

Human-Centered Explainable AI (XAI): From Algorithms to User Experiences

Q. VERA LIAO*, Microsoft Research, Canada

KUSH R. VARSHNEY, IBM Research, United States

(Book Chapter Draft 10/2021) As a technical sub-field of artificial intelligence (AI), explainable AI (XAI) has produced a vast collection of algorithms in recent years. However, explainability is an inherently human-centric property and the field is starting to embrace inter-disciplinary perspectives and human-centered approaches. As researchers and practitioners begin to leverage XAI algorithms to build XAI applications, explainability has moved beyond a demand by data scientists or researchers to comprehend the models they are developing, to become an essential requirement for people to trust and adopt AI deployed in numerous domains. Human-computer interaction (HCI) research and user experience (UX) design in this area are therefore increasingly important. In this chapter, we begin with a high-level overview of the technical landscape of XAI algorithms, then selectively survey recent HCI work that takes human-centered approaches to design, evaluate, provide conceptual and methodological tools for XAI. We ask the question “*what are human-centered approaches doing for XAI*” and highlight three roles that they should play in shaping XAI technologies: to drive technical choices by understanding users’ explainability needs, to uncover pitfalls of existing XAI methods through empirical studies and inform new methods, and to provide conceptual frameworks for human-compatible XAI.

1 INTRODUCTION

In everyday life, people seek explanations when there is a gap of understanding. Explanations are sought for many goals that this understanding is meant to serve, such as predicting future events, diagnosing a problem, resolving cognitive dissonance, assigning blame, and rationalizing one’s action. In interactions with computing technologies, an appropriate understanding of how the system works, often referred to as the user’s “mental model” [78], is the foundation for users to correctly anticipate system behaviors and interact effectively. A user’s understanding is constantly being shaped by what they see and experience with the system, and can be refined by being directly explained how the system works.

With the increasing adoption of AI technologies, especially popular inscrutable “opaque-box” machine learning (ML) models such as neural networks models, understanding becomes increasingly difficult. Meanwhile, the need for stakeholders to understand AI is heightened by the uncertain nature of ML systems and the hazardous consequences they can possibly cause as AI is now frequently deployed in high-stakes domains such as healthcare, finance, transportation, and even criminal justice. Some are concerned that this challenge of understanding will become the bottleneck for people to trust and adopt AI technologies. Others have warned that a lack of human scrutiny will inevitably lead to failures in usability, reliability, safety, fairness, and other moral crises of AI.

It is with this overwhelming challenge of modern AI that the term explainable AI (XAI) has made its way into numerous academic works, industry efforts, as well as public policy and regulatory requirements. For example, the European Union General Data Protection Regulation (GDPR) now requires that “meaningful information about the logic involved” must be provided to people who are affected by automated decision-making systems. However, despite a vast collection of XAI algorithms produced by the AI research community and a recent emergence of off-the-shelf toolkits (e.g. [1–4, 12]) for AI developers to incorporate state-of-the-art XAI techniques in their own models, successful examples of XAI are still relatively scarce in real-world AI applications.

*Reviewed work by the first author was done while working at IBM Research

Authors’ addresses: Q. Vera Liao, Microsoft Research, Montreal, Canada, veraliao@microsoft.com; Kush R. Varshney, IBM Research, Yorktown Heights, United States, krvarshn@us.ibm.com.

Developing XAI applications is challenging because explainability, or the effectiveness of explanation, is not intrinsic to the model but lies in the perception and reception of the person receiving the explanation. Making a model completely transparent to its nuts and bolts does not guarantee that the person at the receiving end can make sense of all the information, or is not overwhelmed. What makes an explanation good—to provide appropriate information that can be understood and utilized—is contingent on the receiver’s current knowledge and their goal for receiving the explanation, among other human factors.

Therefore, developing XAI applications requires human-centered approaches that center the technical choices around people’s explainability needs, and define success by human experience and well-being. It also means that XAI presents as much of a design challenge as an algorithmic challenge. Hence there are rich opportunities for HCI researchers and design practitioners to contribute insights, solutions, and methods to make AI more explainable and responsible. We hope this chapter serves as a call to engagement in this inter-disciplinary endeavor by presenting a selected overview of recent AI and HCI work on the topic of XAI. We will highlight three perspectives that XAI can gain from human-centered approaches:

- There is no one-fits-all solution in the growing collection of XAI techniques. The technical choice should be driven by users’ explainability needs, for which HCI and user research can offer insights and methodological tools (Section 3).
- Empirical studies with users can reveal pitfalls of existing XAI methods. To overcome the pitfalls requires both design efforts to fill the gaps and challenging fundamental assumptions in techno-centric views (Section 4).
- Theories of human cognition and behaviors can offer conceptual tools to inspire new computational and design frameworks for XAI. However, this is still a nascent area and relevant theories in social science, behavioral science, and information science are yet to be explored (Section 5).

Before getting to these points, we will start with a brief overview of the technical landscape of XAI to ground our discussions (Section 2). For interested readers, we suggest several recent papers that provided deeper technical surveys [6, 11, 20, 44].

2 WHAT IS EXPLAINABLE AI AND WHAT ARE THE TECHNIQUES?

The definitions of explainability and related terms such as transparency, interpretability, intelligibility, and comprehensibility, are in a bit of flux. Scholars sometimes disagree on their scopes and how these terminologies intersect. However, XAI work often shares a common goal of *making AI understandable by people*. By adopting this pragmatic, human-centered definition in the chapter, we consider XAI as broadly encompassing all technical means to this end of understanding, including direct interpretability, generating an explanation or justification, providing transparency information, etc. (and avoid the philosophical question “what is or is not an explanation” altogether [80]). We note a distinction between a narrow scope of XAI focusing on explaining the model processes or internals, versus a broad scope that covers all explanatory information about the model, also including the training data, performance, uncertainty, and so on [63, 96]. Since our focus is on XAI applications, we believe a broad view is necessary as users are often interested in a holistic understanding of the system. However, technical challenges are commonly presented in the inscrutability of the model internals, so the XAI techniques we review in this section are within the narrower scope of XAI.

It is worth mentioning that while the majority of XAI work focuses on ML models, and so does this chapter, there are emerging areas of other types of XAI including explainable planning, multi-agent systems, robotics, etc. In fact, the

term XAI was coined half a century ago in the context of expert systems. We are currently in its second wave spurred by the popularity of ML.

At a high level, XAI techniques falls into two camps [44, 66]: 1) choosing a directly interpretable model, such as simpler models like decision trees, rule-based models, and linear regression; 2) choosing an “opaque-box” model such as deep-neural networks and large tree ensembles, and then using a post-hoc technique to generate explanations. The choice between the two is sometimes discussed under the term “performance-interpretability tradeoff”, as complex opaque-box models tend to perform better in many tasks. However, this tradeoff is not always true. Research has shown that in many contexts, especially with well-structured datasets and meaningful features, directly interpretable models can reach comparable performance than opaque-box models [90]. Moreover, an active research area of XAI focuses on developing new algorithms that possess both performance advantages and interpretability properties. For example, decision sets [58], generalized linear rule models [101], GA2Ms [19], and CoFrNets [84] are recent algorithms that have more advanced computational properties than simple rule-based or linear models, but the model behaviors are still represented in meaningful rules or coefficients that can be understood relatively easily.

However, opaque-box models are often chosen in practice because of their performance advantage for a given dataset, often lower requirement for human effort (e.g., feature engineering), or the availability of off-the-shelf solutions. In these cases, one will have to use a post-hoc XAI technique. Based on their purpose, Guidotti et al. [44] categorize post-hoc XAI techniques into *global explanation* on the overall logic of the model, *local explanation* on a particular prediction, and *counterfactual inspection* that supports understanding how the model would behave with alternative input. Within these categories, XAI techniques commonly generate either *feature-based* explanations to elucidate the model internals, or *example-based* explanations to support case-based reasoning.

Note the three categories also apply to directly interpretable models. For example, a shallow decision-tree can be presented directly as a global explanation, highlighted of a particular path to locally explain a prediction, or traced in alternative paths to perform counterfactual inspection. It is, however, much less straightforward with opaque-box models, which require separate post-hoc techniques, as we will give some examples below.

Examples of global explanation. Since it is impossible to understand the complex internals of an opaque-box model, the goal of global explanation is to provide an approximate overview of how the model behaves. This is often done by training a simple directly interpretable model such as decision tree, rule sets or regression with the same training data, and performing optimization to make the simple model behaving more closely to the original model. For example, a technique called distillation changes the learning objective of the interpretable model to matching the original model’s predictions [95]. SRatio reweighs the training data based on the original model’s predictions and then re-trains the interpretable model [27]. With these approaches, depending on the choice of the approximate model, global explanations can take the form of a decision-tree, a set of rules the model follows, or feature weights.

Examples of local explanation. To explain a prediction made on an instance, a number of algorithms can be used to estimate the importance of each feature of this instance for the model’s prediction. For example, LIME (local interpretable model-agnostic explanations) [86] starts by adding a small amount of noise to the instance to create a set of neighbor instances, with them it fits a simple linear model that mimics the original model’s behavior in the local region. The linear model’s weights can then be used as the feature importance to explain the prediction. Another popular algorithm SHAP (Shapley additive explanations) [70] defines feature importance based on Shapley values, inspired by cooperative game theory, to assign credit to each feature. Feature-importance explanations can be shown to users by visualizing the importance, or simply describing the most important features for the prediction. To explain deep neural networks, many other algorithms can be used to identify important parts of input features based on gradient [91], propagation [13],

occlusion [62], etc. They are sometimes referred to as saliency methods and, when applied to image data, generate saliency maps. Example-based methods are useful to explain a prediction as well. For example, with some notion of similarity, finding similar instances in the training data with the same predicted outcome can be used to justify the prediction [45, 55].

Examples of counterfactual inspection. Different from local explanations that describe the model’s prediction process for a given instance, counterfactual explanations—“counter to the facts”—are sought when people are interested in how the model would behave when the current input changes. In other words, people are interested in the “why not a different prediction” or “how to change to get a different prediction” questions rather than a descriptive “why” question. Such explanations are especially sought when seeking recourse to a current, often undesirable, prediction, such as ways to improve a patient’s predicted high risk of disease. Several algorithms can be used to generate counterfactual explanations by identifying changes, often with some notion of minimum changes, needed for an instance to receive a different prediction [26, 67]. They are sometimes referred to as contrastive explanations for a counterfactual outcome (differentiated from other kinds of counterfactuals such as counterfactual causes). Example-based methods can also be used to generate counterfactual examples—instances with minimum difference from the original one but having a different outcome [76, 97]. In other situations, people may want to zoom in on a specific feature and explore how its changes impact the model’s prediction, i.e. asking a “what if” question. For this purpose, feature inspection techniques such as partial dependence plot (PDP) [47] and individual conditional expectation (ICE) can be used [42].

Most post-hoc techniques make some approximations. Distillation and LIME approximate the complex model’s behaviors with a simpler model’s. PDP leaves out interactions between features. Example-based methods explain by samples in the data. There is a long-standing debate regarding the potential risk of using approximate post-hoc techniques to explain instead of a directly interpretable model, as approximation will inevitably leave out some corner cases or even be unfaithful to what the original model computes [90].

However, in addition to the practical reasons mentioned earlier to opt for an opaque-box model, there is a pragmatic argument to be made about the diverse communication devices people use to reach “sufficient understanding” to achieve a given objective. For example, if one is to make precise diagnosis of a problem, they may need explanations that describe a causal chain; whereas if one’s goal is to predict future events, following approximate rules or case-based reasoning could be sufficient and less cognitively demanding. One can also argue that when the model and the person have different epistemic access, approximation can be seen as a form of translation necessary to bridge the two. There is an emerging area of XAI research on generating human-consumable explanations with supervision of human explanations [34, 50, 56], which essentially translates model reasoning into meaningful human explanations applied to the same prediction. This kind of explanation is a complete approximation but could be especially useful for lay people who have difficulty understanding how ML models work, but want to get a sense of the reasonability of a prediction. We further highlight this objective and user dependent nature of the choice of explanation methods in the next sections.

That being said, developers of AI have a responsibility to understand, mitigate, and transparently communicate the limitations of approximate explanations to stakeholders. For example, an explainability metric known as faithfulness can be used to detect faulty post-hoc explanations [7]. This is an actively researched topic and there is still a lack of principled approaches to identify and communicate the limitations of post-hoc explanations.

3 DIVERSE EXPLAINABILITY NEEDS OF AI STAKEHOLDERS

It is easy to see that there is no “one-fits-all” solutions from this vast, and still rapidly growing, collection of XAI algorithms, and the choice should be based on target users’ explainability needs. The challenge here are twofold: First,

users of XAI are far from a uniform group and their explainability needs can vary significantly depending on their goals, backgrounds, usage contexts, etc. Second, XAI algorithms were often not developed with specific usage contexts in mind, or were developed primarily to help model developers or AI researchers inspect the model [74]. Hence their appropriateness to support an end users' explainability needs can be unclear.

A starting point to address these challenges is to map out the design space of XAI and develop frameworks that account for people's diverse explainability needs. Many have summarized common user groups that demand explainability and what they would use AI explanations for [11, 49, 83]:

- **Model developers**, to improve or debug the model.
- **Business owners or administrators**, to assess an AI application's capability, regulatory compliance, etc. to determine its usage.
- **Decision-makers**, who are direct users of AI decision support applications, to form appropriate trust in the AI and make informed decisions.
- **Impacted groups**, whose life could be impacted by the AI, to seek recourse or contest the AI.
- **Regulatory bodies**, to audit for legal or ethical concerns such as fairness, safety, privacy, etc.

While useful for considering different personas interacting with XAI, this kind of categorization lacks granularity to characterize people's explainability needs. For example, a doctor using a patient risk-assessment AI (i.e., a decision-maker) would want to have an overview of the system during the onboarding stage, but delve into AI's reasoning for a particular patient's risk assessment when they treat the patient. Also, people in any of these groups may want to assess model capabilities or biases at certain usage points.

In a recent HCI paper, Suresh et al. define stakeholder's *knowledge* and their *objectives* as two components that cut across to determine one's explainability needs [93]. The authors characterize stakeholders' knowledge by formal, instrumental, and personal knowledge and how it manifests in the contexts of machine learning, data domain, and general milieu. For stakeholders' goals and objectives, the authors propose a multi-level typology, ranging from long-term goals (building trust and understanding the model), immediate objectives (debug and improve, ensure compliance with regulations, take actions based on model output, justify actions influenced by a model, understand data usage, learn about a domain, contest model decisions), and specific tasks to perform with explanations (assess reliability of a prediction, detect mistakes, understand information used by the model, understand feature influence, understand model strengths and limitations).

While these efforts can be seen as top-down approaches to characterize the overall space of explainability needs, a complementary approach is to follow user-centered design and start with user research to identify application or interaction specific explainability needs. For example, Eiband et al. proposed a participatory design method that starts with analyzing users' current mental model and gaps with how the system should be understood, based on an appropriate mental model prescribed by experts, to identify what needs to be explained [37].

In our own research with collaborators, we proposed to identify users' explainability needs by eliciting user questions to understand the AI [63]. This notion is based on prior HCI work using prototypical questions to represent "intelligibility types" [65], and social science literature showing that people's explanatory goals can be expressed in different kinds of questions [48]. By interviewing 20 designers, we collected common questions users ask across 16 ML applications and developed an *XAI Question Bank*, with more than more than 50 detailed user questions organized in 9 categories:

- **How** (global model-wide): asking about the general logic or process the AI follows to have a global view.
- **Why** (a given prediction): asking about the reason behind a specific prediction.

- **Why Not** (a different prediction): asking why the prediction is different from an expected or desired outcome.
- **How to be That** (a different prediction)¹: asking about ways to change the instance to get a different prediction.
- **How to Still Be This** (the current prediction): asking what change is allowed for the instance to still get the same prediction.
- **What if**: asking how the prediction changes if the input changes.
- **Performance**: asking about the performance or of the AI.
- **Data**: asking about the training data.
- **Output**: asking what can be expected or done with the AI’s output.

These questions demonstrate that XAI should be defined broadly, not limited to explaining model internals, as users are also interested in explanatory information about the performance, data, and scope of output, among other dimensions.

This XAI question bank maps out the space of common explainability needs and can be used as a tool to identify applicable questions in user research. In a follow-up work [64], we propose a question-driven user centered design method that starts with identifying key user questions with user research, using them to guide the choice of XAI techniques and iterative design. To facilitate this process and foreground users’ explainability needs, we suggest reframing the technical space of XAI by the user question that each XAI method can address. For example, a feature-importance local explanation technique can answer the *Why* question, while a counterfactual explanation can answer the *How to be That* question. We provide a suggested mapping between the question categories and example XAI methods in Table 1, focusing on post-hoc methods that are available in current open-source XAI toolkits accessible for practitioners [1–4].

In short, a growing collection of XAI techniques offer a rich toolbox for researchers and practitioners to build XAI applications. Making effective and responsible choices from this toolbox should be guided by users’ explainability needs. HCI research not only offers means to understand user needs for specific applications, but also insights about real-world user needs to better frame and organize this toolbox, as well as methodological tools to help navigating the toolbox. In the next section, we discuss how HCI research can also inform limitations of the current technical XAI toolbox.

4 PITFALLS OF XAI: MINDING THE GAPS BETWEEN ALGORITHMIC EXPLANATIONS AND ACTIONABLE UNDERSTANDING

With so many XAI algorithms developed, one must ask: do they work? The answer is complicated because of the diverse contexts that XAI is sought for. The answer is also difficult because it requires understanding how people perceive, process, and use AI explanations. HCI research, and more broadly human-subject studies, are key to evaluating XAI in the context of use [30], identifying where it falls short, and informing human-centered solutions. While many studies showed positive results that XAI techniques can improve people’s understanding of models [46, 59, 68, 87], in this section we draw attention to a few pitfalls of XAI based on recent HCI research.

We start with the position that users’ goal with XAI is not an understanding defined in a vacuum, but an *actionable understanding* that is sufficient to serve the objective that they seek explanations for. As discussed above, these objectives are diverse and dynamic. One common pitfall and overall obstacle for the current XAI field is that, despite growing efforts, there is still *a disconnect between technical XAI approaches and their downstream effectiveness in supporting different user objectives in deployment*. A recent study by Bućinca et al. points out that “proxy tasks” widely used by AI

¹The difference between *Why Not* and *How to Be That* can be subtle and context-dependent. User may ask *Why Not* when seeing an unexpected prediction and interested in comparing what gets the counterfactual outcome. User may ask *How to Be That* when seeking recourse so the explanation should more specifically focus on minimum or actionable changes they can make to the current input.

Question	Ways to explain	Example XAI methods
How (global model-wide)	<ul style="list-style-type: none"> Describe the general model logic as feature impact*, rules† or decision-trees‡ If user is only interested in a high-level view, describe what are the top features or rules considered 	ProfWeight*†‡ [28], Global feature importance* [69, 102], Global feature inspection plots* (e.g. PDP [47]), Tree surrogates‡ [24]
Why (a given prediction)	<ul style="list-style-type: none"> Describe how features of the instance, or what key features, determine the model’s prediction of it* Or describe rules that the instance fits to guarantee the prediction† Or show similar examples with the same predicted outcome to justify the model’s prediction‡ 	LIME* [86], SHAP* [70], LOCO* [61], Anchors† [87], ProtoDash‡ [45]
Why Not (a different prediction)	<ul style="list-style-type: none"> Describe what features of the instance determine the current prediction and/or with what changes the instance would get the alternative prediction* Or show prototypical examples that have the alternative outcome† 	CEM* [26], Counterfactuals* [67], ProtoDash† (on alternative prediction) [45]
How to Be That (a different prediction)	<ul style="list-style-type: none"> Highlight feature(s) that if changed (increased, decreased, absent, or present) could alter the prediction, often ones with minimum effort required* Or show examples with minimum differences but had a different outcome than the prediction† 	CEM* [26], Counterfactuals* [67], Counterfactual instances† [97], DiCE† [76]
How to Still Be This (the current prediction)	<ul style="list-style-type: none"> Describe features/feature ranges or rules that could guarantee the same prediction* Or show examples that are different from the instance but still had the same outcome† 	CEM* [26], Anchors† [87]
What if	<ul style="list-style-type: none"> Show how the prediction changes corresponding to the inquired change of input 	PDP [47], ALE [9], ICE [42]
Performance	<ul style="list-style-type: none"> Provide performance information of the model Provide uncertainty information for each prediction* Describe potential strengths and limitations of the model 	Precision, Recall, Accuracy, F1, AUC; Communicate uncertainty of each prediction* [40]; See examples in FactSheets [10] and Model Cards [75]
Data	<ul style="list-style-type: none"> Provide comprehensive information about the training data, such as the source, provenance, type, size, coverage of population, potential biases, etc. 	See examples in FactSheets [10] and DataSheets [38]
Output	<ul style="list-style-type: none"> Describe the scope of output or system functions. If applicable, suggest how the output should be used for downstream tasks or user workflow 	See examples in FactSheets [10] and Model Cards [75]

Table 1. A mapping between categories of user questions in XAI question bank [63] and example XAI methods to answer these questions, with descriptions of their output in “Ways to explain” column. XAI methods are selected based on what are available in current open-source XAI toolkits [1–4]. The last three rows (in *italic*) are broader XAI needs not limited to explaining model processes

researchers to evaluate their proposed XAI techniques can be misleading [16]. A common example of proxy task is a “simulation test”, which asks people to predict the model’s output based on an input and explanation. Such tests to assess people’s understanding without a specific end goal can fail to predict the success of using XAI in real tasks that people seek explanations for, such as debugging models [54] or improving decision-making [16, 100].

There can be a multitude of reasons for this divide between effectiveness in proxy tasks and deployment. As Bućinca et al. point out, performing proxy tasks in a controlled setting could induce a different cognitive process from a realistic

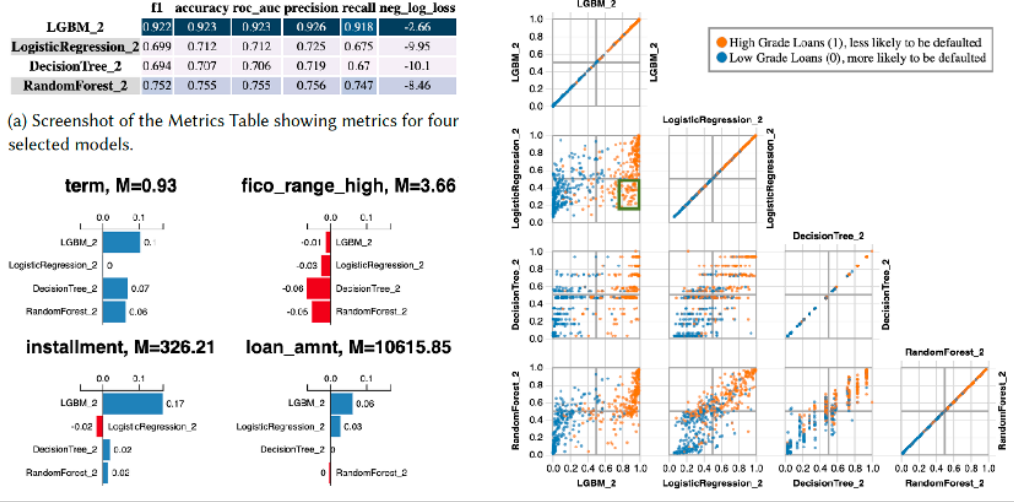


Fig. 1. Model LineUpper, an example XAI tool that supports ML developers to compare multiple candidate models by comparing their feature-importance explanations at multiple levels (by selecting an instance, a region of instances, or viewing all from the right-side Scatterplot Matrix panel) from Narkar et al. [77]

setting, such as granting more attention to the explanations [16]. Moreover, the ability to simulate a model prediction may simply not match the need for a user to perform a realistic task. For instance, in the context of decision-making, the key to success of a human-AI team is appropriate reliance—knowing when to trust the AI’s recommendation and when to be cautious. An actionable understanding for appropriate reliance requires not only knowing how the model makes predictions, but also how to judge if the reasoning is flawed. Filling this gap of understanding may necessitate a different kind of transparency. For example, recent HCI studies repeatedly found that showing uncertainty information of individual predictions is more effective than local explanations to help people achieve the objective of appropriate reliance on AI [14, 105].

To close this gap between technical research of XAI and user experiences will require both studying user interactions with XAI in the contexts of use, and better operationalizing human-centered perspectives in algorithmic work of XAI, including developing evaluation methods that can account for more fine-grained user needs and their context or objective dependent nature.

For the first part, recent HCI research provides useful insights for XAI user experiences in some common usage contexts. For model development or debugging, research suggests that users often need a range of explanations for different levels of the model behaviors to perform comprehensive diagnosis [51, 52, 77]. For example, Figure 1 shows *Model LineUpper* [77], an XAI tool that we designed with our collaborators to support data scientists to compare multiple models (in the context of choosing candidate models generated by AutoML) by comparing their feature-importance explanations. This comparison can happen at different levels: for a global view, for an input region they select on the Scatterplot Matrix on the right, and down to individual instances. For decision-makers, such as users of an AI system supporting medical diagnosis [18, 63, 104], they may need upfront global explanations about the properties of the model during the on-boarding stage, but local explanations particularly when they get unexpected or suspicious model output. When it comes to auditing for model fairness or biases, our work with collaborators compared the effectiveness

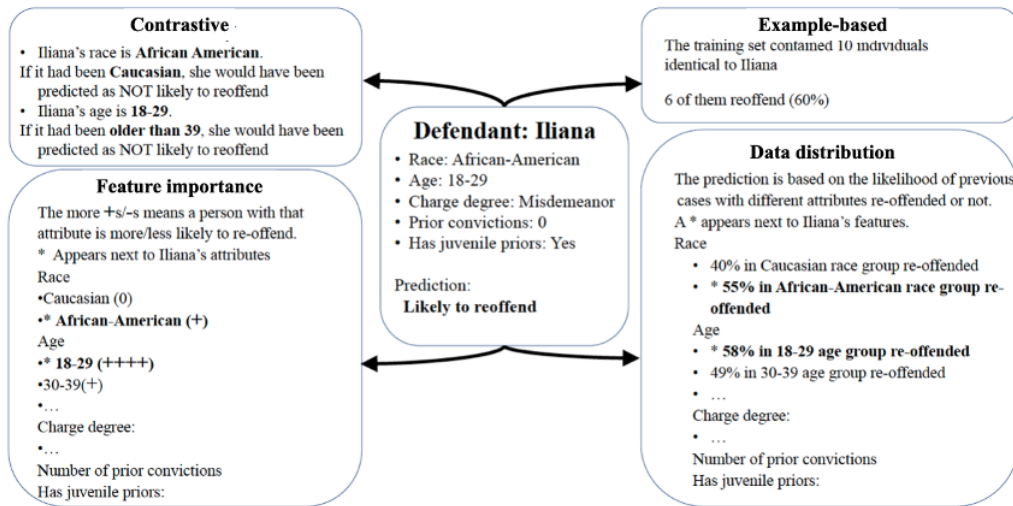


Fig. 2. Four types of XAI features compared in Dodge et al. [29] (with minor updates on the names of explanations from the original paper) to support people's fairness judgment of ML models, with an ML model performing recidivism risk prediction as a use case

of four types of explanation (shown in Figure 2) and found that contrastive explanations can effectively help people identify concerns of individual fairness, where similar individuals are treated differently by the model [29].

By bringing to light how people actually interact with AI explanations, empirical studies can also challenge fundamental assumptions underlying technical approaches to XAI. A pitfall robustly found in recent work is that *explanations can lead to unwarranted trust or confidence in the model*. In a controlled experiment where an ML model was used to assist participants to predict apartment sales prices, Poursabzi-Sangdeh et al. found that, contrary to the hypothesis, showing people an explainable model with feature importance hindered their ability to detect model mistakes [82]. By conducting contextual inquiry with data scientists using popular XAI techniques (e.g. SHAP) during model development, Kauer et al. found that the existence of explanations could mistakenly lead to over-confidence that the model is ready for deployment [54]. In the context of a nutrition recommender, Eiband et al. showed that even placebic explanations, which did not convey useful information, invoked a similar level of trust as real explanations do [36]. In addition, there is the concern of illusory understanding, with which one subjectively over-estimates the understanding they gain from XAI [22]. Explanations can also create information overload and distract people from forming a useful mental model of how a system operates [92].

These observations highlight the danger of deploying technologies without a clear understanding of how people interact with them. One way to move the field forward is to connect with theories and insights about human behaviors and cognition. For example, dual-process theories [53, 81] provide a critical lens to understand how people process XAI and inform new means to make AI more understandable and actionable to users. The central thesis of dual-process theories is that people can engage in two different systems to process information and make decisions. System 1 is intuitive thinking, often following mental shortcuts and heuristics; System 2 is analytical thinking, relying on careful reasoning of information and arguments. Because System 2 is slower and more cognitively demanding, people often resort to System 1 thinking, which, when applied inappropriately, can lead to cognitive biases and sub-optimal decisions. Through this theoretical lens, there is an increasing awareness [16, 32, 79, 85, 99] that while XAI techniques make an

implicit assumption that people can and will attend to every bit of explanations, in reality people are more likely to engage in System 1 thinking.

However, *it remains an open question what kind of heuristics can be triggered by XAI with System 1 thinking*. It is possible that people associate the ability to provide explanations directly with competence, and therefore form unwarranted trust and confidence. Heuristics are developed through past experiences, and can evolve as people experience new technologies or domains. Nourani et al. demonstrated that when interacting with XAI, people were vulnerable to common cognitive biases such as anchoring bias after observing model behaviors early on [79]. A recent study by Ehsan et al. uncovered diverse heuristics people follow in response to AI explanations, such as associating explanations with affirmation, diagnostic support, and social presence, and associating a specific presentation of explanation, such as numerical numbers, with intelligence and algorithmic thinking [32].

Another critical implication with dual-process theories is that people do not equally engage in System 1 or System 2 thinking in all contexts. People are generally inclined to engage in System 1 thinking when they lack either the ability or motivation to perform analytical thinking [81]. This difference can lead to another pitfall of XAI—*potential inequalities of experience* including risks subject to mistrust and misuse of AI. For example, a study found that AI novices, compared to experts, not only had less performance gain from XAI but were also more likely to have illusory satisfaction [94]. Other studies suggest that in time and cognitive resource constraint settings people are less able to process explanations effectively [89, 104]. In our own work with collaborators [39], we showed that adding explanations in an active learning setting (i.e. label instances requested by the model) decreased satisfaction for people scored low in Need for Cognition, a personality trait reflecting one’s general motivation to engage in effortful cognitive activities.

Research has begun to address this mismatch between people’s cognitive processes and current assumptions underlying XAI. One way is to provide interventions to nudge people to engage deeper in System 2 thinking. Bućinca et al. introduced cognitive forcing functions as design interventions for that purpose [17], including asking users to make decisions before seeing the AI’s recommendations, slowing down the process, and letting users choose when to see the AI recommendation. In our own work with collaborators [85], we saw that increasing the time for the users to interact with the ML system mitigated some System 1 biases. Another path is to seek technical and design solutions that reduce the cognitive workload imposed by XAI, by reducing the quantity and improving consumability of information. For example, studies suggest that multi-modalities (text, visual, audio, etc.) can be leveraged to aid attention and understanding of XAI [89, 94]. Progressive disclosure [92], starting with simplified or high-level transparency information and revealing details later or upon user requests, is another effective approach to reduce cognitive workload. Technical approaches that optimize for a balance between explanation accuracy and conciseness have also been explored [5].

We must note that heuristics are an indispensable part of people’s decision-making process. If applied appropriately, they can aid people to make more efficient and optimal decisions. In fact, they may be key to closing the inequality gaps for people with different levels of ability or motivation to process information about AI. For example, we may envision a quality endorsement feature through some authorized third-party inspecting a model with explanations. This could allow lay people to apply a reliable “authority heuristic”. Understanding what heuristics are involved in interactions with XAI and AI in general, and how to leverage reliable heuristics to improve human-AI interaction, are important open questions for the field.

We close this section with an optimistic note that by centering our analysis on people, on how they interact with and process information about AI, and whether they can achieve their objectives, we can move away from a techno-centric focus on generating algorithmic explanations. We can begin to identify opportunities to improve user experiences in the currently under-developed space between algorithmic explanations and actionable understanding, and appreciate

explainable AI as much of a design problem as a technical problem. The design solutions may be concerned with how to communicate algorithmic explanations, such as choosing the right modalities, level of abstraction, work-arounds for privacy or security constraints, and so on. They may also come in the form of interventions to influence how people process XAI, such as providing cognitive forcing functions or checklists that help people better assess information [88].

Furthermore, it is necessary to fill the knowledge or information gaps for users to achieve actionable understanding beyond algorithmic explanations, such as providing necessary domain knowledge (e.g. what a feature means) and general notions of how AI works. In another example, to enable a *socially situated understanding*, with collaborators we proposed the concept of social transparency—making visible the social-organizational factors that govern the use of AI systems [31]. Operationalized in a design framework to present past users’ interactions and reasoning with the AI (e.g., “I am rejecting the AI’s recommendation because this is a long-term profitable customer”), we demonstrated that such information could help users make more informed decisions and improve the collective experience with AI as a sociotechnical system. In short, we must challenge the techno-centric focus on algorithmic explanations, and expand the design space of XAI to support users’ actionable understanding.

5 THEORY-DRIVEN HUMAN-COMPATIBLE XAI

Previously we gave an example of using dual-process theories to retrospectively understand how people interact with XAI. In this section we discuss another important human-centered approach to XAI, by performing theoretical analysis of human explanations, as well as broader cognitive and behavioral processes, to inspire new computational and design frameworks to make XAI more human compatible.

Such work is best represented by Miller’s seminal paper that brings insights from social sciences about fundamental properties of human explanations to common awareness of the AI community [73]. By surveying a large volume of prior work on how people seek, generate, and evaluate explanations in philosophy, psychology, and cognitive science, Miller summarized four major properties of human explanations: 1) Explanations are often contrastive, sought in response to some counterfactual cases. This is because a *Why* question is often triggered by “abnormal or unexpected” events, not asked to understand the cause for an event per se, but the cause of an event relative to some other event that did not occur. In other words, the *Why* question is often an implicit *Why not* question. 2) Explanations are selected, often in a biased manner. People rarely give an actual or complete cause of event, but select a small number of causes based on some criteria or heuristics. 3) Explanations are social, as a transfer of knowledge, often part of a conversation or interaction, and thus presented relative to the explainer’s beliefs about the explainee’s beliefs. 4) Using probabilities or statistical information to explain is often ineffective and unsatisfying. Explicitly referring to causes is more important.

Published in 2019, in just two years, this work has made significant impact on the XAI field. For instance, the point about explanations being contrastive has inspired many to work on counterfactual explanations to answer the *Why Not* or *How to be That* questions, as we reviewed in Section 2. From a user-interaction point of view, the points of explanations being selected and social have profound implications. Miller reviewed several useful theories about how people generate and present explanations to others, which we believe can provide conceptual ground to frame XAI as interaction problems. One of them is Malle’s theory of explanation [72], which breaks the generation of explanations into two distinct and co-influencing groups of psychological processes: 1) Information processes for the explainer (i.e. AI in the case of XAI) to devise explanations, which are determined by what kind of information the explainer has access to. 2) Impression management processes that govern the social interactions with the explainee (i.e., users in the case of XAI), which are driven by the pragmatic goal of the explainer, such as transferring knowledge, generating trust in an explainee, assigning blame, etc.

While currently under-explored, framing communication of explanations as an impression management process can inspire computational and design methods to make XAI effectively selected (and social), which can then mitigate the cognitive load and make XAI more consumable. A useful set of resources to inform XAI work on this topic, as Miller suggested, is to look at the cognitive processes for people to select explanations from available causes. Besides formal models of abductive reasoning, Miller also reviewed common heuristics people follow, such as abnormality (selecting the abnormal cause), intentionality (select intentional actions), necessity, sufficiency, and robustness (selecting causes that would hold in many situations). The choice highly depends on the explainer’s goal, which again highlights the importance of specifying the objective of explaining. Further, we point to broader social science research on impression management [41, 60], on influencing other’s perception by regulating information in social interactions, as well as ethics discussions around it, to draw inspiration from.

The social nature of explanation also maps to an essential requirement for *interactivity* in XAI applications [57]. User interactions do not end at receiving an XAI output, but continue until an actionable understanding is achieved. In other words, as users’ explainability needs are expressed in questions, they will keep asking follow-up questions until satisfied and thus engage in a back-and-forth conversation. Therefore, conversational models of explanation, as well as general principles of conversations and communication (e.g., Grice’s maxims that a speaker follows to optimize for the desired social goal [43]; Theory of grounding in communication [23]), hold promises for informing technologies and design for interactive XAI. Miller reviewed several relevant theories including Hilton’s conversational model of explanations [48], which postulates that a good explanation must be relevant to the focus of a question and present a topology of different causal questions. Antaki and Leudar extended this model to a wider class of argumentative dialogue for the common pattern of claim-backing in explanations [8]. Walton further extended this line of work into a formal dialogue model of explanation [98], including a set of speech act rules. This kind of theories offer appealing grounds to build computational models, and recent XAI has begun to explore dialogue models for interactive XAI [71]. Outside XAI, work in dialogue systems frequently builds on formal models of human conversational and social interactions (e.g. [15]), including systems that generate explanatory dialogues [21].

Theories can also inform design frameworks that guide researchers and practitioners to investigate the design space and make design choices. For example, Wang et al. preformed a comprehensive analysis on the theoretical underpinning of human reasoning and decision-making to derive a conceptual framework that allows linking XAI methods to users’ reasoning needs [99]. This framework includes four dimensions that describe a normative view of how people should reason with explanations, including explanation goals, reasoning process, causal explanation type, and elements in rational choice decisions. It also separately describes people’s natural decision-making, and the errors and limitations they subject to, based on dual-process theories. Designers can use the framework to perform a conceptual analysis to understand, e.g., based on user research, users’ reasoning goals and potential errors, to identify what XAI methods can support their goals, or to investigate gaps in current XAI methods. The authors further provided a mapping between elements under these human-reasoning dimensions and existing XAI approaches, and guidelines on how to use XAI methods to mitigate common cognitive biases.

While hugely promising, theory-driven XAI is still a nascent area. Many areas of cognitive, social, and behavioral theories are yet to be explored. For example, if we center on users of XAI as information seekers to achieve actionable understanding, theories in information science such as models of sense-making [25] and information seeking behaviors [103] (how people’s information needs drive their behaviors and information use) can offer useful theoretical lenses to formalize and anticipate user behaviors. However, the major challenge lies in how to operationalize theoretical

insights and formal behavioral models into computational and design frameworks, which may require, as many have already argued [30, 73, 96], collaboration across the research disciplines of AI, HCI, social sciences and more.

6 SUMMARY

Explainable AI is one of the fastest-growing areas of AI in several directions: a rapidly expanding collection of techniques, substantial industry effort to produce open-source XAI toolkits for practitioners to use, and widespread public awareness and interest in the topic. It is also a fast-growing area for human-centered ML, which can be seen in a proliferation of XAI research published in HCI and social science venues in recent years. Adopting human-centered approaches to XAI is inevitable given that explainability is a human-centric property and XAI must be studied as an interaction problem. However, different from some other topics in this book, HCI work on XAI currently resides in, and often needs to challenge, a techno-centric reality given that the technical AI community have made large strides already. A research community of human-centered XAI [33, 35] has emerged. In this chapter we provide a selected survey on work from this emerging community to encourage future research to continue bridging design practices and state-of-the-art XAI techniques, uncovering pitfalls of and challenging algorithmic assumptions, and building human-compatible XAI from theoretical grounds. We also hope these practiced approaches will inspire work to address broader challenges in human-centered ML.

REFERENCES

- [1] 2017. H2O.ai Machine Learning Interpretability. <https://github.com/h2oai/mli-resources>.
- [2] 2018. Model Interpretation with Skater. <https://oracle.github.io/Skater/>.
- [3] 2019. IBM AIX 360. aix360.mybluemix.net/.
- [4] 2019. Microsoft InterpretML. <https://github.com/interpretml/interpret>.
- [5] Ashraf Abdul, Christian von der Weth, Mohan Kankanhalli, and Brian Y Lim. 2020. COGAM: Measuring and Moderating Cognitive Load in Machine Learning Model Explanations. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–14.
- [6] Amina Adadi and Mohammed Berrada. 2018. Peeking inside the black-box: A survey on Explainable Artificial Intelligence (XAI). *IEEE Access* 6 (2018), 52138–52160.
- [7] David Alvarez-Melis and Tommi S Jaakkola. 2018. Towards robust interpretability with self-explaining neural networks. *arXiv preprint arXiv:1806.07538* (2018).
- [8] Charles Antaki and Ivan Leudar. 1992. Explaining in conversation: Towards an argument model. *European Journal of Social Psychology* 22, 2 (1992), 181–194.
- [9] Daniel W Apley and Jingyu Zhu. 2020. Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 82, 4 (2020), 1059–1086.
- [10] Matthew Arnold, Rachel KE Bellamy, Michael Hind, Stephanie Houde, Sameep Mehta, Aleksandra Mojsilović, Ravi Nair, K Natesan Ramamurthy, Alexandra Olteanu, David Piorkowski, et al. 2019. FactSheets: Increasing trust in AI services through supplier’s declarations of conformity. *IBM Journal of Research and Development* 63, 4/5 (2019), 6–1.
- [11] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bannetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion* 58 (2020), 82–115.
- [12] Vijay Arya, Rachel K. E. Bellamy, Pin-Yu Chen, Amit Dhurandhar, Michael Hind, Samuel C. Hoffman, Stephanie Houde, Q. Vera Liao, Ronny Luss, Aleksandra Mojsilovic, Sami Mourad, Pablo Pedemonte, Ramya Raghavendra, John Richards, Prasanna Sattigeri, Karthikeyan Shanmugam, Moninder Singh, Kush R. Varshney, Dennis Wei, and Yunfeng Zhang. 2020. AI Explainability 360: An Extensible Toolkit for Understanding Data and Machine Learning Models. *J. Mach. Learn. Res.* 21, 130 (2020), 1–6.
- [13] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one* 10, 7 (2015), e0130140.
- [14] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 2021. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–16.
- [15] Timothy Bickmore and Justine Cassell. 2001. Relational agents: a model and implementation of building user trust. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 396–403.

- [16] Zana Bućinca, Phoebe Lin, Krzysztof Z Gajos, and Elena L Glassman. 2020. Proxy tasks and subjective measures can be misleading in evaluating explainable AI systems. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*. 454–464.
- [17] Zana Bućinca, Maja Barbara Malaya, and Krzysztof Z Gajos. 2021. To trust or to think: cognitive forcing functions can reduce overreliance on AI in AI-assisted decision-making. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (2021), 1–21.
- [18] Carrie J Cai, Samantha Winter, David Steiner, Lauren Wilcox, and Michael Terry. 2019. Hello AI: Uncovering the Onboarding Needs of Medical Practitioners for Human-AI Collaborative Decision-Making. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 104.
- [19] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. 2015. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of KDD*.
- [20] Diogo V Carvalho, Eduardo M Pereira, and Jaime S Cardoso. 2019. Machine learning interpretability: A survey on methods and metrics. *Electronics* 8, 8 (2019), 832.
- [21] Alison Cawsey. 1992. *Explanation and interaction: the computer generation of explanatory dialogues*. MIT press.
- [22] Michael Chromik, Malin Eiband, Felicitas Buchner, Adrian Krüger, and Andreas Butz. 2021. I Think I Get Your Point, AI! The Illusion of Explanatory Depth in Explainable AI. In *26th International Conference on Intelligent User Interfaces*. 307–317.
- [23] Herbert H Clark and Susan E Brennan. 1991. Grounding in communication. (1991).
- [24] Mark Craven and Jude Shavlik. 1995. Extracting tree-structured representations of trained networks. *Advances in neural information processing systems* 8 (1995), 24–30.
- [25] Brenda Dervin. 1998. Sense-making theory and practice: An overview of user interests in knowledge seeking and use. *Journal of knowledge management* (1998).
- [26] Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and Payel Das. 2018. Explanations based on the missing: Towards contrastive explanations with pertinent negatives. *arXiv preprint arXiv:1802.07623* (2018).
- [27] Amit Dhurandhar, Karthikeyan Shanmugam, and Ronny Luss. 2020. Enhancing Simple Models by Exploiting What They Already Know. In *International Conference on Machine Learning*. PMLR, 2525–2534.
- [28] Amit Dhurandhar, Karthikeyan Shanmugam, Ronny Luss, and Peder A Olsen. 2018. Improving Simple Models with Confidence Profiles. *Advances in Neural Information Processing Systems* 31 (2018).
- [29] Jonathan Dodge, Q Vera Liao, Yunfeng Zhang, Rachel KE Bellamy, and Casey Dugan. 2019. Explaining models: an empirical study of how explanations impact fairness judgment. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*. 275–285.
- [30] Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* (2017).
- [31] Upol Ehsan, Q Vera Liao, Michael Muller, Mark O Riedl, and Justin D Weisz. 2021. Expanding explainability: Towards social transparency in ai systems. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–19.
- [32] Upol Ehsan, Samir Passi, Q Vera Liao, Larry Chan, I Lee, Michael Muller, Mark O Riedl, et al. 2021. The Who in Explainable AI: How AI Background Shapes Perceptions of AI Explanations. *arXiv preprint arXiv:2107.13509* (2021).
- [33] Upol Ehsan and Mark O Riedl. 2020. Human-centered explainable ai: Towards a reflective sociotechnical approach. In *International Conference on Human-Computer Interaction*. Springer, 449–466.
- [34] Upol Ehsan, Pradyumna Tambwekar, Larry Chan, Brent Harrison, and Mark O Riedl. 2019. Automated rationale generation: a technique for explainable AI and its effects on human perceptions. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*. 263–274.
- [35] Upol Ehsan, Philipp Wintersberger, Q Vera Liao, Martina Mara, Marc Streit, Sandra Wachter, Andreas Riener, and Mark O Riedl. 2021. Operationalizing Human-Centered Perspectives in Explainable AI. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–6.
- [36] Malin Eiband, Daniel Buschek, Alexander Kremer, and Heinrich Hussmann. 2019. The impact of placebic explanations on trust in intelligent systems. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–6.
- [37] Malin Eiband, Hanna Schneider, Mark Bilandzic, Julian Fazekas-Con, Mareike Haug, and Heinrich Hussmann. 2018. Bringing transparency design into practice. In *23rd international conference on intelligent user interfaces*. 211–223.
- [38] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2018. Datasheets for datasets. *arXiv preprint arXiv:1803.09010* (2018).
- [39] Bhavya Ghai, Q Vera Liao, Yunfeng Zhang, Rachel Bellamy, and Klaus Mueller. 2020. Explainable Active Learning (XAL): An Empirical Study of How Local Explanations Impact Annotator Experience. *arXiv preprint arXiv:2001.09219* (2020).
- [40] Soumya Ghosh, Q Vera Liao, Karthikeyan Natesan Ramamurthy, Jiri Navratil, Prasanna Sattigeri, Kush R Varshney, and Yunfeng Zhang. 2021. Uncertainty Quantification 360: A Holistic Toolkit for Quantifying and Communicating the Uncertainty of AI. *arXiv preprint arXiv:2106.01410* (2021).
- [41] Erving Goffman et al. 1978. *The presentation of self in everyday life*. Vol. 21. Harmondsworth London.
- [42] Alex Goldstein, Adam Kapelner, Justin Bleich, and Emil Pitkin. 2015. Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics* 24, 1 (2015), 44–65.
- [43] Herbert P Grice. 1975. Logic and conversation. In *Speech acts*. Brill, 41–58.
- [44] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2019. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)* 51, 5 (2019), 93.

- [45] Karthik S Gurumoorthy, Amit Dhurandhar, Guillermo Cecchi, and Charu Aggarwal. 2019. Efficient data representation by selecting prototypes with importance weights. In *2019 IEEE International Conference on Data Mining (ICDM)*. IEEE, 260–269.
- [46] Peter Hase and Mohit Bansal. 2020. Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior?. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 5540–5552.
- [47] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. 2009. The elements of statistical learnin. *Cited on* (2009), 33.
- [48] Denis J Hilton. 1990. Conversational processes and causal explanation. *Psychological Bulletin* 107, 1 (1990), 65.
- [49] Michael Hind. 2019. Explaining explainable AI. *XRDS: Crossroads, The ACM Magazine for Students* 25, 3 (2019), 16–19.
- [50] Michael Hind, Dennis Wei, Murray Campbell, Noel CF Codella, Amit Dhurandhar, Aleksandra Mojsilović, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. 2019. TED: Teaching AI to explain its decisions. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 123–129.
- [51] Fred Hohman, Andrew Head, Rich Caruana, Robert DeLine, and Steven M Drucker. 2019. Gamut: A design probe to understand how data scientists understand machine learning models. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–13.
- [52] Sungsoo Ray Hong, Jessica Hullman, and Enrico Bertini. 2020. Human factors in model interpretability: Industry practices, challenges, and needs. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW1 (2020), 1–26.
- [53] Daniel Kahneman. 2011. *Thinking, fast and slow*. Macmillan.
- [54] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. 2020. Interpreting Interpretability: Understanding Data Scientists’ Use of Interpretability Tools for Machine Learning. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–14.
- [55] Been Kim, Rajiv Khanna, and Oluwasanmi O Koyejo. 2016. Examples are not enough, learn to criticize! criticism for interpretability. In *Proceedings of NIPS*.
- [56] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. 2018. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*. PMLR, 2668–2677.
- [57] Josua Krause, Adam Perer, and Kenney Ng. 2016. Interacting with predictions: Visual inspection of black-box machine learning models. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. ACM, 5686–5697.
- [58] Himabindu Lakkaraju, Stephen H Bach, and Jure Leskovec. 2016. Interpretable decision sets: A joint framework for description and prediction. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 1675–1684.
- [59] Himabindu Lakkaraju, Ece Kamar, Rich Caruana, and Jure Leskovec. 2017. Interpretable & explorable approximations of black box models. *arXiv preprint arXiv:1707.01154* (2017).
- [60] Mark R Leary and Robin M Kowalski. 1990. Impression management: A literature review and two-component model. *Psychological bulletin* 107, 1 (1990), 34.
- [61] Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. 2018. Distribution-free predictive inference for regression. *J. Amer. Statist. Assoc.* 113, 523 (2018), 1094–1111.
- [62] Jiwei Li, Will Monroe, and Dan Jurafsky. 2016. Understanding neural networks through representation erasure. *arXiv preprint arXiv:1612.08220* (2016).
- [63] Q Vera Liao, Daniel Gruen, and Sarah Miller. 2020. Questioning the AI: informing design practices for explainable AI user experiences. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [64] Q Vera Liao, Milena Pribić, Jaesik Han, Sarah Miller, and Daby Sow. 2021. Question-Driven Design Process for Explainable AI User Experiences. *arXiv preprint arXiv:2104.03483* (2021).
- [65] Brian Y Lim and Anind K Dey. 2009. Assessing demand for intelligibility in context-aware applications. In *Proceedings of the 11th international conference on Ubiquitous computing*. 195–204.
- [66] Zachary C Lipton. 2018. The mythos of model interpretability. *Queue* 16, 3 (2018), 31–57.
- [67] Arnaud Van Looveren and Janis Klaise. 2021. Interpretable counterfactual explanations guided by prototypes. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 650–665.
- [68] Ana Lucic, Hinda Haned, and Maarten de Rijke. 2020. Why does my model fail? contrastive local explanations for retail forecasting. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 90–98.
- [69] Scott M Lundberg, Gabriel Erion, Hugh Chen, Alex DeGrave, Jordan M Prutkin, Bala Nair, Ronit Katz, Jonathan Himmelfarb, Nisha Bansal, and Su-In Lee. 2020. From local explanations to global understanding with explainable AI for trees. *Nature machine intelligence* 2, 1 (2020), 56–67.
- [70] Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Proceedings of the 31st international conference on neural information processing systems*. 4768–4777.
- [71] Prashan Madumal, Tim Miller, Liz Sonenberg, and Frank Vetere. 2019. A grounded interaction protocol for explainable artificial intelligence. *arXiv preprint arXiv:1903.02409* (2019).
- [72] Bertram F Malle. 2006. *How the mind explains behavior: Folk explanations, meaning, and social interaction*. Mit Press.
- [73] Tim Miller. 2018. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence* (2018).
- [74] Tim Miller, Piers Howe, and Liz Sonenberg. 2017. Explainable AI: Beware of inmates running the asylum or: How I learnt to stop worrying and love the social and behavioural sciences. *arXiv preprint arXiv:1712.00547* (2017).

- [75] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*. 220–229.
- [76] Ramaravind Kommiya Mothilal, Amit Sharma, and Chenhao Tan. 2019. Explaining machine learning classifiers through diverse counterfactual explanations. *arXiv preprint arXiv:1905.07697* (2019).
- [77] Shweta Narkar, Yunfeng Zhang, Q Vera Liao, Dakuo Wang, and Justin D Weisz. 2021. Model LineUpper: Supporting Interactive Model Comparison at Multiple Levels for AutoML. In *26th International Conference on Intelligent User Interfaces*. 170–174.
- [78] Don Norman. 2013. *The design of everyday things: Revised and expanded edition*. Basic books.
- [79] Mahsan Nourani, Chiradeep Roy, Jeremy E Block, Donald R Honeycutt, Tahrira Rahman, Eric Ragan, and Vibhav Gogate. 2021. Anchoring Bias Affects Mental Model Formation and User Reliance in Explainable AI Systems. In *26th International Conference on Intelligent User Interfaces*. 340–350.
- [80] Andrés Páez. 2019. The pragmatic turn in explainable artificial intelligence (XAI). *Minds and Machines* 29, 3 (2019), 441–459.
- [81] Richard E Petty and John T Cacioppo. 1986. The elaboration likelihood model of persuasion. In *Communication and persuasion*. Springer, 1–24.
- [82] Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Wortman Vaughan, and Hanna Wallach. 2021. Manipulating and measuring model interpretability. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–52.
- [83] Alun Preece, Dan Harborne, Dave Braines, Richard Tomsett, and Supriyo Chakraborty. 2018. Stakeholders in explainable AI. *arXiv preprint arXiv:1810.00184* (2018).
- [84] Isha Puri, Amit Dhurandhar, Tejaswini Pedapati, Karthikeyan Shanmugam, Dennis Wei, and Kush R. Varshney. 2021. CoFrNets: Interpretable neural architecture inspired by continued fractions. In *Advances in neural information processing systems*.
- [85] Charvi Rastogi, Yunfeng Zhang, Dennis Wei, Kush R. Varshney, Amit Dhurandhar, and Richard Tomsett. 2022. Deciding Fast and Slow: The Role of Cognitive Biases in AI-Assisted Decision-Making. In *Proceedings of the ACM Conference on Computer Supported Cooperative Work and Social Computing*.
- [86] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of KDD*.
- [87] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Anchors: High-precision model-agnostic explanations. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [88] Soo Young Rieh and David R Danielson. 2007. Credibility: A multidisciplinary framework. *Annual review of information science and technology* 41, 1 (2007), 307–364.
- [89] Justus Robertson, Athanasios Vasileios Kokkinakis, Jonathan Hook, Ben Kirman, Florian Block, Marian F Ursu, Sagarika Patra, Simon Demediuk, Anders Drachen, and Oluseyi Olarewaju. 2021. Wait, But Why?: Assessing Behavior Explanation Strategies for Real-Time Strategy Games. In *26th International Conference on Intelligent User Interfaces*. 32–42.
- [90] Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence* 1, 5 (2019), 206–215.
- [91] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*. 618–626.
- [92] Aaron Springer and Steve Whittaker. 2019. Progressive disclosure: empirically motivated approaches to designing effective transparency. In *Proceedings of the 24th international conference on intelligent user interfaces*. 107–120.
- [93] Harini Suresh, Steven R Gomez, Kevin K Nam, and Arvind Satyanarayan. 2021. Beyond Expertise and Roles: A Framework to Characterize the Stakeholders of Interpretable Machine Learning and their Needs. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–16.
- [94] Maxwell Szymanski, Martijn Millecamp, and Katrien Verbert. 2021. Visual, textual or hybrid: the effect of user expertise on different explanations. In *26th International Conference on Intelligent User Interfaces*. 109–119.
- [95] Sarah Tan, Rich Caruana, Giles Hooker, Paul Koch, and Albert Gordo. 2018. Learning global additive explanations for neural nets using model distillation. *arXiv preprint arXiv:1801.08640* (2018).
- [96] Jennifer Wortman Vaughan and Hanna Wallach. 2020. A human-centered agenda for intelligible machine learning. *Machines We Trust: Getting Along with Artificial Intelligence* (2020).
- [97] Sandra Wachter, Brent Mittelstadt, and Chris Russell. 2017. Counterfactual explanations without opening the black box: Automated decisions and the GDPR.
- [98] Douglas Walton. 2004. A new dialectical theory of explanation. *Philosophical Explorations* 7, 1 (2004), 71–89.
- [99] Danding Wang, Qian Yang, Ashraf Abdul, and Brian Y Lim. 2019. Designing Theory-Driven User-Centric Explainable AI. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, 601.
- [100] Xinru Wang and Ming Yin. 2021. Are Explanations Helpful? A Comparative Study of the Effects of Explanations in AI-Assisted Decision-Making. In *26th International Conference on Intelligent User Interfaces*. 318–328.
- [101] Dennis Wei, Sanjeeb Dash, Tian Gao, and Oktay Gunluk. 2019. Generalized linear rule models. In *International Conference on Machine Learning*. PMLR, 6687–6696.
- [102] Pengfei Wei, Zhenzhou Lu, and Jingwen Song. 2015. Variable importance analysis: a comprehensive review. *Reliability Engineering & System Safety* 142 (2015), 399–432.

- [103] Tom D Wilson. 1981. On user studies and information needs. *Journal of documentation* (1981).
- [104] Yao Xie, Melody Chen, David Kao, Ge Gao, and Xiang ‘Anthony’ Chen. 2020. CheXplain: Enabling Physicians to Explore and Understand Data-Driven, AI-Enabled Medical Imaging Analysis. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–13.
- [105] Yunfeng Zhang, Q Vera Liao, and Rachel KE Bellamy. 2020. Effect of confidence and explanation on accuracy and trust calibration in AI-assisted decision making. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 295–305.