Intelligent User Interfaces* (2020).

[38] Chu, E., Roy, D. & Andreas, J. Are visual explanations useful? A case study in model-in-the-loop prediction. Preprint at *arXiv* https://doi.org/10.48550/arXiv.2007.12248 (2020).

[39] Alufaisan, Y., Marusich, L. R., Bakdash, J. Z., Zhou, Y. & Kantarcioglu, M. Does explainable artificial intelligence improve human decision-making? *AAAI Conference on Artificial Intelligence* **35** (2021).

[40] Schemmer, M., Kuehl, N., Benz, C., Bartos, A. & Satzger, G. Appropriate reliance on AI advice: Conceptualization and the effect of explanations. In *International Conference on Intelligent User Interfaces* (2023).

[41] Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N. & Lee, S.-I. From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence* **2**, 56–67 (2020).

[42] Das, N. *et al.* Collaboration between explainable artificial intelligence and pulmonologists improves the accuracy of pulmonary function test interpretation. *European Respiratory Journal* **61**, 2201720 (2023).

[43] Gaube, S., Suresh, H., Raue, M., Lermer, E., Koch, T. K., Hudecek, M. F. C., Ackery, A. D., Grover, S. C., Coughlin, J. F., Frey, D., Kitamura, F. C., Ghassemi, M. & Colak, E. Non-task expert physicians benefit from correct explainable AI advice when reviewing X-rays. *Scientific Reports* **13**, 1383 (2023).

[44] Jesus, S., Belém, C., Balayan, V., Bento, J., Saleiro, P., Bizarro, P. & Gama, J. How can I choose an explainer? In *ACM Conference on Fairness, Accountability, and Transparency* (2021).

[45] Sivaraman, V., Bukowski, L. A., Levin, J., Kahn, J. M. & Perer, A. Ignore, trust, or negotiate: Understanding clinician acceptance of AI-based treatment recommendations in health care. In *CHI Conference on Human Factors in Computing Systems* (2023).

[46] Jabbour, S., Fouhey, D., Shepard, S., Valley, T. S., Kazerooni, E. A., Banovic, N., Wiens, J. & Sjoding, M. W. Measuring the impact of AI in the diagnosis of hospitalized patients: A randomized clinical vignette survey study. *JAMA* **330**, 2275–2284 (2023).

[47] Bergmann, P., Löwe, S., Fauser, M., Sattlegger, D. & Steger, C. Improving unsupervised defect segmentation by applying structural similarity to autoencoders. Preprint at *arXiv* https://doi.org/10.48550/arXiv.1807.02011 (2019).

[48] Bergmann, P., Fauser, M., Sattlegger, D. & Steger, C. MVTec AD — A comprehensive real-world dataset for unsupervised anomaly detection. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019).

[49] Binder, A., Bockmayr, M., Hägele, M., Wienert, S., Heim, D., Hellweg, K., Ishii, M., Stenzinger, A., Hocke, A., Denkert, C., Müller, K.-R. & Klauschen, F. Morphological and molecular breast cancer profiling through explainable machine learning. *Nature Machine Intelligence* **3**, 355–366 (2021).

[50] Saporta, A., Gui, X., Agrawal, A., Pareek, A., Truong, S. Q. H., Nguyen, C. D. T., Ngo, V.-D., Seekins, J., Blankenberg, F. G., Ng, A. Y., Lungren, M. P. & Rajpurkar, P. Benchmarking saliency methods for chest x-ray interpretation. *Nature Machine Intelligence* **4**, 867–878 (2022).

[51] Juran, J., Gryna, F. & Bingham, R. *Quality control handbook* (New York: McGraw-Hill, 1979).

[52] Hoyle, D. *Quality Management Essentials* (Butterworth-Heinemann, 2007).

[53] Baudin, M. & Netland, T. H. *Introduction to manufacturing: An industrial engineering and management perspective* (Taylor & Francis, 2022).

[54] Tschandl, P. *et al.* Human-computer collaboration for skin cancer recognition. *Nature Medicine* **26**, 1229–1234 (2020).

[55] Pantanowitz, L., Valenstein, P. N., Evans, A. J., Kaplan, K. J., Pfeifer, J. D., Wilbur, D. C., Collins, L. C. & Colgan, T. J. Review of the current state of whole slide imaging in pathology. *Journal of Pathology Informatics* **2**, 36 (2011).

[56] Ibrahim, R. & Shafiq, M. O. Explainable convolutional neural networks: A taxonomy, review, and future directions. *ACM Computing Surveys* **55**, 1–37 (2023).

[57] Amershi, S., Weld, D., Vorvoreanu, M., Fourney, A., Nushi, B., Collisson, P., Suh, J., Iqbal, S., Bennett, P. N., Inkpen, K., Teevan, J., Kikin-Gil, R. & Horvitz, E. Guidelines for human-AI interaction. In *CHI Conference on Human Factors in Computing Systems* (2019).

[58] De-Arteaga, M., Fogliato, R. & Chouldechova, A. A case for humans-in-the-loop: Decisions in the presence of erroneous algorithmic scores. In *CHI Conference on Human Factors in Computing Systems* (2020).

28

[59] Fügener, A., Grahl, J., Gupta, A. & Ketter, W. Will humans-in-the-loop become borgs? Merits and pitfalls of working with AI. *MIS Quarterly* **45**, 1527–1556 (2021).

[60] Cadario, R., Longoni, C. & Morewedge, C. K. Understanding, explaining, and utilizing medical artificial intelligence. *Nature Human Behaviour* **5**, 1636–1642 (2021).

[61] Sun, J., Zhang, D. J., Hu, H. & van Mieghem, J. A. Predicting human discretion to adjust algorithmic prescription: A large-scale field experiment in warehouse operations. *Management Science* **68**, 846–865 (2022).

[62] Kawaguchi, K. When will workers follow an algorithm? A field experiment with a retail business. *Management Science* **67**, 1670–1695 (2021).

[63] Castelo, N., Bos, M. W. & Lehmann, D. R. Task-dependent algorithm aversion. *Journal of Marketing Research* **56**, 809–825 (2019).

[64] Yeomans, M., Shah, A., Mullainathan, S. & Kleinberg, J. Making sense of recommendations. *Journal of Behavioral Decision Making* **32**, 403–414 (2019).

[65] Senoner, J., Netland, T. & Feuerriegel, S. Using explainable artificial intelligence to improve process quality: Evidence from semiconductor manufacturing. *Management Science* **68**, 5704–5723 (2022).

[66] Slack, D., Hilgard, S., Jia, E., Singh, S. & Lakkaraju, H. Fooling LIME and SHAP: Adversarial attacks on post hoc explanation methods. In *AAAI/ACM Conference on AI, Ethics, and Society* (2020).

[67] Algorithmic Accountability Act. H.R.2231 - Algorithmic Accountability Act of 2019 (116th Congress) (2019). URL https://www.congress.gov/bill/116th-congress/house-bill/2231.

[68] European Commission. Ethics Guidelines for Trustworthy AI (2019). URL https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai.

[69] Wachter, S., Mittelstadt, B. & Floridi, L. Transparent, explainable, and accountable AI for robotics. *Science Robotics* **2**, eaan6080 (2017).

[70] Muller, H., Mayrhofer, M. T., van Veen, E.-B. & Holzinger, A. The ten commandments of ethical medical AI. *Computer* **54**, 119–123 (2021).

[71] Jobin, A., Ienca, M. & Vayena, E. The global landscape of AI ethics guidelines. *Nature Machine Intelligence* **1**, 389–399 (2019).

[72] UTSouthwestern Medical Center. Pulmonary nodules and lung lesions (Accessed 04/09/24). URL [https://utswmed.org/conditions-treatments/pulmonary-nodules-and-lung-lesions/](https://utswmed.org/conditions-treatments/pulmonary-nodules-and-lung-lesions/).

[73] Oestmann, J. W., Greene, R., Kushner, D. C., Bourgouin, P. M., Linetsky, L. & Llewellyn, H. J. Lung lesions: correlation between viewing time and detection. *Radiology* **166**, 451–453 (1988).

[74] Irvin, J. *et al.* CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison. *AAAI Conference on Artificial Intelligence* **33** (2019).

[75] Hollon, T. *et al.* Artificial-intelligence-based molecular classification of diffuse gliomas using rapid, label-free optical imaging. *Nature Medicine* **29**, 828–832 (2023).

[76] Linardatos, P., Papastefanopoulos, V. & Kotsiantis, S. Explainable AI: A review of machine learning interpretability methods. *Entropy* **23** (2020).

[77] Nelder, J. A. & Wedderburn, R. W. M. Generalized linear models. *Journal of the Royal Statistical Society. Series A (General)* **135**, 370–384 (1972).

[78] Hastie, T. & Tibshirani, R. *Generalized additive models* (Chapman & Hall/CRC, 1990).

[79] Lou, Y., Caruana, R. & Gehrke, J. Intelligible models for classification and regression. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2012).

[80] Kraus, M., Tschernutter, D., Weinzierl, S. & Zschech, P. Interpretable generalized additive neural networks. *European Journal of Operational Research* **317**, 303–316 (2024).

[81] Shapley, L. S. A value for n-person games. In *Contributions to the Theory of Games*, Annals of Mathematics Studies, 307–318 (Princeton University Press, 1953).

[82] Hansen, K., Baehrens, D., Schroeter, T., Rupp, M. & Müller, K.-R. Visual interpretation of kernel-based prediction models. *Molecular Informatics* **30**, 817–826 (2011).

[83] Buhrmester, V., Münch, D. & Arens, M. Analysis of explainers of black box deep neural networks for computer vision: A survey. *Machine Learning and Knowledge Extraction* **3**, 966–989 (2021).

[84] Zeiler, M. D. & Fergus, R. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, 8689 (2014).

[85] Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R. & Samek, W. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS ONE* **10**, e0130140 (2015).

[86] Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, 70 (2017).

[87] Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propagating activation differences. In *International Conference on Machine Learning*, 70 (2017).

[88] Zhou, B., Khosla, A., Lapedriza, A., Oliva, A. & Torralba, A. Learning deep features for discriminative localization. In *IEEE Conference on Computer Vision and Pattern Recognition* (2016).

[89] Chattopadhay, A., Sarkar, A., Howlader, P. & Balasubramanian, V. N. Grad-CAM++: Generalized gradient-based visual explanations for deep convolutional networks. In *IEEE Winter Conference on Applications of Computer Vision*, 839–847 (2018).

[90] Bany Muhammad, M. & Yeasin, M. Eigen-CAM: Visual explanations for deep convolutional neural networks. *SN Computer Science* **2** (2021).

[91] Verma, S., Dickerson, J. P. & Hines, K. E. Counterfactual explanations for machine learning: A review. Preprint at *arXiv* https://doi.org/10.48550/arXiv.2010.10596 (2020).

[92] Zipfel, J., Verworner, F., Fischer, M., Wieland, U., Kraus, M. & Zschech, P. Anomaly detection for industrial quality assurance: A comparative evaluation of unsupervised deep learning models. *Computers & Industrial Engineering* **177**, 109045 (2023).

[93] Liu, J., Xie, G., Wang, J., Li, S., Wang, C., Zheng, F. & Jin, Y. Deep industrial image anomaly detection: A survey. *Machine Intelligence Research* **21**, 104–135 (2024).

[94] Doshi-Velez, F. & Kim, B. Towards a rigorous science of interpretable machine learning. Preprint at *arXiv* https://doi.org/10.48550/arXiv.1702.08608 (2017).

[95] Plumb, G., Molitor, D. & Talwalkar, A. S. Model agnostic supervised local explanations. In *Advances in Neural Information Processing Systems* (2018).

[96] Hooker, S., Erhan, D., Kindermans, P.-J. & Kim, B. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems* (2019).

[97] Alvarez-Melis, D. & Jaakkola, T. S. On the robustness of interpretability methods. Preprint at *arXiv* https://doi.org/10.48550/arXiv.1806.08049 (2018).

[98] Narayanan, M., Chen, E., He, J., Kim, B., Gershman, S. & Doshi-Velez, F. How do humans understand explanations from machine learning systems? An evaluation of the human-interpretability of explanation. Preprint at *arXiv* https://doi.org/10.48550/arXiv.1802.00682 (2018).

[99] Amarasinghe, K., Rodolfa, K. T., Jesus, S., Chen, V., Balayan, V., Saleiro, P., Bizarro, P., Talwalkar, A. & Ghani, R. On the importance of application-grounded experimental design for evaluating explainable ML methods. *AAAI Conference on Artificial Intelligence* **38** (2024).

[100] Liu, M., Shi, J., Li, Z., Li, C., Zhu, J. & Liu, S. Towards better analysis of deep convolutional neural networks. *IEEE Transactions on Visualization and Computer Graphics* **23**, 91–100 (2017).

[101] Ming, Y., Cao, S., Zhang, R., Li, Z., Chen, Y., Song, Y. & Qu, H. Understanding hidden memories of recurrent neural networks. In *IEEE Conference on Visual Analytics Science and Technology* (2017).

[102] Pezzotti, N., Hollt, T., van Gemert, J., Lelieveldt, B. P. F., Eisemann, E. & Vilanova, A. DeepEyes: Progressive visual analytics for designing deep neural networks. *IEEE Transactions on Visualization and Computer Graphics* **24**, 98–108 (2018).

[103] Strobelt, H., Gehrmann, S., Pfister, H. & Rush, A. M. LSTMVis: A tool for visual analysis of hidden state dynamics in recurrent neural networks. *IEEE Transactions on Visualization and Computer Graphics* **24**, 667–676 (2018).

[104] Candrian, C. & Scherer, A. Rise of the machines: Delegating decisions to autonomous AI. *Computers in Human Behavior* **134**, 107308 (2022).

[105] Jessica Ochmann, Leonard Michels, Sandra Zilker, Verena Tiefenbeck & Sven Laumer. The influence of algorithm aversion and anthropomorphic agent design on the acceptance of AI-based job recommendations. In *International Conference on Information Systems* (2020).

[106] Hou, Y. T.-Y. & Jung, M. F. Who is the expert? Reconciling algorithm aversion and algorithm appreciation in AI-supported decision making. *ACM on Human-Computer Interaction* **5**, 477 (2021).

[107] Bogert, E., Schecter, A. & Watson, R. T. Humans rely more on algorithms than social influence as a task becomes more difficult. *Scientific Reports* **11**, 8028 (2021).

[108] Fleiß, J., Bäck, E. & Thalmann, S. Mitigating algorithm aversion in recruiting: A study on explainable AI for conversational agents. *ACM SIGMIS Database: the DATABASE for Advances in Information Systems* **55**, 56–87 (2024).

[109] Cai, C. J., Reif, E., Hegde, N., Hipp, J., Kim, B., Smilkov, D., Wattenberg, M., Viegas, F., Corrado, G. S., Stumpe, M. C. & Terry, M. Human-centered tools for coping with imperfect algorithms during medical decision-making. In *CHI Conference on Human Factors in Computing Systems* (2019).

[110] Branley-Bell, D., Whitworth, R. & Coventry, L. User trust and understanding of explainable AI: Exploring algorithm visualisations and user biases. In *Human-Computer Interaction. Human Values and Quality of Life*, 12183, 382–399 (2020).

[111] Zhang, Y., Liao, Q. V. & Bellamy, R. K. E. Effect of confidence and explanation on accuracy and trust calibration in AI-assisted decision making. In *Conference on Fairness, Accountability, and Transparency*, 295–305 (2020).

[112] Lancaster Farrell, C.-J. Explainability does not improve biochemistry staff trust in artificial intelligence-based decision support. *Annals of Clinical Biochemistry* **59**, 447–449 (2022).

[113] Panigutti, C., Beretta, A., Fadda, D., Giannotti, F., Pedreschi, D., Perotti, A. & Rinzivillo, S. Co-design of human-centered, explainable AI for clinical decision support. *ACM Transactions on Interactive Intelligent Systems* **13** (2023).

[114] Leichtmann, B., Humer, C., Hinterreiter, A., Streit, M. & Mara, M. Effects of explainable artificial intelligence on trust and human behavior in a high-risk decision task. *Computers in Human Behavior* **139**, 107539 (2023).

[115] Buçinca, Z., Malaya, M. B. & Gajos, K. Z. To trust or to think: Cognitive forcing functions can reduce overreliance on AI in AI-assisted decision-making. *ACM on Human-Computer Interaction* **5**, 188 (2021).

[116] Lee, J. D. & See, K. A. Trust in automation: designing for appropriate reliance. *Human factors* **46**, 50–80 (2004).

[117] Green, B. & Chen, Y. The principles and limits of algorithm-in-the-loop decision making. *ACM on Human-Computer Interaction* **3**, 50 (2019).

[118] Lage, I., Chen, E., He, J., Narayanan, M., Kim, B., Gershman, S. J. & Doshi-Velez, F. Human evaluation of models built for interpretability. *AAAI Conference on Human Computation and Crowdsourcing* **7**, 59–67 (2019).

[119] Lai, V. & Tan, C. On human predictions with explanations and predictions of machine learning models. In *Conference on Fairness, Accountability, and Transparency* (2019).

[120] Cai, C. J., Jongejan, J. & Holbrook, J. The effects of example-based explanations in a machine learning interface. In *International Conference on Intelligent User Interfaces* (2019).

[121] Carton, S., Mei, Q. & Resnick, P. Feature-based explanations don't help people detect misclassifications of online toxicity. *International AAAI Conference on Web and Social Media* **14** (2020).

[122] Lai, V., Liu, H. & Tan, C. "Why is 'Chicago' deceptive?" Towards building model-driven tutorials for humans. In *CHI Conference on Human Factors in Computing Systems* (2020).

[123] Yang, F., Huang, Z., Scholtz, J. & Arendt, D. L. How do visual explanations foster end users' appropriate trust in machine learning? In *International Conference on Intelligent User Interfaces* (2020).

[124] Alqaraawi, A., Schuessler, M., Weiß, P., Costanza, E. & Berthouze, N. Evaluating saliency map explanations for convolutional neural networks. In *International Conference on Intelligent User Interfaces* (2020).

[125] Wang, X. & Yin, M. Are explanations helpful? A comparative study of the effects of explanations in AI-assisted decision-making. In *International Conference on Intelligent User Interfaces* (2021).

[126] Poursabzi-Sangdeh, F., Goldstein, D. G., Hofman, J. M., Wortman Vaughan, J. W. & Wallach, H. Manipulating and measuring model interpretability. In *CHI Conference on Human Factors in Computing Systems* (2021).

[127] van der Waa, J., Nieuwburg, E., Cremers, A. & Neerincx, M. Evaluating XAI: A comparison of rule-based and example-based explanations. *Artificial Intelligence* **291**, 103404 (2021).

[128] Kim, S. S. Y., Meister, N., Ramaswamy, V. V., Fong, R. & Russakovsky, O. Hive: Evaluating the human interpretability of visual explanations. In *European Conference on Computer Vision*, 13672 (2022).

[129] Leichtmann, B., Hinterreiter, A., Humer, C., Streit, M. & Mara, M. Explainable artificial intelligence improves human decision-making: Results from a mushroom picking experiment at a public art festival. *International Journal of Human–Computer Interaction* 1–18 (2023).

[130] Müller, R., Reindel, D. F. & Stadtfeld, Y. D. The benefits and costs of explainable artificial intelligence in visual quality control: Evidence from fault detection performance and eye movements. *Human Factors and Ergonomics in Manufacturing & Service Industries* (2024).

[131] Metta, C., Beretta, A., Guidotti, R., Yin, Y., Gallinari, P., Rinzivillo, S. & Giannotti, F. Improving trust and confidence in medical skin lesion diagnosis through explainable deep learning. *International Journal of Data Science and Analytics* (2023).

[132] Nagendran, M., Festor, P., Komorowski, M., Gordon, A. C. & Faisal, A. A. Quantifying the impact of AI recommendations with explanations on prescription decision making. *npj Digital Medicine* **6**, 206 (2023).

[133] Pimentel, M. A., Clifton, D. A., Clifton, L. & Tarassenko, L. A review of novelty detection. *Signal Processing* **99**, 215–249 (2014).

[134] Wang, Z., Bovik, A. C., Sheikh, H. R. & Simoncelli, E. P. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing* **13**, 600–612 (2004).

[135] Huang, G., Liu, Z., van der Maaten, L. & Weinberger, K. Q. Densely connected convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 4700–4708 (2017).

[136] Peer, E., Vosgerau, J. & Acquisti, A. Reputation as a sufficient condition for data quality on Amazon Mechanical Turk. *Behavior Research Methods* **46**, 1023–1031 (2013).

[137] NASA. Nasa Task Load Index (TLX) (1986). URL https://humansystems.arc.nasa.gov/groups/TLX/downloads/TLX.pdf.

[138] Davis, F. D. Perceived usefulness, perceived ease of use, and user acceptance of information technology. *MIS Quarterly* **13**, 319–340 (1989).

[139] Jian, J.-Y., Bisantz, A. M. & Drury, C. G. Foundations for an empirically determined scale of trust in automated systems. *International Journal of Cognitive Ergonomics* **4**, 53–71 (2000).

# Acknowledgements

# Author contributions

All authors contributed to the research design, data analysis, interpretation of results, and writing of the paper. JS and BK performed the manufacturing experiment. SS performed the medical experiment.

# Competing interests

The funding bodies had no control over design, conduct, data, analysis, review, reporting, or interpretation of the research conducted.

# Data and code availability

All analyses were conducted using Python (3.11) with *numpy* (1.24.3), *pandas* (1.5.3), *scipy* (1.11.1), and *statsmodels* (0.14.0). The data visualizations were created with *seaborn* (0.12.2) and *matplotlib* (3.7.1). The data and code to reproduce the results from all studies will be made publicly available at https://osf.io/ upon publication.

# Supplements

## Contents

# Supplement A  Extended literature review

In the following, we provide an extended literature review of explainable artificial intelligence (AI). In particular, we differentiate research on explainable AI in computer science (which is primarily focused on methodological outcomes) from our work (which is focused on behavioral science outcomes). An overview is provided in Table S1.

## A.1  Explainable AI in computer science

In the field of computer science, the primary objective concerning explainable AI is to develop and evaluate new methods to achieve better transparency of AI algorithms. For a general overview of explainable AI, see for example [9, 16, 76, 17]. AI algorithms can be broadly divided into two categories: algorithms that are considered to be inherently interpretable and algorithms that are not due to their complexity [8]. The latter are often referred to as black-box algorithms.

An inherently interpretable model is linear regression, where the decision-making can be directly followed by inspecting the coefficients. Extensions of linear regression are generalized linear models (GLMs) and generalized additive models (GAMs) [77, 78]. GLMs were introduced as a unification of various methods that allow for different distributions of the dependent variable (e.g., a binary dependent variable as in logistic regression). GAMs were introduced to also allow for non-linear relationships between an independent and the dependent variable, which can be modeled for example with decision trees or shallow neural networks (see e.g., [79, 80]). While GLMs and GAMs are still considered to be inherently interpretable, they are not as straightforward to interpret as linear regression.

Decision trees are also considered to be inherently interpretable by simply following the decision rules from the root to the leaf nodes. Other inherently interpretable models are naïve Bayes classifier, k-nearest neighbor algorithm, rule-based learning, etc. (see [18] for an introduction).

In contrast, post-hoc explanation techniques can be applied to better understand black-box algorithms such as neural networks. These explanation techniques are applied after the AI algorithm has been trained. Post-hoc explanation methods can be divided into global and local methods [76], where the former aim at explaining the algorithm's overall decision-making process while the latter provide explanations for a single, specific input. An example of global methods are feature importance rankings, which rank the features based on their importance in predict-

ing a model's outcome, usually measured across the entire model rather than for individual predictions. A further example are partial dependence plots, which show the effect of a single feature on the predicted outcome of a model, averaged over a dataset. Prominent examples of local methods are local interpretable model-agnostic explanations (*LIME*) and SHapley Additive exPlanations (*SHAP*) [19, 20]. *LIME* approximates a black-box model locally around the prediction with an interpretable model (like a linear model) to explain individual predictions. *SHAP* leverages a concept from cooperative game theory (Shapley values [81]) to explain the output of a model by computing the contribution of each feature to the prediction while also considering possible interaction effects. Post-hoc explanation methods can be further categorized into model-specific and model-agnostic methods. Model-specific methods are designed for a specific class of AI algorithms or even for a single AI algorithm. Model-specific methods exist, for example, for convolutional neural networks [22] or for kernel-based AI algorithms such as support vector machines [82]. In contrast, model-agnostic methods can be applied to any AI algorithm; notable examples are *LIME* and *SHAP*. Another dimension to differentiate post-hoc explanation methods is for which data type (tabular, text, audio, images, etc.) the method was developed. In this study, we focus on images, and, in the following, we thus present some of the most relevant post-hoc explanation methods for AI algorithms in computer vision. For general overviews of explainable AI in computer vision, we refer to [83, 56].

One of the earliest attempts to explain convolutional neural networks (CNN), which are nowadays widely used in computer vision, was made by Zeiler and Fergus [84]. Therein, the authors present the deconvolutional network (*DeconvNet*) as a visualization method that maps feature activations back to the input image. Additionally, a simple technique called occlusion was discussed as a method for explaining the predictions of a CNN. For that, different portions of the input image are systematically occluded with a grey square, and the impact on the output of the network is observed. Significant changes in the output probabilities indicate the regions of the image most important for classification.

In [21], a method was introduced for generating saliency maps by computing the gradient of the output category with respect to the input image. This technique highlights the regions of the image that contribute most to the model's classification decision, offering a straightforward visual explanation of where the network is "looking" to make its predictions.

Layer-wise relevance propagation *(LRP)* backtracks the output decision of the network

through the layers to assign relevance scores to individual pixels. This method helps in understanding which parts of the input image were most relevant for the model's decision, emphasizing a layer-by-layer decomposition of the prediction [85].

Integrated gradients attribute the prediction of a neural network to its input features, calculating the gradients of the output prediction with respect to the input image. It integrates these gradients along the path from a baseline (zero input) to the actual input, offering a way to visualize the importance of each pixel [86].

Deep learning important features ($DeepLIFT$) compares the activation of each neuron to its 'reference activation' and assigns contribution scores according to the difference. This method can identify which features of the input contribute to differences in the output from some baseline, offering a more detailed view than simple gradient-based methods [87].

By using the global average pooling layers in CNNs, class activation mapping ($CAM$) generates heatmaps that highlight the discriminative parts of the image used by the network to identify specific classes, facilitating visual explanations of model decisions [88]. Several extensions to this approach exist [22, 89, 90], with $GradCAM$ being one of the most often used approach [22]. In contrast to $CAM$, $GradCAM$ employs a gradient-based approach to generate heatmaps, and, as a result, no changes to the network architecture are required.

Counterfactual explanations provide insights by showing how a small change in the input image could change the classification result. This method helps in understanding model decisions by answering "what-if" scenarios, offering a direct way to comprehend how the model might react to different inputs [91].

A special case of post-hoc explanation in computer vision are anomaly heatmaps. These were developed for unsupervised computer vision AI algorithms, i.e., algorithms that do not require a labeled dataset but rather aim at finding anomalies automatically [47, 92, 93].

A plethora of post-hoc explanation methods exists and often it is not obvious which method to choose. Thus, different evaluation measures for post-hoc explanation methods have been proposed [94]. These include fidelity, which measures how accurately the explanations reflect the decisions of the underlying model [19, 95, 96, 50], and robustness, assessing stability under small changes in the input [97]. Another measure is human-interpretability, which examines how understandable the explanations are to humans [98]. But also application-grounded measures have been proposed, where the evaluation metric is how humans perform in a certain task [44, 99].

In another stream of literature in computer science, tools that help AI engineers with designing and training AI algorithms are developed. Especially, deep neural networks are difficult to train and require a certain amount of experience. Therefore, visualization tools have emerged that facilitate this task (e.g., see [100, 101, 102, 103]).

## A.2 Explainable AI in behavioral science

In behavioral science, different outcomes of human-AI collaboration have been studied. Human delegation of tasks and decisions to AI algorithms has been intensively studied recently [59, 23, 24, 104]. Specifically, it has been examined whether the use of explainable AI can increase the likelihood of humans delegating decisions to AI algorithms [25].

Algorithm aversion refers to the phenomenon where humans are reluctant to use algorithms [26, 27, 28, 29, 105, 106, 107, 60, 61, 62, 63, 64]. To overcome human aversion towards AI algorithms, it has been hypothesized that providing an explanation of an AI's decision may be beneficial. This has been tested with mixed findings [30, 108]. Missing trust in an AI's decision can lead to algorithm aversion [31]. Therefore, previous studies investigated whether explainable AI can increase the trust in AI algorithms [109, 32, 110, 111, 112, 33, 113, 114, 45].

A contrary phenomenon to algorithm aversion is overreliance [115, 35], where humans place too much trust in AI algorithms, potentially overlooking or ignoring their limitations. Previous research has found mixed results on whether explainability of AI leads to decreased overreliance [34, 35, 36].

However, for a good task performance, it is crucial that humans only adhere to correct AI predictions and overrule wrong ones, which is also referred to as appropriate reliance [116]. The effect of explainable AI on task performance has been studied previously. In the majority of studies, however, non-experts (e.g., via Amazon Mechanical Turk or university students) were recruited and those oftentimes performed simplified, non-realistic, or even non-relevant tasks (see, e.g., [117, 118, 119, 120, 37, 121, 122, 123, 38, 111, 124, 125, 126, 127, 39, 34, 115, 128, 36, 129, 130]). It has been hypothesized that non-experts can not fully harness explanations due to a lack of domain knowledge [36]. Thus, an empirical evaluation of the effect of explainable AI on task performance in real job tasks, requires actual domain experts of those tasks [94].

Previous studies that recruited domain experts to perform real-job tasks have other drawbacks. For example, prior work has compared the effect of explainable AI against humans

alone [41, 42]. Others have used expert annotations as a proxy for explainable AI [43] or research designs that prevent isolating the treatment effect of explainable AI on task performance [44, 46, 45]. Finally, also other outcomes have been studied such as trust, confidence, and perceived usefulness [131, 132], while, as our novelty, we add by focusing on task performance with domain experts.

**Table S1: Overview of key literature on explainable AI.**

| Domain | Concept | Research summary | Dependent variable | References (examples) |
|---|---|---|---|---|
| Computer science | New explanation methods | Derivations of new methods where the focus is on mathematical / algorithmic contributions | n/a | [77, 78, 79, 80, 19, 20, 82, 84, 21, 85, 86, 87, 88, 22, 89, 90, 91, 47, 92] |
| | Benchmarking methods/ datasets | Proposing new methods or datasets to benchmark the performance of explainable AI | n/a | [95, 96, 50, 97, 98, 44, 99] |
| | New visualization tools | Visualization tools that facilitate the development of complex machine learning models | n/a | [100, 101, 102, 103] |
| Behavioral science | Delegation between humans and AI | Humans avoid delegation to algorithms | Delegation frequency | [59, 23, 24, 104, 25] |
| | Algorithm aversion | Humans reject advice from algorithm | Adherence | [26, 27, 28, 29, 105, 106, 107, 60, 61, 62, 63, 64, 30, 108] |
| | Trust in AI | Humans do not trust AI algorithms | Trust | [109, 32, 110, 111, 112, 33, 113, 114, 45] |
| | Overreliance on AI | Humans follow advice from algorithms blindly | Overreliance | [34, 35, 36] |
| | Task performance in response to explainable AI | Comparison of black-box AI vs explainable AI for task performance using unrealistic tasks, non-experts, or non-causal research designs | Task performance | [117, 118, 119, 120, 37, 121, 122, 123, 38, 111, 124, 125, 126, 127, 39, 34, 115, 128, 36, 129, 130, 41, 42, 43, 44, 46, 45] |
| | | Real-world job tasks with domain experts for estimating treatment effects of explainable AI vs black-box AI | Task performance | **ours** |

# Supplement B  Research setting

## B.1  Manufacturing setting

Poor quality generates 10% to 15% of the operating expenses in manufacturing.[1]  Identifying defective products before they move downstream in the value chain is essential to maintain a high operational performance. For this purpose, manufacturers conduct visual quality inspections to assess whether products have defects (e.g., assembly errors or surface damages) [47]. In manufacturing operations, many quality inspections are still conducted manually, which is often a tedious, tiring, and error-prone task.  AI offers promising opportunities to overcome these drawbacks by supporting factory workers in automatically detecting quality defects before products are sold to customers. Specifically, AI can assist workers in detecting the location and type of error so rework can be conducted more effectively and efficiently.  Therefore, AI algorithms enable factory workers to be more productive by focusing on their key value-creating work tasks.

Our research was carried out at *Siemens* Smart Infrastructure in Zug, Switzerland. To test our hypotheses, the company provided us with real-world product images (each with $1920 \times 1080$ pixels) from their factory. The images comprise four different types of electronic products, all of which are printed circuit boards. Figure S1 shows example images of the four types of electronic products that were inspected during the experiment. Figure S2 shows three examples of quality defects, which include products with wrong components, products with assembly errors, and products with faulty components.

Overall, we received two datasets. The first dataset comprised 200 images, including 43 correct products and 7 defective products for each of the four product types. All experiments (and thus the empirical results in the main analysis) are based on the first dataset. The second dataset comprised 200 additional faultless images (50 for each of the four product types). These images were used to train the AI algorithm that was used to compute the quality scores in the experiment (see Supplement C).

---

[1]American Society for Quality.  *Cost of Quality (COQ)*. URL: https://asq.org/quality-resources/cost-of-quality, last accessed on June 13, 2024.

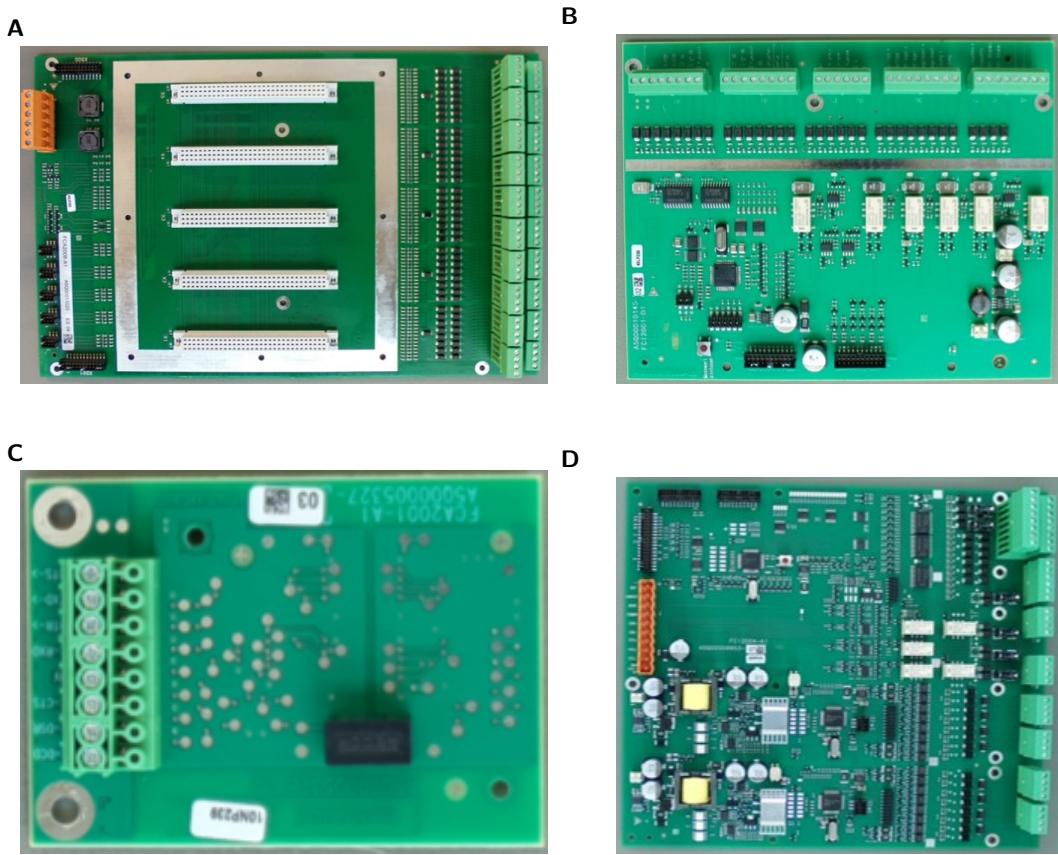**Figure S1: Four types of electronic products (printed circuit boards).** (**A-D**) Exemplary images of faultless products that were inspected during the experiment.
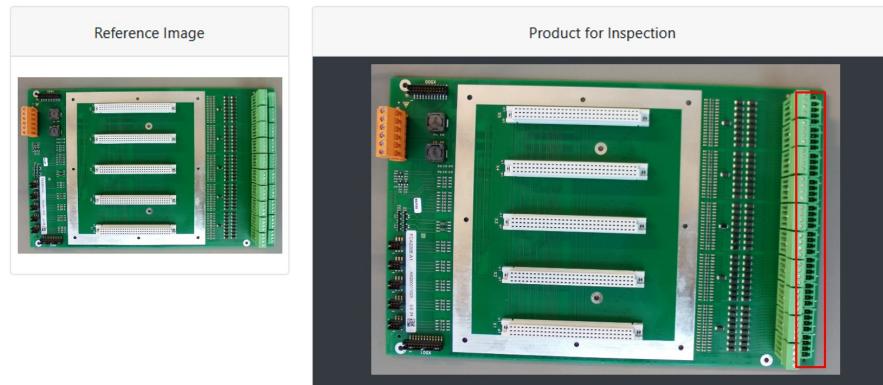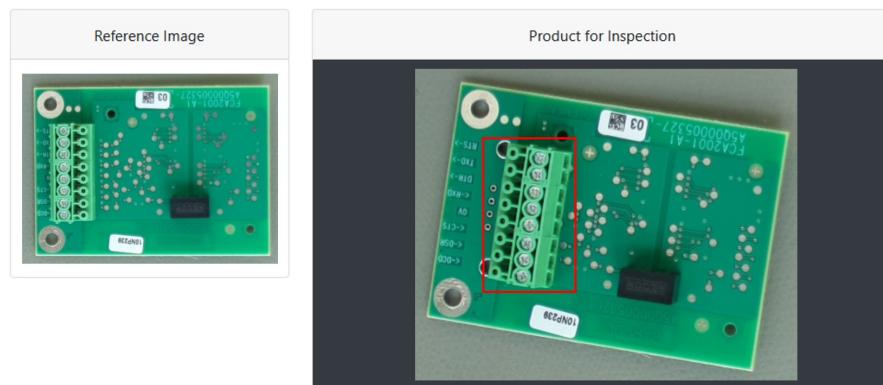
**Figure S2: Examples of quality defects.** (**A**) Example of a defective product with wrong components. (**B**) Example of a defective product with a component assembled in the wrong orientation. (**C**) Example of a defective product with a faulty component.

## B.2    Medical setting

Chest radiography (capturing X-ray images) is a widely performed diagnostic imaging test across the world and plays a crucial role in the screening, diagnosis, and management of numerous diseases that pose a threat to life [74]. One of such diseases are lung lesions, which include lung nodules and masses in our experiment. Lung nodules are common and are encountered on roughly one out of 500 chest X-ray images [72].

Overlooking a lung lesion on a chest X-ray can have serious, potentially life-threatening consequences for patients. The failure to detect a lesion at an early stage can lead to a delay in diagnosis and treatment, allowing diseases to progress to more advanced stages. This can significantly worsen the prognosis for conditions such as lung cancer, tuberculosis, and pneumonia, where early intervention can often lead to better outcomes. Beyond the immediate health risks, there are also implications for patient care, including increased medical costs due to more complex and prolonged treatment that may become necessary as a disease progresses. However, subtle lung lesions can be easily overlooked even without any constraints on how long radiologists are allowed to inspect the chest X-ray image [73]. Given these reasons, identifying lung lesions on chest X-ray images is an important, non-trivial task in daily, medical care. To that end, giving physicians a decision aid for this task is crucial.

Example chest X-ray images including the corresponding heatmaps are provided in Figure S3.

**Figure S3: Example chest X-ray images with corresponding heatmaps**. (**A**) Chest X-ray image without lung lesions. (**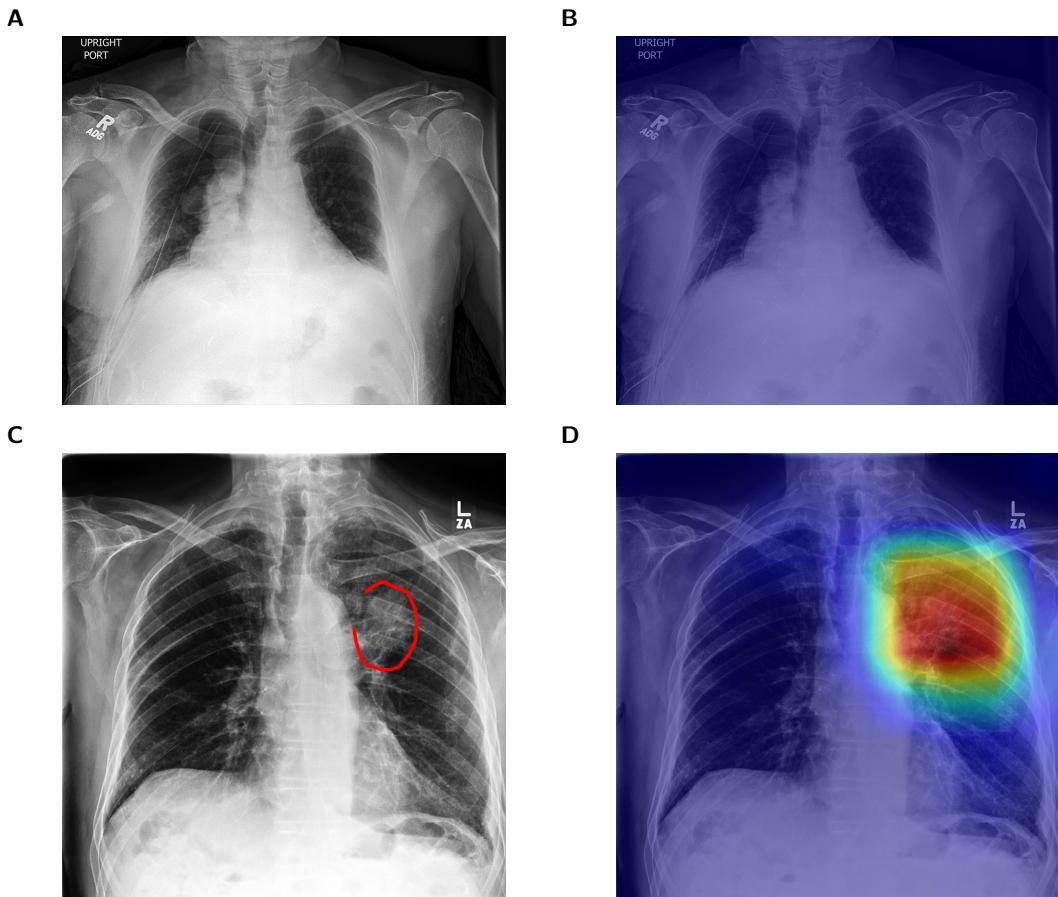B**) Heatmap overlaid over the chest X-ray image from **A**. (**C**) Chest X-ray image with a lung lesion annotated in red by an experienced radiologist. (**D**) Heatmap overlaid over the chest X-ray image from **C**.

# Supplement C   Implementation of AI algorithm

## C.1   Manufacturing setting

As part of this research, we implemented an AI algorithm that provided the predictions (i.e., quality scores) that were shown to the participants during the manufacturing experiment. Our AI algorithm builds upon unsupervised anomaly detection [133] and, as such, follows common standards in industry for visually analyzing the quality of product images [47, 48]. Algorithms based on unsupervised anomaly detection are particularly suitable for industrial settings because they only require a set of faultless product images to be trained. Therefore, there is no need to specify defect types beforehand, which also allows identifying quality defects that have never been observed before. This is reflected in anomaly detection where product images with sufficiently large deviations from "normal" products are labeled as defective.

**AI algorithm.** In our case, the AI algorithm performs unsupervised anomaly detection as follows [133]. First, we are given an existing training set $T$ with images of faultless products. Upon deployment, an out-of-sample product image $x$ is subject to assessment; that is, whether it is similar to any of the images $t \in T$ and thus likely faultless or whether it is highly dissimilar and thus likely defective. Here, anomaly detection compares the similarity (with regard to some similarity function $d$) between the new image $x$ and the existing images $t \in T$. For each, a similarity $d(x, t)$, for all $t \in T$ is computed. If the similarity $d$ falls below a certain threshold $\theta^*$, an image is labeled as defective.

**AI predictions.** The similarity of product images is computed by following best practice in computer vision. As such, we refrain from simply computing the L2-norm (or some other norm) between $x$ and $t$. The reason is that such distance would give equal weight to all pixels and cannot properly account for the semantic similarity in images. Rather, we follow established practice and compute the similarity via the so-called structural similarity index [134]. The structural similarity index is a standard computer vision method for quantifying the similarity of images between 0 (i.e., no similarity at all) and 1 (i.e., perfect similarity). Product images with a low structural similarity indicate an increased probability of a quality defect because they are less similar to the training data (i.e., images of faultless products). For details on the computation of the structural similarity index, we refer to [134]. Eventually, we scaled the structural similarity

index of all images between 0 and 100 and rounded the values to the nearest integer to enhance readability. The resulting similarity measure corresponds to the quality scores that were shown to the participants in the experiment.

**Prediction performance.** We evaluated the out-of-sample prediction accuracy of the AI algorithm as follows. In a first step, we mapped the quality scores onto a binary faultless/defective label. For this, we introduced a quality score of $\theta^* = 90$ as a cutoff (i.e., predicting that a product is defective if the quality score is below 90 and faultless otherwise). We then compared the predictions of the algorithm against the ground-truth quality labels provided by *Siemens*. Table S2 gives the confusion matrix for the 200 out-of-sample images used in the experiment. We measure the prediction performance via the balanced accuracy (i.e., average sensitivity across faultless and defective products). We choose the balanced accuracy as our main performance metric because it accounts for the unbalanced distribution between faultless and defective products (i.e., 172 products are faultless and 28 products are defective). The standalone AI algorithm achieves a balanced accuracy of 95.6% (i.e., $0.5 \times [169/172 + 26/28]$). Additionally, we evaluated the defect detection rate (true negative rate) of the AI algorithm, i.e., how many of the defective products were identified as such. The defect detection rate of the standalone AI algorithm was 92.9% (26/28).

**Table S2: Confusion matrix comparing AI predictions with ground-truth labels in the manufacturing setting**

|  |  | Predicted label | |
|---|---|---|---|
|  |  | Faultless | Defective |
| *Actual label* | Faultless | 169 | 3 |
|  | Defective | 2 | 26 |

**Explainable AI.** We extended the above AI algorithm to produce explanations for each prediction as follows. We followed other research in computer vision that generates so-called "anomaly heatmaps" [47]. Anomaly heatmaps visualize in what area of an image a quality defect is predicted to be. Formally, in an anomaly heatmap, each pixel $x_i$ is associated with a score measuring the likelihood of a defect at that location. Pixels that receive a bright color (yellow, orange, red, etc.) correspond to "anomalous" regions because they have a large distance to the training data (i.e., the pixel is dissimilar to the one in a faultless product). In contrast,

pixels that are colored in blue have a small distance to the training data and should thus be considered as "normal." For better usability, we overlay the anomaly heatmap over the actual product image (with partially transparent colors). Examples of two anomaly heatmaps are shown in Figure S4. In the experiment, the heatmaps were shown to the participants in the explainable AI treatment arm as an additional decision aid.



**Figure S4: Anomaly heatmaps for AI predictions**. (**A**) Example anomaly heatmap for a faultless product. (**B**) Example anomaly heatmap for a defective product.

## C.2   Medical setting

**AI algorithm.** We used an already trained DenseNet121 from [50]. DenseNet121 is a convolutional neural network that is part of the DenseNet family, known for its dense connectivity pattern where each layer is connected to every other layer in a feed-forward fashion [135]. The "121" in DenseNet121 stands for the total number of layers in the network, including convolutional layers, pooling layers, and fully connected layers, summing up to 121. The DenseNet121 was set up as a multi-label classifier, which takes a chest X-ray image as input and outputs probabilities for the following 10 labels: airspace opacity, atelectasis, cardiomegaly, consolidation, edema, enlarged

cardiomediastinum, lung lesion, pleural effusion, pneumothorax, and support devices. It was trained on 224,316 chest X-ray images from 65,240 patients.

**AI predictions.** The probabilities returned by the DenseNet121 were mapped onto binary yes/no labels by finding the probability threshold that maximized the balanced accuracy on a validation set of chest X-ray images, which were not used during training. In order to have a score identical to the quality score from the manufacturing setting, where a smaller score indicates a greater likelihood of showing a defect and with a cutoff of 90 that divides the quality scores into defective and faultless, the following transformations were performed: (i) the probabilities were inverted, (ii) the threshold that divides the two classes was set to 90, (iii) the inverted probabilities larger than that threshold were rescaled using min-max scaling on a scale from 90 to 100, and (iv) the inverted probabilities smaller than that threshold were rescaled using min-max scaling on a scale from 0 to 90.

**Prediction performance.** As in the manufacturing setting, we evaluate performance of the AI algorithm on the 50 chest X-ray images by calculating the balanced accuracy and the disease detection rate. Those 50 images were neither used for training nor for finding the class dividing threshold. The standalone AI algorithm achieved a balanced accuracy of 82.2% (i.e., $0.5 \times [40/43 + 5/7]$) and a disease detection rate of 71.4% (i.e., 5/7). Additionally, the confusion matrix for the 50 images is shown in Table S3.

**Table S3: Confusion matrix comparing AI predictions with ground-truth labels in the medical setting**

|  | Predicted label | |
|---|---|---|
|  | No lung lesion | At least one lung lesion |
| No lung lesion | 40 | 3 |
| At least one lung lesion | 2 | 5 |

*Actual label* (row group label)

**Explainable AI.** As explanation technique for the above AI algorithm, we used *GradCAM* [22]. *GradCAM* outputs heatmaps similar to the anomaly heatmaps from the manufacturing setting and showed state-of-the-art localization performance on chest X-ray images across a variety of diagnoses [50]. Analogous to the manufacturing setting, pixels with bright colors (yellow, orange, red, etc.) correspond to regions that were most relevant for predicting lung lesions, whereas blue pixels were least relevant. To increase usability, heatmaps were overlaid

over the raw chest X-ray images with partially transparent colors. Examples of two heatmaps next to the original chest X-ray images are shown in Figure S3. Heatmaps were only provided to radiologists in the explainable AI treatment arm.

# Supplement D   Experimental interface

## D.1   Manufacturing setting

The experiment was carried out via a computer interface that was analogously designed to the real-world quality inspection setup at *Siemens*. The experiment comprises the following steps: (1) the study description and study consent, (2) a tutorial on how to use the application, (3) the visual inspection task involving 200 images, (4) a post-experimental questionnaire.

Depending on the randomly assigned treatment, different versions of the quality inspection interface were shown to participants (Figure S5). Similar to the real-world setting at *Siemens*, all participants had access to a reference image, which showed a faultless product. The participants were asked to evaluate each of the 200 images individually and to make an "approve" (faultless product) or "reject" (defective product) decision by clicking the respective buttons. This represents the quality assessments that we use for all analyses. The participants were allowed to change their quality assessment before submitting their decision and proceeding to the next image. Once a decision was submitted, participants could no longer return to the previous image. Overall, the participants were given 35 minutes to solve the task, which corresponds to realistic field conditions. The remaining time was always shown on the top of the interface.

We tracked several metrics during the experiment. In the tutorial, we tracked whether participants were following the steps correctly and screened out those that did not complete the tutorial successfully. During the visual inspection task, we tracked the final quality assessment (i.e., faultless or defective) and the decision speed of the users. In the post-experimental questionnaire, we saved the answers to individual questions. The aggregated user data were stored in a database and later converted into a comma-separated values (CSV) file.

**Figure S5: Different interfaces depending on treatment arm.** (**A**) The interface for the black-box AI. (**B**) The interface for the explainble AI. (**C**) The interface for the human without AI treatment. (**D**) The interface for the post-experimental questionnaire.

## D.2 Medical setting

The experiment was conducted via Qualtrics. The experiment was divided in the following steps: (1) physician confirmation and study consent, (2) a tutorial on how to perform the experiment, (3) the visual inspection task consisting of 50 chest X-ray images, and (4) a post-experimental questionnaire.

A different version of the chest X-ray inspection interface was shown to radiologists depending on the randomly assigned treatment (Figure S6). For both treatment arms, the chest X-ray image to inspect was shown on the left. To the right of it, an enlarged view was shown, which could be altered by moving the mouse over the chest X-ray image. Radiologists in the treatment arm with explainable AI additionally received a heatmap, which was displayed right of the enlarged view. The radiologists were asked to inspect 50 chest X-ray images and to answer the question "Is at least one lung lesion visible in the chest X-ray image below?" with either "YES" or "NO" for each image. The radiologists were allowed to change their assessment before submitting their decision (clicking the blue arrow button to proceed to the next page). Each chest X-ray image was shown on a separate page and radiologists were not allowed to go back to a previous image

once a decision was submitted. The radiologists had 35 minutes to complete the task and the remaining time was always shown on the top-left of the page.

Several metrics were recorded during the experiment: In the tutorial, we tracked whether radiologists understood how to perform the task. In the visual inspection task, the final assessment for each image as well as the corresponding decision speed were recorded. In the post-hoc questionnaire, we saved the answers to the individual questions. The data was stored on Qualtrics and exported as a CSV file.
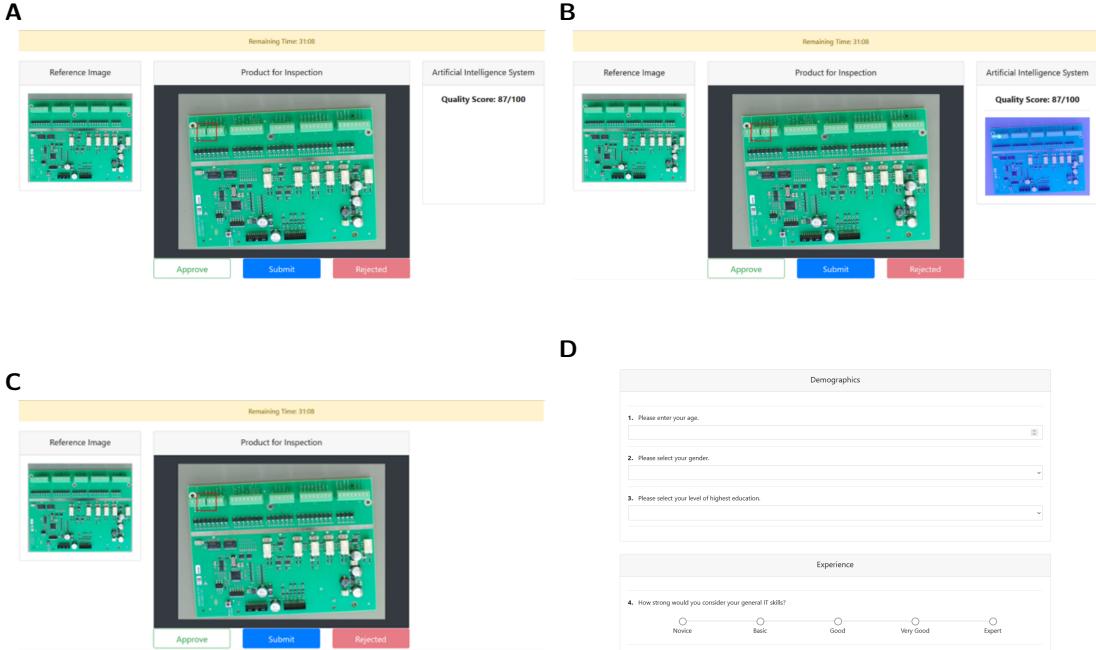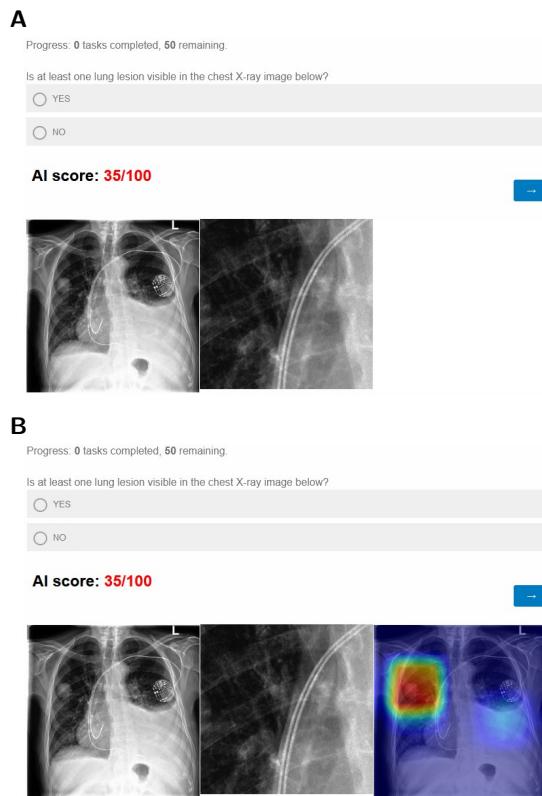


**Figure S6: Different interfaces depending on treatment arm.** (**A**) The interface for the black-box AI. (**B**) The interface for the explainble AI.

# Supplement E   Randomization checks

## E.1   Study 1: Manufacturing experiment

We performed randomization checks to confirm that the distribution of workers in the two treatment arms of the manufacturing experiment was unbiased. The following demographic variables were collected: age bracket [<20, 20–30, 30–40, 40–50, 50–60, 60–70, >70], gender [male, female, not listed], and highest level of education [ISCED1, ISCED2, ISCED3, ISCED4, ISCED5, ISCED6, ISCED7].[2] We further collected the participant-specific tenure at *Siemens* measured in years since start of employment. Table S4 reports the observed frequencies and the mean tenure (with standard deviation in parentheses) for both treatment arms. The randomization checks for age, gender, and education are based on $\mathcal{X}^2$-tests of independence. The randomization check for tenure is based on a two-sided Welch's $t$-test. The results suggest no statistically significant differences between the participants in the two treatment arms.

**Table S4: Randomization checks for manufacturing experiment**

|  | Human with black-box AI | Human with explainable AI | *P*-value |
|---|:---:|:---:|:---:|
| Age | 0 \| 1 \| 7 \| 8 \| 4 \| 2 \| 0 | 0 \| 1 \| 3 \| 9 \| 12 \| 1 \| 0 | 0.223 |
| Gender | 16 \| 6 \| 0 | 17 \| 9 \| 0 | 0.815 |
| Education | 0 \| 9 \| 4 \| 2 \| 3 \| 3 \| 1 | 0 \| 14 \| 2 \| 3 \| 3 \| 3 \| 1 | 0.897 |
| Tenure | 11.91 (8.83) | 15.38 (10.42) | 0.217 |
| Observations | 22 | 26 | – |

Notes: The table reports the frequency of participants that fall in the specific subgroups of age, gender, and education (separated by vertical bars) and the average tenure per treatment arm (standard deviation in parentheses). The *P*-values for the randomization checks are computed based on $\mathcal{X}^2$-tests of independence (age, gender, education) and a two-sided Welch's $t$-test (tenure).

## E.2   Study 2: Medical experiment

We performed a randomization check to confirm that the distribution of radiologists with respect to tenure in the two treatment arms of the medical experiment was unbiased. The randomization check is based on a two-sided Welch's $t$-test. The result suggests no statistically significant differences between the radiologists in the two treatment arms.

---

[2]UNESCO.   *International   Standard   Classification   of   Education   (ISCED)*.   URL: http://uis.unesco.org/en/topic/international-standard-classification-education-isced, last accessed on June 13, 2024.

**Table S5: Randomization checks for medical experiment**

| | Human with black-box AI | Human with explainable AI | *P*-value |
|---|:---:|:---:|:---:|
| Tenure | 11.89 (8.96) | 15.4 (11.71) | 0.08 |
| Observations | 61 | 52 | – |

Notes: The table reports the average tenure per treatment arm (standard deviation in parentheses). The *P*-value is computed based on a two-sided Welch's *t*-test.

## Supplement F   Robustness of the heatmap

We applied two additional algorithms to generate the heatmaps in the medical setting in order to show that different algorithms lead to similar heatmaps. In particular, we compared our heatmaps generated by *GradCAM* to heatmaps generated by *DeepLIFT* and *LRP* (for an introduction to these two methods see Supplement A) [87, 85]. Figure S7 shows the heatmaps generated by the three distinct algorithms for the eight chest X-ray images, where the AI algorithm predicted that lung lesions are visible.

Additionally, we calculated Pearson correlation coefficients between the heatmaps generated by different algorithms to quantify whether they highlight similar regions. The Pearson correlation coefficient is calculated via $r = \frac{\text{cov}(x_i, x_j)}{\sqrt{\text{Var}(x_i)\,\text{Var}(x_j)}}$, where $x_i$ and $x_j$ denote the flattened array of heatmaps generated by algorithm $i$ and $j$. The average Pearson correlation coefficient between *GradCAM* and *DeepLIFT* is $r = 0.92$, and the average Pearson correlation coefficient between *GradCAM* and *LRP* is $r = 0.63$. The exact Pearson correlation coefficient for each heatmap pair is reported in Figure S7. The Pearson correlation coefficient was statistically significant for each pair of heatmaps ($P < 0.001$). In general, we observe that all three algorithms produce similar heatmaps. The heatmaps we used in our medical setting (generated by *GradCAM*) are more similar to the heatmaps generated by *DeepLIFT* than to those generated by *LRP*.

**Figure S7: Heatmaps generated by three different algorithms.** The left column shows the heatmaps generated by *GradCAM* (the algorithm we used for our medical experiment). The middle columns shows the heatmaps generated by *DeepLIFT*. The right column shows the heatmaps generated by *LRP*. *r* denotes the Pearson correlation coefficient between the heatmaps generated by *GradCAM* (the algorithm we used) and *DeepLIFT/LRP*, respectively.

# Supplement G   Results with precision as task performance metric

In addition to the balanced accuracy and defect detection rate, we also report precision as a metric for task performance of the participants in combination with the defect/disease detection rate. Formally, precision is computed via $TN/PN$ with true negatives $TN$ and predicted negatives $PN$. We again compare the effect of augmenting humans with explainable AI versus black-box AI. Figure S8 reports the results for the manufacturing experiment (Study 1) and Figure S9 for the medical experiment (Study 2)



**Figure S8: Results of manufacturing experiment.** The boxplots compare the task performance between the two treatments: black-box AI and explainable AI. The task performance is measured by the defect detection rate (**A**) and the precision (**B**) based on the quality assessment of workers and the ground-truth labels of the product images. The standalone AI algorithm attains a defect detection rate of 92.9% and a precision of 89.7% (orange dashed lines). Statistical significance is based on a one-sided Welch's $t$-test ($^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$). In the boxplots, the center line denotes the median; box limits are upper and lower quartiles; whiskers are defined as the 1.5x interquartile range.

**Figure S9: Results of medical experiment.** The boxplots compare the task performance between the two treatments: black-box AI and explainable AI. The task performance is measured by the disease detection rate (**A**) and the precision (**B**) based on the quality assessment of radiologists and the ground-truth labels of the chest X-ray images. The standalone AI algorithm attains a disease detection rate of 71.4% and a precision of 62.5% (orange dashed lines). Statistical significance is based on a one-sided Welch's $t$-test ($^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$). In the boxplots, the center line denotes the median; box limits are upper and lower quartiles; whiskers are defined as the 1.5x interquartile range.

## Supplement H    Regression models

This section reports various regression models estimating the treatment effect of augmenting humans with explainable AI. The models are estimated via

$$Y_i = \beta_0 + \beta_1\, Treatment_i + \beta_2\, X_i + \varepsilon_i, \tag{S1}$$

where $Y_i$ is the observed task performance (i.e., balanced accuracy or defect/disease detection rate), $Treatment_i$ is a binary variable which equals 0 if participant $i$ received the black-box AI treatment and 1 if participant $i$ received the explainable AI treatment, and $X_i$ is the vector of participant-specific control variables. The above regression is estimated via ordinary least squares (OLS).

We acknowledge that the balanced accuracy is only defined between 0 and 100. Because OLS regression models could return values below 0 and above 100, we additionally estimate quasi-binomial regression models with a logit link function. For this, we s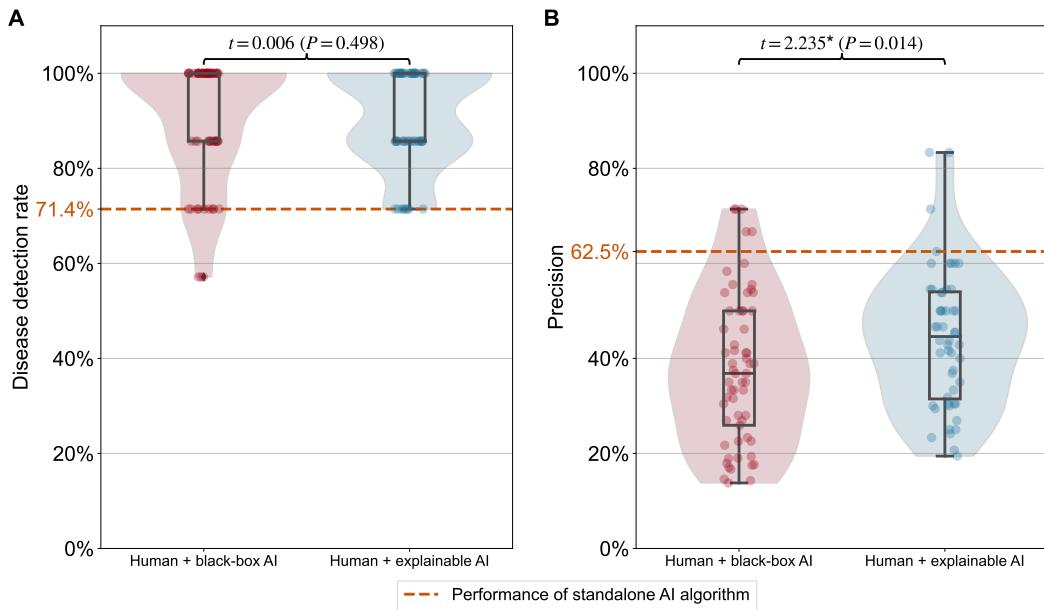et the scale parameter of the regression models to the Pearson $\mathcal{X}^2$-statistic divided by the residual degrees of freedom.

### H.1    Study 1: Manufacturing experiment

Table S6 reports three OLS regression models estimating the treatment effect with different control variables. Model (1) estimates the treatment effect for explainable AI with demographic controls (age, gender, and highest level of education) and the tenure at *Siemens* measured in years from start of employment. Model (2) estimates the treatment effect for explainable AI with demographic controls, tenure, and self-reported IT skills (ranging from 1: "novice" to 5: "expert"). Model (3) estimates the treatment effect for explainable AI with demographic controls, tenure, self-reported IT skills, and the decision speed (median across the 200 images). All three models return a significant treatment effect for both metrics (balanced accuracy and defect detection rate) as dependent variables.

**Table S6: OLS regression results for treatment effect (manufacturing experiment)**

| | Balanced accuracy | | | Defect detection rate | | |
|---|---|---|---|---|---|---|
| | Model (1) | Model (2) | Model (3) | Model (1) | Model (2) | Model (3) |
| Treatment | 8.131*** | 7.513*** | 7.508** | 11.783** | 10.914** | 10.888** |
| (explainable AI) | (2.087) | (2.098) | (2.117) | (3.732) | (3.790) | (3.717) |
| Demographics | Yes | Yes | Yes | Yes | Yes | Yes |
| Tenure | Yes | Yes | Yes | Yes | Yes | Yes |
| IT skills | No | Yes | Yes | No | Yes | Yes |
| Decision speed | No | No | Yes | No | No | Yes |
| Observations | 48 | 48 | 48 | 48 | 48 | 48 |

Notes: The table reports three OLS regression models with different sets of control variables and two different metrics as dependent variables. The standard errors of the treatment effect are reported in parentheses. Statistical significance: ***$P < 0.001$, **$P < 0.01$, *$P < 0.05$.

Table S7 reports three quasi-binomial regression models estimating the treatment effect with the same control variables as before. Again, all three models return a significant treatment effect for both metrics as dependent variables.

**Table S7: Quasi-binomial regression results for treatment effect (manufacturing experiment)**

| | Balanced accuracy | | | Defect detection rate | | |
|---|---|---|---|---|---|---|
| | Model (1) | Model (2) | Model (3) | Model (1) | Model (2) | Model (3) |
| Treatment | 1.295*** | 1.175** | 1.184*** | 1.157** | 1.060** | 1.089** |
| (explainable AI) | (0.340) | (0.358) | (0.357) | (0.369) | (0.386) | (0.375) |
| Demographics | Yes | Yes | Yes | Yes | Yes | Yes |
| Tenure | Yes | Yes | Yes | Yes | Yes | Yes |
| IT skills | No | Yes | Yes | No | Yes | Yes |
| Decision speed | No | No | Yes | No | No | Yes |
| Observations | 48 | 48 | 48 | 48 | 48 | 48 |

Notes: The table reports three quasi-binomial regression models with different sets of control variables and two different metrics as dependent variables. The standard errors of the treatment effect are reported in parentheses. Statistical significance: ***$P < 0.001$, **$P < 0.01$, *$P < 0.05$.

## H.2 Study 2: Medical experiment

Table S8 reports three OLS regression models estimating the treatment effect with different control variables. Model (1) estimates the treatment effect for explainable AI with tenure measured in years as a control variable. Model (2) estimates the treatment effect for explainable AI with tenure and self-reported IT skills (ranging from 1: "novice" to 5: "expert"). Model (3) estimates the treatment effect for explainable AI with tenure, self-reported IT skills, and the

decision speed (median across the 50 images). All three models return a significant treatment effect for balanced accuracy as a dependent variable.

**Table S8: OLS regression results for treatment effect (medical experiment)**

| | Balanced accuracy | | | Defect detection rate | | |
|---|---|---|---|---|---|---|
| | Model (1) | Model (2) | Model (3) | Model (1) | Model (2) | Model (3) |
| Treatment | 4.637* | 4.452* | 4.473* | 0.129 | 0.304 | 0.645 |
| (explainable AI) | (1.834) | (1.853) | (1.863) | (2.286) | (2.312) | (2.215) |
| Tenure | Yes | Yes | Yes | Yes | Yes | Yes |
| IT skills | No | Yes | Yes | No | Yes | Yes |
| Decision speed | No | No | Yes | No | No | Yes |
| Observations | 113 | 113 | 113 | 113 | 113 | 113 |

Notes: The table reports three OLS regression models with different sets of control variables and two different metrics as dependent variables. The standard errors of the treatment effect are reported in parentheses. Statistical significance: $^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$.

Table S9 reports three quasi-binomial regression models estimating the treatment effect with the same control variables as before. Again, all three models return a significant treatment effect for balanced accuracy as dependent variable.

**Table S9: Quasi-binomial regression results for treatment effect (medical experiment)**

| | Balanced accuracy | | | Defect detection rate | | |
|---|---|---|---|---|---|---|
| | Model (1) | Model (2) | Model (3) | Model (1) | Model (2) | Model (3) |
| Treatment | 0.308* | 0.296* | 0.297* | 0.015 | 0.035 | 0.073 |
| (explainable AI) | (0.120) | (0.121) | (0.122) | (0.264) | (0.268) | (0.259) |
| Tenure | Yes | Yes | Yes | Yes | Yes | Yes |
| IT skills | No | Yes | Yes | No | Yes | Yes |
| Decision speed | No | No | Yes | No | No | Yes |
| Observations | 113 | 113 | 113 | 113 | 113 | 113 |

Notes: The table reports three quasi-binomial regression models with different sets of control variables and two different metrics as dependent variables. The standard errors of the treatment effect are reported in parentheses. Statistical significance: $^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$.

# Supplement I   Analysis with excluded participants

In our data analyses, we followed our preregistration and excluded participants who did not finish the task in time or participants with obvious misbehavior. Specifically, six and two participants were excluded from Study 1 and Study 2, respectively, because they did not finish the task in time. Further, in Study 1, we excluded participants who did not label a single product as defective (which corresponds to one participant) and in Study 2, radiologists were excluded if they assigned only one label to all chest X-ray images (which corresponds to 1 radiologist). Further, participants whose performance was more than three standard deviations worse than the mean of their respective treatment arm were excluded (which corresponds to one worker in Study 1 and two radiologists in Study 2). This section repeats the OLS regression from the main paper (without control variables) with participants who were excluded due to obvious misbehavior. Overall, we arrive at consistent findings.

**Table S10: Excluded participants across treatment arms**

|  | Study 1: Manufacturing | | Study 2: Medical | |
|---|---|---|---|---|
|  | Black-box AI | Explainable AI | Black-box AI | Explainable AI |
| Time-out | 5 | 1 | 0 | 2 |
| No defective | 1 | 0 | – | – |
| Single label | – | – | 1 | 0 |
| Worse than $3\sigma$ | 0 | 1 | 1 | 1 |

## I.1   Study 1: Manufacturing experiment

Table S11 reports the OLS regression model estimating the treatment effect for all different combinations of exclusion criteria. As in our main analysis, the effect of explainable AI is statistically significant for both metrics, balanced accuracy and defect detection rate. The only exception is for the defect detection rate when workers that timed-out and did not label a single product as defective were excluded while including those that were worse than three standard deviations than the mean.

**Table S11: OLS regression results with excluded participants (manufacturing experiment)**

| Excluded: | | | Observations | Balanced accuracy | Defect detection rate |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Time-out | No defective | Worse than $3\sigma$ | | | |
| ✗ | ✗ | ✗ | 56 | 8.143** (2.754) | 11.590* (4.965) |
| ✓ | ✗ | ✗ | 50 | 7.944* (3.016) | 12.192* (5.472) |
| ✗ | ✓ | ✗ | 55 | 6.769** (2.432) | 8.651* (4.077) |
| ✗ | ✗ | ✓ | 54 | 8.120*** (2.042) | 10.961** (3.391) |
| ✓ | ✓ | ✗ | 49 | 6.253* (2.642) | 8.628 (4.477) |
| ✓ | ✗ | ✓ | 48 | 7.653** (2.178) | 11.014** (3.680) |
| ✗ | ✓ | ✓ | 54 | 8.120*** (2.042) | 10.961** (3.391) |

Notes: The table reports the OLS regression model with two different metrics as dependent variables for all combinations of exclusion criteria. The standard errors of the treatment effect are reported in parentheses. Statistical significance: $^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$.

## I.2 Study 2: Medical experiment

Table S12 reports the OLS regression model estimating the treatment effect for all different combinations of exclusion criteria. As in our main analysis, the effect of explainable AI is statistically significant for balanced accuracy, irrespective of the exclusion criteria. Only when radiologists that assigned one label to all chest X-ray images are excluded the treatment effect for balanced accuracy was not statistically significant. For the disease detection rate, the treatment effect was not statistically significant in the main analysis. This did not change when different exclusion criteria are considered.

**Table S12: OLS regression results with excluded participants (medical experiment)**

| Excluded: | | | Observations | Balanced accuracy | Disease detection rate |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Time-out | Single label | Worse than $3\sigma$ | | | |
| ✗ | ✗ | ✗ | 118 | 4.418* (2.067) | -0.282 (2.337) |
| ✓ | ✗ | ✗ | 116 | 5.266** (1.996) | 0.543 (2.274) |
| ✗ | ✓ | ✗ | 117 | 3.967 (2.026) | -0.121 (2.349) |
| ✗ | ✗ | ✓ | 115 | 4.988** (1.846) | -0.256 (2.212) |
| ✓ | ✓ | ✗ | 115 | 4.815* (1.950) | 0.704 (2.286) |
| ✓ | ✗ | ✓ | 114 | 5.162** (1.857) | -0.168 (2.232) |
| ✗ | ✓ | ✓ | 114 | 4.520* (1.790) | -0.102 (2.224) |

Notes: The table reports the OLS regression model with two different metrics as dependent variables for all combinations of exclusion criteria. The standard errors of the treatment effect are reported in parentheses. Statistical significance: $^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$.

## Supplement J  Experiment with non-experts

To extend our findings to non-experts, we conducted a third experiment with participants recruited via Amazon MTurk to perform the visual inspection task in the manufacturing setting. We chose the manufacturing task for this because, in principle, non-experts could compare electronic products against a reference image (a faultless product looks always identical). For chest X-ray images, this is hardly possible since these can look very different across different healthy patients. We followed common practice by only admitting MTurk workers with an approval rating above 95% [136]. We prevented double participation by tracking the IP of participants. The participants received a base compensation ($5) and had the opportunity to earn a performance-dependent bonus proportional to the correctly labeled quality defects ($3). Participants were randomly assigned to one of the following three treatments: (a) human with black-box AI, (b) human with explainable AI, and (c) human without AI. Following the preregistration, we aimed to include approximately 600 participants excluding dropouts. We thus recruited 861 participants (U.S. residents) who started the study between July 19 and July 21, 2021. Out of them, 117 participants did not complete the study; 152 failed the tutorial; 92 did not finish on time; and 70 participants were excluded due to obvious misbehavior. The final sample consisted of 430 participants, out of which 288 were assigned to treatment arms (a) or (b), performing $N = 57,600$ assessments of electronic products.

We found that participants supported by explainable AI reached a higher task performance than the participants supported by black-box AI across both metrics (Figure S10). Participants with black-box AI treatment only achieved a balanced accuracy with a mean of 81.4%, whereas participants with explainable AI treatment achieved a balanced accuracy with a mean of 87.6%. We then estimated the overall treatment effect on the task performance by regressing the balanced accuracy on the treatment (black-box AI = 0, explainable AI = 1). The regression results suggest that the treatment effect of explainable AI is statistically significant and large ($\beta = 6.252$, $SE = 1.733$, $P < 0.001$); that is, an improvement of 6.3 percentage points. Accordingly, participants equipped with explainable AI achieved a higher defect detection rate with mean of 77.7% compared to participants with black-box AI with a mean of 66.4%. Again, the regression results showed a large and statistically significant treatment effect of explainable AI ($\beta = 11.271$, $SE = 3.276$, $P = 0.001$). The regression results remain statistically significant

for both metrics when including relevant control variables (demographics, self-reported IT skills, and decision speed) in the regression model (Supplement J.3).

We additionally compared how humans without AI support performed relative to humans with black-box AI or explainable AI. For this, we further recruited 142 participants and assigned them to a third treatment: human without AI. Here, participants only got images of the to-be-inspected products and the corresponding reference images of faultless products, but not the AI-based quality scores or the heatmaps. We found that participants without AI support only achieved a balanced accuracy with a mean of 72.4% (Figure S10) and were significantly outperformed by participants with both black-box AI ($t = 5.507$, $P < 0.001$) and explainable AI ($t = 9.017$, $P < 0.001$). Similar results were found for the defect detection rate, where participants without AI achived a mean of 53.6% and were outperformed by participants with both black-box AI ($t = 5.202$, $P < 0.001$) and explainable AI ($t = 8.733$, $P < 0.001$).

We further explored whether the performance difference between the treatments (black-box AI versus explainable AI) was associated with adherence to AI predictions. For this, we compared how likely participants were to follow quality scores that were accurate (i.e., the AI prediction for the inspected product was correct). The results suggest that participants with explainable AI were more likely to adhere to accurate quality scores than participants with black-box AI (mean = 92.9% for black-box AI, mean = 95.2% for explainable AI). Overall, participants supported by black-box AI were 47.9% more likely to erroneously overrule an AI prediction, despite the prediction being accurate ($t = 2.377$, $P = 0.009$). We also analyzed whether participants were able to identify and overrule AI predictions that were wrong. Here, we found that participants supported by black-box AI only overruled 65.8% of the wrong AI predictions, whereas participants supported by explainable AI overruled 79.1% of the wrong AI predictions. The difference between both treatments is statistically significant ($t = 4.563$, $P < 0.001$). Evidently, explainable AI gives a powerful decision aid: it made participants not only less averse to following accurate AI predictions but also helped them overrule wrong AI predictions.
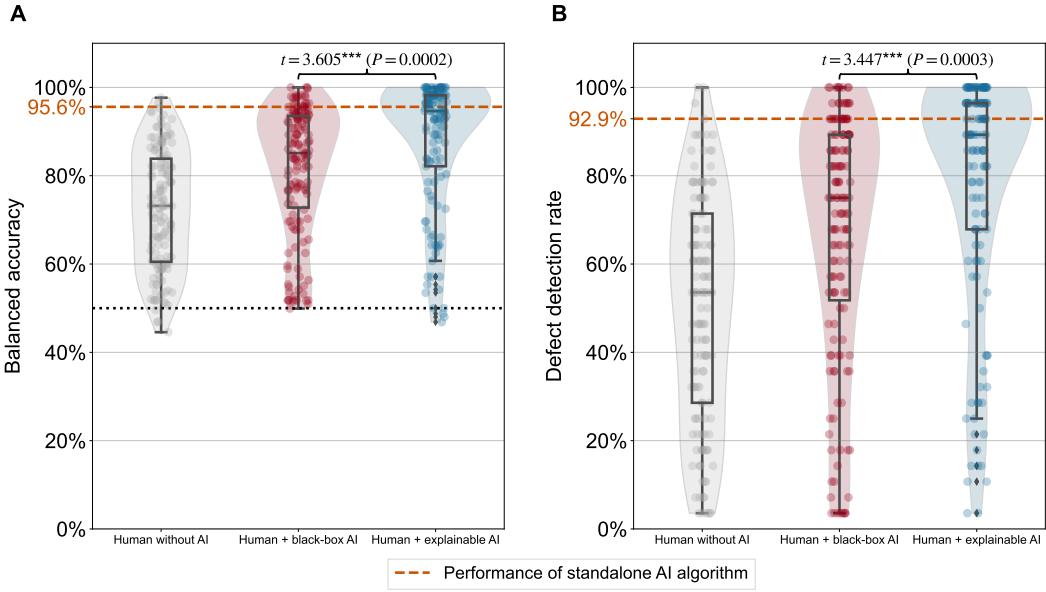
**Figure S10: Results of non-experts experiment.** The boxplots compare the task performance between humans without AI, with black-box AI, and with explainable AI. The task performance is measured by the balanced accuracy (**A**) and the defect detection rate (**B**) based on the quality assessment of participants and the ground-truth labels of the product images. A balanced accuracy of 50% provides a naïve baseline corresponding to a random guess (black dotted line). The standalone AI algorithm attains a balanced accuracy of 95.6% and a defect detection rate of 92.9% (orange dashed lines). Statistical significance is based on a one-sided Welch's $t$-test ($^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$). In the boxplots, the center line denotes the median; box limits are upper and lower quartiles; whiskers are defined as the 1.5x interquartile range.

We also assessed whether participants with explainable AI invested more time for the visual inspection task. For this, we compared participants' median decision speeds across the 200 product images. No significant differences ($t = 0.584$, $P = 0.280$) between both treatments (mean $= 4.61$ s for black-box AI, mean $= 4.50$ s for explainable AI) were observed. Hence, explainable AI improved task performance, but not at the cost of decision speed.

## J.1 Results with precision as task performance metric

In Figure S11, the results with precision as task performance metric are shown. We find that non-experts augmented by explainable AI are more precise in identifying defective electronic products in comparison to peers supported by black-box AI.

71

**Figure S11: Results of non-expert experiment.** The boxplots compare the task performance between the two treatments: black-box AI and explainable AI. The task performance is measured by the defect detection rate (**A**) and the precision (**B**) based on the quality assessment of radiologists and the ground-truth labels of the chest X-ray images. The standalone AI algorithm attains a defect detection rate of 92.9% and a precision of 89.7% (orange dashed lines). Statistical significance is based on a one-sided Welch's $t$-test ($^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$). In the boxplots, the center line denotes the median; box limits are upper and lower quartiles; whiskers are defined as the 1.5x interquartile range.

## J.2 Randomization checks

We performed randomization checks to confirm that the distribution of participants in the three treatment arms of the non-experts experiment was unbiased. The following demographic variables were collected: age bracket [$<20$, 20–30, 30–40, 40–50, 50–60, 60–70, $>70$], gender [male, female, not listed], and highest level of education [no schooling, primary school, some high-school; no degree, high school degree, Bachelor's degree, Master's degree, doctorate]. Table S13 reports the observed frequencies in the three treatment arms. The randomization checks are based on $\mathcal{X}^2$-tests of independence. The results suggest no statistically significant differences between the participants in the three treatment arms.

## Table S13: Randomization checks for non-experts experiment

|  | Human with black-box AI | Human with explainable AI | Human without AI | *P*-value |
|---|---|---|---|---|
| Age | 0 \| 39 \| 58 \| 36 \| 16 \| 1 \| 1 | 0 \| 26 \| 56 \| 34 \| 13 \| 8 \| 0 | 1 \| 46 \| 49 \| 28 \| 13 \| 5 \| 0 | 0.168 |
| Gender | 104 \| 47 \| 0 | 86 \| 51 \| 0 | 89 \| 53 \| 0 | 0.443 |
| Education | 0 \| 0 \| 1 \| 30 \| 88 \| 32 \| 0 | 0 \| 0 \| 1 \| 24 \| 91 \| 21 \| 0 | 0 \| 0 \| 2 \| 16 \| 90 \| 34 \| 0 | 0.278 |
| Observations | 151 | 137 | 142 | – |

Notes: The table reports the frequency of participants that fall in the specific subgroups of age, gender, and education (separated by vertical bars). The *P*-values for the randomization checks are computed based on $\mathcal{X}^2$-tests of independence.

## J.3 Regression models

Table S14 reports three OLS regression models estimating the treatment effect with different control variables. Model (1) estimates the treatment effect for explainable AI with demographic controls (age, gender, and highest level of education). Model (2) estimates the treatment effect for explainable AI with demographic controls and self-reported IT skills (ranging from 1: "novice" to 5: "expert"). Model (3) estimates the treatment effect for explainable AI with demographic controls, self-reported IT skills, and the decision speed (median across the 200 images). All three models return a significant treatment effect for both metrics (balanced accuracy and defect detection rate) as dependent variables.

## Table S14: OLS regression results for treatment effect (non-experts experiment)

|  | Balanced accuracy | | | Defect detection rate | | |
|---|---|---|---|---|---|---|
|  | Model (1) | Model (2) | Model (3) | Model (1) | Model (2) | Model (3) |
| Treatment | 5.792*** | 5.832*** | 5.570** | 10.435** | 10.509** | 10.299** |
| (explainable AI) | (1.729) | (1.720) | (1.707) | (3.274) | (3.258) | (3.263) |
| Demographics | Yes | Yes | Yes | Yes | Yes | Yes |
| IT skills | No | Yes | Yes | No | Yes | Yes |
| Decision speed | No | No | Yes | No | No | Yes |
| Observations | 288 | 288 | 288 | 288 | 288 | 288 |

Notes: The table reports three OLS regression models with different sets of control variables and two different metrics as dependent variables. The standard errors of the treatment effect are reported in parentheses. Statistical significance: ***$P < 0.001$, **$P < 0.01$, *$P < 0.05$.

Table S15 reports three quasi-binomial regression models estimating the treatment effect with

the same control variables as before. Again, all three models return a significant treatment effect for both metrics as dependent variables.

**Table S15: Quasi-binomial regression results for treatment effect (non-experts experiment)**

|  | Balanced accuracy | | | Defect detection rate | | |
| --- | --- | --- | --- | --- | --- | --- |
|  | **Model (1)** | **Model (2)** | **Model (3)** | **Model (1)** | **Model (2)** | **Model (3)** |
| Treatment | 0.449** | 0.455*** | 0.442** | 0.528** | 0.536** | 0.528** |
| (explainable AI) | (0.137) | (0.136) | (0.136) | (0.168) | (0.167) | (0.168) |
| Demographics | Yes | Yes | Yes | Yes | Yes | Yes |
| IT skills | No | Yes | Yes | No | Yes | Yes |
| Decision speed | No | No | Yes | No | No | Yes |
| Observations | 288 | 288 | 288 | 288 | 288 | 288 |

Notes: The table reports three quasi-binomial regression models with different sets of control variables and two different metrics as dependent variables. The standard errors of the treatment effect are reported in parentheses. Statistical significance: $^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$.

## J.4 Analysis with excluded participants

Table S16 reports the number of patients that were excluded according to the three different criteria that we have preregistered.

**Table S16: Excluded participants across treatment arms**

|  | Black-box AI | Explain-able AI | Without AI |
| --- | --- | --- | --- |
| Time-out | 26 | 31 | 35 |
| No defective | 21 | 21 | 27 |
| Worse than $3\sigma$ | 0 | 1 | 0 |

Table S17 reports the OLS regression model estimating the treatment effect for all different combinations of exclusion criteria. As in our main analysis, the effect of explainable AI is statistically significant for both metrics, balanced accuracy and defect detection rate, irrespective of the exclusion criteria.

**Table S17: OLS regression results with excluded participants (non-experts experiment)**

| Excluded: | | | Observations | Balanced accuracy | Defect detection rate |
|---|---|---|---|---|---|
| Time-out | No defective | Worse than $3\sigma$ | | | |
| ✗ | ✗ | ✗ | 388 | 6.344*** (1.825) | 11.539** (3.548) |
| ✓ | ✗ | ✗ | 331 | 4.843* (1.990) | 8.765* (3.902) |
| ✗ | ✓ | ✗ | 342 | 6.966*** (1.648) | 12.652*** (3.023) |
| ✗ | ✗ | ✓ | 387 | 6.561*** (1.817) | 11.788*** (3.549) |
| ✓ | ✓ | ✗ | 289 | 5.918*** (1.756) | 10.863** (3.288) |
| ✓ | ✗ | ✓ | 330 | 5.102* (1.980) | 9.054* (3.904) |
| ✗ | ✓ | ✓ | 341 | 7.238*** (1.631) | 12.988*** (3.014) |

Notes: The table reports the OLS regression model with two different metrics as dependent variables for all combinations of exclusion criteria. The standard errors of the treatment effect are reported in parentheses. Statistical significance: $^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$.

# Supplement K  Preregistered hypotheses

The following hypotheses were preregistered at https://osf.io/7djxb (Study 1) and https://osf.io/69yqt (Study 2):

> **Hypothesis 1 (H1):** *Explainable AI improves the overall decision performance (measured by the balanced accuracy and defect detection rate) compared to humans without AI (i.e., manual inspection) ($\alpha = 0.05$).*

> **Hypothesis 2 (H2):** *Explainable AI improves the overall decision performance (measured by the balanced accuracy and defect detection rate) compared to black-box AI ($\alpha = 0.05$).*

> **Hypothesis 3 (H3):** *Explainable AI reduces variation in decision performance (measured by the variance in the balanced accuracy and defect detection rate) compared to black-box AI ($\alpha = 0.05$).*

> **Hypothesis 4 (H4):** *Explainable AI increases the trust in model decisions (measured by the rate of correct model decisions that are not overruled by the user) compared to black-box AI ($\alpha = 0.05$).*

Table S18 summarizes the results from all three studies: (1) the manufacturing experiment at *Siemens*, (2) the medical experiment, and (3) the manufacturing experiments with non-experts from Amazon MTurk. We report the $P$-values for both the balanced accuracy and defect detection rate for hypotheses H1, H2, and H3. The statistical testing for hypotheses H1, H2, and H4 are based on one-sided Welch's $t$-tests. The statistical testing for Hypothesis H3 is based on Levene's test for equality of variances. As specified in our preregistration, we refrained from testing Hypothesis H1 in our manufacturing field experiment and our medical experiment. The reason is that we wanted sufficient power in our main treatment arms of interest (i.e., black-box AI versus explainable AI). All hypotheses except for Hypothesis H3 in the non-experts experiment and Hypotheses H2 and H3 for the disease detection rate in the medical experiment were confirmed at a significance level of $\alpha = 0.05$. The latter can be expected since missing a lung lesion has more serious consequences than erroneously believing a lung lesion is visible; thus, leading to conservative decision-making of radiologists. Therefore, we additionally inspected precision as a task performance metric. We find that radiologists augmented with explainable AI were significantly more precise in identifying lung lesions compared to radiologists with black-box AI (see Supplement G).

**Table S18: Comparison of results against preregistered hypotheses.**

|  | Study 1: Manufacturing | Study 2: Medical | Study 3: Non-experts |
|---|---|---|---|
| **H1** (BACC) | *not part of preregistration* | *not part of preregistration* | ✓ ($P < 0.001$) |
| **H1** (DDR) | *not part of preregistration* | *not part of preregistration* | ✓ ($P < 0.001$) |
| **H2** (BACC) | ✓ ($P = 0.001$) | ✓ ($P = 0.004$) | ✓ ($P < 0.001$) |
| **H2** (DDR) | ✓ ($P = 0.004$) | ✗ ($P = 0.498$) | ✓ ($P < 0.001$) |
| **H3** (BACC) | ✓ ($P = 0.002$) | ✓ ($P = 0.033$) | ✗ ($P = 0.356$) |
| **H3** (DDR) | ✓ ($P = 0.023$) | ✗ ($P = 0.790$) | ✗ ($P = 0.217$) |
| **H4** | ✓ ($P = 0.011$) | ✓ ($P = 0.001$) | ✓ ($P = 0.009$) |

Notes: The significance level was preregistered at $\alpha = 0.05$ and the marks denote whether the corresponding $P$-value was significant at this level. BACC refers to balanced accuracy and DDR to defect/disease detection rate.

# Supplement L  Post-experimental questionnaire

For post-hoc exploratory analyses, we asked participants to complete a questionnaire. The questions involved established constructs, such as self-reported task load [137], perceived usefulness [138], perceived ease of use [138], and self-reported trust [139]. We further asked participants about their previous experience and the perceived performance of the AI algorithm. In the manufacturing experiment, the questions were translated to German. In the medical experiment, the questions were adapted for the medical setting (for the exact wording, see our preregistration https://osf.io/69yqt). Participants in the non-experts experiment were asked to answer the questions from the viewpoint of a factory worker ("Imagine you work in a factory with a similar job task as you just did."). The results for all three studies are provided in Tables S19 to S23.

## Table S19: Self-reported task load

| Question | Study 1: Manufacturing | | Study 2: Medical | | Study 3: Non-experts | |
|---|---|---|---|---|---|---|
| | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI |
| How mentally demanding was the task? (1 = very low, 7 = very high) | 3.41 (1.37) | 3.54 (1.17) | 3.80 (1.34) | 3.73 (1.60) | 4.72 (1.58) | 4.87 (1.62) |
| How physically demanding was the task? (1 = very low, 7 = very high) | 3.09 (1.54) | 2.85 (1.38) | 2.18 (1.32) | 2.52 (1.42) | 4.09 (2.10) | 3.95 (2.06) |
| How hurried or rushed was the pace of the task? (1 = very low, 7 = very high) | 3.59 (1.10) | 3.92 (1.23) | 2.76 (1.35) | 3.16 (1.66) | 4.56 (1.55) | 4.44 (1.68) |
| How successful were you in accomplishing what you were asked to do? (1 = very poor, 7 = very good) | 5.27 (0.94) | 5.81 (0.85) | 5.64 (1.09) | 5.57 (0.93) | 5.68 (1.08) | 5.87 (1.02) |
| How hard did you have to work to accomplish your level of performance? (1 = very low, 7 = very high) | 4.00 (1.35) | 4.00 (0.94) | 3.33 (1.28) | 3.36 (1.30) | 5.18 (1.38) | 5.29 (1.46) |
| How insecure, discouraged, irritated, stressed, and annoyed were you? (1 = very low, 7 = very high) | 2.77 (1.31) | 2.77 (1.58) | 2.76 (1.28) | 2.91 (1.43) | 3.52 (1.98) | 3.31 (1.95) |

Notes: The table reports the average scores for the self-reported task load. Standard deviations are reported in parentheses.

## Table S20: Perceived usefulness

| Question | Study 1: Manufacturing | | Study 2: Medical | | Study 3: Non-experts | |
| --- | --- | --- | --- | --- | --- | --- |
| | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI |
| Using the Artificial Intelligence System in my job would enable me to accomplish tasks more quickly. (1 = very unlikely, 7 = extremely likely) | 4.95 (1.70) | 5.12 (1.37) | 5.00 (1.40) | 5.25 (1.06) | 5.68 (1.02) | 5.86 (1.08) |
| Using the Artificial Intelligence System would improve my job performance. (1 = very unlikely, 7 = extremely likely) | 4.91 (1.54) | 5.23 (0.99) | 5.04 (1.24) | 5.09 (1.14) | 5.68 (1.04) | 5.98 (1.06) |
| Using the Artificial Intelligence System in my job would increase my productivity. (1 = very unlikely, 7 = extremely likely) | 4.95 (1.43) | 4.96 (1.22) | 5.04 (1.38) | 5.39 (1.15) | 5.71 (1.11) | 5.98 (1.01) |
| Using the Artificial Intelligence System in my job would enhance my effectiveness on the job. (1 = very poor, 7 = very good) | 5.05 (1.53) | 4.88 (1.14) | 4.98 (1.34) | 5.20 (1.25) | 5.66 (1.12) | 5.91 (0.97) |
| Using the Artificial Intelligence System would make it easier to do my job. (1 = very unlikely, 7 = extremely likely) | 5.18 (1.26) | 5.08 (1.32) | 4.93 (1.45) | 5.23 (1.08) | 5.65 (1.23) | 6.01 (1.00) |
| I would find the Artificial Intelligence System useful in my job. (1 = very unlikely, 7 = extremely likely) | 5.27 (1.42) | 5.23 (1.24) | 5.00 (1.41) | 5.09 (1.22) | 5.79 (1.09) | 6.07 (1.04) |

Notes: The table reports the average scores for the perceived usefulness of the AI algorithm. Standard deviations are reported in parentheses.

## Table S21: Perceived ease of use

| Question | Study 1: Manufacturing | | Study 2: Medical | | Study 3: Non-experts | |
| --- | --- | --- | --- | --- | --- | --- |
| | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI |
| Learning to operate the Artificial Intelligence System would be easy for me. (1 = very unlikely, 7 = extremely likely) | 5.00 (1.45) | 5.35 (1.20) | 5.78 (0.88) | 6.00 (0.86) | 5.69 (1.11) | 5.90 (0.95) |
| I would find it easy to get the Artificial Intelligence System to do what I want it to do. (1 = very unlikely, 7 = extremely likely) | 4.05 (1.36) | 4.46 (0.90) | 4.87 (1.22) | 5.09 (1.03) | 5.61 (1.02) | 5.78 (0.97) |
| My interaction with the Artificial Intelligence System would be clear and understandable. (1 = very unlikely, 7 = extremely likely) | 5.09 (0.97) | 5.27 (0.96) | 5.16 (1.07) | 5.20 (1.21) | 5.65 (1.08) | 5.99 (0.93) |
| I would find the Artificial Intelligence System to be flexible to interact with. (1 = very poor, 7 = very good) | 4.82 (1.22) | 5.00 (1.06) | 4.82 (1.28) | 4.84 (1.27) | 5.36 (1.24) | 5.55 (1.10) |
| It would be easy for me to become skillful at using the Artificial Intelligence System. (1 = very unlikely, 7 = extremely likely) | 5.14 (1.21) | 5.38 (0.85) | 5.38 (1.21) | 5.61 (0.84) | 5.73 (0.97) | 5.93 (0.99) |
| I would find the Artificial Intelligence System easy to use. (1 = very unlikely, 7 = extremely likely) | 5.14 (0.94) | 5.42 (1.03) | 5.16 (1.17) | 5.66 (0.83) | 5.71 (1.03) | 6.05 (0.96) |

Notes: The table reports the average scores for the perceived ease of use of the AI algorithm. Standard deviations are reported in parentheses.

## Table S22: Previous experience and perceived performance

| Question | Study 1: Manufacturing | | Study 2: Medical | | Study 3: Non-experts | |
|---|---|---|---|---|---|---|
| | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI |
| How strong would you consider your general IT skills? (1 = novice, 5 = expert) | 2.86 (0.94) | 3.08 (0.84) | 3.47 (1.01) | 3.77 (1.01) | 3.48 (0.99) | 3.48 (1.08) |
| How often do you interact with Artificial Intelligence in your job? (1 = very little, 5 = very much) | 2.32 (1.13) | 2.23 (1.11) | *not part of preregistration* | | 3.18 (1.31) | 3.12 (1.32) |
| How familiar do you feel with Artificial Intelligence in general? (1 = very little, 5 = very much) | 2.41 (1.10) | 2.88 (1.07) | 3.00 (1.07) | 3.07 (0.95) | 3.56 (0.98) | 3.48 (0.93) |
| How well did the Artificial Intelligence System perform in comparison to your expectations? (1 = very poor, 7 = very good) | 5.00 (1.57) | 6.08 (1.06) | 4.49 (1.41) | 4.48 (1.47) | 5.72 (1.01) | 6.15 (0.85) |
| How likely is the Artificial Intelligence System to make a bad estimate? (1 = very unlikely, 7 = very likely) | 3.55 (1.34) | 3.04 (1.08) | 3.93 (1.14) | 4.02 (1.13) | 3.25 (1.61) | 2.93 (1.68) |
| Completing the quality inspections task has changed my opinion about Artificial Intelligence. (1 = strongly disagree, 7 = strongly agree) | 4.64 (1.36) | 4.08 (1.38) | *not part of preregistration* | | 4.68 (1.65) | 4.83 (1.61) |
| How much did you rely on the Artificial Intelligence System? (1 = very little, 7 = very much) | 4.32 (1.46) | 4.77 (1.37) | *not part of preregistration* | | 4.96 (1.43) | 5.54 (1.38) |
| The Artificial Intelligence System provides clear explanations for its outputs. (1 = strongly disagree, 7 = strongly agree) | 4.59 (1.33) | 5.04 (1.00) | *not part of preregistration* | | 4.99 (1.54) | 5.71 (1.20) |

Notes: The table reports the average scores for previous experience and the perceived performance of the AI algorithm. Standard deviations are reported in parentheses.

## Table S23: Self-reported trust in AI algorithm

| Question | Study 1: Manufacturing | | Study 2: Medical | | Study 3: Non-experts | |
|---|---|---|---|---|---|---|
| | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI | Human with black-box AI | Human with explainable AI |
| The Artificial Intelligence System is deceptive. (1 = strongly disagree, 7 = strongly agree) | 3.27 (1.20) | 2.88 (1.18) | 3.31 (1.41) | 3.39 (1.28) | 3.79 (2.04) | 3.72 (2.16) |
| I am suspicious of the Artificial Intelligence System's intent, action, or outputs. (1 = strongly disagree, 7 = strongly agree) | 2.91 (1.27) | 2.62 (0.98) | 3.09 (1.69) | 2.98 (1.56) | 3.88 (2.04) | 3.69 (2.15) |
| The Artificial Intelligence System's actions will have a harmful outcome. (1 = strongly disagree, 7 = strongly agree) | 3.09 (1.41) | 2.65 (1.09) | 3.38 (1.51) | 3.34 (1.43) | 3.73 (2.04) | 3.48 (2.08) |
| I am confident in the Artificial Intelligence System. (1 = very poor, 7 = very good) | 4.91 (1.06) | 4.96 (1.15) | 4.04 (1.31) | 4.34 (1.22) | 5.52 (1.20) | 5.89 (0.97) |
| The Artificial Intelligence System is reliable. (1 = strongly disagree, 7 = strongly agree) | 4.95 (0.84) | 5.27 (1.04) | 4.07 (1.39) | 4.27 (1.09) | 5.63 (1.10) | 5.93 (0.85) |
| I can trust the Artificial Intelligence System. (1 = strongly disagree, 7 = strongly agree) | 4.86 (0.77) | 5.04 (0.92) | 3.87 (1.39) | 4.07 (1.11) | 5.66 (1.06) | 5.80 (0.94) |

Notes: The table reports the average scores for the self-reported trust in the AI algorithm. Standard deviations are reported in parentheses.