# Towards a Comprehensive Human-Centred Evaluation Framework for Explainable AI

Ivania Donoso-Guzmán[1,2][0000−0002−2427−9128], Jeroen Ooge[1][0000−0001−9820−7656], Denis Parra[2][0000−0001−9878−8761], and Katrien Verbert[1][0000−0001−6699−7710]

[1] KU Leuven, Department of Computer Science
[2] Pontificia Universidad Católica de Chile

**Abstract.** While research on explainable AI (XAI) is booming and explanation techniques have proven promising in many application domains, standardised human-centred evaluation procedures are still missing. In addition, current evaluation procedures do not assess XAI methods holistically in the sense that they do not treat explanations' effects on humans as a complex user experience. To tackle this challenge, we propose to adapt the User-Centric Evaluation Framework used in recommender systems: we integrate explanation aspects, summarise explanation properties, indicate relations between them, and categorise metrics that measure these properties. With this comprehensive evaluation framework, we hope to contribute to the human-centred standardisation of XAI evaluation.

**Keywords:** XAI Evaluation · Human-centred evaluation · Evaluation framework

## 1 Introduction

*Explainable AI* (XAI) is advancing fast: between 2017 and 2021 alone, the number of XAI papers increased eight-fold [39] and researchers have proposed XAI methods for virtually all existing media types and families of AI models. However, it is still unclear to what extent explanations are effective in practice [34] because full-fledged standardised evaluation procedures are missing. This is partly due to lacking consensus on which explanation properties should be assessed and which measurements should be used [8, 34, 35, 39, 49].

To better assess XAI methods, researchers have tried to disentangle explanation's characteristics into simple measurable properties such as completeness [5, 39, 49], novelty [8, 30, 32, 43], and interactivity [21, 39, 49]. However, there is little evidence on how these properties relate to explanations being appropriate in real scenarios [29]. In addition, while many researchers stress the importance of context, we are unaware of XAI evaluation methods that treat explanations' effects on humans as a **complex user experience** involving factors such as user perception and system interaction.

To evaluate explanations holistically, we are working towards a human-centred evaluation framework for XAI, which extends pioneering work on developing and

evaluating user experience [25] and explanations [46] for recommender systems. We categorise explanation properties according to this framework and indicate their relations reported in the literature. Additionally, we present the *explanation elements* that help to classify metrics to simplify the choice of measurements. This adapted user-centric framework will allow researchers and practitioners to evaluate explanations of AI-based systems and potentially increase deployment of such systems in their respective domains [34, 36].

The contributions of this paper are three-fold: first, we present an extensive analysis of existing definitions of explanation properties and methods, as well as their interrelationships. Our analysis aligns different properties and methods as defined by different research communities. Second, based on this analysis, we define a human-centred evaluation framework for XAI that presents an integrative approach and combines user-centric evaluation and functional metrics. Third, we present an example of the use of this framework.

## 2    Background and Related work

### 2.1    Human-Centred Explainable AI

The XAI area of research has been led mostly by the AI community, even though it is a multidisciplinary area of research. For this reason, XAI methods have been criticised for being developed with the AI researchers' intuition of what constitutes a good explanation [35]. In particular, the design and evaluation of XAI methods are often conducted without considering the final users' needs and their cognitive processes [29].

More recently, the HCI community started proposing ideas for tackling the XAI design, considering how the users reason about explanations: Wang et al. [51] proposed a framework to design explanations based on how humans reason; Chen et al. [10] characterised how explanations affect human understanding of task decision boundary, model decision boundary and model error; Most recently, Chen et al. [12] conducted a study to investigate the decision-making process users follow when faced with AI predictions and their explanations.

Another line of work has been understanding the wants and needs of different shareholders and ensuring they are considered in the design. Mohseni et al. [36] categorised the goals of target user groups and developed design guidelines to iteratively design and evaluate Explainable AI systems; Suresh et al. [44] proposed a framework to characterise users with two multidimensional criteria: knowledge and interpretability needs, that together help to understand the system's users; Langer et al. [26] review the main types of users of XAI systems and their wants and needs, to propose a model for designing XAI systems according to these desiderata; Liao et al. [28] proposed a question-driven design process to fulfil the Explainable AI user's needs; Rong et al. [41] analysed human-based XAI evaluations and provided guidelines for conducting user studies in the area.

Overall, these studies have emphasised the importance of users' characteristics and the tasks they perform during the design phase of XAI experiences. Although

it has been stated as an important aspect of the final adoption of XAI systems [36, 41], to the best of our knowledge, evaluation procedures that capture the complexity of the human-AI interaction have not yet been proposed. We contribute by adapting a widely accepted procedure in recommender systems to evaluate explanations generated by XAI methods holistically.

## 2.2 Evaluating Explanations

Even though AI/ML models have standard evaluation metrics, there is still no consensus on the strategy to evaluate XAI methods. Doshi-Velez et al. [18] proposed the first standardisation of XAI evaluation. According to their work, the evaluation could be performed in three levels: application-grounded, with real tasks and users; human-grounded, with real users and proxy tasks; and functionality-grounded, with proxy tasks and no users. Currently, application or human-grounded approaches have been criticized for their lack of rigour [22, 23], and for using proxy tasks [6].

To conduct functionally-grounded evaluations, i.e. proxy tasks and no users, some studies have focused on grouping concepts and defining properties [5, 8, 34, 49] and their corresponding metrics [39]. These works aggregate existing literature that defines properties or presents metrics to assess them. The proposed properties try to measure the quality of the explanations without context so that they can be used in functionality-grounded evaluation. Similarly, Hoffman et al. [20] proposed to evaluate explanations using the 'goodness criteria' that assess the explanation quality without context. Most recently, Agarwal et al. [2] presented a framework to benchmark different XAI methods using automatic metrics. Still, it is limited to particular methods and only works with specific datasets created for the benchmark.

Little work has been conducted to present the connections between these properties. Most papers state that trade-offs exist [8, 30, 34, 39], but they have not quantified them. To the best of our knowledge, only the study by Balog et al. [4] uncovered conflicting relationships between some of the proposed properties, but they did not evaluate XAI-generated explanations.
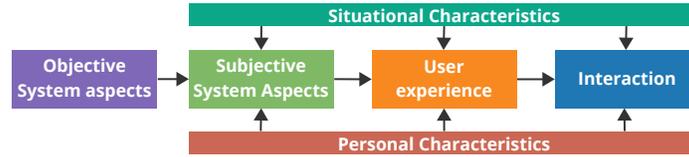
Given the number of properties to evaluate, selecting the aspects to consider in the evaluation is becoming an important topic. According to Liao et al. [30], this selection depends on the tasks the system has to support because the user accomplishment of these tasks determines the overall system's success. Knijnenburg et al. [25] indicate the selection is made according to theoretical models, i.e., it results from previous studies or from the hypothesis that is tested. Recently, Liao et al. [30] presented a study that connects tasks with evaluation criteria to provide general guidelines for the field. In this study, experts and end-users selected the most appropriate properties to evaluate diverse XAI tasks. They found that XAI tasks obtained different property rankings regardless of the application domain (loan application, medical diagnosis, among others).

Our work builds upon these previous studies by proposing a unified framework that integrates previously proposed definitions and measurements by making the relations between them explicit and grounded in previous work. Additionally,

we analysed measurement procedures and classified them by which explanation element they measure according to Miller's definition of explanation [35], which declares that explanations are composed of a cognitive process, a product and a social process. This new criteria to classify measurements provides researchers and practitioners with a new understanding of how to measure properties of explanations.

### 2.3   User Centric Evaluation of Recommender systems

The User Centric Evaluation Framework for recommender systems in Figure 1 was proposed by Knijnenburg et al. [25] to explain how users experience the interaction with a recommendation system and to predict how users behave under similar circumstances. This framework has six *conceptual components* encompassing different *constructs* that can be measured during a user study. For example, the conceptual component *Subjective system aspects* groups constructs such as *Perceived recommendation quality* or *Interaction adequacy*, while *User experience* contains *Choice difficulty* and *Choice satisfaction* among others. The constructs and the causal relations between them found with Structural Equation Modelling (SEM) [24] help explain how different aspects of the experience affect each other and influence the outcomes.



**Fig. 1.** The User-Centric Evaluation Framework by Knijnenburg et al. [25]. Each box represents a *conceptual component* that groups related *constructs*.

This evaluation framework has been used and appreciated in recommender systems because of its capacity to provide relations between different user experience aspects. By capturing the causal relations between different measurements, researchers can not only report and compare these measurements but also explain why differences do or do not occur. This provides a better understanding of what makes a system more adapted to the users and, ultimately, predicts whether it will be successful and why.

In this work, we expand this successful framework for XAI evaluation. We believe our comprehensive work sheds light on which explanation aspects are more important and relevant to users and their circumstances. Furthermore, since the framework provides causal relations between different properties, we believe it can provide better guidelines for XAI design.

# 3   Methods

To adapt the user-centric evaluation framework by Knijnenburg et al. [25], we analysed current literature on the topic with a grounded theory approach. This section describes how we collected papers and categorised them along two axes (conceptual components and explanation elements), to build the foundation for our XAI framework.

## 3.1   Paper Collection

Finding relevant literature on XAI evaluation requires searching several research disciplines. Evaluation, in particular, has been published in several types of venues (workshops, posters, surveys), presenting concrete methods and execution procedures but also proposals and blue-sky ideas. To include as much relevant literature as possible, we consulted Google Scholar with this query:

```
intitle:properties OR intitle:evaluation OR intitle:metrics
OR intitle:property OR intitle:metric
("explainable" OR "interpretable")
("artificial intelligence" OR "machine learning" OR XAI OR AI)
```

The search was conducted at the end of October 2022 and was limited to the years 2017 and onwards because Doshi-Velez et al. [18] then proposed one of the first XAI evaluation procedures. This query returned approximately 5970 results. As a first step, only the titles were reviewed to check whether the result was related to AI or XAI. We checked all result pages until the first page where no papers related to XAI or AI appeared. This occurred on page 25, similar to the results of Vilone et al. [49]. This first screening yielded 80 research works.

These works were analysed by looking at the abstract and, in doubt, at the full paper. The aim of this second screening was to remove duplicate works and keep only works that describe properties, relations between them and measurements. The exclusion criteria were the following:

- The research did not use or propose properties or measurements for XAI explanations.
- The study considered only non-XAI-generated explanations.
- The research compared different XAI methods using different metrics, but said metrics were not grounded on explanation quality aspects.
- The evaluation of the explanations was performed with a ground truth explanation.
- The search result was a master's or PhD thesis, and one or more papers were already published based on the same research, making it redundant.

After this screening process, only 19 results were kept. From their references, other related papers were found. We also included [46] because it is a comprehensive review of the evaluation of explanations in the context of recommender systems. The final number of papers included was 29.

### 3.2   Classification Axis 1: Conceptual Components

A Grounded Theory [9] approach was followed to analyse the collected works in three steps: Initial Coding, aimed at finding quotes that related to properties of explanation; Focused Coding, which consisted of labelling the passages according to a set of concepts; and finally Axial Coding, which connects and groups the different concepts.

The Initial Coding step was conducted in-vivo. Definitions of explanation properties, definitions of metrics to measure aspects of explanations, and relations between properties were searched for. Some of the papers had definitions of properties based on multiple previous works. In those cases, we kept the summarised definition and did not look for primary sources. In contrast, if the definition made in the survey paper did not fully explain metrics, we added the primary source to the group of papers.
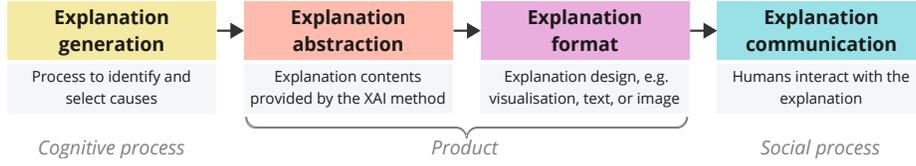
The Focused Coding Step consisted of labelling the different definitions with the most appropriate concept, independently of the name the authors had coined. This iterative process aimed to group the definitions that point to the same desiderata of an explanation while avoiding overlapping concepts. The definition of each property was created at this step. In addition, passages that described a procedure to measure the property were marked as such. The procedure to analyse those quotations is described in Section 3.3.

The Axial Coding phase was conducted by first collecting the relations that were described in the selected papers. After these relations were captured, new relations that emerged from the definitions were investigated and added to the model. Additionally, relations were added based on evidence of other papers the researchers were aware of.

Finally, each of the found properties was matched to a conceptual component as defined in Knijnenburg's framework [25]. Our analysis yielded very few and general properties for the situational and personal characteristics components, so it was decided to leave those properties out of the current analysis. During this phase, it was noted that some properties belonged to a new category that captured the abstract quality of the explanation. This idea aligns with the nature of XAI methods: the original framework was made for recommender systems, i.e., an AI model that selects objects, but XAI methods **generate** an object. To evaluate the quality of generated objects, it was decided to add the conceptual component *Explanation Aspects* (see Figure 3), which groups properties that evaluate the explanation quality.

### 3.3   Classification Axis 2: Explanation Elements

Previous analysis of properties had classified measurement and metrics depending on their user dependency [5], the nature of the procedure (objective, subjective)[16, 21] or according to umbrella properties [38, 39]. However, during the analysis of the conceptual components and the properties of explanations, it was found that similar properties are often named differently because of the ways in which they are measured. For example, Carvalho et al. [8] defined two similar concepts that were

| Explanation generation | Explanation abstraction | Explanation format | Explanation communication |
|---|---|---|---|
| Process to identify and select causes | Explanation contents provided by the XAI method | Explanation design, e.g. visualisation, text, or image | Humans interact with the explanation |

*Cognitive process* — *Product* — *Social process*

**Fig. 2.** The four elements of explanations. We use the same ideas as [35] but changed the names of the elements.Additionally, we further divide the explanation product into *abstraction* and *format*.

applied in two types of evaluation. They used the name *Representativeness* for the evaluation without users and the concept *General and probable* for evaluation with user studies, even though both refer to the number of instances that can be explained with the same causes. We argue this inconsistency occurs because explanations are made of different elements. Miller [35] states that *explanations* are both processes and products: the *Cognitive process* selects a subset of the causes; the *Product* is the resulting outcome; and the *Social process* consists of transferring the knowledge from explainer to explainee.

   With these ideas in mind, a focused coding was conducted only of the passages marked as describing a procedure to measure a property. Each passage was labelled as *generation*, *product* or *communication*. It was found that many metrics that were labelled *product* were very format dependent: for example, BLEU (BiLingual Evaluation Understudy)[14], which evaluates machine-translation quality, cannot be applied to visual-based explanations, but Covariate Homogeneity [39] could be applied to both text and visual-based explanations. For this reason, the metrics under the *product* label were further categorised between *abstraction* and *format*. Figure 2 displays the new definitions and the relation to Miller's definitions.
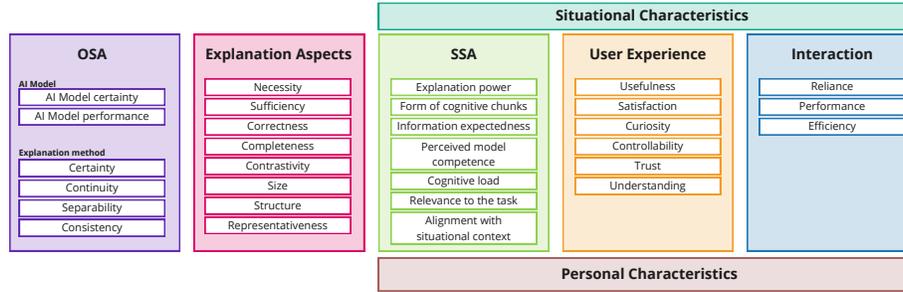
   This categorisation allows classifying measurement procedures under three criteria: property they measure, element of explanation and type of procedure (questionnaire, metrics, etc). Different measurements can be applied to evaluate the properties along the four explanation elements. Some properties can only be assessed by measuring one element, while others can be measured in more than one. These new criteria are explained and justified in Section 4.3.

## 4   A User-Centric Evaluation Framework for XAI

This section presents an adapted version of the *User-Centric Evaluation Framework*. To describe it, we use the following terminology: **conceptual components** group **explanation properties**, which in turn can be measured with **measurements**. While each measurement applies to only one **explanation element**, a single property can be measured by several measurements.

   This section is organised as follows: in Section 4.1, the choice of properties for each conceptual component is justified and explained, and the properties are defined; then, in Section 4.2, the connections between properties are presented;

finally in Section 4.3 the classification criteria for measurements is presented and justified, as well as the existing measurements for each property.



**Fig. 3.** The User-Centric Evaluation Framework by Knijnenburg et al. [25] extended with a new conceptual component: *Explanation Aspects*. Each conceptual component displays its properties. The box of Objective system aspects (OSA) marks the properties that apply to the AI model and the ones that apply to the XAI method.

### 4.1   Explanation Properties

**Objective system aspects.** Objective systems aspects (OSAs) are 'the aspects of the system that are currently being evaluated' [25]. It was found from the analysis that characteristics from the particular instance of the XAI method and AI model can affect the explanation. For instance, the AI model performance will affect the level of Trust users can achieve. Making these characteristics explicit in the framework can help to understand the specific aspects of the XAI method and AI model that affect the user experience.

The analysis yielded six properties: AI model performance, AI model certainty, Certainty, Continuity, Separability and Consistency. The first two properties measure the AI model, and the last four are applied to the XAI method. Continuity was described in several works as the desired 'smoothness' of the XAI function. In the beginning, Separability and Continuity were one concept, but it was noted that providing similar explanations to similar instances does not guarantee that different instances will get different explanations. Consistency evaluates the randomness of the XAI method: if different runs of the XAI method algorithm return different functions, the model will be highly inconsistent.

AI model certainty and XAI method certainty were complicated properties. Uncertainty quantification is a very active field of research within AI, and several approximation methods have been proposed. However, the problem is still being investigated due to its high computational cost [1]. Papers' definitions for these concepts emphasised the fact that the models needed to tell the users when to trust their outputs. For this reason, we decided to keep them, even though there are no proven ways to compute them yet.

**Explanation Aspects.** The Explanation aspects component was added to the original framework (see Section 3.2). This component groups the properties that measure the quality of the generated explanation. These concepts have been generally associated with Functionality-Grounded evaluation because these properties can be measured with metrics at the abstraction level, that is, without the need for users.
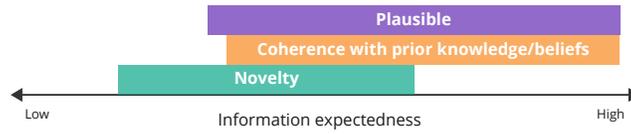
From the analysis, eight properties were found. Necessity, Sufficiency and Contrastivity specifically measure the quality of the selected causes. Their goal is to evaluate whether the reasons the XAI method is providing clearly inform the prediction that was made. Correctness and Completeness are analogous to precision and recall in AI performance metrics. Correctness describes whether the XAI method selected the causes that the AI model used to make a prediction. For explanations generated using the AI model parameters, such as linear regression, the correctness will always be high. Explanations generated by surrogate models will have lower correctness. Completeness quantifies if all the causes that the model used to generate the prediction are present in the explanation. Representativeness determines whether the explanations are unique to each instance or they generalise over multiple instances. This property helps to estimate the Cognitive Load the users will face when using the system. Size and Structure evaluate the explanations' length and organisation, which affects how easy it will be for users to understand them.

**Subjective System Aspects.** Subjective System Aspects (SSA) are "users' perceptions of the Objective System Aspects" [25]. These properties provide evidence that the users perceive the Objective System Aspects. In this modified framework, they help to establish whether the users perceive the OSAs and the Explanation Aspects. Additionally, this component helps us to understand the pertinence of the generated explanations to the users' situational context. These properties are mostly measured at the communication level, but some of them have measures at the abstraction and format level that can be used as proxies of the real value.

The analysis yielded seven properties for this component. Explanation power measures the perceived quality of the selected causes. Explanations with high power provide valuable justifications for the AI model behaviour. Form of cognitive chunks estimates the semantics of the information provided by the explanation. This concept was coined by Doshi-Velez et al. [18] and it has been widely used in the XAI domain. Information expectedness measures whether the explanation provides new knowledge to the user. The analysed works used three concepts for this notion: plausibility, coherence with prior knowledge/beliefs, and novelty. We decided to keep these notions under one umbrella term because we found that they are part of the same scale (see Figure 4). The relation of each concept with information expectedness is the following:

- Plausibility [5, 8, 38]: if the information is expected, the user will think it is plausible. However, the contrary does not necessarily holds. The information can be new but still plausible in the user's mind.

– Coherence with prior knowledge/beliefs [8, 39, 43]: the information provided by the explanation should have some level of connection to the user's background. If that relation does not exist, it will be hard for the user to understand the explanation.
– Novelty [8, 30, 32, 35, 43]: explanations should focus on abnormal causes [35] and provide information the user does not expect to increase her engagement with the system. However, if the reasons are too unexpected, the user will probably dismiss them and ignore the system.



**Fig. 4.** Relation between plausibility, coherence with prior knowledge/beliefs and Novelty with Information Expectedness. Each bar represents the amount of new information that according to the concepts relates to user acceptance.

Perceived model competence evaluates whether the user thinks the AI model can perform as expected. The Cognitive load measures the cognitive effort the user makes to understand the explanations.

The last two properties measure the fit between the explanation and the situational context. Relevance to the task measures whether the explanation provides insights that help to perform the task better. An explanation has to be relevant to be useful for the task the user has to perform; otherwise, she will not exploit it. For example, in a medical context, this would measure whether the explanations are actionable in the patient's state. Alignment with situational context evaluates whether the provided explanation is appropriate for the usage context. For instance, a complex visualisation cannot be used correctly in a time-constrained context.

Table 1: Table of all explanation properties and their definitions based on the reviewed literature.

| Property | Definition | References |
|---|---|---|
| *Objective system aspects* | | |
| AI Model performance | The accomplishment level the AI model has with respect to the task for which it was trained. | |
| AI Model certainty | The confidence the AI model has in its prediction. | [8, 30, 39, 49] |
| Certainty | The confidence the XAI method has in the explanation. | [5, 8, 20, 30, 39, 46] |
| Continuity | The function should provide similar explanations for similar instances. | [8, 16, 20, 21, 34] |

Table 1 – continued from previous page

| Property | Definition | References |
|---|---|---|
| Separability | The XAI method should return different explanations for different instances. | [8] |
| Consistency | The degree to which different runs of the XAI method yield similar XAI functions. | [5, 8, 21, 30, 32, 39, 47, 49] |

*Explanation aspects*

| Property | Definition | References |
|---|---|---|
| Necessity | Measures whether the explanation method selected the causes that are responsible for the prediction. If the necessary causes change, then the prediction will also change. | [8, 30, 35, 39] |
| Sufficiency | Measures whether the explanation method did not select causes that do not affect the prediction. If non-selected causes change, the prediction would still hold, and thus the explanation should not change. | [35] |
| Correctness | Quantifies the extent to which the selected causes are correct with respect to the model reasoning | [8, 30, 34, 39, 43, 49] |
| Completeness | Quantifies if all the causes that the model used to generate the prediction are present in the explanation. | [8, 30, 34, 39, 49] |
| Contrastivity | Measures whether the explanation contains reasons that highlight differences with respect to other possible outcomes. Low contrastivity will provide the same reasons for instances in which the model predicts different classes. | [8, 35, 39] |
| Size | Refers to the amount of information present in the explanation. | [8, 18, 30, 34, 39, 49] |
| Structure | The information should be displayed in a way that allows the users to understand the hierarchy of the information quickly | [8, 21, 34, 39, 43, 46, 47, 49] |
| Representativeness | An explanation is representative if it holds for many distinct but similar instances. | [8, 16, 43, 49] |

*Subjective system aspects*

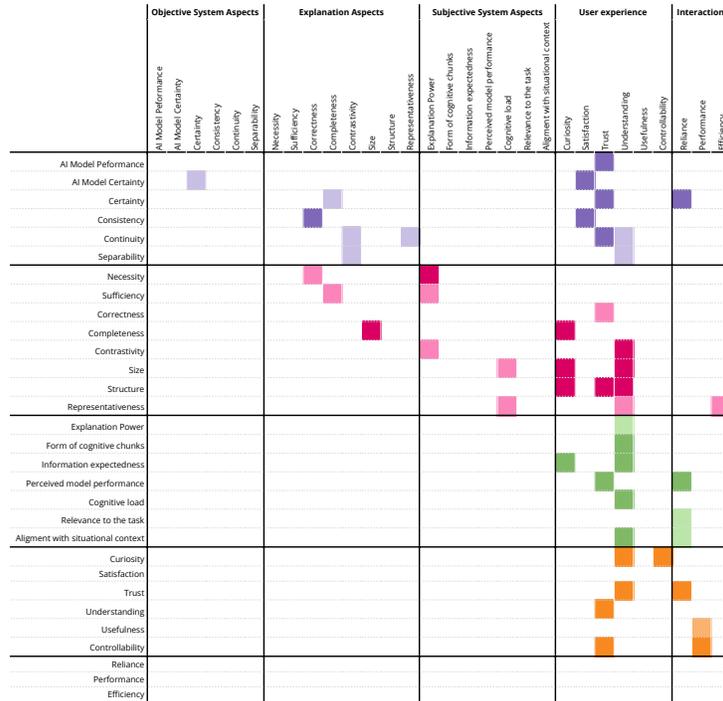| Property | Definition | References |
|---|---|---|
| Explanation power | Measures whether the selected causes make the user understand the reasons the model considered when making a decision | [8, 30, 49, 52] |
| Form of cognitive chunks | Refers to the semantics and structure of the pieces of information the user will receive. | [8, 18, 39, 49] |
| Information expectedness | Level of surprise of the information revealed by the explanation | [5, 8, 21, 30, 32, 35, 37–39, 43, 47, 49] |
| Perceived model competence | Measures the user's impression of the model competence for the task at hand | [11] |
| Cognitive Load | Refers to the cognitive effort the user has to do to achieve the task. | [8, 30, 49] |
| Relevance to the task | Level of explanation usefulness to the user's task. | [21, 30, 35, 38, 39, 47, 49, 52] |

Table 1 – continued from previous page

| Property | Definition | References |
|---|---|---|
| Alignment with situational context | Level of appropriateness of the explanation to the usage context | [8, 21, 30, 43, 47] |
| *User experience* | | |
| Curiosity | Measures whether the user is intrinsically motivated to understand the explanation. If the user is curious, she will be more attentive to the task and, therefore, more engaged with the system. | [20, 21, 49] |
| Satisfaction | Refers to the level of fulfilment the user gets while interacting with the system. This satisfaction is always measured at the communication level because it is for the users to decide whether they feel good about the overall system interaction. | [20, 21, 46, 49, 52] |
| Trust | We use the definition by Tintarev et al. [46]: "perceived confidence in a system's competence" | [3, 11, 20, 21, 30, 46, 47, 52] |
| Understanding | Refers to the ability of the user to interpret the system's output correctly. The user fails to understand when she cannot interpret or incorrectly interprets the system's explanation and prediction. This involves the creation of the user's mental model and how that aligns with the system's functionality. | [8, 20, 21, 32, 34, 49, 52] |
| Usefulness | Measures whether the explanation helps the user to understand the AI prediction. | [5, 46, 49, 52] |
| Controllability | Measures whether the user perceives she has some level of control over the system. This could manifest as the ability to reverse actions, correct the system, filter or zoom the explanation, or ask questions to clarify the explanation or prediction. | [20, 21, 30, 39, 43, 46, 49] |
| *Interaction* | | |
| Efficiency | Measures the speed at which a task can be performed. | [46, 49] |
| Performance | Measures how well the user can do the task while using the system (prediction+explanations). | [20, 21, 32] |
| Reliance | Measures whether the user is willing to provide control to the machine for the given task. | [20, 21, 46, 49] |

**User experience.** The User experience factors evaluate what the user encounters when interacting with the system [25]. The analysis did not find surprising aspects because all but Curiosity, have been studied in recommender systems. The aspects found are Satisfaction, Trust, Understanding, Usefulness, Controllability and Curiosity. This last one was highlighted as extremely relevant in the case of explanations because the search for an explanation is modulated by the user's curiosity [21]. Moreover, the motivation to ask or explore an explanation is determined by the user's curiosity [20].

**Interaction.** Interaction factors measure aspects related to the possible adoption of the system. Three properties were found to be relevant: Efficiency, Performance and Reliance. Efficiency measures how fast the user can perform the task. Performance evaluates the level of achievement the user reaches while using the system. Finally, Reliance measures to which extent the user is willing to provide control to the AI model to perform the task.

## 4.2    Relations Between Properties



**Fig. 5.** Relations between properties. The relations are directed: horizontal properties are the source, and vertical properties receive the effect. High saturation squares indicate that the relationship has been described in the literature, and low saturation squares indicate the relation was inferred from the definitions.

Explanation properties are related to each other in intricate ways [25]. As stated in Section 3.2, by scanning past research, we identified such relationships and linked the properties in our framework as described in the literature; for instance, explanation size affects user curiosity [20]. In this way, we mapped out relations proposed in the literature and relations inferred from the properties' definitions. Table 2 describes the relations found for each property and Figure 5 displays a visual summary of the interactions. These relations help theorise the

expected causal effects between the properties. In practical terms, they serve as hypotheses for the Structural Equation Model.

Table 2: Relations between properties. Relations without reference were hypothesised based on the properties' definitions.
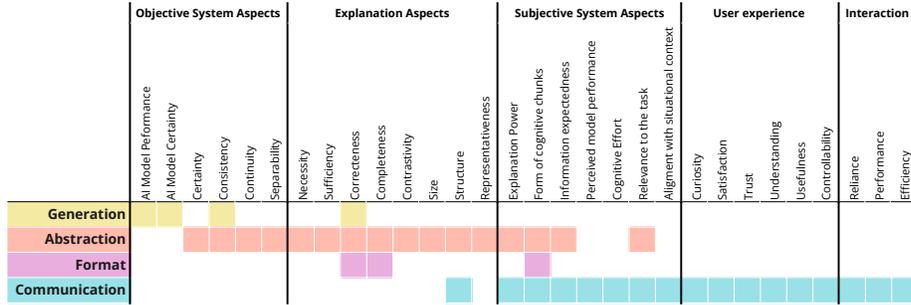
| Explanation Property | Relations with other properties |
| --- | --- |
| *Objective system aspects* | |
| Model performance | The performance of the AI will affect the level of Trust the users can achieve [17, 46] |
| Model certainty | If the AI model is uncertain of the predictions, the XAI method will have more difficulties obtaining consistent explanations, which will affect the XAI method's certainty. This confidence will also affect the satisfaction with the system because, as Tintarev et al. [46] explains, a user might be more forgiving if the system admits it is not confident about a prediction. |
| Certainty | If the explanation shows its limitations, the user may not relay or trust the system [20, 30]. Low certainty will affect the correctness of the explanation. |
| Continuity | Higher continuity increases the understanding of the model because the similarity of explanations helps to learn from the model. It also helps to produce contrastive and representative explanations. Ultimately, high continuity can also increase Trust [47]. |
| Separability | Higher separability increases the understanding of the model and the contrastivity of the explanations. |
| Consistency | Low consistency may decrease user satisfaction [21] and correctness [30] |
| *Explanation Aspects* | |
| Necessity | Affects the explanation power [30]. Additionally, if the necessary causes are selected, then the correctness will be high. |
| Sufficiency | Affects the explanation power. Moreover, if the sufficient causes are selected the completeness of the explanation will increase. |
| Correctness | An explanation with high correctness will faithfully reflect the decision process of the AI model. This could increase Trust in the explanation and AI model |
| Completeness | The explanation size is related to completeness: the bigger the explanation, the more complete it will be [8, 30]. However, bigger explanations might decrease curiosity [21]. |
| Contrastivity | High contrastivity will increase the explanation power. Additionally, this property will affect understanding because the people expect explanations to be contrastive [35]. |
| Size | The amount of information affects curiosity in an inverted U-shaped pattern: little or excessive information reduces curiosity [21]. The size of the explanation also affects how easily a user can understand the explanation [30, 34]. This last effect could be mediated by the Cognitive Load. |

Table 2 – continued from previous page

| Explanation Property | Relations with other properties |
| --- | --- |
| Structure | The design of the information that will be shown affects its trustworthiness [46], curiosity [21] and ultimately how easy they can be understood [34, 43]. |
| Representativeness | This property affects understanding, cognitive load and efficiency because the user can understand an explanation more quickly if it is similar to those she has seen before. |

*Subjective System Aspects*

| | |
| --- | --- |
| Explanation power | The quality of the selected causes will help increase understanding. |
| Form of cognitive chunks | It affects understanding because this property measures how interpretable are the information pieces the user receives [8] |
| Information expectedness | If the information is coherent with the user's beliefs, they will be more likely to understand it [43]. However, if the information does not add anything new to their existing knowledge, they are less likely to be curious [43]. |
| Perceived model competence | When the users perceive the AI model can perform, they are more likely to trust it and eventually to rely on it [11]. |
| Cognitive Load | An explanation with low cognitive load will be easier to understand [30]. |
| Relevance to the task | If the explanations help the development of the task, the user is more likely to rely on the AI advice. |
| Alignment with situational context | Trust in the system is context-dependent. If the system is aligned with the situation the user has to perform, she will be more likely to trust it. [20]. Additionally, this could build up until the user starts to rely on the AI agent. |

*User experience*

| | |
| --- | --- |
| Curiosity | Mental model formation, which is the final goal of understanding, is modulated by Curiosity [21]. Additionally, Curiosity encourages users to explore and interact with the system [21]. |
| Satisfaction | |
| Trust | Reliance is an outcome of appropriate trust [11, 20, 21, 46]. Mental model formation is also modulated by Trust in the system [21] |
| Understanding | If users cannot understand the behaviour, Trust will be lost [21] |
| Usefulness | If the user finds the explanations helpful, they are more likely to increase the user performance with the system. |
| Controllability | The possibility of interaction increases Trust in the system [46]. Good interaction with the system can increase the performance of the users [21] |

### 4.3   Measurements

As explained in Section 3.3, we classified measurements of properties with three criteria: property they measure, explanation element in which they are applied and

**Fig. 6.** Existence of a measurement procedure for each property in each explanation element. Each coloured rectangle indicates that a measurement has been defined for the tuple (property, explanation element)

type of procedure. The explanations elements allow us to capture the complexity of the explanations: they are not simply an object we show to users; a model has generated them and then transformed them to be shown to users in a specific situational context. The four elements are:

- **Generation element**. Refers to the process that was conducted to select the causes that will be displayed in the explanation for a specific object. The measures of this element are applied to the XAI function and AI model. They check the function's parameters to obtain indicators.
- **Abstraction element**. Represents the selected causes of the explanation without considering the format in which they will be displayed. For example, for feature importance, this could be a table with the features and their corresponding importance values. Measurements that are applied at this level look at the data that was selected by the XAI method as an explanation.
- **Format element**. Refers to the manner in which the causes will be presented to the user. This could be as example-based, text, visualisation, etc. In this study, few measurements were found to be applicable at this level. However, each specific media type has its own measurements that could be modified to be applied. For instance, for visual explanations, the data-ink ratio could be applied to analyse whether the most important features use more ink in the visualisation.
- **Communication element**. It refers to the process of interacting with the formatted explanation. During this process, information can be captured as interaction measurements, as well as self-reported information.

In Figure 6, a coloured square is present if at least one measurement exists for that (property, element) tuple. It is noted that the User Experience and Interaction conceptual components are only measured at the communication level, i.e. only when the explanation is displayed to users. What stands out in this figure is the number of measurements at the abstraction level for the Subjective System Aspects component. In Knijnenburg's framework [25], these aspects were recommended to be measured with self-reporting questionnaires.

However, our analysis found that for some of them, metrics have been proposed at the Abstraction element, which means that some computational metric is applied to the abstract explanation to obtain a value.

The main advantage of decoupling the properties from the ways to measure them is that it allows researchers to select measurements considering the study constraints. For instance, if the study is conducted with users that do not have much time to answer questionnaires, and the researchers want to measure *Explanation power* , *Form of cognitive chunks* , *Curiosity* and *Understanding* , they may choose to measure the first two properties at the abstract level of the explanation and the last two at the communication level with questionnaires. In this way, they do not overwhelm the users with questions but still measure the required properties.

As pointed out before, we also classified the measurements by the type of procedure. In this analysis, only procedures that produce a quantitative value were considered. This means that qualitative interviews were not considered, nor were experiment tasks. The four types are:

- **Quantitative interviews**: closed-ended questions, usually in Likert scale.
- **Computational metrics**: mathematical functions that are applied to the explanations or XAI methods.
- **Behaviour metrics**: indicators of user behaviour and interaction with a system. For example, the number of interactions within the system and the time to complete a task.
- **Objective Body Measurements**: measurements taken from the user body. The most common is eye-tracking.

Table 3: Measurements of properties

| Explanation Property | Measurement |
| --- | --- |
| *Generation* | |
| AI Model performance | Measured by the model type appropriate metrics: accuracy, f-score, precision, recall and others. |
| AI Model certainty | This property can be measured for each individual prediction and the global model. If it is measured globally, it should be measured over a dataset similar to the data the real system will face. [16] |
| Consistency | Implementation invariance: check whether the XAI function parameters are the same after different runs of the XAI method creation [8, 49] |
| Correctness | Translucency [8] |
| *Abstraction* | |
| Certainty | Confidence Accuracy [39] |
| Continuity | Connectedness [39] also in [7, 19, 37, 38, 42, 49]; Stability for Slight Variations [39]; Fidelity for Slight Variations [39] |

Table 3 – continued from previous page

| Explanation Property | Measurement |
| --- | --- |
| Separability | Separability [8] |
| Consistency | Stability of explanation: check whether the explanations for a single object change for different instances of the XAI method [8, 16] |
| Necessity | Responsability of an outcome [35]; Sparsity and Sparsity rate [38]; Deletion Check [27, 39] |
| Sufficiency | Count whether the AI model prediction changes when the non-selected causes change [35] |
| Correctness | Model Parameter Randomization Check, Explanation Randomization, White Box Check, Controlled Synthetic Data Check, Predictive Performance [39]; Fidelity [16, 39, 43]; Alignment between AI model features and explanation features [27, 48] |
| Completeness | Preservation Check [39]; Completeness [16, 37]; Recall [48] |
| Contrastivity | Data Randomization, Target Sensitivity, Target Dicriminativeness [39]; Sensitivity [49] |
| Size | Total size or sparsity [39] |
| Structure | Incremental Deletion [16, 39]; Covariate Regularity [16, 39]; Chronology [43]; Single Deletion [39] |
| Representativeness | Explanation support (number of instances to which the explanation applies over the number of instances) [8, 16, 43, 49] |
| Explanation power | Sensitivity Axiom [8] |
| Form of cognitive chunks | Covariate Homogeneity [39] |
| Information expectedness | Alignment with Domain Knowledge [39] |
| Relevance to the task | Pragmatism [16, 38, 39]; Attribute costs [49] |

*Format*

| | |
| --- | --- |
| Correctness | Percentage of invalid rules [49] |
| Completeness | Rules redundancy [49] |
| Form of cognitive chunks | BLEU and METEOR [49]; Perceptual Realism [39] |

*Communication*

| | |
| --- | --- |
| Structure | Questionnaire [11, 49] |
| Explanation power | Questionnaire [50] |
| Form of cognitive chunks | Perceived Homogeneity [39] |
| Information expectedness | Questionnaire [50] |
| Perceived model competence | Questionnaire [3] |
| Cognitive Load | NASA TLX [21] |
| Relevance to the task | Questionnaire [50] |
| Alignment with situational context | Goodness explanation [20] |
| Curiosity | Curiosity Checklist [20]; Eye Movement Pattern [21] |
| Satisfaction | Explanation Satisfaction Scale [20]; Eye Movement Pattern [21]; Loyalty [46]; Questionnaire [3] |
| Trust | Trust Scale [20]; Questionnaire [3, 11, 40, 50] |

Table 3 – continued from previous page

| Explanation Property | Measurement |
|---|---|
| Understanding | Questionnaires [3, 46, 50] |
| Helpfulness | Questionnaires [46, 50]; Evaluate user action before and after explanation [46] |
| Controllability | Concept-level feedback Satisfaction Ratio [13]; The extent to which a user can produce certain outcomes [20] |
| Efficiency | Interaction time and number of interactions to perform a task [46] |
| Performance | Performance metrics with respect to the primary goal [20] |
| Reliance | Questionnaires [3, 11]; Willingness to accept AI agent advice [21] |

The metrics for each (property, element) are listed in Table 3. This table was built under the following rules:

- Several measurements have been defined in multiple works. To avoid naming all of them in the tables, we built upon existing work by using the name proposed by Nauta et al. [39] to summarise metrics every time a similar metric was defined in another work. This new work was added as a reference under the same name.
- Some procedures were not described with a specific name in the paper. In those cases, an explanatory sentence was used to name them.
- If asking questions was proposed as a procedure, but no measurement model or questions were provided, the measurement was not considered.
- For a given property, questions proposed in different papers were joined together under the *Questionnaire* term.

## 5    Illustrative Example

In this section, we provide an illustrative example of how our framework can be used. Researchers have an AI model that predicts whether a patient will be readmitted to the emergency department within 30 days. *SHAP* [33] is used to determine the *feature importance* on a patient level and this information is then visualised in a *force plot* [33]. Finally, medical staff *analyses* the prediction and visual explanation to decide whether they discharge a patient.

In this context, assessing the explanation requires several steps. First, researchers have to decide which properties to measure. This is a decision support system, so according to [30], the most relevant properties would be *Trust*, *Controllability*, and *Understanding*. The researchers conjecture that *Reliance* and *Performance* will be good indicators of adoption. Second, they have to select explanation properties that relate to these five properties. Following the theoretical causal relations in Table 2 and Figure 5, such properties are:

- *AI Model performance*, *Certainty*, *Continuity*

- *Size* , *Structure* , *Representativeness*
- *Form of cognitive chunks* , *Information expecteness* , *Cogntive load* , *Perceived model competence* , *Aligment with situational context*
- *Curiosity*

Finally, to assess all selected properties, researchers pick appropriate metrics from Table 3. The metrics' scores applied to the elements *abstraction* and *format* are averaged over the single explanations, and the questionnaires are applied at the end of the experience. This data is then analysed using structural equation modelling.

## 6  Conclusion and Future Work

In this work, we have presented a user-centric evaluation framework for XAI inspired by research on recommender systems allowing researchers to conduct systematic user experience evaluations in the context of XAI-based systems. Our proposal integrates the current state of the art in XAI evaluation but also allows to easily incorporate new properties or metrics that might become relevant for new applications. By decoupling the aspects of explanations and the procedures to measure them, this framework provides researchers with more tools to choose what and how to measure, and why it is necessary to do it, with the ultimate goal of evaluating the user experience under these new XAI scenarios.

For future work, we plan to validate the framework with user studies. We aim at validating metrics, properties, as well as mediation and causal effects between them. Additionally, we could include experimental designs that compare different explanations, for instance, by comparing the user experience under two different visualisations for explanations generated with SHAP. Furthermore, we did not analyse how specific situational and personal characteristics affect the properties. This area has been explored [15, 31, 45], but more work is needed to connect those findings to explanation properties. Another area of improvement is proposing a standardised report of results to increase fair comparison with previous studies. Lastly, there is no comprehensive survey on the maturity of each of the measurements and on the relations between the properties. Such a survey would help researchers and practitioners to understand the maturity of each property and measurement to help them plan their studies based on current evidence.

### Acknowledgements

# References

1. Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U.R., Makarenkov, V., Nahavandi, S.: A review of uncertainty quantification in deep learning: Techniques, applications and challenges (12 2021). https://doi.org/10.1016/j.inffus.2021.05.008

2. Agarwal, C., Krishna, S., Saxena, E., Pawelczyk, M., Johnson, N., Puri, I., Zitnik, M., Lakkaraju, H.: OpenXAI: Towards a Transparent Evaluation of Model Explanations (6 2022). https://doi.org/10.48550/arxiv.2206.11104, https://arxiv.org/abs/2206.11104v2

3. Ashoori, M., Weisz, J.D.: In AI We Trust? Factors That Influence Trustworthiness of AI-infused Decision-Making Processes (12 2019), http://arxiv.org/abs/1912.02675

4. Balog, K., Radlinski, F.: Measuring Recommendation Explanation Quality. In: Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 329–338. ACM, New York, NY, USA (7 2020). https://doi.org/10.1145/3397271.3401032, https://dl.acm.org/doi/10.1145/3397271.3401032

5. Beckh, K., Müller, S., Rüping, S.: A Quantitative Human-Grounded Evaluation Process for Explainable Machine Learning. Tech. rep. (2022), http://ceur-ws.org

6. Buçinca, Z., Lin, P., Gajos, K.Z., Glassman, E.L.: Proxy tasks and subjective measures can be misleading in evaluating explainable AI systems. International Conference on Intelligent User Interfaces, Proceedings IUI pp. 454–464 (2020). https://doi.org/10.1145/3377325.3377498

7. Carlevaro, A., Lenatti, M., Paglialonga, A., Mongelli, M.: Counterfactual Building and Evaluation via eXplainable Support Vector Data Description. IEEE Access **10**, 60849–60861 (2022). https://doi.org/10.1109/ACCESS.2022.3180026

8. Carvalho, D.V., Pereira, E.M., Cardoso, J.S.: Machine Learning Interpretability: A Survey on Methods and Metrics. Electronics **8**(8), 832 (7 2019). https://doi.org/10.3390/electronics8080832, https://www.mdpi.com/2079-9292/8/8/832

9. Charmaz, K.: Constructing Grounded Theory: A Practical Guide through Qualitative Analysis. No. 4, Sage, Londo, second edn. (2014)

10. Chen, C., Feng, S., Sharma, A., Tan, C.: Machine Explanations and Human Understanding (2 2022), http://arxiv.org/abs/2202.04092

11. Chen, L., Kong, H., Pu, P.: Trust building in recommender agents. Tech. rep. (2005), https://www.researchgate.net/publication/229020498

12. Chen, V., Liao, Q.V., Vaughan, J.W., Bansal, G.: Understanding the Role of Human Intuition on Reliance in Human-AI Decision-Making with Explanations (1 2023), http://arxiv.org/abs/2301.07255

13. Chen, Z., Wang, X., Xie, X., Parsana, M., Soni, A., Ao, X., Chen, E.: Towards Explainable Conversational Recommendation. Tech. rep. (2020), https://concept.research.microsoft.com/

14. Clinciu, M.A., Eshghi, A., Hastie, H.: A Study of Automatic Metrics for the Evaluation of Natural Language Explanations. In: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume. pp. 2376–2387. Association for Computational Linguistics, Stroudsburg, PA, USA (2021). https://doi.org/10.18653/v1/2021.eacl-main.202, https://aclanthology.org/2021.eacl-main.202

15. Conati, C., Barral, O., Putnam, V., Rieger, L.: Toward personalized XAI: A case study in intelligent tutoring systems. Artificial Intelligence **298**, 103503 (9 2021). https://doi.org/10.1016/J.ARTINT.2021.103503

16. Coroama, L., Groza, A.: Evaluation Metrics in Explainable Artificial Intelligence (XAI). In: Communications in Computer and Information Science. vol. 1675 CCIS, pp. 401–413. Springer Science and Business Media Deutschland GmbH (2022). https://doi.org/10.1007/978-3-031-20319-0_30

17. Dominguez, V., Donoso-Guzmán, I., Messina, P., Parra, D.: The effect of explanations and algorithmic accuracy on visual recommender systems of artistic images. In: International Conference on Intelligent User Interfaces, Proceedings IUI. vol. Part F1476 (2019). https://doi.org/10.1145/3301275.3302274

18. Doshi-Velez, F., Kim, B.: Towards A Rigorous Science of Interpretable Machine Learning. Arxiv pp. 1–13 (2 2017), http://arxiv.org/abs/1702.08608

19. Ge, Y., Liu, S., Li, Z., Xu, S., Geng, S., Li, Y., Tan, J., Sun, F., Zhang, Y.: Counterfactual Evaluation for Explainable AI (9 2021), http://arxiv.org/abs/2109.01962

20. Hoffman, R.R., Mueller, S.T., Klein, G., Litman, J.: Metrics for Explainable AI: Challenges and Prospects pp. 1–50 (12 2018), http://arxiv.org/abs/1812.04608

21. Hsiao, J.H.w., Ngai, H.H.T., Qiu, L., Yang, Y., Cao, C.C.: Roadmap of Designing Cognitive Metrics for Explainable Artificial Intelligence (XAI) (7 2021). https://doi.org/10.48550/arxiv.2108.01737, https://arxiv.org/abs/2108.01737v1

22. Johs, A.J., Agosto, D.E., Weber, R.O.: Qualitative Investigation in Explainable Artificial Intelligence: A Bit More Insight from Social Science. In: Association for the Advancement of Artificial Intelligence (11 2020), http://arxiv.org/abs/2011.07130

23. Johs, A.J., Agosto, D.E., Weber, R.O.: Explainable artificial intelligence and social science: Further insights for qualitative investigation. Applied AI Letters **3**(1) (2 2022). https://doi.org/10.1002/ail2.64

24. Kline, R.B.: Principles and practice of structural equation modeling. Guilford publications, 5th edition edn. (2023)

25. Knijnenburg, B.P., Willemsen, M.C.: Evaluating Recommender Systems with User Experiments. In: Recommender Systems Handbook, pp. 309–352. Springer US, Boston, MA (1 2015). https://doi.org/10.1007/978-1-4899-7637-6_9, https://link.springer.com/10.1007/978-1-4899-7637-6_9

26. Langer, M., Oster, D., Speith, T., Hermanns, H., Kästner, L., Schmidt, E., Sesing, A., Baum, K.: What do we want from Explainable Artificial Intelligence (XAI)? – A stakeholder perspective on XAI and a conceptual model guiding interdisciplinary XAI research. Artificial Intelligence **296**, 103473 (7 2021). https://doi.org/10.1016/J.ARTINT.2021.103473

27. Li, Y., Zhou, J., Verma, S., Chen, F.: A Survey of Explainable Graph Neural Networks: Taxonomy and Evaluation Metrics (7 2022), http://arxiv.org/abs/2207.12599

28. Liao, Q.V., Pribić, M., Han, J., Miller, S., Sow, D.: Question-Driven Design Process for Explainable AI User Experiences **1**(1), 1–23 (2021), http://arxiv.org/abs/2104.03483

29. Liao, Q.V., Varshney, K.R.: Human-Centered Explainable AI (XAI): From Algorithms to User Experiences (10 2021), http://arxiv.org/abs/2110.10790

30. Liao, Q.V., Zhang, Y., Luss, R., Doshi-Velez, F., Dhurandhar, A.: Connecting Algorithmic Research and Usage Contexts: A Perspective of Contextualized Evaluation for Explainable AI. Proceedings of the AAAI Conference on Human Computation and Crowdsourcing **10**(1), 147–159 (10 2022). https://doi.org/10.1609/hcomp.v10i1.21995, https://ojs.aaai.org/index.php/HCOMP/article/view/21995

31. Lim, B.Y., Dey, A.K., Avrahami, D.: Why and why not explanations improve the intelligibility of context-aware intelligent systems. Conference on Human Factors in Computing Systems - Proceedings pp. 2119–2128 (2009). https://doi.org/10.1145/1518701.1519023, https://dl.acm.org/doi/10.1145/1518701.1519023

32. Löfström, H., Hammar, K., Johansson, U.: A Meta Survey of Quality Evaluation Criteria in Explanation Methods (3 2022), http://arxiv.org/abs/2203.13929

33. Lundberg, S.M., Nair, B., Vavilala, M.S., Horibe, M., Eisses, M.J., Adams, T., Liston, D.E., Low, D.K.W., Newman, S.F., Kim, J., Lee, S.I.: Explainable machine-learning predictions for the prevention of hypoxaemia during surgery. Nature Biomedical Engineering **2**(10), 749–760 (10 2018). https://doi.org/10.1038/s41551-018-0304-0

34. Markus, A.F., Kors, J.A., Rijnbeek, P.R.: The role of explainability in creating trustworthy artificial intelligence for health care: a comprehensive survey of the terminology, design choices, and evaluation strategies. Journal of Biomedical Informatics **113**, 103655 (7 2020). https://doi.org/10.1016/j.jbi.2020.103655, http://arxiv.org/abs/2007.15911http://dx.doi.org/10.1016/j.jbi.2020.103655

35. Miller, T.: Explanation in artificial intelligence: Insights from the social sciences. Artificial Intelligence **267**, 1–38 (2 2019). https://doi.org/10.1016/j.artint.2018.07.007, https://linkinghub.elsevier.com/retrieve/pii/S0004370218305988

36. Mohseni, S., Zarei, N., Ragan, E.D.: A Multidisciplinary Survey and Framework for Design and Evaluation of Explainable AI Systems. ACM Transactions on Interactive Intelligent Systems **1**(3-4), 1–45 (2021). https://doi.org/10.1145/3387166, http://arxiv.org/abs/1811.11839

37. Moraffah, R., Karami, M., Guo, R., Raglin, A., Liu, H.: Causal Interpretability for Machine Learning-Problems, Methods and Evaluation. Tech. rep.

38. Moreira, C., Chou, Y.L., Hsieh, C., Ouyang, C., Jorge, J., Pereira, J.M.: Benchmarking Counterfactual Algorithms for XAI: From White Box to Black Box (3 2022), http://arxiv.org/abs/2203.02399

39. Nauta, M., Trienes, J., Pathak, S., Nguyen, E., Peters, M., Schmitt, Y., Schlötterer, J., van Keulen, M., Seifert, C.: From Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review on Evaluating Explainable AI (1 2022), http://arxiv.org/abs/2201.08164

40. Pu, P., Chen, L.: Trust-inspiring explanation interfaces for recommender systems. Knowledge-Based Systems **20**(6), 542–556 (8 2007). https://doi.org/10.1016/j.knosys.2007.04.004

41. Rong, Y., Leemann, T., Nguyen, T.t., Fiedler, L., Qian, P., Unhelkar, V., Seidel, T., Kasneci, G., Kasneci, E.: Towards Human-centered Explainable AI: User Studies for Model Explanations (10 2022), http://arxiv.org/abs/2210.11584

42. Singh, V., Cyras, K., Inam, R.: Explainability Metrics and Properties for Counterfactual Explanation Methods. In: Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics). vol. 13283 LNAI, pp. 155–172. Springer Science and Business Media Deutschland GmbH (2022). https://doi.org/10.1007/978-3-031-15565-9_10

43. Sokol, K., Flach, P.: Explainability fact sheets: A framework for systematic assessment of explainable approaches. In: FAT* 2020 - Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency. pp. 56–67. Association for Computing Machinery, Inc (1 2020). https://doi.org/10.1145/3351095.3372870

44. Suresh, H., Gomez, S.R., Nam, K.K., Satyanarayan, A.: 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. vol. 16, pp. 1–16. ACM, New York, NY, USA (5 2021). https://doi.org/10.1145/3411764.3445088, https://dl.acm.org/doi/10.1145/3411764.3445088

45. Szymanski, M., Abeele, V.V., Verbert, K.: Explaining health recommendations to lay users: The dos and don'ts. Tech. rep. (2022), http://ceur-ws.org

46. Tintarev, N., Masthoff, J.: Explaining Recommendations: Design and Evaluation. In: Recommender Systems Handbook, pp. 353–382. Springer US, Boston, MA (1 2015). https://doi.org/10.1007/978-1-4899-7637-6_10, https://link.springer.com/10.1007/978-1-4899-7637-6_10

47. Tonekaboni, S., Joshi, S., McCradden, M.D., Goldenberg, A.: What Clinicians Want: Contextualizing Explainable Machine Learning for Clinical End Use. Proceedings of Machine Learning Research (5 2019), http://arxiv.org/abs/1905.05134

48. Velmurugan, M., Ouyang, C., Moreira, C., Sindhgatta, R.: Developing a Fidelity Evaluation Approach for Interpretable Machine Learning (6 2021), http://arxiv.org/abs/2106.08492

49. Vilone, G., Longo, L.: Notions of explainability and evaluation approaches for explainable artificial intelligence. Information Fusion **76**, 89–106 (12 2021). https://doi.org/10.1016/J.INFFUS.2021.05.009

50. Vilone, G., Longo, L.: A Novel Human-Centred Evaluation Approach and an Argument-Based Method for Explainable Artificial Intelligence. In: IFIP Advances in Information and Communication Technology. vol. 646 IFIP, pp. 447–460. Springer Science and Business Media Deutschland GmbH (2022). https://doi.org/10.1007/978-3-031-08333-4_36

51. Wang, D., Yang, Q., Abdul, A., Lim, B.Y.: Designing theory-driven user-centric explainable AI. In: Conference on Human Factors in Computing Systems - Proceedings. Association for Computing Machinery (5 2019). https://doi.org/10.1145/3290605.3300831

52. Wanner, J., Herm, L.V., Heinrich, K., Janiesch, C.: A social evaluation of the perceived goodness of explainability in machine learning. Journal of Business Analytics **5**(1), 29–50 (2022). https://doi.org/10.1080/2573234X.2021.1952913