# Designing and Personalising Hybrid Multi-Modal Health Explanations for Lay Users

Maxwell Szymanski[1], Cristina Conati[2], Vero Vanden Abeele[1] and Katrien Verbert[1]

[1]*Department of Computer Science, KU Leuven, Leuven, Belgium*

[2]*Department of Computer Science, University of British Columbia, Vancouver, Canada*

## Abstract

Recommender systems are increasingly used in mobile health applications. While researchers have highlighted the importance of explaining these recommendations to lay users, with benefits such as increased trust and a higher tendency to follow up on these recommendations, how to design explanations for lay users in critical contexts such as health remains largely unexplored. This paper explores and evaluates a unimodal visual and textual explanation, as well as a combined hybrid explanation modality for chronic pain related recommendations through a qualitative and quantitative lens via a repeated measures study ($n = 262$), and links user's preference towards these modalities to differences in user perception such as trust and acceptance. Additionally, we explore the effects of how personal characteristics such as need for cognition and ease-of-satisfaction affect the user's preference and reception of the different explanations. Results indicate a strong preference towards the combined hybrid explanations. We also found interaction effects of ease-of-satisfaction and need for cognition on the perception of different explanation designs, indicating that users with a higher need for cognition tend to trust unimodal explanations more compared to hybrid explanations.

## 1. Introduction

In recent years, recommender systems (RS) have gained popularity across a multitude of domains, including domains with high stakes such as health. In such domains, explanations to make systems more scrutable have been identified as a core requirement [1, 2]. Various factors need to be taken into account when designing such explanations, such as the goals of the explanations [2] and the type of end user [3]. Ribera et al. indicate that users can be differentiated with respect to their expertise - AI-experts, domain-experts and lay users - and that goals of explanations differ across these groups [3]. Personality traits are another characteristic amongst which users can be distinguished. One such trait that has often been studied in the context of explanations is *need for cognition* (NFC - the tendency to engage in and enjoy activities that require thinking). It has been previously linked to the effectiveness of explanations, with studies indicating that

users with a low NFC tend to benefit more from explanations and hints as they help them make more confident decisions [4, 5]. Additionally, Kouki et al. investigated the effect of a user's *ease-of-satisfaction* (EOS - the user's natural propensity to be satisfied) on their reception of explanations [6, 7], and found that it has an influence on the subjective persuasiveness of certain explanation types.

Despite the many interesting studies on the effectiveness of explanations in relation to user characteristics, very little is known on how explanations need to be designed for lay users in high stakes domains such as health: only 10% of explanations in the health domain are designed with lay users in mind [8]. There is a strong need to design and evaluate effective explanation designs specifically for these lay users, as it has been shown that lay users are more susceptible to a plethora of biases when interacting with explanations [9, 10]. Recent work focusing on lay users suggests that hybrid explanations, that combine two or more modalities, can be beneficial for lay users in terms of understanding and bias mitigation [6, 9]. However, these studies have only contextualised their findings and benefits in low-stakes contexts, such as music and news recommender systems. To the best of our knowledge, lay user evaluations of hybrid explanations in high-stakes domains such as health remain largely unexplored. It is therefore important to evaluate whether findings from studies in other domains generalise to domains with higher stakes.

In this paper, we explore whether combining a unimodal textual and visual explanation into a hybrid explanation produces benefits for lay users in the higher-stakes health domain, or introduce any potential information overload or negative perception. As previous findings also highlight the role of NFC and EOS in differences of perception, we also investigate any potential interaction effects of these traits on the user's perception of explanations [4, 6]. This leads us to answer the following research questions:

**RQ1** Which health explanation modality (a unimodal TEXTUAL or VISUAL, or a HYBRID modality) do lay users prefer and why?

**RQ2** In which way do personal characteristics (need for cognition, ease-of-satisfaction) influence lay user perception of unimodal and hybrid explanations?

To address these research questions, we developed a mobile health application that is able to coach and inform users experiencing chronic musculoskeletal pain. We developed this application in collaboration with IDEWE, the largest occupational health service provider in Belgium. This app contains a conversational RS that provides knowledge-based recommendations regarding how to effectively manage pain flare-ups. These recommendations were developed based on input of six ergonomists and prevention advisors and refined based on data collected through an initial longitudinal study ($N = 249$). In a follow-up mixed-methods between-subject study ($N = 262$), we compared UNIMODAL textual and visual explanations to a hybridised combination of the two. Results indicate that although subjective perception (apart from usefulness) remains the same across the UNIMODAL and HYBRID explanations, lay users strongly tend to prefer HYBRID health explanations. Through a thematic analysis, we distil themes as to why HYBRID explanations are preferred, and additionally find that extending VISUAL explanations with a TEXTUAL modality is more beneficial for users than extending existing TEXTUAL explanations with VISUALS. Secondly, we see an interaction effect of NFC on user

perception of different explanations modalities: users with a higher NFC have a more positive reception of UNIMODAL explanations, compared to HYBRID explanations, nuancing previous findings in favour of HYBRID explanations. Ease-of-satisfaction seems to be a strong predictor of attitudes towards explanation modalities in general, with high EOS users having a more positive perception compared to users with low EOS indicating a need for further research to explore effective explanations for users with low ease-of-satisfaction.

## 2. Related work

In this section, we will discuss the use of recommender systems in the health domain, as well as recent advances of incorporating explanations into recommender systems.

### 2.1. RecSys in Health

Recommender systems (RS) have become prominent in health applications, where they help retrieve relevant information or recommend possible next actions tailored to the needs of the end user. These health recommender systems (HRS) are used both in clinical settings as well as in personal contexts where health applications aid users in their daily lives. A recent systematic review of HRS for lay users shows that the majority of HRS that used a graphical user interface focus on mobile applications [11].

However, the increased use of HRS is also paralleled with certain barriers. One such issue is a mismatch in recommendations to the user's expectations. Such mismatch can lead to a decrease in system effectiveness [12] and a decrease in trust towards the system [11], potentially steering the user away from future use. Early research mainly focused on increasing the accuracy of RS in order to mitigate this issue. However, recent research increasingly explores the effects of human factors, including research on explanations to increase transparency, human-in-the-loop feedback to correct misunderstandings, and using conversational RS to increase familiarity towards the system's interface [13]. This broader approach in reasoning about RS should allow researchers to improve RS effectiveness beyond quantitative algorithmic capability.

### 2.2. Explaining health recommendations

As highlighted earlier, adding explanations to recommendations can improve their overall effectiveness. These explanations make the system interpretable and transparent, which in turn can improve trust towards the system [14, 15]. There exist HRS that explain their rationale to the end user, such as the food recommender system of Wayman et al. that explains why certain recipes are recommended based on the user's nutritional intake [16], or a visualisation for medical experts that is able to explain breast cancer similarities [17]. However, the systematic review of De Croon et al. states that only 10% of HRS that focus on lay users provide explanations [11]. Additionally, a study of Bussone et al. points out that providing overly detailed explanations for health recommenders can create unforeseen effects, such as creating over-reliance on explanations [18], indicating that health recommendation explanations should be designed with sufficient care. As such, designing explanations with lay users in mind, and evaluating them with these users, is paramount.

## 2.3. Personalising explanations for end users

Research has explored the influence of different characteristics on the effect of the type, amount and content of explanations to show. Naveed et al. explored the effects of a person's thinking style and found that participants with a fast and intuitive thinking style depended more on explanations, whereas participants with logical and rational thinking acted more independently from the given explanations [19]. Millecamp et al. repeatedly studied the effects of said thinking style (by measuring need for cognition) in music RS explanations, but additionally investigated the effects of a person's musical sophistication and openness on the attitude towards explanations [20, 4].

Additionally, an increasing amount of research has indicated that the *expertise* of end users should be taken into account when designing explanations. Ribera et al. highlight differences in the needs, goals and limitations of different user groups, including AI-experts, domain experts and lay users. AI expert users, for example, use explanations to verify or improve the underlying AI system, whereas domain experts can leverage explanations to gain additional insights and learn from the system. Lay users have their own set of goals, but also their own array of limitations. Wang et al. have highlighted several shortcomings of lay users that relate to cognitive biases, such as confirmation and anchoring biases, due to a backward-oriented, hypothesis-driven reasoning process [10]. Tsai et al. also noticed a *reinforcing effect*, where users avoid interacting with content they are not familiar with [21]. Szymanski et al. additionally pointed out that lay users, despite having these biases and incorrectly interpreting certain complex explanations, can still have a preference for said complex explanations over other, simpler explanation modalities due to these cognitive biases [9].

Thus we see that interpretability through explanations has multiple benefits and can result in an increased trust towards the system. However, as previously mentioned, the adoption of explanations in HRS is still low. Furthermore, most health-related AI explanations are being researched with AI and domain expert users in mind [8], which leaves a big gap for explanations w.r.t. lay users. Keeping the aforementioned biases in mind that lay users are prone to, it is therefore important to assess whether explanations are indeed interpretable to make sure no misalignment in trust is created.

## 2.4. Hybrid explanations

There has been some previous work looking into hybrid explanations, which essentially combine multiple explanations to enhance understanding and transparency [22, 23, 9, 24]. Hybrid explanations can serve a twofold purpose: combining two or more *explanation styles* in the same representation, to provide more information and context to end users, or combining two or more *explanation representations* that offer the same information in different representations, to help alleviate shortcomings of using only one representation.

However, when designing hybrid explanations, the trade-off between information overload and completeness needs to be considered. Kouki et al. investigated the optimal number of explanation styles that users want to see and found that this number in part depends on ones personal characteristics [6]. The authors of TeleGam additionally focused their research on the concept of user-specified resolutions of textual explanations, where users could choose the
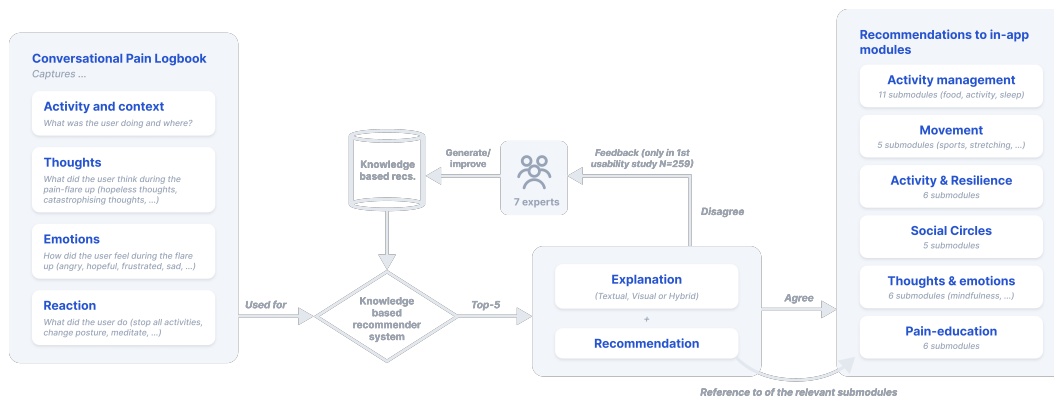
**Figure 1:** High-level overview of the conversational pain logbook, together with the knowledge-based rule-engine RS and possible recommendation topics available in the app

level of detail of verbalisations that accompanied the visual explanations [23]. Considering the benefits of hybrid explanations, as well as the pitfalls due to limited understanding and potential biases related to lay users, more research needs to be done with regards to hybrid explanations [25]. Given that previous work only situated their benefits in low-stakes domains, future research, including this work, should explore if any benefits generalise to higher-stakes domains.

## 3. Explanation design

In this section, we present the health recommender system that was developed, as well as the different explanation modalities to explain said recommendations.

### 3.1. Designing a health recommender system

For the purpose of this research, we have developed a mobile health application that is able to coach and inform users experiencing chronic musculoskeletal pain. The main function of the application is to monitor and educate users through a variety of information, questionnaires and interactive exercises, by guiding them through content modules regarding various topics, such as activity management, pain education, mindfulness, etc. The content of the app, including the modules, were developed in collaboration with researchers from IDEWE, and a team of 6 ergonomists and prevention advisors [26].

To help users with specific episodes of a pain flare up, we have developed a so called pain logbook that users can fill in in order to receive specific and tailored recommendations. The logbook includes a conversational RS to increase engagement and ease-of-use [27]. It asks the user about situational data (how intense the pain is and in which situation it occurred), as well as the user's reaction and thoughts they had at the time (e.g. stopping all activities, being frustrated or scared, feeling helpless). Based on these inputs, recommendations are generated using a knowledge-based recommender system on how to better cope with the pain flareup next time or how to better handle the given reactions and emotions. Domain knowledge of a
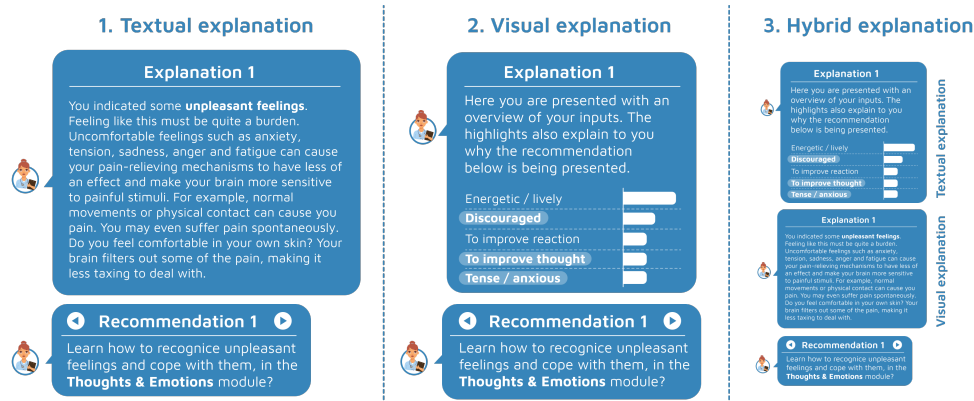
**Figure 2:** Designs of the TEXTUAL (left), VISUAL (centre) and HYBRID (right) explanation

team of six ergonomists and prevention advisors, together with a physiotherapy researcher, was translated into a set of rules that guide users towards one of the 39 in-app submodules related to mindfulness, resilience, activity, etc. Note that more than one relevant recommendation can be given at once if relevant. In that case, recommendations are sorted based on their feature importances (which themes were most prevalent in their input) as well as past pain logbook entries. For example, if a user indicates being really frustrated when the flare up occurs, and ceases all activities, two recommendations will be given. A first one will be related to the frustrated emotions, as negative emotions sensitise our neurons causing a higher pain perception, and will recommend a specific mindfulness exercise within the *thoughts and emotions* module. A second recommendation will be related to the user stopping all physical activities, as frequently abstaining from activities will cause the muscles to become even weaker over time, which in turn will increase the amount of pain episodes. Therefore, a recommendation regarding *activity management* will also be shown, which will tell the users to take a short break, and adapt or lighten their current activities so they can continue. A high-level overview of how the pain logbook, as well as the knowledge-based RS work, is given in Figure 1.

### 3.2. Designing the explanations

To explain the pain logbook recommendations, we designed different explanation modalities, including TEXTUAL and VISUAL explanations presented in Figure 2. Based on data from a think aloud study with these different explanation modalities, we found that feature importance (FI) explanations were favoured by most users for their ability to give an overview as a comprehensive and complete explanation, despite a possible information overload. Additionally, TEXTUAL explanations were favoured by users for whom the recommendation strongly aligned with their expectations, as the explanation was able to give more in-depth information into the topic that the users already agree with. Users that preferred less (overwhelming) amounts of information were also more in favour of TEXTUAL explanations, and less in favour of FI and other visually more elaborate designs. To illustrate how these explanations will be presented, we give the following example. When a user has indicated being scared during a pain flare up, and having

thoughts such as *"Why is this happening to me?"*, a recommendation towards the mindfulness submodule is given. The VISUAL explanation gives an overview of all themes that where captured from the user's input, and indicates negative emotions (e.g. frustrated, scared) with high feature attributions as they contribute most towards said recommendation. The TEXTUAL explanation on the other hand highlights that negative emotions were detected in the input, and explains that negative emotions increases pain-perception and can be minimised through mindfulness. When the user clicks through to the next recommendation, different topics that are relevant to the new recommendation will be highlighted in the VISUAL explanation, and a different text will be displayed in the TEXTUAL explanation.

We evaluated these unimodal TEXTUAL and VISUAL explanations, as well as our knowledge-based recommendations, through a longitudinal study with 249 participants who used the application for 4 months. Participants could interact with the coaching application and pain logbook that either had the VISUAL or TEXTUAL explanation (randomly assigned), and were able to give suggestions through an in-app feedback module. The data of this iteration was then used to refine the knowledge-based recommender system, as well as the explanations, for the following iteration. Participants mentioned that some recommendations and their respective textual explanations were too general in some cases. For instance, some participants found that explanations related to activity management did not take their personal context into account, such as which activity they were doing. To address this limitations, our collaborative team of 6 ergonomists further diversified the textual explanations to take specific nuances in inputs into account.

In addition to the two UNIMODAL explanations, we also add a new HYBRID explanation design (Figure 2, right) that essentially combines both the TEXTUAL and the VISUAL explanation. As discussed in Section 2.4, some previous work has shown positive effects of using HYBRID explanations in terms of increased understanding and preference [9, 23, 6]. However, previous work regarding HYBRID health explanations focuses more on multiple complex visuals that suit domain experts, or states that a combination of complex feature attributions and class attributions should be incorporated, which in practice might be overwhelming for lay users [10]. Our HYBRID design therefore draws inspiration from simpler combinations of TEXTUAL and VISUAL explanations that have been proven to be beneficial for lay users, albeit not yet tested in high-stakes domains such as health [9, 23].

## 4. Methodology

### 4.1. Study design

Through a within-subject study with $N = 262$ participants (ethically approved by the Ethics Committee Research of the UZ Leuven universitary hospital (EC Research) with application number S-65610), we explore if, and how extending a unimodal TEXTUAL explanation with a VISUAL explanation might benefit the user, and vice-versa. We also investigate whether extending TEXTUAL explanations with VISUAL explanations has different benefits for lay users than doing so the other way around. During the within-subject study, participants are first presented with contextual information regarding the user study, and are required to consent to their (questionnaire) data being anonymously recorded in order to continue. After consenting,
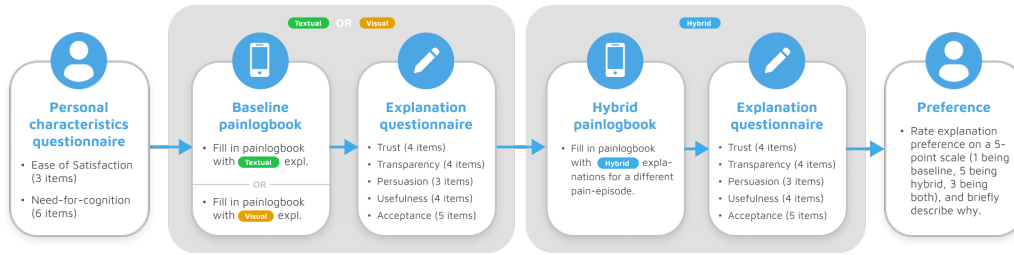
**Figure 3:** Methodology of the repeated measures within-subject user study

the participants receive a questionnaire regarding their *personal characteristics* (PC). This questionnaire relates to their ease-of-satisfaction (EOS, from Kouki et al., slightly adapted to fit the health recommendation setting [6]), and the abbreviated 6-item version of the need for cognition scale (NFC) [28]). Next, participants interacted with the pain logbook with the UNIMODAL explanation (being either TEXTUAL or VISUAL depending on the randomly assigned group). The study asks them to fill in the logbook according to a (chronic) pain scenario they recent experienced, after which they receive several relevant recommendations, accompanied by either a TEXTUAL or VISUAL explanation. Participants are free to go through the recommendations and their respective explanation, and choose one that suits them, or select the 'none' option if no recommendation is relevant. This step is followed by a questionnaire regarding the *perception of explanations*, inquiring the participant about the perceived *trust, transparency, persuasiveness, usefulness* and *satisfaction*, taken from [29, 30]. After filling in the questionnaire, participants are presented with the same pain logbook, but with an explanation that extends the first UNIMODAL explanation they were presented with (by a VISUAL explanation if TEXTUAL was shown first, or vice versa). Participants are asked to fill in the logbook according to a different pain episode they recently experienced. This step is followed by the same questionnaire relating to the perceived *trust, transparency, persuasiveness, usefulness* and *satisfaction* of the HYBRID explanations. In the final step, participants are asked to give their preference towards either of the explanation designs on a 5-point Likert scale (1 - strong preference for UNIMODAL  5 - strong preference for HYBRID  3 - both), and give a reason as to why they like their preferred explanation design, and dislike the other.

## 4.2. Recruitment

Figure 3 shows an overview of the within-subject study design, in which users first interacted with either the TEXTUAL or VISUAL explanations as the UNIMODAL design, followed by interacting with the HYBRID explanations. We recruited a total of 291 English-speaking participants through the online Prolific platform, where users receive $7.29 for participating in the study (approx. 25 minutes). The inclusion criterion was to have experienced chronic pain for a period of at least three months in the past three years. After filtering out participants that did not complete all questionnaires or incorrectly answered the 'alertness' question, the resulting number of participants equates to 262 (143 participants in the TEXTUAL↔HYBRID setting, 119 in the VISUAL↔HYBRID setting). The demographics (age, gender and scores on NFC and EOS) of the participants can be seen in Figure 4. We also find that the medians of the personal characteristics
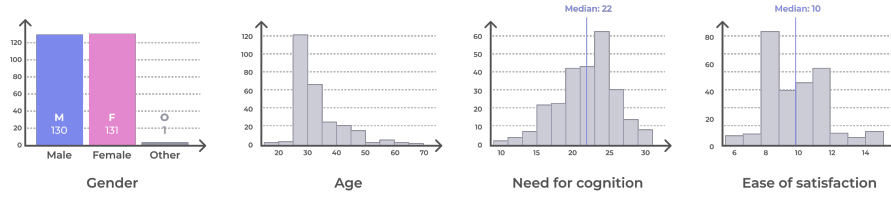
**Figure 4:** Demographics of participants, including gender, age, scores on the NFC and EOS scale
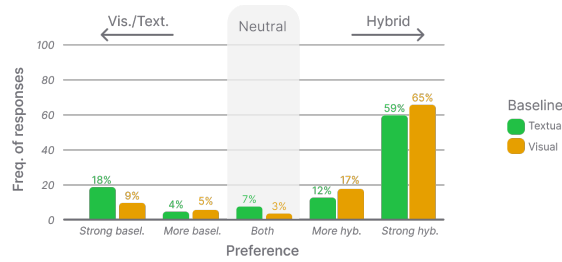


**Figure 5:** Preference scores for each within-subject study

are in line with other studies [4].

# 5. Results RQ1 - Which explanation do users prefer and why?

To answer RQ1 regarding user preferences, we need to analyse the preference data. As it is non-normal (Shapiro-Wilk test, TEXTUAL↔HYBRID: $W = 0.65, p < 0.001$, VISUAL↔HYBRID: $W = 0.62, p < 0.001$), we cannot perform a (M)ANOVA analysis and need to opt for a non-parametric univariate analysis using a one-sided Wilcoxon signed rank test. Additionally, we do a two-coder thematic analysis to gain insights into the reasoning behind their preference.

## 5.1. Main effect for preference

We asked users to give a preference towards either the UNIMODAL explanation (either VISUAL or TEXTUAL), and the hybrid explanation on a 5-point Likert scale, with 1 being a strong preference towards the UNIMODAL explanation, 5 a strong preference towards HYBRID, and 3 being a preference for both. To test whether there is a significant preference for the novel hybrid explanations or not, we test the sample median against a median of $\mu_{MED} = 3$. Since preference is non-normally distributed (Shapiro-Wilk results: $W = 0.668, p < .001$), we perform a one-sided one-sample Wilcoxon rank test. Results indicate strong evidence to assume the alternative hypothesis (i.e. the preference median being greater than 3) ($V = 420, p < .001, d = .543$ (LARGE EFFECT)), indicating that HYBRID explanations were largely preferred compared to UNIMODAL explanation (Figure 5).

## 5.2. How does extending a unimodal explanation affect user preference?

To delve deeper into the reasons behind user preferences, we extend the previous results with results from the thematic analysis to gain more insight in to *why* users prefer certain modalities over others. We do this through a two-coder iterative thematic analysis, with an agreement percentage of 89.2% and Cohen's kappa $\kappa = 0.68$, resulting in a substantial inter-coder agreement [31]. For brevity, we only report the overarching themes rather than individual codes.

### 5.2.1. Discussion

The thematic analysis points to an interesting synergy between the TEXTUAL and VISUAL explanations. When used in a standalone setting, TEXTUAL explanations were found to be good at describing specific factors for recommendation in a detailed way, helping end users to understand both the explanation, as well as the recommendation. However, most users also found them to be less engaging and harder to scan through. On the other hand, VISUAL explanations gave most users a concise and global overview of all their inputs, making it easy to scan for information at first glance. Yet we also notice a dislike for the lack of depth and detail, with some users not even perceiving the visuals and highlights as an explanation, but instead as a mere summary of their inputs. **Combining both modalities has proven to integrate the strengths of both modalities, whilst at the same time alleviating the shortcomings of using them separately** (i.e. TEXTUAL being less engaging and hard to scan, and the VISUAL being too general or not offering a 'real' explanation).

Whilst most users did express a positive sentiment towards the HYBRID explanations, we do have to keep in mind that the preferred explanation modality also has its shortcomings. Some users found the combination of both explanations to be either too overwhelming at first glance due to the addition of text, or had no interest in seeing visuals due to a general dislike or lack of perceived added value. This relates closely to the well-known dilemma known as the completeness-conciseness (or simplicity-power) trade-off, for which several solutions have been proposed. One of the solutions is *progressive disclosure*, where we seek to offer information on demand by providing users with initial essential information, and deferring detailed or more advanced information to a pop-up or secondary screen [32]. A second method consists of making TEXTUAL explanations more readable by applying methods such as *chunking*, *highlighting* or applying *brevity*, [33]. These methods can be combined with progressive disclosure to gradually expose users to information in the HYBRID explanation when needed, i.e. when the recommendation doesn't align with their expectations or when the user shows high interest in the topic and wants more information. Further research could investigate if these proposed solutions have a potential to further decrease the limited shortcomings of the hybrid explanations.

We also investigate whether extending textual explanations with visuals has different benefits compared to adding visual explanations with a textual modality. We see a stronger preference for textual modalities (26,3% standalone textual compared to 16,3% standalone visual, or 83,7% hybrid preference when textuals are added compared to 73,7% when visuals are added). When looking at overarching themes when unimodal explanations are extended, we see insightfulness as a big theme when textual explanations are added ($N = 82$ unique users), with reasons being a more

**Table 1**

Wilcox sign rank tests to compare either textual or visual with hybrid explanation (w/ Bonferroni corr.)

| Comparison | Trust | Transparency | Persuasion | Satisfaction | Usefulness |
|---|---|---|---|---|---|
| TEXTUAL↔HYBRID | $p = 1.00$ | $p = 1.00$ | $p = 1.00$ | $p = 1.00$ | $p < 0.001$ *** |
| VISUAL↔HYBRID | $p = 1.00$ | $p = 1.00$ | $p = 1.00$ | $p = 0.17$ | $p < 0.001$ *** |

detailed, better to understand and more informative design. When visual explanations are added, we still see insightfulness emerge, albeit to a lesser extent ($N = 50$ unique users). In addition to insightfulness, users mention visual engagement as a factor to prefer the addition of visual explanation ($N = 35$ unique users), with users finding the visuals to make the explanation more engaging, appealing and less boring. The stronger preference for textual modalities contradicts previous findings that indicate that lay users tend to have a strong preference towards visual explanations over textual ones due to their visual engagement [9]. However, previous findings were contextualised in non-critical settings. In high-stakes contexts such as health, users saw a higher need for more more informative explanations (a need-to-have), rather than "visual engagement" (which is a nice-to-have).

## 6. Results RQ2 - How does the user's perception of explanations differ?

For RQ2 regarding perception, we find our independent variables to be normally distributed and thus are able to do a (M)ANOVA analysis.

### 6.1. Main effects of perception

We can compare the HYBRID explanation with both the VISUAL and TEXTUAL explanation to see whether they perform better or worse in any of these categories. We perform a non-parametric Wilcox signed rank test in each setting to see if the results of the hybrid explanations differ from the other explanation. Using a Bonferroni correction (adjusting for $n = 5$ tests in each setting, accounting for the 5 dimensions we test on), we find that only *usefulness* differs with the HYBRID explanation, compared to both the TEXTUAL explanation, as well as the VISUAL explanation. Using a (Bonferroni corrected) one-way Wilcox signed rank test, there is a strong statistical significance ($p = < 0.001$***) that indicates **participants found the hybrid explanations to be more useful compared to both the textual and visual explanations** (Table 1). Research by Tsai et al. [30] mentions that when explanations that are being combined are *complementary*, they are perceived to be more *useful* for users. This leads us to conclude that our design of the textual and simplified feature importances are complementary to each other, and in turn are perceived to be more useful.

### 6.2. Interaction effects of personal characteristics on perception

We now explore the effects of the user's personal characteristics (*PC*), consisting of NFC and EOS, on the perception of explanations. We divide this research question into two parts: how

**Table 2**

MANOVA results of personal characteristics influencing attitude towards explanations (both baseline and hybrid)

| | NFC | | | EOS | | |
|---|---|---|---|---|---|---|
| | LOW | HIGH | $F_{1,233}/p$ | LOW | HIGH | $F_{1,220}/p$ |
| Trust | $63.48 \pm 20.45$ | $64.66 \pm 21.74$ | $0.371/0.543$ | $56.38 \pm 22.12$ | $71.53 \pm 16.46$ | $68.397/ < 0.001$ *** |
| Trans | $64.72 \pm 19.12$ | $64.72 \pm 19.71$ | $0.113/0.737$ | $59.41 \pm 20.43$ | $69.62 \pm 15.68$ | $35.543/ < 0.001$ *** |
| Persuasion | $71.19 \pm 17.40$ | $70.90 \pm 18.55$ | $0.029/0.866$ | $65.55 \pm 18.26$ | $76.33 \pm 14.78$ | $47.276/ < 0.001$ *** |
| Satisfaction | $71.88 + 19.99$ | $69.95 \pm 21.13$ | $1.035/0.310$ | $64.08 \pm 21.20$ | $75.97 \pm 18.48$ | $39.778/ < 0.001$ *** |
| Usefulness | $69.70 + 22.18$ | $68.14 \pm 22.41$ | $0.578/0.448$ | $60.64 \pm 23.02$ | $75.30 \pm 19.49$ | $52.696/ < 0.001$ *** |
| Preference | $3.99 \pm 1.47$ | $4.07 \pm 1.52$ | $0.156/0.694$ | $4.01 \pm 1.46$ | $4.14 \pm 1.42$ | $0.620/0.432$ |

do *PC* affect a user's perception of explanations *in general*, and how do *PC* affect the perception of *different unimodal and hybrid explanations*. For each *PC* (NFC and EOS), we do a median split (leave-median-out) to create a group that has a low score (below median) and a high score (above median). This allows us to perform a MANOVA with either NFC or EOS as the predictor, with trust, transparency, persuasion, acceptance, usefulness and preference being the outcome, to try to answer the aforementioned research questions. The scores regarding trust, transparency, persuasion, satisfaction and usefulness have been rescaled between 0 and 100 for easier interpretation and comparison.

### 6.2.1. Effects of NCS and EOS on general explanation perception

Table 2 shows the means of the scores of each subgroup of participants (low and high NFC and EOS), as well as the MANOVA $p$ that indicates whether the difference between the low and the high subgroup is significant. We find that users with a **higher ease-of-satisfaction have a more positive perception of explanations in general**, with their average scores for not only the perceived *satisfaction*, but also *trust*, *transparency* and *persuasiveness* being higher by 10 to 15% compared to users with a lower ease-of-satisfaction. For NFC, we do not find any significance w.r.t. explanation perceptions in general.

### 6.2.2. Effects of NCS and EOS on differences in perception between unimodal and hybrid explanations

Now we use NFC and EOS as the predictor on the *differences* in perception between the hybrid explanation and the UNIMODAL explanation (e.g. $\Delta$Trust represents how much more/less participants trust the HYBRID explanations compared to the UNIMODAL explanation, with a positive score of $x$ meaning that they trust hybrid explanations more by $x$% compared to the UNIMODAL design). The results are shown in Table 3. We find that **users with a higher NFC tend to score the hybrid explanations lower in terms of trust, transparency and usefulness compared to the unimodal explanation.** To contextualise this finding, we look at related work regarding NFC and music recommendation explanations by Millecamp et al. [4]. They state that users with a lower NFC tend to benefit more from additional explanations

**Table 3**
MANOVA results of personal characteristics influencing differences in perception between hybrid compared to baseline

| | NFC | | | EOS | | |
|---|---|---|---|---|---|---|
| | LOW | HIGH | $F_{1,233}$ / $p$ | LOW | HIGH | $F_{1,220}$ / $p$ |
| $\Delta$Trust | 1.85 + 18.94 | -4.01 $\pm$ 20.66 | 5.140/0.024* | 2.84 $\pm$ 20.77 | -2.25 $\pm$ 18.09 | 3.790/0.053 |
| $\Delta$Trans | 1.85 $\pm$ 17.82 | -3.69 $\pm$ 20.22 | 4.967/0.027* | 4.64 $\pm$ 18.69 | -1.35 $\pm$ 18.51 | 5.671/0.018* |
| $\Delta$Persuasion | 0.00 $\pm$ 15.33 | -3.92 $\pm$ 18.26 | 3.174/0.076 | -0.52 $\pm$ 16.40 | -1.73 $\pm$ 16.51 | 0.299/0.585 |
| $\Delta$Satisfaction | -1.38 $\pm$ 16.79 | -5.18 $\pm$ 19.43 | 2,581/0.110 | -0.97 $\pm$ 17.08 | -4.45 $\pm$ 20.65 | 1.803/0.181 |
| $\Delta$Usefulness | 1.61 $\pm$ 19.07 | -3.89 $\pm$ 21.51 | 4.308/0.039* | 0.98 $\pm$ 19.52 | -2.20 $\pm$ 20.62 | 1.3601/0.244 |

in general, whereas users with a high NFC only do so when the recommendations are not in line with their expectation and have an explicit need for explanations. When looking at the metadata in our study of whether or not users agreed with the recommendations, we see that the majority of users did agree with the recommendations, which explains why the need for more information is lower, and consequently a slightly lower need for and perceived usefulness of the hybrid, more informative explanations. However, the slightly lower *trust*, *transparency*, *persuasion* and *usefulness* does not translate to their preference, as most users (low and high NFC alike) still opt to choose for the hybrid explanations. Thus we can conclude that it is still feasible to present hybrid explanations by default, albeit with the modification of showing either the textual or visual part on-demand when users request to have more information.

For EOS, we see a similar trend where **participants with a high EOS score the hybrid explanations slightly lower compared to participants with a low EOS, but only in terms of transparency.** Similar results have been reported by Kouki et al, who found that users with a high conscientiousness score also high on EOS, and prefer to see a lower amount of explanations [34]. While we didn't measure conscientiousness directly, we can speculate that the high conscientiousness - high EOS finding is generalisable (as it relates to personal characteristics, not the study setup), which might explain why our participants with a higher EOS scored the HYBRID modality with more explanations slightly lower in terms of perception.

## 7. Conclusion and future work

In this work, we have designed and compared a TEXTUAL, VISUAL and HYBRID explanation modality that shows non-expert users why they receive certain recommendations regarding their chronic pain. Through an initial longitudinal study with 249 participants, we found certain shortcomings when using previously proposed feature importances as visual explanations for lay users. Using their feedback, we adapted the feature importance designs to fit their mental model, yet still convey largely the same information. Afterwards, we performed an online within-subject study with 262 participants to compare the UNIMODAL explanation (either TEXTUAL or newly adapted VISUAL) to a HYBRID explanation in terms of preference and user reception. We find that most users prefer the HYBRID explanation, i.e. a combination of different explanation modalities that is able to give them more insight into why the recommendation is

given. However, we found this preference to be higher with users that were first presented with visual explanations, indicating a higher need for accompanying textual explanations to existing visual explanations. Through a thematic analysis, we explored themes to gain insight into *why* users liked said explanation modalities, and extracted general guidelines for designing health recommendations for non-expert users. We found that using the adapted VISUAL explanations proved to be a good fit for conveying a summary regarding their input, as it allowed users to see the general picture at a glance. However, VISUAL explanations alone often lacked depth and specific information, and could be confusing for non-expert users to interpret. Adding TEXTUAL explanations allowed users to gain more insight into both the explanation and recommendation, and is often regarded as easy to understand by most non-expert users. However, using text alone also had it's shortcomings, such as being perceived as less engaging. Interestingly, both the thematic analysis and quantitative data have shown that both modalities in tandem as a HYBRID explanation alleviated most of the downsides of using VISUAL and TEXTUAL explanations separately, and proved to be the most opted for by the majority of participants. Additionally, we found that the combination of modalities didn't introduce any information overload effect, since the negative themes regarding HYBRID explanations mainly related to a lack of need for more information, and not due to an overwhelming amount of information. These results add to the limited corpus of previous findings within other domains, which state that complementarily designed hybrid explanations aid end users rather than overwhelm them [9, 30], especially in the underexplored area of designing hybrid explanations for lay users in high-stakes domain such as health [8, 11].

We also explored the effects of a user's personal characteristics to see whether their ease-of-satisfaction and need for cognition have an influence on the reception of the explanations. We find that ease-of-satisfaction seems to be a good predictor on the user's general perception of explanations, as users with a high EOS tend to find explanations to be not only more satisfactory, but also more trustworthy, transparent, persuasive, and useful compared to users with low EOS. Need for cognition did not seem to affect a user's general attitude towards explanations, but it did however highlight differences *between* explanation modalities (HYBRID vs. UNIMODAL). We found that users with a low NFC tend to find the HYBRID explanations more useful, transparent and trustworthy compared to the UNIMODAL explanations, and find the opposite effect in users with a high NFC. This nuances the overall benefits of HYBRID explanations by stating that, although users with a high NFC still tend to prefer HYBRID explanations, they find them to be slightly less useful, transparent and trustworthy compared to unimodal designs. This partly relates to previous work by Szymanski et al. [9] where lay users preferred explanation modalities that they unknowingly performed worse with, and calls for more in-depth research as to where these discrepancies come from and how they might potentially affect lay users.

**Limitations and future work**

While careful consideration has been put into the study design process, certain decisions inherently lead to excluding other trajectories to explore. While our within-study design allowed us to gain insightful qualitative data comparing standalone explanations to underexplored hybrid ones (by seeing how extending either a unimodal textual or visual explanation with the other modality affects user perception and preference and if extending textual first differs form visual

first), it limits us from directly comparing all three explanation modalities with each other. A follow-up between-subject study could prove to be useful for evaluating and comparing the visual, textual and hybrid explanation modalities through a fully quantitative lens. Additionally, the choice of representation of the visual and textual explanations naturally have their impact on how users perceive those explanations. We based our explanations on the previous guidelines for designing health explanations for lay users, but exploring the effects of different visual representations on how they influence user perception, and potentially exploring multiple types of visual explanations side by side, could also prove be interesting [35]. Lastly, many surveys and papers note that besides measuring aspects such as trust, transparency, usefulness and so on, user understanding should also be measured [36, 1]. While there are tools and measures to inquire about AI and explanation understanding in an offline way, previous research has shown that due to biases present with non-expert users, there is a possibility that these users unknowingly have an incorrect understanding of the system [10, 9].

Future work could also focus on exploring the effects of other personal characteristics on how non-expert users perceive health explanations. Similar research by Kouki et al. has shown that other characteristics, such as a user's dependability, neuroticism and agreeableness can also impact the way users perceive explanations [6]. Capturing the user's previous knowledge or experience with health RS could also prove to be useful. Other aspects of user perception, such as the quality of the recommendation (e.g. perceived accuracy and novelty), also seem to be influenced by both the explanation modality, as well as user characteristics, and could prove to be a useful path to explore in the context of health recommendations [34]. And lastly, as mentioned in Section 7, performing an in-person study to assess user understanding of explanations has a lot of potential, as non-expert users often suffer from cognitive biases that can lead to incorrect understanding unbeknownst to them.

## Acknowledgments

## References

[1] S. Mohseni, N. Zarei, E. D. Ragan, A multidisciplinary survey and framework for design and evaluation of explainable ai systems, ACM Trans. Interact. Intell. Syst. 11 (2021). URL: https://doi.org/10.1145/3387166. doi:10.1145/3387166.

[2] N. Tintarev, Explanations of recommendations, in: Proceedings of the 2007 ACM Conference on Recommender Systems, RecSys '07, Association for Computing Machinery, New York, NY, USA, 2007, pp. 203–206. URL: https://doi.org/10.1145/1297231.1297275. doi:10.1145/1297231.1297275.

[3] M. Ribera, A. Lapedriza, Can we do better explanations? a proposal of user-centered explainable ai, CEUR Workshop Proceedings 2327 (2019).

[4] M. Millecamp, N. N. Htun, C. Conati, K. Verbert, To explain or not to explain: The effects of personal characteristics when explaining music recommendations, in: Proceedings of the 24th International Conference on Intelligent User Interfaces, IUI '19, Association for Computing Machinery, New York, NY, USA, 2019, pp. 397–407. URL: https://doi.org/10.1145/3301275.3302313. doi:10.1145/3301275.3302313.

[5] C. Conati, O. Barral, V. Putnam, L. Rieger, Toward personalized xai: A case study in intelligent tutoring systems, Artif. Intell. 298 (2021). URL: https://doi.org/10.1016/j.artint.2021.103503. doi:10.1016/j.artint.2021.103503.

[6] P. Kouki, J. Schaffer, J. Pujara, J. O'Donovan, L. Getoor, Generating and Understanding Personalized Explanations in Hybrid Recommender Systems, ACM Trans. Interact. Intell. Syst. 10 (2020) 1–40. URL: https://doi.org/10.1145/3365843. doi:10.1145/3365843.

[7] J. Schaffer, J. O'Donovan, T. Höllerer, Easy to please: Separating user experience from choice satisfaction, in: Proceedings of the 26th Conference on User Modeling, Adaptation and Personalization, UMAP '18, Association for Computing Machinery, New York, NY, USA, 2018, pp. 177–185. URL: https://doi.org/10.1145/3209219.3209222. doi:10.1145/3209219.3209222.

[8] J. Ooge, G. Stiglic, K. Verbert, Explaining artificial intelligence with visual analytics in healthcare, WIREs Data Mining and Knowledge Discovery 12 (2021). URL: https://wires.onlinelibrary.wiley.com/doi/abs/10.1002/widm.1427. doi:https://doi.org/10.1002/widm.1427.

[9] M. Szymanski, M. Millecamp, K. Verbert, Visual, textual or hybrid: The effect of user expertise on different explanations, in: 26th International Conference on Intelligent User Interfaces, IUI '21, Association for Computing Machinery, New York, NY, USA, 2021, pp. 109–119. URL: https://doi.org/10.1145/3397481.3450662. doi:10.1145/3397481.3450662.

[10] D. Wang, Q. Yang, A. Abdul, B. Y. Lim, Designing Theory-Driven User-Centric Explainable AI, Association for Computing Machinery, New York, NY, USA, 2019, pp. 1–15. URL: https://doi.org/10.1145/3290605.3300831.

[11] R. D. Croon, L. V. Houdt, N. N. Htun, G. Štiglic, V. V. Abeele, K. Verbert, Health recommender systems: Systematic review, Journal of medical Internet research 23 (2021). URL: https://pubmed.ncbi.nlm.nih.gov/34185014/. doi:10.2196/18035.

[12] K. Balog, F. Radlinski, Measuring Recommendation Explanation Quality: The Conflicting Goals of Explanations, in: Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '20, Association for Computing Machinery, New York, NY, USA, 2020, pp. 329–338. URL: https://doi.org/10.1145/3397271.3401032. doi:10.1145/3397271.3401032.

[13] A. Calero Valdez, M. Ziefle, K. Verbert, Hci for recommender systems: The past, the present and the future, in: Proceedings of the 10th ACM Conference on Recommender Systems, RecSys '16, Association for Computing Machinery, New York, NY, USA, 2016, pp. 123–126. URL: https://doi.org/10.1145/2959100.2959158. doi:10.1145/2959100.2959158.

[14] D. V. Carvalho, E. M. Pereira, J. S. Cardoso, Machine learning interpretability: A survey on methods and metrics, Electronics 8 (2019). URL: https://www.mdpi.com/2079-9292/8/8/832.

doi:`10.3390/electronics8080832`.

[15] S. Stumpf, V. Rajaram, L. Li, W.-K. Wong, M. Burnett, T. G. Dietterich, E. Sullivan, J. Herlocker, Interacting meaningfully with machine learning systems: Three experiments, International Journal of Human Computer Studies 67 (2009) 639–662. URL: https://openaccess.city.ac.uk/id/eprint/12417/. doi:`10.1016/j.ijhcs.2009.03.004`, © 2009, Elsevier. Licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International http://creativecommons.org/licenses/by-nc-nd/4.0/.

[16] E. Wayman, S. Madhvanath, Nudging Grocery Shoppers to Make Healthier Choices, in: Proceedings of the Ninth Conference on Recommender Systems, ACM, 2015, pp. 289–292. doi:`10.1145/2792838.2799669`.

[17] J.-B. Lamy, B. Sekar, G. Guezennec, J. Bouaud, B. Séroussi, Explainable artificial intelligence for breast cancer: A visual case-based reasoning approach, Artificial Intelligence in Medicine 94 (2019) 42–53. URL: https://www.sciencedirect.com/science/article/pii/S0933365718304846. doi:`https://doi.org/10.1016/j.artmed.2019.01.001`.

[18] A. Bussone, S. Stumpf, D. M. O'Sullivan, The role of explanations on trust and reliance in clinical decision support systems, 2015 International Conference on Healthcare Informatics (2015) 160–169.

[19] S. Naveed, T. Donkers, J. Ziegler, Argumentation-based explanations in recommender systems: Conceptual framework and empirical results, in: Adjunct Publication of the 26th Conference on User Modeling, Adaptation and Personalization, UMAP '18, Association for Computing Machinery, New York, NY, USA, 2018, pp. 293–298. URL: https://doi.org/10.1145/3213586.3225240. doi:`10.1145/3213586.3225240`.

[20] M. Millecamp, R. Haveneers, K. Verbert, Cogito Ergo Quid? The Effect of Cognitive Style in a Transparent Mobile Music Recommender System, UMAP '20 (2020) 323–327. URL: https://doi.org/10.1145/3340631.3394871.

[21] C.-H. Tsai, P. Brusilovsky, Beyond the ranked list: User-driven exploration and diversification of social recommendation, in: 23rd International Conference on Intelligent User Interfaces, IUI '18, Association for Computing Machinery, New York, NY, USA, 2018, pp. 239–250. URL: https://doi.org/10.1145/3172944.3172959. doi:`10.1145/3172944.3172959`.

[22] D. Parra, P. Brusilovsky, User-controllable personalization: A case study with setfusion, International Journal of Human-Computer Studies 78 (2015) 43–67. URL: https://www.sciencedirect.com/science/article/pii/S1071581915000208. doi:`https://doi.org/10.1016/j.ijhcs.2015.01.007`.

[23] F. Hohman, A. Srinivasan, S. M. Drucker, Telegam: Combining visualization and verbalization for interpretable machine learning, in: 2019 IEEE Visualization Conference (VIS), 2019, pp. 151–155. doi:`10.1109/VISUAL.2019.8933695`.

[24] K. Verbert, D. Parra, P. Brusilovsky, E. Duval, Visualizing recommendations to support exploration, transparency and controllability, in: Proceedings of the 2013 International Conference on Intelligent User Interfaces, IUI '13, Association for Computing Machinery, New York, NY, USA, 2013, pp. 351–362. URL: https://doi.org/10.1145/2449396.2449442. doi:`10.1145/2449396.2449442`.

[25] M. Szymanski, K. Verbert, V. Vanden Abeele, Designing and evaluating explainable ai for non-ai experts: Challenges and opportunities, in: Proceedings of the 16th ACM Conference on Recommender Systems, RecSys '22, Association for Computing Machinery,

New York, NY, USA, 2022, pp. 735–736. URL: https://doi.org/10.1145/3523227.3547427. doi:10.1145/3523227.3547427.

[26] C. Puri, S. Keyaerts, M. Szymanski, L. Godderis, K. Verbert, S. Luca, B. Vanrumste, Daily pain prediction in workplace using gaussian processes, Proceedings of the 16th International Joint Conference on Biomedical Engineering Systems and Technologies (2023). doi:10.5220/0011611200003414.

[27] D. Jannach, A. Manzoor, W. Cai, L. Chen, A survey on conversational recommender systems, 2020. arXiv:2004.00646.

[28] G. L. de Holanda Coelho, P. H. P. Hanel, L. J. Wolf, The very efficient assessment of need for cognition: Developing a six-item version, Assessment 27 (2020) 1870–1885. URL: https://doi.org/10.1177/1073191118793208. doi:10.1177/1073191118793208. arXiv:https://doi.org/10.1177/1073191118793208, pMID: 30095000.

[29] J. D. Van Der Laan, A. Heino, D. De Waard, A simple procedure for the assessment of acceptance of advanced transport telematics, Transportation Research Part C: Emerging Technologies 5 (1997) 1–10. URL: https://www.sciencedirect.com/science/article/pii/S0968090X96000253. doi:https://doi.org/10.1016/S0968-090X(96)00025-3.

[30] C.-H. Tsai, P. Brusilovsky, Evaluating Visual Explanations for Similarity-Based Recommendations: User Perception and Performance, in: Proceedings of the 27th ACM Conference on User Modeling, Adaptation and Personalization, UMAP '19, Association for Computing Machinery, New York, NY, USA, 2019, pp. 22–30. URL: https://doi.org/10.1145/3320435.3320465. doi:10.1145/3320435.3320465.

[31] N. J.-M. Blackman, J. J. Koval, Interval estimation for cohen's kappa as a measure of agreement, Statistics in Medicine 19 (2000) 723–741. doi:https://doi.org/10.1002/(SICI)1097-0258(20000315)19:5<723::AID-SIM379>3.0.CO;2-A.

[32] A. Springer, S. Whittaker, Progressive disclosure: When, why, and how do users want algorithmic transparency information?, ACM Trans. Interact. Intell. Syst. 10 (2020). URL: https://doi.org/10.1145/3374218. doi:10.1145/3374218.

[33] N. Wichman, Speaking of sentences: Chunking, Teaching English in the Two Year College 36 (2009) 281–290. URL: https://www.proquest.com/scholarly-journals/speaking-sentences-chunking/docview/220969438/se-2, copyright - Copyright National Council of Teachers of English Mar 2009; Document feature - Tables; ; Last updated - 2019-11-22.

[34] P. Kouki, J. Schaffer, J. Pujara, J. O'Donovan, L. Getoor, Personalized Explanations for Hybrid Recommender Systems, in: Proceedings of the 24th International Conference on Intelligent User Interfaces, IUI '19, Association for Computing Machinery, New York, NY, USA, 2019, pp. 379–390. URL: https://doi.org/10.1145/3301275.3302306. doi:10.1145/3301275.3302306.

[35] M. Szymanski, V. V. Abeele, K. Verbert, Explaining health recommendations to lay users: The dos and don'ts, Technical Report, 2022. URL: http://ceur-ws.org.

[36] V. Lai, C. Chen, Q. V. Liao, A. Smith-Renner, C. Tan, Towards a science of human-ai decision making: A survey of empirical studies, CoRR abs/2112.11471 (2021). URL: https://arxiv.org/abs/2112.11471. arXiv:2112.11471.