

# XEQ Scale for Evaluating XAI Experience Quality

Anjana Wijekoon<sup>1,\*†</sup>, Nirmalie Wiratunga<sup>2</sup>, David Corsar<sup>2</sup>, Kyle Martin<sup>2</sup>,  
Ikechukwu Nkisi-Orji<sup>2</sup>, Belen Díaz-Agudo<sup>3</sup> and Derek Bridge<sup>4</sup>

<sup>1</sup>University College London, London, England

<sup>2</sup>Robert Gordon University, Aberdeen, Scotland

<sup>3</sup>Universidad Complutense de Madrid, Spain

<sup>4</sup>University College Cork, Ireland

## Abstract

Explainable Artificial Intelligence (XAI) aims to improve the transparency of autonomous decision-making through explanations. Recent literature has emphasised users' need for holistic "multi-shot" explanations and personalised engagement with XAI systems. We refer to this user-centred interaction as an XAI Experience. Despite advances in creating XAI experiences, evaluating them in a user-centred manner has remained challenging. In response, we developed the XAI Experience Quality (XEQ) Scale. XEQ quantifies the quality of experiences across four dimensions: learning, utility, fulfilment and engagement. These contributions extend the state-of-the-art of XAI evaluation, moving beyond the one-dimensional metrics frequently developed to assess single-shot explanations. This paper presents the XEQ scale development and validation process, including content validation with XAI experts, and discriminant and construct validation through a large-scale pilot study. Our pilot study results offer strong evidence that establishes the XEQ Scale as a comprehensive framework for evaluating user-centred XAI experiences.

## Keywords

Explanation Experience Quality, Multi-shot Explanation, Psychometric Theory, Scale Development

## 1. Introduction

Explainable Artificial Intelligence (XAI) describes a range of techniques to elucidate autonomous decision-making and the data that informed that AI system [1, 2, 3]. Each technique typically provides explanations focusing on a specific aspect of the system and its decisions, often answering a singular question fulfilling a specific intent [4, 5]. Accordingly, the utility of employing multiple techniques for a holistic explanation of a system becomes increasingly clear [3, 6]. The collection of explanations, provided by different techniques and describing different components of the system, forms what we describe as "multi-shot" explanations. Multi-shot explanation treats the elucidation of an AI system as a process which can invoke a variety of knowledge sources to construct explanations that comprehensively answer a user's explanation needs. This is unlike single-shot, which provides a lone explanation that targets a specific aspect of an autonomous decision to resolve an individual query (see Figure 1). Furthermore, viewing multi-shot explanations as an interactive process allows users to engage through graphical interfaces [3] or conversations [7, 6], enabling personalised experiences tailored to individual needs [8].

While the utility of user-centred interactive multi-shot explanations is evident in recent literature, evaluating them remains a key research challenge. Current works primarily target the development of objective metrics for single-shot techniques [9, 10], emphasising the need for reproducible benchmarks on public datasets [11]. Such metrics are system-centred and model-agnostic, meaning they are compatible with a variety of XAI systems and giving the advantage of generalisability. However,

---

*Preprint*

\*Corresponding author.

†This work was done while affiliated with Robert Gordon University, Aberdeen, Scotland.

✉ a.wijekoon@ucl.ac.uk (A. Wijekoon)

ORCID 0000-0003-3848-3100 (A. Wijekoon); 0000-0003-4040-2496 (N. Wiratunga); 0000-0001-7059-4594 (D. Corsar);  
00000-0003-0941-3111 (K. Martin); 0000-0001-9734-9978 (I. Nkisi-Orji); 0000-0003-2818-027X (B. Díaz-Agudo);  
0000-0002-8720-3876 (D. Bridge)



© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

Multi-shot Explanation	
<p>Why was my application rejected?</p> <p>Our recent credit inquiry suggests you are a high-risk candidate. If the most recent inquiry was more than a month ago, your application would fall in the approved category.</p>	Single-shot Explanation
<p>Apart from credit rating, what else can help when I apply again?</p> <p>Your consolidated risk marker was at 77, it would be helpful to get it above 88 when you apply next time.</p>	Single-shot Explanation
<p>How did you calculate the consolidated risk marker ?</p> <p>The consolidated risk marker is calculated as an aggregation of your payment history, credit utilisation ratio, length of credit history and types of credit accounts you already own.</p>	Single-shot Explanation
<p>Can you give me an example of what someone in a similar situation did to get their credit approved?</p> <p>Yes of course, here are a couple of examples....</p>	Single-shot Explanation

**Figure 1:** A conversational interaction demonstrating a single-shot explanation (adapted from IBM Explainability 360 platform [3]) extended to a multi-shot XAI experience. Each single-shot explanation is represented by individual turns in the conversation, while a multi-shot explanation is represented by the conversation in its entirety.

objective metrics fail to acknowledge the requirements of different stakeholder groups. A satisfactory explanation is reliant on the recipient’s expertise within that domain [12] and their previous experiences with AI [13, 14]. Subjective metrics, such as those described in [15, 16], allow an evaluation which is personalised to the individual and domain. However, existing subjective evaluations lack the capacity to measure the interactive process that underpins multi-shot explanations and how they impact user experience. There is a clear need for a tool that facilitates the evaluation of XAI experiences. In turn, this would empower the iterative development of XAI systems to ensure they address the diverse needs of various stakeholder groups.

**Table 1**

Glossary of term and definitions used throughout this paper.

Term	Definition
XAI System	An automated decision-making system that is designed and developed to provide information about its reasoning.
Stakeholder	An individual or group with a vested interest in the XAI system. Stakeholders are a diverse group, encompassing system designers and developers, who hold an interest in the system’s technical functionality, the end consumers relying on its decisions, and regulatory authorities responsible for ensuring fair and ethical use.
Single-Shot Explanation	An explanatory process where a lone explanation is provided to answer a specific query.
Multi-Shot Explanation	An explanatory process where multiple techniques are used to provide two or more alternative and/or complementary insights into a single XAI system or output.
XAI Experience (XE)	A user-centred process of a stakeholder interacting with an XAI system to gain knowledge and/or improve comprehension.
XE Quality (XEQ)	The extent to which a stakeholder’s explanation needs are satisfied by their XE.

We address this challenge by introducing the XAI Experience Quality (XEQ) Scale. We define an XAI Experience as the user-centred process of a stakeholder interacting with an XAI system to gain knowledge and/or improve comprehension. XAI Experience Quality (XEQ) is defined as *the extent to which a stakeholder’s explanation needs are satisfied by their XAI Experience*. A glossary of all related terminology used throughout this paper is included in Table 1. Specifically, we ask the research question: “How to evaluate an XAI experience, in contrast to assessing single-shot (non-interactive) explanations?”. To address this, we follow a formal psychometric scale development process [17] and outline the following objectives:

1. conduct a literature review to compile a collection of XAI evaluation questionnaire items;

2. conduct a content validity study with XAI experts to develop the XEQ scale; and
3. perform a pilot study to refine and validate the XEQ scale for internal consistency, discriminant and construct validity and test-retest reliability.

The rest of this paper expands on each objective. We discuss related work in Section 2. Section 3 presents the creation of the initial items bank. The Content Validity study details and results are presented in Section 4 followed by Section 5 presenting the pilot study for the refinement and validation of the XEQ Scale. Key implications and the practical use of the Scale are discussed in Section 6. Finally, we offer conclusions in Section 7.

## 2. Related Work

In the literature, there are several methodologies for developing evaluation metrics or instruments for user-centred XAI.

Hoffman et al. [15] employed Psychometric Theory to construct the Satisfaction Scale, evaluating both content validity and discriminant validity. A similar methodology was adopted in [18] to develop the Madsen-Gregor human-machine trust scale, relying on pre-existing item lists from previous scales in conjunction with expert insights [19]. Jian et al. [20] pursued a factor analysis approach involving non-expert users to formulate a human-machine trust scale. They compiled words and phrases associated with trust and its variants, organising them based on their relevance to trust and distrust, which were then clustered to identify underlying factors and formulate corresponding statements. This methodology is particularly suitable in cases where no prior items exist for initial compilation. While these methodologies are robust to produce reliable scales they are resource and knowledge-intensive processes.

A more frequent approach to scale development is deriving them from existing scales in psychology research. For instance, the System Causability Scale [16] draws inspiration from the widely used System Usability Scale [21], while the Hoffman Curiosity Checklist originates from scales designed to assess human curiosity [15]. Similarly, the Cahour-Forzy Trust Scale [22] selected questions from research on human trust, and the Hoffman Trust Scale incorporates items from previous trust scales [22, 20]. Notably, these derived scales were not evaluated for reliability or other desirable factors, they rely on the quality of the original scales for validity. In this paper, we opt for the psychometric theory approach to establish the content, construct and discriminant validity of the resulting scale. While this approach is resource-intensive, the complexity and the novelty of the evaluation task necessitate a rigorous approach to scale development.

## 3. Initial Items Bank Compilation

This section presents the literature review findings that led to the compilation of the initial items bank for the XEQ Scale.

### 3.1. Methodology

Netemeyer [23] outlines a framework for developing an item pool, via domain sampling grounded literature review. Following this approach, we conducted a targeted literature review to examine existing research, identify key dimensions, and develop the initial item bank. The reasoning for a targeted review instead of a systematic review is twofold: 1) the purpose of the review is to form the initial item bank which involves in-depth analysis of selected literature (depth over breadth); and 2) literature under this topic is significantly limited. The initial findings highlighted that while many publications discuss and emphasise the importance of evaluation dimensions (what should be or is evaluated), only a few propose and systematically develop metrics for XAI evaluation.

## 3.2. Findings: Evaluation metrics and dimensions

Hoffman et al. [15] are one of the leading contributors and their work has been widely utilised in many user-centred XAI research. They conceptually modelled the “process of explaining in XAI” outlining dimensions and metrics for evaluating single-shot explanations from stakeholders’ perspectives. They considered six evaluation dimensions: goodness, satisfaction, mental model, curiosity, trust and performance. For each dimension, they either systematically developed an evaluation metric or critiqued metrics available in literature offering a comprehensive evaluation methodology for XAI practitioners. System Causability Scale [16] is the other most prominent work in XAI evaluation. Notably, all the above scales are designed to be application-domain agnostic, with each item adaptable to reference the specific AI system under evaluation. We discuss each scale briefly below.

**Hoffman’s Goodness Checklist** is utilised to objectively evaluate explanations with an independent XAI expert to improve the “goodness”. It consists of 7 items answered by either selecting ‘yes’ or ‘no’. It was developed by referring to literature that proposes “goodness” properties of explanations.

**Hoffman Satisfaction Scale** was designed using psychometric theory to evaluate the subjective “goodness” of explanations with target users. It consists of 8 items responded in a 5-step Likert Scale. It is viewed as the user-centred variant of the Goodness Checklist with many shared items. This scale has been evaluated for content validity with XAI experts as well as construct and discriminant validity in pilot studies.

**Hoffman Curiosity Checklist** is designed to elicit stakeholder explanation needs, i.e. which aspects of the system pique their curiosity. This metric consists of one question *Why have you asked for an explanation? Check all that apply.* and the responses inform the design and implementation of the XAI system.

**Hoffman Trust Scale** measures the development of trust when exploring a system’s explainability. The authors derived this trust scale by considering the overlaps and cross-use of scales from trust scales in literature for measuring trust in autonomous systems (not in the presence of explainability, e.g. trust between human and a robot) [20, 24, 25, 22].

**System Causability Scale** measures the effectiveness, efficiency and satisfaction of the explainability process in systems involving multi-shot explanations [16]. Derived from the widely-used System Usability Scale [21], this scale comprises 10 items rated on a 5-step Likert scale. Notably, it includes items that measure stakeholder engagement, addressing a gap in previous scales designed for one-shot explainability settings. However, the validation of the scale is limited to one small-scale pilot study in the medical domain.

### 3.2.1. Other Dimensions

Many other publications emphasised the need for user-centred XAI evaluations and explored evaluation dimensions. Two other dimensions considered by Hoffman et al. [15] are mental model and performance concerning task completion. They recommended eliciting the mental model of stakeholders in think-aloud problem-solving and question-answering sessions. Performance is measured by observing the change in productivity and change in system usage. The evaluation of these dimensions requires metrics beyond questionnaire-based techniques. Another domain-agnostic survey finds many overlaps with Hoffman et al., defining 4 user-centred evaluation dimensions: mental model, usefulness and satisfaction, trust and reliance and human-task performance [12]. Zhou et al., [26] summarise previous literature, emphasising three subjective dimensions - trust, confidence and preference that overlap with dimensions identified in [15]. Conversely to Hoffman et al., they consider task completion to be an objective dimension in user-centred XAI evaluation.

Carvalho et al., delineate characteristics of a human-friendly explanation in the medical domain, including some subjective or user-centred properties such as comprehensibility, novelty, and consistency with stakeholders' prior beliefs [27]. Notably, consistency with stakeholders' prior beliefs aligns with the mental model from Hoffman et al. [15], while novelty can influence stakeholder engagement [16]. Nauta and Seifert [28] recognise 12 properties of explanation quality for image classification applications. They identify three user-centred properties: context - how relevant the explanation is to the user; coherence - how accordant the explanation is with prior knowledge and beliefs; and controllability - how interactive and controllable the explanation is. In comparison to other literature, controllability aligns with engagement [16] and coherence aligns with the mental model [15, 27]. Context can be associated with several properties such as curiosity, satisfaction, and preference [15, 26].

These findings highlighted that there are many overlaps between evaluation dimensions identified in recent years. However, we highlight two main gaps in this current work: 1) there is no consensus in previous literature regarding the applicable metrics to measure these evaluation dimensions; and 2) the majority of the existing dimensions and metrics focus on evaluating individual explanations (i.e. single-shot), not multi-shot XAI experiences.

### 3.3. Initial item bank

Informed item sampling from literature resulted in 33 items, including 7 from the Goodness Checklist [15]; 8 from the Satisfaction Scale [15]; 8 from the Trust Scale [15]; and 10 from the System Causability Scale [16]. As evidenced in the literature review, while these items covered constructs related to user-centred evaluation of single-shot explanations, they lacked explicit representativeness of multi-shot explanations. Recognising that *interaction* is crucial to creating effective explanation experiences, the research team generated seven additional items to capture this construct. These items were designed to capture stakeholder perspectives on the interactive experience, a dimension explicitly not covered in previous XAI evaluation literature. This expanded the item pool to a total of 40 items.

The item pool was then standardised to a symmetric 5-point Likert Scale, ranging from "I Strongly Agree" to "I Strongly Disagree". The existing user-centred XAI scales identified in the literature review commonly adopt the 5-point format, offering a balance between granularity and minimising ambiguity between response options. Broader research in psychology has also shown via mathematical modelling that using a scale with more than five points offers no significant improvement in internal consistency reliability measured by Cronbach's alpha [29].

Finally, following similar approaches from recent literature in scale development [30, 31], the item pool then underwent a 2-stage rigorous review and revision process. Each stage of the revision process was led by two authors, who proposed revisions, and subsequently validated by the remaining authors, who either accepted, rejected, or further revised them. In the first stage, lexically identical items were removed, and items with significant semantic overlap were consolidated, resulting in 32 items. During the second stage, expressions were made more concise, non-suggestive, and aligned with the intended construct (i.e., measuring XAI experiences rather than explanations). Rephrased items were further validated by iterating wording until unanimous agreement between authors to ensure clarity for stakeholders with varying knowledge levels and to eliminate suggestive wording to minimise potential bias. The resulting 32 items formed the initial XEQ Scale (included in Supplementary Material).

#### 3.3.1. Evaluation Dimensions

We reviewed evaluation dimensions from previous literature and consolidated XEQ items into four evaluation dimensions representing XAI experience quality: learning, utility, fulfilment, and engagement. These dimensions are relevant to capturing personalised experiences for a given stakeholder. We define them as follows:

**Learning:** the extent to which the experience develops knowledge or competence;

**Utility:** the contribution of the experience towards task completion;

**Fulfilment:** the degree to which the experience supports the achievement of XAI goals; and

**Engagement:** the quality of the interaction between the user and the XAI system.

In the next sections, we describe the development, refinement and validation of the XEQ Scale following Psychometric Theory [32].

## 4. XEQ Scale Development

This section presents the details and results of the expert-led evaluation of the initial set of 32 items to develop the XEQ scale, following the Content Validity Ratio (CVR) method [33, 34].

### 4.1. Study instrument

The study instrument consisted of three components: 1) an introduction to terminology; 2) three XAI experience examples; and 3) a structured survey for data collection. The first component introduced the participants to the terminology used throughout the study (also included in Table 1). This is intended to establish a consistent mental model between experts from diverse XAI backgrounds. Next, participants explored three XAI experiences curated to capture different application domains, end-user explanation needs and interaction modalities. Lastly, participants responded to the survey with three sub-sections: a) rating the 32 items regarding their relevance for measuring XAI experience quality on a 5-step Likert scale ranging from *Not Relevant at All* to *Extremely Relevant*; b) rating the 32 items regarding their clarity on a similar Likert scale ranging from *Not Clear at All* to *Extremely Clear*; and c) rephrasing items that they found relevant but not clear.

#### 4.1.1. XAI experiences development and selection

The development of sample XAI experiences was conducted over three stages. First, we established several desiderata for sample experiences: a) reflect diverse, real-world AI systems where explaining AI outcomes is critical; b) possess the level of complexity that is suited for multi-shot explanation; and c) capture varying levels of domain expertise and AI familiarity.

The second stage builds on these desiderata, to identify relevant application domains for sample XAI experiences considering previous literature on XAI design and development, and feedback from industry and academic partners. We selected radiograph fracture detection from the medical domain, welfare application approval auditing from the governance domain, and student support for course selection from academia. Regulations governing the use of AI decision-support in medical imaging underscore the need for XAI, particularly through comprehensive explanation methodologies [35, 36, 37, 38]. Similar regulations are being implemented for AI in civic and governmental decision-making [39] to ensure that automated decisions are interpretable, contestable, auditable, and fair [40]. Finally, student support represents an area with significantly lower legal/ethical requirements for explainability, but with high user experience impact [41].

In the third and final stage, we designed each of the sample experiences. For each experience, the target user was identified to ensure a variety of AI knowledge and domain expertise levels. The interaction modality was selected to suit the target user's explanation and interaction needs. The sample interactions articulated different types of explanations addressing the target user's explanation needs. The three sample experiences are summarised in Table 2. A sample dialogue with the interactive CourseAssist chatbot is shown in Figure 2 positive XAI experience. This dialogue, in its entirety, captures an XAI experience that we wish to evaluate. The experiences constructed for medical imaging and welfare application approval auditing are available in the Supplementary Material.

### 4.2. Participant recruitment and selection

The CVR method recommends a small group of domain experts (5-10) for quantifying the strength of psychometric scale items [33]. We invited 38 XAI experts with diverse expertise within XAI, considering

**Table 2**

Sample XAI experiences; for each experience, we summarise the AI system, attributes of the target stakeholder and the explanations presented in their XAI experience. Knowledge levels are adapted from the Dreyfus 6-stage skill acquisition model [42].

Identifier	Attributes	Name	Domain knowledge	AI familiarity	Interaction Modality
CourseAssist	<b>Stakeholder</b>	Student	Novice	Novice	Conversation
	<b>AI System</b>	A rule-based agent that recommends courses given a student's chosen career path			
	<b>Explanations</b>	Distribution and statistics of the data used in the recommendations; Factual explanations; Nearest-neighbour cases			
RadioAssist	<b>Stakeholder</b>	Trainee Radiologist	Proficient	Competent	Graphical User Interface
	<b>AI System</b>	A neural network black-box that predicts the presence of fracture given a radiograph in PNG format			
	<b>Explanations</b>	Distribution and statistics of the data used to develop the model; Feature attributions; Nearest-neighbour cases; Performance metrics of the AI model			
AssistHub	<b>Stakeholder</b>	Auditor	Expert	Proficient	Graphical User Interface
	<b>AI System</b>	A decision tree algorithm that recommends the outcome given a welfare application			
	<b>Explanations</b>	Distribution and statistics of the data used to create the AI system; Prototypical approved and rejected cases; Outlier cases; Simplified decision tree			

those who are actively publishing in the XAI domain since 2020. The study recruited a diverse group of experts, 3 from industry and 10 from academia. Their specific areas of expertise were identified as follows: 3 experts in implementing XAI in industrial engineering applications (Energy and Telecommunications); 1 expert in XAI evaluation scale development; 5 experts in example-based and counterfactual XAI methods; 2 experts in human-centric XAI design and development; and 2 experts in the role of causality in XAI. They collectively represent the diverse aspects of crafting effective XAI experiences.

### 4.3. Data collection and analysis

The study was hosted on the Jisc Online Surveys platform for 3 weeks and collected responses from 13 XAI experts. Their responses were collated and analysed using the following metrics.

#### Content Validity

The Content Validity Index (CVI) assesses item validity based on responses to the relevance property. Lower scores indicate items that may need modification or removal. Given scale  $\mathbf{S}$  with  $M$  items where  $i$  indicates an item,  $r_j^i$  denotes the response of participant  $j$  to item  $i$ . For analysis, each response ( $r_j^i$ ) is modified as follows.

$$r_j^i = \begin{cases} 1, & \text{if } r_j^i \in [\text{Extremely Relevant or Somewhat Relevant}] \\ 0, & \text{otherwise} \end{cases}$$

We calculate the following two forms of the Content Validity Index (CVI) scores.

**Item-Level CVI:** measures the validity of each item independently; the number of responses is  $N$  and the expected score is  $\geq 0.78$ .

$$I-CVI_i = \frac{\sum_{j=1}^N (r_j^i)}{N}$$

**Scale-Level CVI:** measures the overall scale validity using a) Average method i.e. the mean Item-Level CVI score where the expected score is  $\geq 0.90$ ; and b) Universal Agreement method i.e. the percentage of items experts always found relevant with the expected value of  $\geq 0.80$ .

$$S-CVI(a) = \frac{\sum_{i=1}^M (I-CVI_i)}{M}$$

$$S-CVI(b) = \frac{\sum_{i=1}^M 1_{[I-CVI_i=1]}}{M}$$

Here, once the average of the I-CVIs is calculated for all items with  $S-CVI(a)$ ,  $S-CVI(b)$  counts the number of items with an I-CVI of 1 (indicating complete agreement among experts that the item is relevant) and divides this by the total number of items.

#### 4.4. Results

We refer to columns Item and I-CVI from Table 4 for the results of the Content Validity study. We first removed items with low validity ( $I-CVI_i \leq 0.75$ ) and thereafter S-CVI scores were used to establish the content validity of the resulting scale. Here we marginally divert from the established baseline of 0.78 for I-CVI to further investigate items with  $0.75 \leq I-CVI \leq 0.78$  during the pilot study. The Likert responses to the clarity property and free text feedback influenced the re-wording of 7 items to improve clarity (indicated by †). The item selection and rephrasing were done based on the suggestions from the XAI experts and the consensus of the research team. The resulting scale comprised 18 items, which we refer to as the *XEQ Scale*. In Table 4, items are ordered by their I-CVI scores.

$S-CVI(a)$  and  $S-CVI(b)$  of the scale were 0.8846 and 0.2222. While  $S-CVI(a)$  is comparable to the baseline of 0.9,  $S-CVI(b)$  indicate universal agreement is not achieved. However, existing literature suggests that meeting one of the baseline criteria is sufficient to proceed to pilot studies. Notably, the 14 items with  $I-CVI \geq 0.78$  also only achieve average agreement ( $S-CVI(a) = 0.9179$ ) and not universal agreement ( $S-CVI(b) = 0.2667$ ).

Following the item selection, each item was assigned an evaluation dimension based on the consensus of the research team (see Column “Evaluation Dimension” in Table 4). These will be used in further investigations using confirmatory factor analysis in Section 5 to establish the construct validity.

## 5. XEQ Scale Refinement and Validation

In this section, we present the pilot study conducted to refine the XEQ Scale for internal consistency and construct validity, discriminant validity and test-retest reliability.

### 5.1. Study instrument

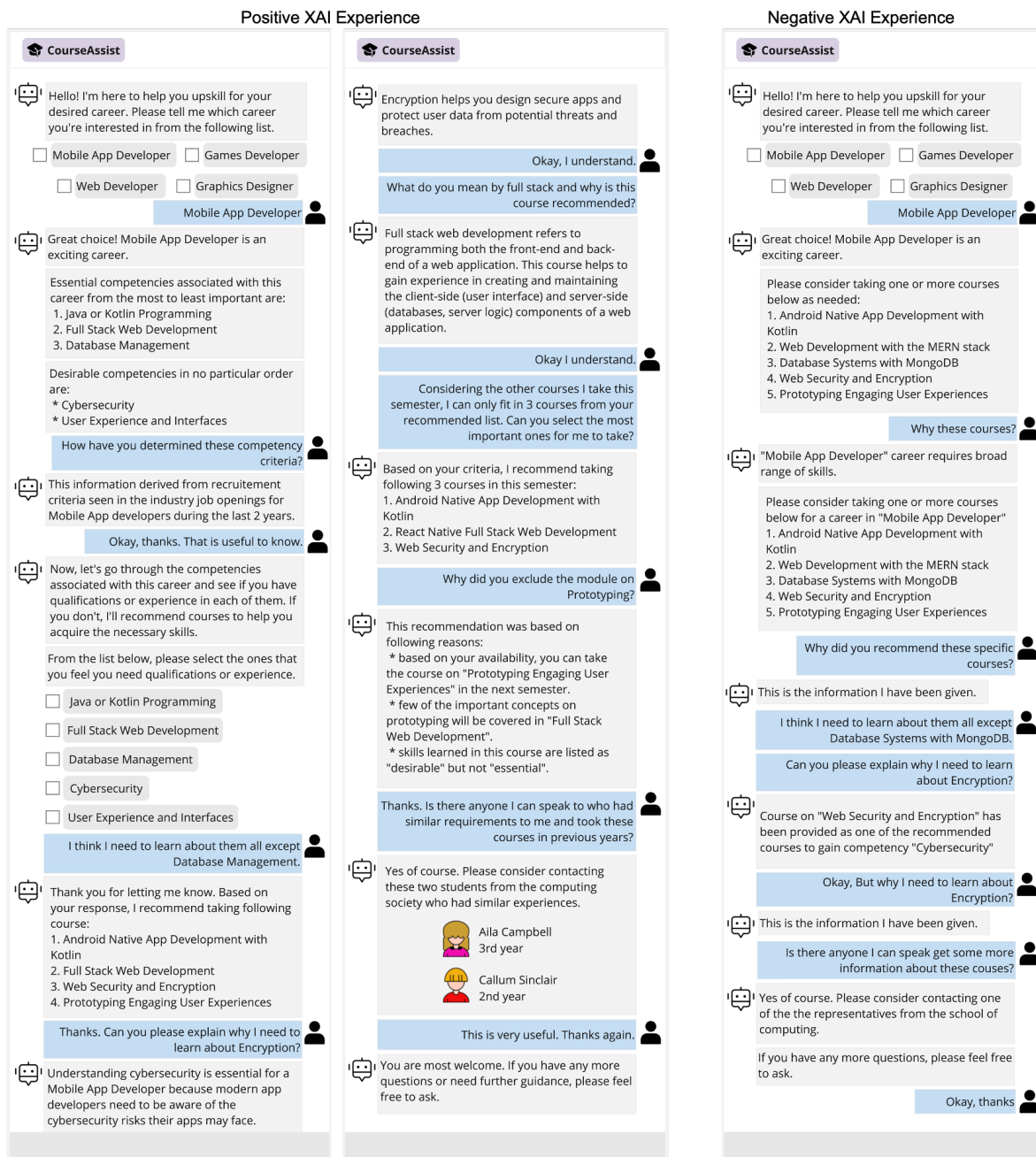
The pilot study instrument consisted of three components: 1) an introduction to the terminology; 2) the explanation experience; and 3) a structured survey for data collection. Similar to the CVR study, introducing the terminology is intended to establish a common understanding among participants. Next participants explored a sample XAI experience stepping into the role of the target user. Finally, participants evaluated the sample XAI experience by responding to the 18 items each on a 5-step Likert scale ranging from *I strongly agree* to *I strongly disagree*.

#### 5.1.1. XAI experiences development and selection

The pilot study is large-scale (100+ participants) and aims, among other goals, to establish discriminant validity. To achieve this, we identified the following practical considerations for designing sample experiences: a) ensuring access to a large pool of target users via survey recruitment platforms or within the research institute; b) curating a relatively negative experience counterpart for each sample for validating the scale’s discriminant properties in a controlled setting; and c) ensuring access to the same sample over longitudinal follow-up to validate the scale’s test-retest reliability.

Based on these considerations, we selected the CourseAssist and AssistHub systems for sample experiences from the point of view of a student and a welfare applicant respectively. Negative experiences, in contrast to positive ones, were deliberately designed to lack best practices recommended in recent human-centred XAI literature [43, 44, 7, 45]. Specifically, they lacked: a) diversity in explanations within the experience; b) alignment between the target user’s intent and the explanations provided; c)





**Figure 2:** The static previews of the relatively positive (left 2 columns) and negative (right column) XAI Experiences with the CourseAssist Chatbot designed for the Pilot Study.

alignment between the user's knowledge level and the explanations; and d) mechanisms to encourage engagement with the XAI agent. The resulting four sample experiences are identified as CourseAssist-P, CourseAssist-N, AssistHub-P, AssistHub-N respectively where P and N indicate positive and negative experiences. Figure 2 presents the static views of CourseAssist-P and CourseAssist-N. AssistHub samples are included in Supplementary Material.

## 5.2. Participant recruitment and selection

The inclusion and exclusion criteria for participant recruitment to evaluate the CourseAssist and AssistHub experiences are summarized in Table 3. After getting through the inclusion criteria, each participant was randomly assigned to assess either the positive or negative sample experience in

their respective application. After preliminary review, we excluded 5, 1 and 1 responses from groups CourseAssist-P, AssistHub-P and AssistHub-N who failed the following attention checks: 1) spend less than half of the allocated time; and/or 2) responded to the questionnaire in a pattern. This yielded 68, 70, 50, and 50 responses for the CourseAssist-P, CourseAssist-N, AssistHub-P, and AssistHub-N cohorts.

**Table 3**

Recruitment details for the pilot study; The only exclusion criterion was that participants of the CourseAssist study were not eligible for AssistHub, and vice versa.

Identifier	Stakeholder	Inclusion criteria	Source	Sample
CourseAssist	Student	Age – between 18 and 30; Current education level – Undergraduate degree;	Institute	68
		Degree subjects – Mathematics and statistics, Information and Communication Technologies, Natural sciences; Language of instruction – English	Prolific	70
AssistHub	Welfare applicant	Age – above 30; Household size – 3 or larger; Property ownership – social housing or affordable-rented accommodation; Employment status – part-time, due to start a new job within the next month, unemployed, or not in paid work (e.g. homemaker or retired)	Prolific	100

### 5.3. Data collection and analysis

The study was hosted on the Jisc Online Surveys platform. Internal consistency, discriminant validity and construct validity were evaluated with both CourseAssist and AssistHub experiences with 238 participants. Among the 68 participants from the research institute who evaluated CourseAssist experiences, 35 took part in the test-retest validity study. AssistHub experiences were excluded from the test-retest validity study due to the practical challenges of following up and re-identifying participants for retesting on the Prolific platform. The duration between the test and retest was 4 weeks.

For analysis, we introduce the following notations. Given  $r_j^i$  is the participant  $j$ 's response to item  $i$ , the participant's total is  $r_j$  and the item total is  $r^i$ . We transform 5-step Likert responses to numbers as follows: Strongly Disagree-1, Somewhat Disagree-2, Neutral-3, Somewhat Agree-4, and Strongly Agree-5. Accordingly, for the 18-item XEQ Scale,  $r_j \leq 90$  ( $5 \times 18$ ).

#### 5.3.1. Internal Consistency

Internal consistency refers to the degree of inter-relatedness among items within a scale. We employ the following metrics from psychometric theory to assess the XEQ Scale items.

**Item-Total Correlation** calculates the Pearson correlation coefficient between the item score and the total score the expected value per item is  $\geq 0.50$ . The Item-Total Correlation of item  $i$ ,  $iT$  is calculated as follows.

$$iT = \frac{\sum_{j=1}^N (r_j^i - \bar{r}^i)(r_j - \bar{r})}{\sqrt{\sum_{j=1}^N (r_j^i - \bar{r}^i)^2 \sum_{j=1}^N (r_j - \bar{r})^2}}$$

Here  $\bar{r}^i$  is the average response score for the  $i$ th item, and  $\bar{r}$  is the overall response average.

**Inter-Item Correlation** is a measure of the correlation between different items within a scale and values between 0.2 and 0.8 are considered expected since  $\geq 0.80$  indicate redundancy and  $\leq 0.20$  indicate poor item homogeneity. The calculation is similar to the previous but is between two items.

**Cronbach's alpha** measures the extent to which all items in a scale are measuring the same underlying construct [46]. High internal consistency is indicated by  $\alpha \geq 0.7$ . If  $s^i$  is the standard deviation of responses to item  $i$ , and  $s$  is the standard deviation of response totals,  $\alpha$  is calculated as follows.

$$\alpha = \frac{M}{M-1} \left( 1 - \frac{\sum_{i=1}^M (s^i)^2}{s^2} \right)$$

As such this helps to quantify how much of the total variance is due to the shared variance among items, which reflects their consistency in measuring the same underlying construct.

### 5.3.2. Discriminant Validity

Discriminant validity measures the ability of the scale to discern between positive and negative experiences and we used the following two methods.

**Discriminant Analysis** treats the pilot study responses as a labelled dataset to train a classification model with a linear decision boundary. The items are considered as features and the group (A or B) is considered as the label. A holdout set then evaluates the model's ability to distinguish between groups A and B.

**Parametric Statistical Test** uses a mixed-model ANOVA test to measure if there is a statistically significant difference between the two groups A and B (agnostic of the domain). Our null hypothesis is "no significant difference is observed in the mean participant total between groups A and B". Our sample sizes meet the requirements for a parametric test determined by an a priori power analysis using G\*Power [47].

### 5.3.3. Construct Validity

Construct validity evaluates the degree to which the scale assesses the characteristic of interest, i.e. factors [48]. Via factor analysis we aim to uncover underlying factors (i.e. dimensions) and validate them. Two forms of factor analysis applied are:

**Exploratory Factor Analysis (EFA)** finds the number of underlying factors in the scale by assessing the variance explained through the Principal Component Analysis (PCA) coefficients (i.e. eigenvalues).

**Confirmatory Factor Analysis (CFA)** tests a pre-defined factor model, with factor loadings expected to be  $>0.5$  to indicate strong support. For the XEQ scale, the evaluation dimensions assigned in Table 4 are validated as a 4-factor model.

### 5.3.4. Test-retest Reliability

The consistency and stability of the XEQ scale over time are measured with the test-retest reliability. We provide quantitative evidence by applying the following metrics.

**Pearson Correlation Coefficient** measures the correlation between scale scores at test and retest. Let  $r_j$  and  $r'_j$  represent participant  $j$ 's scale scores at the first instance (test) and the second instance (retest), respectively, with  $\bar{r}$  and  $\bar{r}'$  being the mean scale scores for all participants at each time point. A  $\rho > 0.7$  indicates strong test-retest reliability, suggesting high consistency between the two measurement instances.

$$\rho = \frac{\sum_{j=1}^N (r_j - \bar{r})(r'_j - \bar{r}')}{\sqrt{\sum_{j=1}^N (r_j - \bar{r})^2 \sum_{j=1}^N (r'_j - \bar{r}')^2}}$$

**Table 4**

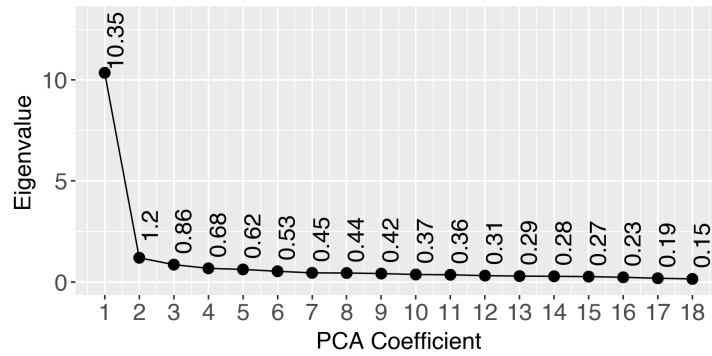
Results: 18 items were selected from the item bank based on the Item-Level Content Validity Index(I-CVI>0.75) and were assigned an evaluation dimension. From the pilot studies, we find Item-Total Correlation (iT)>0.5 for all 18 items. One-factor loadings all over 0.5 confirm the overarching single factor within the XEQ Scale. Confirmatory (C.) Factor Analysis shows significant loading (>0.5) from each item for their respective evaluation dimension.

#	Item	I-CVI	iT	One-Factor Loading	Evaluation Dimension	C. Factor Loading
1	The explanations received throughout the experience were consistent <sup>†</sup> .	1.0000	0.6274	0.6076	Engagement	0.7988
2	The experience helped me understand the reliability of the AI system.	1.0000	0.6416	0.6300	Learning	0.9409
3	I am confident about using the AI system.	1.0000	0.7960	0.7790	Utility	1.0096
4	The information presented during the experience was clear.	1.0000	0.7666	0.7605	Learning	1.1634
5	The experience was consistent with my expectations <sup>†</sup> .	0.9231	0.7959	0.7831	Fulfilment	1.0822
6	The presentation of the experience was appropriate for my requirements <sup>†</sup> .	0.9231	0.8192	0.8083	Fulfilment	0.9300
7	The experience has improved my understanding of how the AI system works.	0.9231	0.6169	0.5859	Learning	0.8280
8	The experience helped me build trust in the AI system.	0.9231	0.7160	0.7018	Learning	0.9787
9	The experience helped me make more informed decisions.	0.9231	0.7460	0.7279	Utility	0.7909
10	I received the explanations in a timely and efficient manner.	0.8462	0.7015	0.6841	Engagement	0.8592
11	The information presented was personalised to the requirements of my role <sup>†</sup> .	0.8462	0.7057	0.6801	Utility	0.7938
12	The information presented was understandable within the requirements of my role <sup>†</sup> .	0.8462	0.7876	0.7803	Utility	1.1452
13	The information presented showed me that the AI system performs well <sup>†</sup> .	0.8462	0.8112	0.8016	Fulfilment	0.9071
14	The experience helped to complete the intended task using the AI system.	0.8462	0.8299	0.8241	Utility	1.0787
15	The experience progressed sensibly <sup>†</sup> .	0.7692	0.8004	0.7912	Engagement	1.1588
16	The experience was satisfying.	0.7692	0.7673	0.7529	Fulfilment	1.0309
17	The information presented during the experience was sufficiently detailed.	0.7692	0.8168	0.8035	Utility	1.0299
18	The experience provided answers to all of my explanation needs.	0.7692	0.8472	0.8444	Fulfilment	0.9041

**Intra-class Correlation Coefficient (ICC)** assesses the consistency of scale scores from the same group of participants between test and retest, using a two-way mixed-effects model (i.e. ICC3). Given that each participant  $j$  has a mean score of  $\bar{r}_j$  between test and retest, and the overall mean across all participants is  $\bar{r}$  the ICC is computed as:

$$ICC = \frac{ms - ms_e}{ms + (k - 1) \times ms_e}$$

, where  $ms$  is the mean square between participants,  $ms = \sum_{j=1}^N (\bar{r}_j - \bar{r})^2 / (N - 1)$  and  $ms_e$  is the mean square error  $ms_e = \sum_{j=1}^N [(r_j - \bar{r}_j)^2 + (r'_j - \bar{r}_j)^2] / [N \times (k - 1)]$ . An ICC value above 0.75 indicates good reliability, and above 0.90 reflects excellent reliability, indicating strong consistency in scale scores between test and retest.



**Figure 3:** Scree plot of eigenvalues derived from the Principal Component Analysis (PCA) performed for the exploratory factor analysis (EFA); a sharp decline in eigenvalues indicates the presence of a single overarching factor.

## 5.4. Results

### 5.4.1. Internal Consistency

Table 4 column  $iT$  reports the Item-Total Correlation. All items met the baseline criteria of  $iT \geq 0.5$  and baseline criteria for Inter-Item correlation. Cronbach’s alpha is 0.9562 which also indicates strong internal consistency.

### 5.4.2. Discriminant Validity

We performed discriminant analysis over 100 trials where at each trial a different train-test split of the responses was used. Each trial used a stratified split, with 70% of the responses for training and 30% for testing. Over the 100 trials, we observed accuracy of  $0.63 \pm 0.05$  and a macro F1-score of  $0.63 \pm 0.05$  which is significantly over the baseline accuracy of 0.50 for a binary classification task. Mixed-model ANOVA test showed a statistically significant difference between groups P and N with a p-value of  $1.63e - 12$  where the mean participant total for groups P and N were  $70.96 \pm 0.47$  and  $57.97 \pm 1.84$ . Also, it revealed a substantial variability within groups indicated by the group variance of 104.86, which we account for including responses from two application domains. Furthermore, Cohen’s d was 1.7639 which indicates a large effect size confirming a significant difference between groups A and B. A standard t-test also obtained a p-value of  $1.13e - 09$  further verifying the statistically significant difference. Based on this evidence we reject the null hypothesis and confirm the discriminant validity of the scale.

### 5.4.3. Construct Validity

We first explore the number underlying factors in the XEQ Scale using EFA. Figure 3 presents the eigenvalues for PCA coefficients derived from scale responses which shows a sharp drop and plateau of eigenvalues for coefficient 2 onwards. This signifies a single overarching factor evident throughout the scale which we refer to as “XAI Experience Quality”. We further confirm this single factor via one-factor loadings (i.e. CFA with a single factor model) given in Table 4 where all item supports are  $> 0.5$ .

An additional CFA was conducted to confirm that evaluation dimensions had substantial factor loadings from their respective items. The rightmost column in Table 4 presents the CFA factor loadings. Each item meets the required factor loading of  $\geq 0.5$  for their evaluation dimension. These observations reinforce that while there is an over-arching factor on “XAI Experience Quality”, it is substantially underpinned by the four evaluation dimensions *Learning*, *Utility*, *Fulfilment* and *Engagement*.

#### 5.4.4. Test-retest Reliability

The test-retest reliability of the scale was assessed with a 4-week interval using the two CourseAssist experiences. Pearson correlation coefficient  $\rho = 0.7561$  with a p-value of  $1.5e - 07$ , indicated a strong linear relationship between the test and retest scores. Additionally,  $ICC = 0.8609$  with a p-value of  $3.17e - 08$ , falls within the range that indicates strong consistency. Both measures establish the test-retest reliability of the XEQ scale.

This concludes our multi-faceted evaluation and refinement of the XEQ Scale based on pilot study results.

## 6. Discussion

### 6.1. Implications and limitations

In psychometric theory, conducting a pilot study involves administering both the scale under development and existing scales to participants. The objective is to assess the correlation between the new scale and those found in existing literature, particularly in shared dimensions. However, our pilot studies did not incorporate this, since to the best of our knowledge there are no previous scales that measured the XAI experience quality or multi-shot XAI experiences. While the System Causability Scale [16] is the closest match in the literature, it was not applicable as it featured in the initial items bank. Also, the current pilot studies had limited variability in application domains. To address this limitation, we are currently planning pilot studies with two medical applications: 1) fracture prediction in Radiographs; and 2) liver disease prediction from CT scans. In the future, we will further validate and refine the scale as necessary.

### 6.2. XEQ scale in practice

We propose the XEQ Scale as a tool to support the development and evaluation of interactive user-centric XAI systems. We anticipate calculating 3 aggregate metrics based upon responses to the 18 items in the scale:

**Stakeholder XEQ Score** quantifies individual experiences and is calculated as the mean of stakeholder's responses to all items.

**Factor Score** quantifies the quality along each evaluation dimension and is calculated as the mean of all responses to the respective subset of items. For XAI system designers, a lower factor score indicates the dimensions that need improvement.

**System XEQ Score** quantifies the XAI experience quality of the system as a whole and is calculated as the aggregate of Factor scores. The System XEQ Score helps the XAI system designers to iteratively develop a well-rounded system grounded in user experience. The designers can also choose to assign a higher weight to one or a subset of dimensions that they find important at any iteration. System XEQ Score can also be utilised by external parties (e.g. regulators, government), either to evaluate or benchmark XAI systems.

We have created the following user story to demonstrate the usefulness of the XEQ Scale for the iterative development and evaluation of an applied XAI system.

*“Jane is the systems manager for a motor insurance company. The company have recently introduced an AI system for calculating an insurance quote using applicants' immediate driving history and previous claims. As part of the company's regulatory obligations, the AI system is supported by several explanation algorithms. However, the systems management team are having complaints from sales staff that they struggle to justify system outcomes using these explanations. Jane therefore decides to use the XEQ Scale to measure the XE for 2*

*stakeholder groups: customers and sales staff. They start by recruiting a representative user group for each role. Subsequently, each participant used the XAI system and scored their experience with the XEQ Scale.”*

*“ Jane collates the responses and calculates the stakeholder scores, factor scores and system scores for each role. When comparing system scores, Jane finds that sales staff consistently score the XAI system lower compared to customers. Looking at granular factor scores, they identify that the Utility factor score is particularly low for sales staff, specifically, item 14 (‘The experience helped me to complete the intended task using the AI system’) received consistently negative feedback. From these insights, Jane recognises that an additional explanation algorithm should be added to support sales staff in task completion. Jane updated the explanation algorithms and trials once again with a sample of sales staff, finding that the new system increased the Utility factor score and overall, led to an improved system XEQ score.”*

From the above user story, it can be observed that the XEQ Scale can support the iterative development of XAI systems by facilitating the targeted update via user-centred evaluation.

### **6.3. XEQ Benchmark Development**

The next phase for the XEQ scale entails developing a benchmark for XAI experience quality. This process includes administering the XEQ scale to over 100 real-world AI systems that provide explainability to stakeholders and establishing a classification system. We are currently following the established benchmark maintenance policy of the User Experience Questionnaire [49] where we develop and release an XEQ Analysis tool with the benchmark updated regularly.

We envision when the scale is administered to stakeholders of a new XAI system, the benchmark will categorise the new system based on the mean participant total in each evaluation dimension as follows - *Excellent*: Within the top 10% of XAI systems considered in the benchmark; *Good*: Worse than the top 10% and better than the lower 75%; *Above average*: Worse than the top 25% and better than the lower 50%; *Below average*: Worse than the top 50% and better than the lower 25%; and *Bad*: Within the 25% worst XAI systems. Accordingly, the XEQ benchmark will enable XAI system owners to iteratively enhance the XAI experience offered to their stakeholders.

## **7. Conclusions**

In this paper, we presented the XEQ scale. The XEQ scale provides a framework for the comprehensive evaluation of user-centred XAI experiences. It fills a novel gap in the evaluation of multi-shot explanations which is currently not adequately fulfilled by any other evaluation metric(s). Throughout this paper, we have described the development and validation of the scale following psychometric theory. We make this scale available as a public resource for evaluating the quality of XAI experiences. In future work, we plan to investigate the generalisability of the XEQ scale on additional domains, AI systems and stakeholder groups. Beyond this, we propose to establish a benchmark using the XEQ scale. Our goal is to facilitate the user-centred evaluation of XAI and support the emerging development of best practices in the explainability of autonomous decision-making.

## **Acknowledgments**

The authors thank all participants in the CVR and pilot studies. iSee is an EU CHIST-ERA project that received funding for the UK from EPSRC under grant number EP/V061755/1, for Ireland from the Irish Research Council under grant number CHIST-ERA-2019-iSee, for France from the French National Research Agency under grant number ANR-21-CHR4-0004, and for Spain from the MCIN/AEI and European Union “NextGenerationEU/PRTR” under grant number PCI2020-120720-2.

## References

- [1] T. Miller, Explanation in artificial intelligence: Insights from the social sciences, *Artificial intelligence* 267 (2019) 1–38.
- [2] B. Hu, P. Tunison, B. Vasu, N. Menon, R. Collins, A. Hoogs, Xaitk: The explainable ai toolkit, *Applied AI Letters* 2 (2021) e40.
- [3] V. Arya, R. K. Bellamy, P.-Y. Chen, A. Dhurandhar, M. Hind, S. C. Hoffman, S. Houde, Q. V. Liao, R. Luss, A. Mojsilović, et al., Ai explainability 360 toolkit, in: *Proceedings of the 3rd ACM India Joint International Conference on Data Science & Management of Data (8th ACM IKDD CODS & 26th COMAD)*, 2021, pp. 376–379.
- [4] Q. V. Liao, D. Gruen, S. Miller, Questioning the ai: Informing design practices for explainable ai user experiences, in: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI '20, Association for Computing Machinery, New York, NY, USA, 2020, p. 1–15. URL: <https://doi.org/10.1145/3313831.3376590>. doi:10.1145/3313831.3376590.
- [5] A. Wijekoon, D. Corsar, N. Wiratunga, K. Martin, P. Salimi, Tell me more: Intent fulfilment framework for enhancing user experiences in conversational xai, *arXiv preprint arXiv:2405.10446* (2024).
- [6] A. Wijekoon, N. Wiratunga, K. Martin, D. Corsar, I. Nkisi-Orji, C. Palihawadana, D. Bridge, P. Pradeep, B. D. Agudo, M. Caro-Martínez, Cbr driven interactive explainable ai, in: *International Conference on Case-Based Reasoning*, Springer, 2023, pp. 169–184.
- [7] L. Malandri, F. Mercorio, M. Mezzanzanica, N. Nobani, Convxai: a system for multimodal interaction with any black-box explainer, *Cognitive Computation* 15 (2023) 613–644.
- [8] B. Finzel, D. E. Tafler, S. Scheele, U. Schmid, Explanation as a process: user-centric construction of multi-level and multi-modal explanations, in: *KI 2021: Advances in Artificial Intelligence: 44th German Conference on AI*, Virtual Event, September 27–October 1, 2021, *Proceedings 44*, Springer, 2021, pp. 80–94.
- [9] A. Rosenfeld, Better metrics for evaluating explainable artificial intelligence, in: *Proceedings of the 20th international conference on autonomous agents and multiagent systems*, 2021, pp. 45–50.
- [10] C.-K. Yeh, C.-Y. Hsieh, A. Suggala, D. I. Inouye, P. K. Ravikumar, On the (in) fidelity and sensitivity of explanations, *Advances in Neural Information Processing Systems* 32 (2019).
- [11] P. Q. Le, M. Nauta, V. B. Nguyen, S. Pathak, J. Schlötterer, C. Seifert, Benchmarking explainable ai: a survey on available toolkits and open challenges, in: *International Joint Conference on Artificial Intelligence*, 2023.
- [12] S. Mohseni, N. Zarei, E. D. Ragan, A multidisciplinary survey and framework for design and evaluation of explainable ai systems, *ACM Transactions on Interactive Intelligent Systems (TiiS)* 11 (2021) 1–45.
- [13] S. S. Y. Kim, E. A. Watkins, O. Russakovsky, R. Fong, A. Monroy-Hernández, "help me help the ai": Understanding how explainability can support human-ai interaction, in: *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, Association for Computing Machinery, New York, NY, USA, 2023. URL: <https://doi.org/10.1145/3544548.3581001>. doi:10.1145/3544548.3581001.
- [14] U. Ehsan, S. Passi, Q. V. Liao, L. Chan, I.-H. Lee, M. Muller, M. O. Riedl, The who in xai: How ai background shapes perceptions of ai explanations, in: *Proceedings of the CHI Conference on Human Factors in Computing Systems*, CHI '24, Association for Computing Machinery, New York, NY, USA, 2024. URL: <https://doi.org/10.1145/3613904.3642474>. doi:10.1145/3613904.3642474.
- [15] R. R. Hoffman, S. T. Mueller, G. Klein, J. Litman, Measures for explainable ai: Explanation goodness, user satisfaction, mental models, curiosity, trust, and human-ai performance, *Frontiers in Computer Science* 5 (2023) 1096257.
- [16] A. Holzinger, A. Carrington, H. Müller, Measuring the quality of explanations: the system causability scale (scs) comparing human and machine explanations, *KI-Künstliche Intelligenz* 34 (2020) 193–198.
- [17] G. O. Boateng, T. B. Neilands, E. A. Frongillo, S. L. Young, Best practices for developing and



- validating scales for health, social, and behavioral research: a primer, *Frontiers in public health* 6 (2018) 366616.
- [18] M. Madsen, S. Gregor, Measuring human-computer trust, in: 11th australasian conference on information systems, volume 53, Citeseer, 2000, pp. 6–8.
- [19] G. C. Moore, I. Benbasat, Development of an instrument to measure the perceptions of adopting an information technology innovation, *Information systems research* 2 (1991) 192–222.
- [20] J.-Y. Jian, A. M. Bisantz, C. G. Drury, Foundations for an empirically determined scale of trust in automated systems, *International journal of cognitive ergonomics* 4 (2000) 53–71.
- [21] J. Brooke, et al., Sus-a quick and dirty usability scale, *Usability evaluation in industry* 189 (1996) 4–7.
- [22] B. Cahour, J.-F. Forzy, Does projection into use improve trust and exploration? an example with a cruise control system, *Safety science* 47 (2009) 1260–1270.
- [23] R. G. Netemeyer, *Scaling procedures: Issues and applications*, Sage Publications, 2003.
- [24] B. D. Adams, L. E. Bruyn, S. Houde, P. Angelopoulos, K. Iwasa-Madge, C. McCann, *Trust in automated systems*, Ministry of National Defence (2003).
- [25] K. Schaefer, *The perception and measurement of human-robot trust* (2013).
- [26] J. Zhou, A. H. Gandomi, F. Chen, A. Holzinger, Evaluating the quality of machine learning explanations: A survey on methods and metrics, *Electronics* 10 (2021) 593.
- [27] D. V. Carvalho, E. M. Pereira, J. S. Cardoso, Machine learning interpretability: A survey on methods and metrics, *Electronics* 8 (2019) 832.
- [28] M. Nauta, C. Seifert, The co-12 recipe for evaluating interpretable part-prototype image classifiers, in: *World Conference on Explainable Artificial Intelligence*, Springer, 2023, pp. 397–420.
- [29] R. W. Lissitz, S. B. Green, Effect of the number of scale points on reliability: A monte carlo approach., *Journal of applied psychology* 60 (1975) 10.
- [30] R. Hasan, R. Weil, R. Siegel, K. Krombholz, A psychometric scale to measure individuals' value of other people's privacy (vopp), in: *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 2023, pp. 1–14.
- [31] D. Votipka, D. Abrokwa, M. L. Mazurek, Building and validating a scale for secure software development self-efficacy, in: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, 2020, pp. 1–20.
- [32] J. C. Nunnally, I. H. Bernstein, *Psychometric Theory: Nunnally and Bernstein*, McGraw Hill, 3rd edition, 1993.
- [33] C. H. Lawshe, et al., A quantitative approach to content validity, *Personnel psychology* 28 (1975) 563–575.
- [34] H. Gehlbach, M. E. Brinkworth, Measure twice, cut down error: A process for enhancing the validity of survey scales, *Review of general psychology* 15 (2011) 380–387.
- [35] W. Jin, X. Li, M. Fatehi, G. Hamarneh, Guidelines and evaluation of clinical explainable ai in medical image analysis, *Medical image analysis* 84 (2023) 102684.
- [36] M. Miró-Nicolau, G. Moyà-Alcover, A. Jaume-i Capó, Evaluating explainable artificial intelligence for x-ray image analysis, *Applied Sciences* 12 (2022) 4459.
- [37] H. Müller, A. Holzinger, M. Plass, L. Brcic, C. Stumptner, K. Zatloukal, Explainability and causability for artificial intelligence-supported medical image analysis in the context of the european in vitro diagnostic regulation, *New Biotechnology* 70 (2022) 67–72.
- [38] M. Lang, A. Bernier, B. M. Knoppers, Artificial intelligence in cardiovascular imaging: “unexplainable” legal and ethical challenges?, *Canadian Journal of Cardiology* 38 (2022) 225–233.
- [39] W. G. De Sousa, E. R. P. de Melo, P. H. D. S. Bermejo, R. A. S. Farias, A. O. Gomes, How and where is artificial intelligence in the public sector going? a literature review and research agenda, *Government Information Quarterly* 36 (2019) 101392.
- [40] H. de Bruijn, M. Warnier, M. Janssen, The perils and pitfalls of explainable ai: Strategies for explaining algorithmic decision-making, *Government information quarterly* 39 (2022) 101666.
- [41] S. Brdnik, V. Podgorelec, B. Šumak, Assessing perceived trust and satisfaction with multiple

- explanation techniques in xai-enhanced learning analytics, *Electronics* 12 (2023) 2594.
- [42] B. S. Rousse, S. E. Dreyfus, Revisiting the six stages of skill acquisition, *Teaching and learning for adult skill acquisition: applying the Dreyfus & Dreyfus model in different fields* (2021) 3–28.
- [43] Q. V. Liao, D. Gruen, S. Miller, Questioning the ai: informing design practices for explainable ai user experiences, in: *Proceedings of the 2020 CHI conference on human factors in computing systems*, 2020, pp. 1–15.
- [44] H. Suresh, S. R. Gomez, K. K. Nam, A. Satyanarayan, 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*, 2021, pp. 1–16.
- [45] R. Nimmo, M. Constantinides, K. Zhou, D. Quercia, S. Stumpf, User characteristics in explainable ai: The rabbit hole of personalization?, in: *Proceedings of the CHI Conference on Human Factors in Computing Systems*, 2024, pp. 1–13.
- [46] L. J. Cronbach, *Essentials of psychological testing*. (1949).
- [47] F. Faul, E. Erdfelder, A.-G. Lang, A. Buchner, G\* power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences, *Behavior research methods* 39 (2007) 175–191.
- [48] A. E. Kazdin, *Research design in clinical psychology*, Cambridge University Press, 2021.
- [49] M. Schrepp, J. Thomaschewski, A. Hinderks, Construction of a benchmark for the user experience questionnaire (ueq) (2017).

## Appendix

### A: Initial Items Bank

Table 5 presents the initial item bank compiled from literature and reviewed by the research team.

### B: Content Validity Study

This study aimed to establish the content validity of the XEQ scale with XAI experts. Since XAI Experiences are a novel concept, we included three example XAI experiences that capture a variety of stakeholder types and application domains. In addition to the CourseAssist chatbot example included in the paper, they were presented with the following experiences in video format.

- The AssistHub AI platform is a website for processing welfare applications and is used by a local council to accelerate the application process. An auditor is exploring the website and its XAI features to understand the fairness and bias of the AI system being used in the decision-making process. A non-interactive preview of the experience is presented in Figure 4.
- RadioAssist AI platform is a desktop application, used by the local hospital to support clinicians in their clinical decision-making processes. An AI system predicts the presence of fracture in Radiographs and explains its decisions to the clinicians. A non-interactive preview of the experience is presented in Figure 5.

Both Figures 5 and 4 are annotated with notes that describe the XAI features that were available to the stakeholder. Finally, Figure 6 presents a preview of the Study page.

### C: Pilot Study

A pilot study was conducted with 238 participants over two application domains where they evaluated either a positive or negative XAI experience. In addition to the CourseAssist chatbot examples provided in the paper, we included two XAI experiences of welfare applicants interacting with the AssistHub AI platform (see Figures 7 and 8). Notes refer to how different aspects of the explanations can lead to a positive or negative XAI experience. Similar to the previous study, all XAI experiences were available

**Table 5**  
Initial Items Bank

---

Item
I like using the system for decision-making.
The information presented during the experience was clear.
The explanations received throughout the experience did not contain inconsistencies.
I could adjust the level of detail on demand.
The experience helped me make more informed decisions.
The experience helped me establish the reliability of the system.
I received the explanations in a timely and efficient manner.
The experience was satisfying.
The experience was suitable for the intended purpose of the system.
I was able to express all my explanation needs.
The experience revealed whether the system is fair.
The experience helped me complete my task using the system.
The experience was consistent with my expectations within the context of my role.
The presentation of the experience was appropriate.
The experience has improved my understanding of how the system works.
The experience helped me understand how to use the system.
The experience was understandable in the context of my role.
The experience helped me build trust in the system.
The experience was personalised to the context of my role.
I could request more detail on demand if needed.
I did not need external support to understand the explanations.
The experience was helpful to achieve my goals.
The experience progressed logically.
The experience was consistent with my understanding of the system.
The duration of the experience was appropriate within the context of my role.
The experience improved my engagement with the system.
The experience was personalised to my explanation needs.
Throughout the experience, all of my explanation needs were resolved.
The experience showed me how accurate the system is.
All parts of the experience were suitable and necessary.
The information presented during the experience was sufficiently detailed for my understanding of the domain.
I am confident in the system.

---

to participants in video format. Finally Figure 9 presents a preview of the Pilot study where pages 1 and 2 were customised based on the application participants were assigned to.

**AssistHub**

**Applicant 34635**

**Personal Information**  
 Full Name: John Smith  
 Date of Birth: 12/09/1980  
 National Insurance Number: AB1234635  
 Postcode: AB24 7GF

**Financial Information**  
 Income: £1,800/month (wages), £90/month (child benefit)  
 Assets and Resources: £4,000 in a current account  
 Expenses: Rent: £1,200/month, Utilities: £100/month, Childcare: £150/month

**Employment**  
 Employment Status: Employed  
 Employer Information: DPO Ltd, AB11 1CD  
 Work Hours and Income: 35 hours/week, £10/hour

**AI Prediction: Declined**

The AI model has made a decision to decline your application with a confidence level of 65%. This means that the AI is 65% certain that after evaluation, the council benefits officer is going to decline the request for benefits. There is 35% likelihood that the AI prediction is incorrect.

**AI INSIGHTS**

Provide descriptive information to the stakeholder about the AI system, features of system and role of the AI system in their day-to-day life.

**AI INSIGHTS**

Explore data used by the AI system

Explore prototypical application scenarios

Explore outlier application scenarios

Explore the high-level decision tree

Provide unbiased information about the data that has been used to implement the AI system.

The AI system has been developed based on 10,500 diverse set of applications containing a wide range of scenarios from past benefit applications processed in the UK. This data is anonymised, and system does not consider personal information, and demographic information related race and religion in its decision making process. The past applications can be categorised to 7 benefit types. Some statistics about the data is presented below.

This system is not applicable to process Veteran benefits, healthcare and disability benefits, social security and legal benefits or any other benefit type.

**AI INSIGHTS**

Explore data used by the AI system

Explore prototypical application scenarios

Explore outlier application scenarios

Explore the high-level decision tree

Provide prototypical cases that explains the AI system decision making process.

Here are two past applications where the AI system predictions were "Approved" with over 95% confidence level.

**Applicant Data 18293**

**Applicant Data 00037**

Key similarities in past approved applications are Household Size larger than or equal to 4, Highest education qualification is A-levels and Citizenship Status is British citizen.

Here are two past applications where the AI system predictions were "Declined" with over 95% confidence level.

**Applicant Data 43573**

**Applicant Data 10182**

Key similarities in past declined applications are Household Size less than or equal to 3, Income is greater than £1,500/month and Assets include house ownership.

**AI INSIGHTS**

Explore data used by the AI system

Explore prototypical application scenarios

Explore outlier application scenarios

Explore the high-level decision tree

Provide outlier cases that explains the weaknesses in the AI system decision making process.

Here are two applications where the AI system predictions were "Approved" with only 51-52% confidence level.

**Applicant Data 26172**

**Applicant Data 00233**

In these applications, the details are very similar to "prototypical-declined" applications seen in past data and we have highlighted the one or two details that were different and led to an "approved" outcome.

Here are two past applications where the AI system predictions were "Declined" with only 51-52% confidence level.

**Applicant Data 00081**

**Applicant Data 00001**

In these applications, the details are very similar to "prototypical-accepted" applications seen in past data and we have highlighted the one or two details that were different and led to a "Declined" outcome.

**AI INSIGHTS**

Explore data used by the AI system

Explore prototypical application scenarios

Explore outlier application scenarios

Explore the high-level decision tree

Provide simplified view of the high-level decision-making process.

The past applications were used by the AI system to learn a prediction model using the XGBOOST technique. Here is a visualisation of the high-level decision tree that approximates the rules learned by the prediction model.

Eligibility criteria include citizenship = British, residence = the UK and the consent to verify information is given. Other information that are considered in the decision making process are household and financial details. At each node, a score is calculated to decide the next steps. To find more about score calculations at each node please [click here](#).


```

graph TD
    Root[Eligibility Criteria] -- Yes --> Node1[Household details]
    Root -- No --> Declined[Declined]
    Node1 -- score <= 0.7 --> Node2[Financial details]
    Node1 -- score > 0.7 --> Accepted[Accepted]
    Node2 -- score > 0.64 --> Accepted
    Node2 -- score <= 0.64 --> Declined
  
```

**Figure 4:** Positive XAI Experience of a Regulation Officer exploring the AssistHub AI platform; A pop-up is shown when clicked on the *AI INSIGHTS* button, with four navigation pages providing different types of explanations. The help description pop-up appears when clicking on the question mark button.

RadioAssist
 Dr. A. Byrne

**Patient 93023**



Mc

✓ More about this prediction

✓ More about the AI system

Welcome to our AI-powered radiograph analysis tool designed to assist in fracture detection. This tool leverages artificial intelligence to provide accurate insights into radiographs, assisting in diagnostic capabilities. It also offers a unique feature that allows you to explore explanations for AI decisions.

**Key Features:**


- Fracture Detection:** Our AI model is trained on a diverse dataset of radiographs to quickly and accurately identify fractures in X-rays.
- User-Friendly Interface:** Easily access radiographs and receive AI assisted decisions
- Explanations for AI Decisions:** Gain transparency in the decision-making process. Explore detailed explanations to better understand the results.

**Clinical Integration:**  
This tool is designed to complement your expertise, not replace it. It can improve diagnostic accuracy, reduce oversights, and expedite the diagnosis and treatment of fractures. The explanations feature promotes transparency and aids in making informed clinical decisions.

---

RadioAssist
 Dr. A. Byrne

**Patient 93023**



Mc

More information

**AI Prediction: Fracture**

The AI model has made a diagnosis indicating a potential humerus fracture in the radiograph with a confidence level of 81%. This means that the AI is 81% certain that a fracture is present based on its analysis of the image. There is a 19% likelihood that the AI prediction is incorrect.

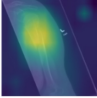
More information

More about this prediction

**Which regions does the AI system identify as potential fracture?**

This explanation is generated using the [Class Activation Map](#) technique.

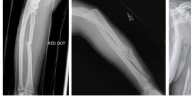
The regions of elevated color intensity within the radiograph signify areas of particular interest that led the AI to predict the presence of a possible humerus fracture.



**Is this the same outcome for similar instances?**

The [nearest neighbour algorithm](#) has chosen three visually similar radiographs.

These radiographs of the upper arm have the same diagnostic prediction of a humerus fracture similar to the current patient's radiograph.



More about the AI system

**How accurate is the AI system?**

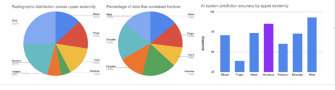
The AI system has an accuracy of 78% in predicting upper extremity fractures. This means that out of 100 predictions, the AI system gets 78 correct, and 22 wrong. This metric was observed when the AI system was tested on a public collection of 3,321 radiographs.

Accuracy: 78.0%  
F1 score: 77.8%

**What evidence has been considered in developing the AI system?**

This AI system was developed based on 40,561 radiographic images that captured upper extremities of 12,173 patients ([MURA Dataset](#)). In this dataset, 45% of the patients were female, and the age range of the patients spanned from 18 to 78 years. The figures display image distribution across upper extremities (left), the percentage of images with fractures in each extremity (middle), and accuracy by upper extremity (right).

Despite constituting just 5% of images, the system has a high prediction accuracy for humerus radiographs compared to Finger or Forearm.



Provide descriptive information to the stakeholder about the AI system, features of system and role of the AI system in their day-to-day work.

Provide information that explains the AI system decision in two forms: 1) feature attributions that led to the current decision; and 2) similar past cases.

Provide unbiased information about the performance of the AI system and the data that has been used to implement the AI system.

**Figure 5:** Positive XAI Experience of a Clinician using the RadioAssist AI platform; The clinician can click on the two minimised pages to expand and view explanations about the AI system and the decision. The question mark button shows the help description.

**XEQ Scale**

0% complete

**Page 1: XAI Experience Quality (XEQ) Scale development and refinement study with XAI practitioners.**

The purpose of this study is to design the XEQ Scale, creating a scale to measure the quality of explanation experiences.

First, we introduce the terminology and definitions related to the study. Please read each of these and make sure you understand them, as they are integral to the study.

Next, you will be presented with three fictitious examples of XAI Experiences. Please explore these carefully before proceeding to the next section.

Lastly, you will take part in the Content Validity Ratio (CVR) study. This study assesses the relevance and authenticity of the criteria that form the basis of the XEQ scale. You will be presented with a list of statements and asked to rate them in terms of relevance and clarity for evaluating the XAI Experience Quality.

The study is conducted anonymously and is estimated to take 15-20 minutes.

[Next >](#)

**XEQ Scale**

25% complete

**Page 2: Preliminaries**

An **XAI System** is an automated decision-making system that has the capacity to provide information about its reasoning.

A **Stakeholder** is an individual or group with a vested interest in the XAI system. Stakeholders encompass a diverse group, ranging from the system designers and developers, who hold an interest in the system's technical functionality, the end consumers relying on its decisions, and regulatory authorities responsible for ensuring fair and ethical use.

**What is an XAI Experience?**

An XAI Experience (XE) is a user-centric process of a stakeholder interacting with an XAI system to gain knowledge and/or improve comprehension.

We envision such experiences involve exploratory interactions with the stakeholder to provide contextual background information about the AI system, as well as deliver explanatory and actionable information about the AI system's decisions. The interactions must be tailored to match stakeholder needs and their level of expertise.

**What is XE Quality?**

The XAI Experience Quality (XEQ) is the extent to which a stakeholder's explanation needs are satisfied by their XE.

**What is the purpose of the XEQ Scale?**

The purpose of the XEQ Scale is to assess the quality of XEs provided by an XAI system.

The XEQ Scale enables a structured evaluation of experiences, ensuring they follow best practices, i.e. fairness, accountability, and transparency.

[< Previous](#) [Next >](#)

**XEQ Scale**

50% complete

**Page 3: XAI Experience Examples**

In this section, we present 3 fictitious scenarios to illustrate different examples of XAI Experiences. Each example presents the interactions between a stakeholder and an XAI system which forms an XAI experience. Please take some time to review each of the examples, as this may help to inform your answers in the following section.

For an optimal viewing experience, expand the video to full screen, adjust the video quality as needed, and feel free to pause or navigate if you require more time during interactions.

[< Previous](#) [Next >](#)

**XEQ Scale**

75% complete

**Page 4: Content Validity Study**

In this section, you will take part in the Content Validity Study as an XAI practitioner.

You are given a set of statements and asked to rate them in terms of **relevance** and **clarity** for evaluating XAI Experience Quality (XEQ). These statements are compiled from existing literature [1,2], with several new additions. They are not presented in any specific order.

Please feel free to refer to the scenarios in the previous section. These scenarios are meant to serve as examples and are *not an exhaustive list of XAI experiences*.

After the completion of this study, the finalised list of statements will form the XEQ Scale. In practice, users will respond to each statement on a 5-point Likert Scale from "Strongly Disagree" to "Strongly Agree" to evaluate their experience.

[1] Holzinger, A., Carrington, A., & Müller, H. (2020). Measuring the quality of explanations: the system causability scale (SCS) comparing human and machine explanations. KI-Künstliche Intelligenz, 34(2), 193-198.

[2] Hoffman, R. R., Mueller, S. T., Klein, G., & Litman, J. (2018). Metrics for explainable AI: Challenges and prospects. arXiv preprint arXiv:1812.04608.

This part of the survey uses a table of questions, [view as separate questions instead?](#)

**1. Relevance for evaluating XAI Experience Quality**

	* Required				
	Not Relevant at all	Somewhat not Relevant	Neutral - can keep it or leave it	Somewhat Relevant	Extremely Relevant
I like using the system for decision-making.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
...					
I am confident in the system.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

This part of the survey uses a table of questions, [view as separate questions instead?](#)

**2. Clarity of the current wording**

	* Required				
	Not Clear at all	Somewhat not Clear	Neutral	Somewhat Clear	Extremely Clear
I like using the system for decision-making.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
The information presented during the experience was clear.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
...					
I am confident in the system.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

**3. If you selected any of the statements as relevant (by selecting Extremely Relevant or Somewhat Relevant) but not clear (by selecting Extremely not Clear or Somewhat not Clear) please propose improved wording. \* Required**

**4. Let us know if you have any other comments. \* Required**

[< Previous](#) [Finish ✓](#)

**Figure 6: Study Preview; Examples removed from Page 3 and list of items shortened in Page 4.**

**AssistHub**

**< Applicant 34635**

Our AI system is here to assist you through the council benefit application process. It leverages the knowledge gained from thousands of past applications to provide insights into potential outcomes. When using our council web application, the AI system offers real-time guidance, helping you complete your application accurately and efficiently. It aims to enhance transparency and streamline the process, making it more user-friendly. While it's not a guarantee of approval, this AI tool is designed to help you make informed decisions throughout your benefit application, ensuring you have the best chance of receiving the support you need.

**AI Prediction: Declined**

The AI model has made a decision to decline your application with a confidence level of 65%. This means that the AI is 65% certain that after evaluation, the council benefits officer is going to decline the request for benefits. There is 35% likelihood that the AI prediction is incorrect.

**AI INSIGHTS**

Provide descriptive information to the stakeholder about the AI system, features of the AI system and role of the AI system in their day-to-day life.

**AI INSIGHTS**

Explore data used by the AI system

Explore the high-level decision tree

Explaining the AI decision

Proposed changes to change the AI decision

**Provide unbiased information about the data that has been used to implement the AI system.**

The AI system has been developed based on 10,500 diverse set of applications containing a wide range of scenarios from past benefit applications processed in the UK. This data is anonymized, and system does not consider personal information, and demographic information related race and religion in its' decision making process. The past applications can be categorised to 7 benefit types. Some statistics about the data is presented below.

Accuracy by benefit category

Benefit Category	Accuracy
Employment	85%
Family Assistance	75%
Financial Assistance	65%
Housing Assistance	55%
Council Tax	45%
Pension	35%
Other	25%

This system is not applicable to process Veteran benefits, healthcare and disability benefits, social security and legal benefits or any other benefit type.

**AI INSIGHTS**

Explore data used by the AI system

Explore the high-level decision tree

Explaining the AI decision

Proposed changes to change the AI decision

**Provide information about the underlying algorithm used by the AI system for decision making**

The past applications were used by the AI system to learn a prediction model using the **XGBoost** technique. Here is a visualisation of the high-level decision tree that approximates the rules learned by the prediction model. Eligibility criteria include citizenship = British, residence = the UK and the consent to verify information is given. Other information that are considered in the decision making process are household and financial details. At each node, a score is calculated to decide the next steps. To find more about score calculations at each node please [click here](#).

**AI INSIGHTS**

Explore data used by the AI system

Explore the high-level decision tree

Explaining the AI decision

Proposed changes to change the AI decision

**Provide information that explains the current AI system decision.**

Following is an illustration of how the application resulted with the current AI decision.

**AI INSIGHTS**

Explore data used by the AI system

Explore the high-level decision tree

Explaining the AI decision

Proposed changes to change the AI decision

**Provide information to the stakeholder on how to get a more desirable decision.**

Making following changes to the current application will result in a change in AI system decision.

- Household details
- Household members who are employed should be None
- Financial details
- Healthcare expenses higher than £500

Figure 7: Relatively positive XAI Experience of a welfare applicant using the AssistHub AI platform

**AssistHub**

**Applicant 34635** Application History

**Personal Information**

Full Name: John Smith  
 Date of Birth: 12/08/1980  
 National Insurance Number: AB1234635  
 Postcode: M20 4PZ

**Household Information**

Household Size: 4  
 Relationship to Household members: Self, Partner (Sarah Smith), Child 1 (Eminia Smith - 5y), Child 2 (Daniel Smith - 7y)

**Financial Information**

Income: £1,800/month (wage), £200/month (child benefit)  
 Assets and Resources: £4,000 in a current account  
 Expenses: Rent - £1,200/month, Utilities - £100/month, Childcare - £150/month

**Employment Information**

Employment Status: Employed  
 Employer Information: DEF Ltd, AB11 1CD  
 Work Hours and Income: 35 hours/week, £10/hour

**Education and Training**

**Health and Disability**

**Citizenship and Residency**

Childcare and Child Custody

**Criminal History**

Documentation and Additional Comments

**AI Prediction: Declined**

The AI model has made a decision to decline your application with a confidence level of 65%. This means that the AI is 65% certain that after evaluation, the council benefits officer is going to decline the request for benefits. There is 35% likelihood that the AI prediction is incorrect.

**AI INSIGHTS**

Stakeholders do not have information about the AI system, features of system or role of the AI system in their day-to-day life.

**AI INSIGHTS**

Explore data used by the AI system >>> The AI system has been developed based on 10,500 applications containing a wide range of scenarios from past benefit applications processed in the UK. Some statistics about the data is presented below.

Explore the high-level decision tree >>>

Proposed changes to change the AI decision >>>

**Information about the data used to implement the AI system is not well explained.**

This system is not applicable to process Veteran benefits, healthcare and disability benefits, social security and legal benefits or any other benefit type.

**AI INSIGHTS**

Explore data used by the AI system >>> The past applications were used by the AI system to learn a prediction model using the XGBoost technique. Here is a visualisation of the high-level decision tree that approximates the rules learned by the prediction model.

Explore the high-level decision tree >>>

Proposed changes to change the AI decision >>>

**Insufficient details about the algorithm and the high-level decision-making process.**

**AI INSIGHTS**

Explore data used by the AI system >>> Making following changes to the current application will result in a change in AI system decision.

Explore the high-level decision tree >>>

Proposed changes to change the AI decision >>>

- Household details
  - Household size increased to 6
  - Household members who are employed should be None
- Financial details
  - Healthcare expenses higher than £1500

**Recommend infeasible changes to the stakeholder to get a more desirable decision.**

Figure 8: Relatively negative XAI Experience of a welfare applicant using the AssistHub AI platform



## Assist Hub Group A

0% complete

---

### Page 1: AssistHub Welfare Application

In this survey, you'll step into the shoes of a welfare benefits applicant exploring AssistHub AI system decision on their application.

Your feedback on its ability to make decisions, provide recommendations, explain them clearly, and assist you in making informed decisions will be valuable for our research.

**Terminology**

AssistHub uses an AI algorithm to predict the decision of welfare benefits applications. This algorithm is referred to as the "AI system".

The applicant can request additional information about the AI system and the decision. In response, the AssistHub can provide additional information, clarifications, explanations and/or reasoning. These are collectively referred to as "explanations".

The complete interaction (i.e. getting the AI decision and explanations from the system) is referred to as an "experience".

1 If you are happy to continue, Please provide your Prolific ID. \* Required

[Next >](#)

---

## Assist Hub Group A

33% complete

---

### Page 2: AssistHub Welfare Applicant's Experience

The video below is a recorded experience of a welfare benefits applicant interacting with the AssistHub system. Carefully watch the video, and try to imagine that you are the applicant who is seeking help with their application.

Use the full-screen view and increase the resolution of the video as needed. Feel free to scroll and control the playback speed to understand the interaction and information. Please note that there is no audio.

Please click next when you are ready to proceed to the last step.

[< Previous](#) [Next >](#)

## Assist Hub Group A

66% complete

---

### Page 3: Welfare Applicant's Experience Evaluation

Now that you have watched the demonstration, consider the statements below. Please rate how well each of the statements describes the experience, from (1 Strongly Disagree) to (1 Strongly Agree).

This part of the survey uses a table of questions, [view as separate questions instead?](#)

2 Evaluating the welfare applicant's experience.

	* Required				
	1 strongly agree	1 somewhat agree	1 am neutral	1 somewhat disagree	1 strongly disagree
The explanations received throughout the experience were consistent.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
...					
The experience was consistent with my expectations.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

This part of the survey uses a table of questions, [view as separate questions instead?](#)

3 Please provide any additional comments about any of the statements you were asked to rate (such as clarity of the statement, relevance to evaluating the experience, etc).

The explanations received throughout the experience were consistent.

...

The experience was consistent with my expectations.

4 Provide any general feedback that you may have.

[< Previous](#) [Finish ✓](#)

Figure 9: Pilot Study Preview; Example removed from Page 2 and list of items shortened in Page 3.