# Towards Human-Compatible XAI:
# Explaining Data Differentials with Concept Induction over Background Knowledge<sup>☆</sup>

Cara Widmer[a], Md Kamruzzaman Sarker[b], Srikanth Nadella[a], Joshua Fiechter[a], Ion Juvina[a,d], Brandon Minnery[a], Pascal Hitzler[c], Joshua Schwartz[c], Michael Raymer[a,d]

[a]*Kairos Research, LLC*
[b]*University of Hartford, USA*
[c]*Kansas State University, USA*
[d]*Wright State University*

## Abstract

Concept induction, which is based on formal logical reasoning over description logics, has been used in ontology engineering in order to create ontology (TBox) axioms from the base data (ABox) graph. In this paper, we show that it can also be used to explain data differentials, for example in the context of Explainable AI (XAI), and we show that it can in fact be done in a way that is meaningful to a human observer. Our approach utilizes a large class hierarchy, curated from the Wikipedia category hierarchy, as background knowledge.

*Key words:* concept induction, explainable AI, class hierarchy

## 1. Introduction

Lack of explainability and understandability are pervasive problems in modern artificial intelligence. Many modern machine learning (ML) methods in particular involve little to no human intervention once training data has been provided, and the resulting models are often "black boxes" whose internal representations and decision mechanisms remain opaque to human users [11]. In order for human operators to be able to trust an ML algorithm's outputs, make actionable decisions based on those outputs, and detect and correct biases, it is important that the algorithm's reasoning processes be understandable. Thus there has been an increased focus on explainable artificial intelligence (XAI) in recent years [3].

Two major approaches for improving explainability are implementing transparency in the ML models themselves and inferring explanations via post-hoc techniques [3]. For example, transparency can be directly introduced to ML models by incorporating simulatability (i.e., developing the model in such a way that the model can be simulated in the thoughts of a human analyzing the model output) [50] or decomposability (i.e., enabling understandability of each part of the model individually) [33].

In contrast, post-hoc explainability techniques do not directly alter a model to make it more explainable but rather infer an explanation based on analysis of the model's behaviors. Examples of post-hoc techniques include generating human-understandable text explanations and visualizations of a model's behavior, as well as analyses that compute quantitative measures of the influence of various input variables on the model's output. While post-hoc methods explain a model's decisions indirectly, they avoid the performance trade-offs that often accompany attempts to engineer transparency directly into the model [15].

In the current work we implement a different post-hoc explainability strategy by applying concept induction [30] to provide human-interpretable explanations of machine learning classifications. We demonstrate that concept induction, together with a suitable class hierarchy as background knowledge, can be used to generate explanations for data differences that are meaningful for a human observer. A class hierarchy consisting of a curated form of the Wikipedia category hierarchy [46] was used as background ontology for concept induction, and the ECII concept induction system [45] was used for explanation generation. We report on experiments that we have conducted with Amazon Mechanical Turk to assess how meaningful the generated explanations are for humans. The first experiment assesses the quality of explanations generated by concept induction in a specific setting related to scene classification, and supports our hypothesis that these explanations are judged by human participants to be more accurate than semi-random explanation, but less accurate than human-authored explanations. The second experiment shows that concept induction can be used to generate explanations for errors made by a ML classifier; the study also shows that these explanations are inferior to human-generated expla-

nations.

Data and other resources for the work reported herein are available from `https://osf.io/xjkcw/`.

The paper is structured as follows. In Section 2 we review the concept induction technique and discuss the background knowledge and data mapping process underlying our work. In Sections 3 and 4 we detail our two experiments. In Section 5 we present related efforts, and in Section 6 we discuss future work and conclude.

## 2. Concept Induction and Utilized Background Knowledge

Concept Induction [30] in its general form can be defined as follows. As background knowledge we assume that we have a given description logic ontology[1] $O$ with TBox $T$ and ABox $A$. We further assume that we are given two sets $P$ and $N$ each consisting of individuals from $O$: $P$ is called the set of *positive* examples and $N$ is called the set of *negative* examples. The concept induction task then consists of identifying complex description logic expressions $C$ such that $O \models C(p)$ for all $p \in P$ and $O \not\models C(n)$ for all $n \in N$. If no such $C$ exists, then a concept induction system is expected to return approximate solutions, together with some accuracy value.[2]

Generally speaking, concept induction is computationally very expensive. Known provably correct algorithms, such as the one underlying the well-known DL-Learner system [31, 7], proceed by *concept refinement*, which is an adaptation from inductive logic programming [35]: A candidate concept expression $C_0$ is assessed as to whether $O \models C(p)$ for all $p \in P$ and $O \not\models C(n)$ for all $n \in N$ hold, and if not, then several candidate *refinements* of $C_0$, meaning concept expressions that are somewhat more or less general than $C_0$, are assessed, and the best of these becomes the new candidate concept expression $C_1$. Then this process is repeated until either a perfect solution is found or some termination criterion is met. Note that for each candidate concept, several calls to a description logic reasoner have to be made for assessment, and in particular with large sets $P, S$ and a large expressive ontology $O$, each of these reasoner calls can take considerable time.

In this paper, we are thus basing our analysis on a heuristic algorithm and system, called ECII (Efficient Concept Induction from Individuals) [45]. The heuristic underlying ECII rests primarily on (a) using only a class hierarchy as ontology, rather than a more complex ontology, (b) a limited search space for candidate concept expressions, and (c) a priori materialization of logical consequences. With ECII, it is possible to perform concept induction over data sets that cannot otherwise be dealt with, and experiments have shown that output accuracy is still rather high [45].

Previous uses of concept induction have been in the context of ontology engineering [29, 7, 40]. We have previously proposed concept induction for explainability, and preliminary studies have been published [47, 46]. In particular, [45, 46] presented preparatory methods that enabled the work we discuss in this paper: the ECII system [45], and our curated Wikipedia category hierarchy [46].

The latter was produced by first retrieving the complete category hierarchy from Wikipedia. Since it is crowdsourced, it actually contains cycles which, if interpreted as a proper class hierarchy (i.e., in the sense of rdfs:subClassOf), would make parts of it collapse. We thus devised an algorithm for breaking these cycles heuristically. The resulting class hierarchy consists of 1,860,342 concepts. In addition, 6,079,748 individuals – which correspond to Wikipedia page URLs – are assigned to the classes according to the categorizations found on Wikipedia.

Using the Wikipedia category hierarchy also enabled us to establish an easy mapping process of examples into the hierarchy: Text labels (more details about the dataset we used can be found in Sections 3 and 4) were mapped to DPbedia [32] entities using DBpedia Spotlight;[3] these entities in turn usually correspond to Wikipedia pages, from which we can obtain the Wikipedia categories the page belongs to.

## 3. First Experiment: Assessing the Quality of Explanations Generated by Concept Induction

In this first experiment we aimed to assess the perceived quality of machine-generated explanations (i.e., ECII's explanations) compared to both human-generated ("gold standard") explanations and semi-random explanations. Explanations in this study consist of brief textual descriptions of the key conceptual differences between two sets of natural images corresponding to two distinct scene categories (A and B). Note that the two image sets notionally represent the decision outputs of a binary scene classification algorithm; however, for this experiment the classifier was merely hypothetical, since the purpose of the experiment was to assess the ability of ECII to explain data differentials *per se*. The study had been reviewed and approved by New England IRB (now part of WCG IRB; protocol #17-1325335-1) and pre-registered[4] before it was conducted.

### 3.1. Hypothesis
We hypothesize that human-generated explanations will be judged most accurate by participants, followed by ECII explanations, followed by semi-random explanations.

### 3.2. Method
#### 3.2.1. Participants
300 participants were recruited through Amazon Mechanical Turk using the Cloud Research platform. Participants were recruited from a listing describing the task and the compensation

---

[1] See [23] for background on description logics and their relation to ontologies, in particular the Web Ontology Language OWL [22].

[2] There are different ways to assess accuracy, including precision, recall, and F1 scores.

[3] See `https://www.dbpedia-spotlight.org/`.

[4] See `https://aspredicted.org/blind.php?x=QYD_JT4` for complete study preregistration information.

structure. Participants were compensated $5 through Mechanical Turk ($7.50 per hour prorated to the estimated time of 40 minutes required to complete the task). Based on an estimated medium effect size of $f2 = 0.15$ and 95% power, we needed a sample size of at least 89 observations (i.e., unique participant judgments) per trial for estimating the parameters of the Bradley-Terry model [5] that was used to model the data. To reach this target we aimed to collect data from 300 participants, which equates to 100 total observations for each trial. This ensured we reached the required total and still allowed for potential exclusions.

### 3.2.2. Materials

*Image sets.* A total of 45 image set pairs (A and B) were created for this study. Images were taken from the ADE20K image dataset [53, 54], which includes approximately 20,000 human-curated images with scene category and object tagging. The tags (annotations) do not only indicate presence of an object but also number of such objects, occlusions, etc., but we ignored these detailed annotations for our purposes and used only the object labels, which were mapped to the Wikipedia categories as described at the end of Section 2. Each image set pair included images from two scene categories (90 scene categories in total). Each set within a pair consisted of eight images selected at random from a particular category (e.g., eight images selected randomly from the *Gazebo* category).

Also included were five "catch trial" image sets that were used to verify that participants were paying adequate attention. The images for these trials were selected in the same manner as the target trial image sets but included a different type of explanation (described below).

*Explanations.* Explanations in this study were defined as lists of up to seven concepts (i.e., strings corresponding to Wikipedia category names) that aimed to describe what was present in category A that was not present in category B (where categories A and B are the two distinct scene categories included in the image set). Concepts could be anything physically present in the images (like a computer or window), or an abstract category that fits the theme of the images (like science or entertainment).

Three types of explanations were created for each of the 45 target image sets: 1) ECII machine-generated explanations, 2) human "gold standard" explanations, and 3) semi-random machine-generated explanations.

ECII explanations were created by providing the image sets' object tags (but not scene category tags) to the ECII algorithm. The algorithm then assessed the images and returned a rating of how well concepts matched images in category A but not category B. ECII explanations were then created by taking the seven highest rated unique concepts. ECII provides several different methods of ranking explanations. Explanations were created using rankings based on F1 scores, recall scores, precision scores, and a hybrid score. Pilot testing revealed that explanations created using F1 scores were seen as most accurate by participants, and so F1-based ECII explanations were selected for use.

Human "gold standard" explanations were created by providing the image sets (but not the object tags or scene category tags) to three human raters. Each rater independently assessed each image set and generated a list of 7 to 10 concepts that (on average) matched the images in category A but not category B. Once all raters had completed this task, gold standard explanations were created by first selecting any concept mentioned by all three raters, then concepts mentioned by two of the raters, and then finally filling the explanation with concepts randomly selected from the remaining concepts across the three raters until seven unique concepts were obtained.

Semi-random explanations were created by randomly shuffling the images in each image set to create new groupings of category A and category B that were not based on the scene category but instead contained a mixture of images from both categories. These new image sets were then provided to ECII, and explanations were generated in the same method as described above. The logic of this approach is that these explanations were likely to be viewed as more plausible than completely random explanations since they would still include concepts that are present in the images participants saw, but they would not (on average) align with one group over the other. Thus, semi-random explanations (hereafter referred to as "random" explanations) constitute a more reasonable baseline against which to assess ECII's performance than fully random explanations, because the latter could be easily outperformed by even a low-quality explanation engine.

Catch trials included their own set of two types of explanations. One set were human explanations which were generated in the same way as the other human gold standard explanations. The other set of explanations for catch trials consisted of completely random concepts taken from a random word generator. This resulted in explanations that were very obviously inaccurate, allowing these trials to serve as an assessment of how well participants were paying attention to the task.

To standardize the presentation of explanations, all concepts in an explanation were presented in alphabetical order.

### 3.2.3. Design

Experiment 1 employed a within-subjects design in which each participant was presented with two explanations per trial and asked to choose the more accurate explanation (two-alternative forced choice design). Each participant saw all three explanation types across all trials, although only two explanation types were compared in any one trial. For each pair of image sets (A and B), the participant completed three trials comparing (1) the ECII versus human explanation; (2) the ECII versus semi-random explanation; and (3) the human versus semi-random explanation. For any given pair of image sets, a given participant completed all three comparisons (i.e., a within-subjects design).

The 45 total pairs of image sets in this study led to a total of 135 unique target trials. Because this was greater than the number of trials an individual participant was likely to be able to complete, participants were randomly assigned to 15 image sets (a total of

45 trials). Image sets were counterbalanced across participants such that all image sets received the same number of responses.

### 3.2.4. Procedure

After providing consent, participants completed brief training on the task, including instructions about how concepts and explanations were defined in this study. Participants then began completing trials. The 50 trials (45 assigned targets and 5 catch trials) were presented in a random order. Figure 1 shows what the stimuli presentation and response options looked like to participants.

### 3.3. Results

Prior to analysis, participant responses to the catch trials were assessed. Participants who failed more than one catch trial were excluded from analysis. While the vast majority of participants did not fail any catch trials (N=295) and one participant failed exactly one trial, four participants failed more (two failed two trials, and two failed three trials). These four participants were excluded from all analyses, leaving a total of 296 participants included in analyses. Across all image sets, human explanations were overwhelmingly chosen over ECII explanations (3856 vs 580; i.e., 87% of the time) and over random explanations (4287 vs 153; i.e., 97% of the time), and ECII explanations were chosen overwhelmingly over random explanations (3862 vs 578; i.e., 87% of the time). See Figure 2.

Participants' pairwise judgments were used in a Bradley-Terry analysis to obtain "ability scores" for each type of explanation, where ability scores in this analysis provide a metric of how much an explanation type was preferred by participants. Ability scores were calculated for each of the 45 image sets. Analysis of these ability scores revealed that human explanations had the highest ability scores (M = 5.32, SD = 3.87), followed by ECII explanations (M = 3.08, SD = 3.99), F(2) = 30.37, $p < 0.001$, $\eta^2 = 0.315$). Random explanations were used as the comparison point in the Bradley-Terry analysis, and thus were set to 0 (with the ability scores for human and ECII explanations indicating how much they were preferred compared to the random explanations). A Tukey's HSD test revealed that differences in the ability scores of human vs. random explanations and ECII vs. random explanations were both significant at $p < 0.001$, while differences in the ability scores of human vs. ECII explanations were significant at $p = 0.004$. See Table 1 for the individual human and ECII ability scores for each image set pair in Experiment 1.

### 3.4. Discussion

Analysis of the results of Experiment 1 provide evidence to support our main hypothesis. Participants found the human-generated explanations to be the most accurate at describing the difference in the image sets, followed by the ECII explanations, and the semi-random explanations proving to be the least accurate. This provides support for the value of ECII explanations. It is not surprising that explanations produced by the ECII algorithm are of lower quality than human-generated explanations at this stage of development. However, the ECII explanations

do contain notable explanatory power, suggesting that ECII explanations can be useful. It should be noted also that there was some variability in ECII's performance across the image sets. For some image sets the ECII explanations were chosen relatively more often than in others, in one case even being chosen more frequently than the human explanation, while in other image sets it was chosen less often than the random explanation. This suggests there is certainly still room for improvement in ECII explanations, but that on average there is promising evidence that ECII can produce explanations that accurately describe the differences between two groups of data.

## 4. Second Experiment: Identifying Errors Made by a Machine Learning Classifier

In Experiment 1 the quality of ECII-generated explanations was (subjectively) judged by human participants via comparison with human-generated ("gold standard") and random-generated explanations. Our second experiment (Experiment 2) sought to obtain an objective measure of the utility of ECII's explanations in helping human users evaluate the decisions of an actual ML system (more precisely, a logistic regression algorithm that classified images into scene categories based on semantic tags of objects present in each image). Specifically, the goal of this second experiment was to test how well ECII explanations (compared to human "gold standard" explanations) helped human participants identify the errors made by the AI. Note that a key difference between this experiment and Experiment 1 is that the explanations in Experiment 1 were of a (hypothetical) ML's *decisions*, whereas the explanations in Experiment 2 were of a (real) ML's decision *errors*. The study had been reviewed and approved by New England IRB (now part of WCG IRB; protocol #17-1325335-1) and pre-registered[5] before it was conducted.

### 4.1. Hypothesis

We hypothesized that participants would be able to match human generated explanations to the correct image set more frequently than for ECII explanations, and that participants would be better than chance for both types of explanations.

### 4.2. Method
### 4.2.1. Participants

Amazon Mechanical Turk participants were recruited using the Cloud Research platform. Participants were recruited from a listing on Mechanical Turk describing the task and the compensation structure. Participants were compensated $5 through Mechanical Turk ($7.50 per hour prorated to the estimated time of 40 minutes required to complete the task). Following recommendations by [44], we initially planned to collect data from 100 participants and compute Bayes Factors for our effects of interest. If the Bayes Factors was not conclusive, then we would continue to collect data in groups of 50 participants, stopping to

| Image Set | H. Ability | E. Ability | H. v E. Wins | H. v R Wins | E v R Wins |
|---|---|---|---|---|---|
| Set 1: Bedroom v Park | 4.58 | 2.04 | 90 – 6 | 94 – 2 | 86 – 10 |
| Set 2: Living Room v Parking Lot | 7.40 | 3.91 | 97 – 2 | 98 – 1 | 98 – 1 |
| Set 3: Office v Playground | 6.09 | 4.80 | 76 – 20 | 95 – 1 | 96 – 0 |
| Set 4: Airport v Amusement Park | 2.60 | 0.09 | 93 – 4 | 87 – 10 | 54 – 43 |
| Set 5: Bathroom v Art Studio | 5.10 | 3.66 | 78 – 18 | 95 – 1 | 94 – 2 |
| Set 6: Beauty Salon v Forest Path | 4.25 | 3.60 | 65 – 33 | 96 – 2 | 96 – 2 |
| Set 7: Bookstore v Child Room | 5.42 | 2.80 | 91 – 6 | 96 – 1 | 92 – 5 |
| Set 8: Hotel Room v Cockpit | 5.35 | 4.17 | 78 – 22 | 98 – 2 | 100 – 0 |
| Set 9: Shoe Store v Alcove | 4.12 | 3.09 | 75 – 26 | 99 – 2 | 97 – 4 |
| Set 10: Alley v Wet Bar | 4.94 | 1.19 | 98 – 2 | 99 – 1 | 77 – 23 |
| Set 11: Closet v Construction Site | 7.76 | 4.61 | 93 – 4 | 97 – 0 | 96 – 1 |
| Set 12: Gazebo v Bowling Alley | 5.31 | 2.49 | 93 – 5 | 97 – 1 | 91 – 7 |
| Set 13: Garage v Hallway | 28.74 | 27.86 | 70 – 29 | 99 – 0 | 99 – 0 |
| Set 14: Laundromat v Pantry | 6.77 | 3.54 | 98 – 4 | 102 – 0 | 99 – 3 |
| Set 15: Conference Room v Waterfall | 6.97 | 3.49 | 95 – 3 | 98 – 0 | 95 – 3 |
| Set 16: Home Office v Bow Window | 4.29 | 1.93 | 93 – 7 | 97 – 3 | 89 – 11 |
| Set 17: Dining Room v Kitchen | 1.77 | 1.64 | 54 – 44 | 81 – 16 | 83 – 14 |
| Set 18: Fast Food v Office Building | 4.09 | 1.26 | 91 – 6 | 96 – 1 | 75 – 22 |
| Set 19: Jacuzzi v Greenhouse | 5.83 | 2.24 | 98 – 3 | 101 – 0 | 91 – 10 |
| Set 20: Gymnasium v Corridor | 3.23 | -0.12 | 93 – 4 | 94 – 3 | 45 – 52 |
| Set 21: Bus v Broadleaf Forest | 6.75 | 3.22 | 99 – 2 | 100 – 1 | 98 – 3 |
| Set 22: Casino v Arrival Gate | 6.17 | 3.55 | 95 – 5 | 98 – 2 | 99 – 1 |
| Set 23: Library v Gas Station | 4.61 | 3.40 | 77 – 23 | 94 – 1 | 92 – 3 |
| Set 24: Valley v Yard | 3.74 | -1.17 | 98 – 0 | 95 – 3 | 24 – 74 |
| Set 25: Mountain v Coast | 0.64 | 0.32 | 58 – 40 | 63 – 35 | 58 – 40 |
| Set 26: Dinette Vehicle v Farm Field | 4.81 | 2.92 | 90 – 11 | 98 – 3 | 98 – 3 |
| Set 27: Poolroom v Driveway | 1.61 | 2.78 | 23 – 75 | 82 – 16 | 92 – 6 |
| Set 28: Bridge v Auditorium | 5.15 | 3.04 | 87 – 10 | 96 – 1 | 94 – 3 |
| Set 29: Museum v Youth Hostel | 1.85 | -0.10 | 89 – 11 | 85 – 15 | 49 – 51 |
| Set 30: Supermarket v Restaurant | 6.16 | 3.97 | 88 – 9 | 96 – 1 | 96 – 1 |
| Set 31: Classroom v Archive | 3.97 | 2.27 | 84 – 14 | 95 – 3 | 90 – 8 |
| Set 32: Dentist Office v Ballroom | 3.79 | 1.25 | 88 – 6 | 95 – 3 | 77 – 21 |
| Set 33: Lighthouse v River | 4.51 | 2.08 | 91 – 7 | 97 – 1 | 88 – 10 |
| Set 34: Creek v Basement | 6.87 | 2.94 | 97 – 1 | 97 – 1 | 94 – 4 |
| Set 35: Building Facade v Ocean | 4.32 | 2.00 | 90 – 8 | 96 – 2 | 87 – 11 |
| Set 36: Courthouse v Parking Garage | 2.28 | -0.45 | 90 – 8 | 91 – 7 | 36 – 62 |
| Set 37: Balcony v Skyscraper | 5.39 | 2.69 | 96 – 7 | 103 – 0 | 96 – 7 |
| Set 38: Game Room v Waiting Room | 4.70 | 3.83 | 69 – 30 | 99 – 0 | 96 – 3 |
| Set 39: Landing Deck v Window Seat | 4.82 | 2.28 | 93 – 6 | 97 – 2 | 91 – 8 |
| Set 40: Bar v Warehouse | 4.81 | 2.87 | 86 – 11 | 95 – 2 | 93 – 4 |
| Set 41: Bakery v Apartment Building | 5.41 | 3.64 | 86 – 14 | 99 – 1 | 98 – 2 |
| Set 42: Needleleaf Forest v Playroom | 7.03 | 4.70 | 92 – 9 | 101 – 0 | 100 – 1 |
| Set 43: Outdoor Window v Roundabout | 4.58 | 1.10 | 98 – 2 | 98 – 2 | 76 – 24 |
| Set 44: Reception v Golf Course | 4.03 | 1.73 | 94 – 7 | 97 – 4 | 88 – 13 |
| Set 45: Staircase v Plaza | 6.46 | 4.78 | 85 – 16 | 101 – 0 | 100 – 1 |

Table 1: Ability Scores and Number of Wins for Human (H.), ECII (E.), and Random (R.) Explanations. Note that random explanations were set as the reference point in the Bradley-Terry analysis and so their ability scores were always equal to 0, and thus are not displayed here.

Which of these better represents what the images in group A have that the images in group B do not?

Bake, Bakery, Bread, Indoor, Product, Store, Woman

Basket, Bread, Cake, Ceiling, Floor, Person, Wall

Figure 1: Experiment 1 task interface, with human explanation presented on the left and ECII explanation on the right.
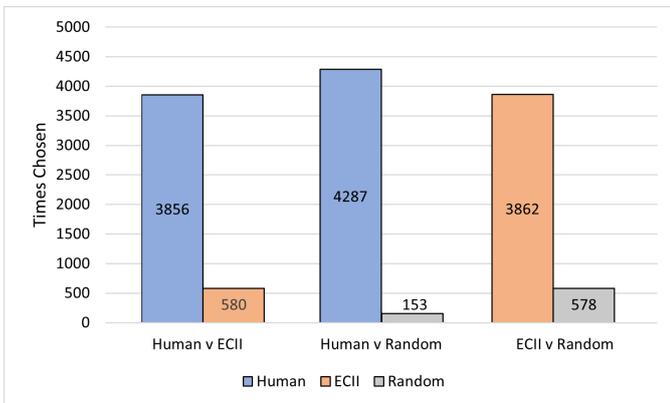


Figure 2: Number of times participants chose different explanation types in Experiment 1.

analyze the data after each new cohort, ceasing data collection if we reached a sample size of 250. We did not need to collect more than the initial 100 responses.

### 4.2.2. Materials

*Image sets.* A total of 16 image sets were created for this study. Images were taken from the Google Open Images database, which includes object tagging and scene category tags. The annotated tags in these images are crowd-generated and thus noisy. Limited manual curation was therefore performed as part of the image selection process: for example, images with very few object tags (5 or fewer) were removed from the dataset,

as images with too few labels would not provide the classifier with enough context to make a classification. Each training set included images from both the target category and images from multiple non-target categories. All images from the target category with sufficient object tags were included. An equal number of images were drawn randomly from non-target categories to provide a balanced dataset for each image set. We utilized the following 16 scene categories: bakeries, bathrooms, bedrooms, bridges, cafes, classrooms, dining rooms, gazebos, greenhouses, kitchens, lobbies, offices, pantries, parks, parking lots, and libraries.

Selected images were fed to a logistic regression classifier to classify scene images into target / non-target categories based on their tags. We utilized 10-fold cross validation to train and test the classifier. The input stimuli were represented as binary object vectors indicating the presence or absence of each object tag in that image. The vector space was generated by taking a set of all the objects present in the image dataset for a target scene category (e.g., all object tags present in the kitchen vs. non-kitchen image set). The classifier then outputs a binary decision as to whether each image is part of the target set or not. These classifications were then grouped into four categories based on ground truth scene categories: true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). Up to nine images from each grouping were included in final image sets displayed to participants, randomly selected from the full set of images in each error group. The same five catch trial image sets from Experiment 1 were also included in Experiment 2.

*Explanations.* Explanation generation in Experiment 2 followed the same process as in Experiment 1, with the exception that only two types of explanations were used for Experiment 2: 1) ECII (machine-generated) explanations and 2) human "gold standard" explanations. Due to the particular design of the study, the semi-random explanations would not have provided additional information and therefore were not used.

ECII explanations were again created by providing the image set's object tags (but not scene category tags) to the ECII algorithm, which provided ratings of concepts in the same manner as in Experiment 1. Object tags from all images in a set were provided to ECII, not only those in the subset of images displayed to participants. Explanations were generated for four image sets for each scene category: FP vs TN, TP vs FN, TP vs FP, and FN vs TN. These four comparisons were chosen because they were deemed most likely to have explainable differences. For example, we reasoned that FP scenes might contain a restricted set of defining elements that distinguish them from TN scenes – e.g., FP kitchen images might be more likely to contain kitchen-relevant features that TN kitchens do not have, such as images of offices that contain microwave ovens). In contrast, FP scenes could contain a wide variety of elements that distinguish them from TP scenes (e.g., a FP kitchen could contain an oven in a warehouse filled with myriad other objects) or from FN scenes (e.g., a FP kitchen and FN kitchen could both contain a large set of kitchen-irrelevant features that are hard to capture concisely).

Explanations in Experiment 1 tended to be more abstract, featuring thematic or high-level category concepts, unlike the very concrete explanations created by human raters, which tended to feature specific objects in the image sets. In order to equate the concreteness of the ECII and human explanations in Experiment 2, we conducted an analysis of the concreteness of the terms used in human explanations (described below), using concreteness ratings provided by [6]. We assessed the concreteness of each concept provided by human raters and in the ECII explanations and confirmed that ECII concepts were noticeably less concrete than the concepts provided by all three human raters. In order to adjust the concreteness of ECII explanations, we eliminated all concepts provided by ECII with a concreteness score of less than 3.5 (See Figure 3 for an example of the concreteness of ECII explanations before and after adjustment, compared to the concreteness of the concepts provided by the three human raters). Explanations were then created by taking the seven highest rated concepts from the filtered set returned by ECII.

Explanations generated by ECII can be ranked with respect to several criteria; to determine which criteria would be most helpful for human participants, we ran a pilot study in which we compared the effectiveness of explanations maximized with respect to F1, precision, and recall. Discriminability from this pilot study is plotted in Figure 4. Overall, we found that all three explanation types generated similar discriminability for our "TP vs. FP" and "TP vs. TN" comparisons, and that explanations maximized with respect to F1 and precision yielded



Figure 3: Average concreteness in three image groups for ECII and human concepts before and after adjustment.
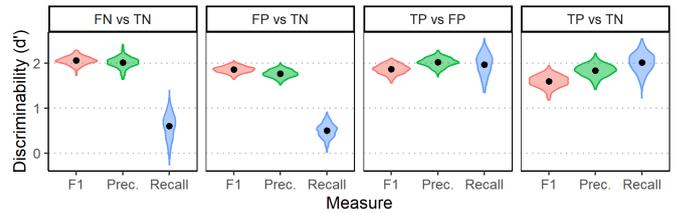


Figure 4: Discriminability as a function of comparison (noted in the panel titles) and explanation type (maximized with respect to F1, precision, or recall) for our pilot study. Violins indicate the distribution of participants' d' values. Black dots indicate group means.

similar discriminability for the remaining two comparisons. However, only F1 generated higher discriminability than recall-based explanations for the "FN vs. TN" comparison, whereas precision discriminability was not convincingly better than that of recall (i.e., BF < 3). We therefore elected to use explanations that were generated with respect to F1 scores in Experiment 2.

Human "gold standard" explanations were again created by providing the image sets (but not object tags or scene category tags) to three human raters. Raters again independently generated concepts for the same four comparisons for each scene category that ECII did. Gold standard explanations were created by pooling the concepts humans provided following the same procedure as in Experiment 1.

Catch trials utilized the same explanations as in Experiment 1, although the randomly generated explanations were not included due to the design of Experiment 2.

Because participants did not need to compare two different explanations in Experiment 2, concepts were not presented in alphabetical order, unlike in Experiment 1.

*4.2.3. Design*
Experiment 2 utilized a fully crossed within-subjects design with 2 (explanation type: ECII, human) × 4 (comparison sets: FP vs TN, TP vs FN, TP vs FP, FN vs TN) conditions for a total of 8 conditions. Participants were randomly assigned to 8 (out of 16) scene categories. For each scene category, partici-

pants saw all four comparisons. For half of the categories, they saw the human explanations, and for the other half they saw the ECII explanation. This resulted in a total of 32 target trials for each participant (4 in each explanation type × comparison set condition).

### 4.2.4. Procedure

After providing consent, participants completed brief training on the task, including instructions about how concepts and explanations were defined in this study. Participants then began completing trials. The 37 trials (32 assigned targets and 5 catch trials) were presented in a random order. Figure 5 shows what the stimuli presentation and response options looked like to participants.

### 4.3. Results

We analyzed choice data from Experiment 2 with a Bayesian hierarchical signal-detection model (SDT; [21]). In its simplest form, SDT parameterizes binary choices as a function of discriminability (i.e., $d'$) and bias (i.e., $c$). Discriminability is the distance between a posited mental Gaussian distribution of noise and a Gaussian distribution of signal. Bias is the degree to which an individual's decision threshold varies from ideal placement (i.e., the point at which the signal and noise distributions are equally likely); negative or positive bias indicates a tendency to respond with one choice over another. This separation of discriminability from bias ensures that individuals are not misjudged as accurate when they generate one response more often than the other (e.g., consider a student who achieves 50% accuracy on a balanced true-false exam by always responding "true").

Our SDT model was estimated in Stan probabilistic language [9] via the "brms" package [8] in R statistical software. We estimated population-level coefficients for discriminability and bias for each comparison and explanation type (i.e., ECII versus humans), as well as participant- and scene-level effects for discriminability and bias. Population-level coefficients were assessed with Bayes factors, which were estimated via Savage-Dickey ratios [52]. Following recommendations from [25], we considered evidence as convincing once it supported one hypothesis over another by a 3:1 ratio.

Because our analyses were rooted in Bayesian estimation, we took advantage of optional stopping during data collection [44]. Our primary interest was in evaluating whether (a) human- and ECII-generated explanations facilitated set discrimination to different degrees and (b) whether human- and ECII-generated explanations yielded non-zero discriminability. We therefore elected to initially collect data from 100 individuals, and then collect data in cohorts of 50 participants either until we (a) found convincing evidence in favor of either the null or alternative hypothesis (i.e., a Bayes factor either $\geq 3$ or $\leq 0.33$) or else (b) collected data from 250 people total. We ultimately stopped data collection after 100 participants.

Estimates of participants' discriminability are plotted in Figure 6. Population-level estimates of discriminability and bias

| Source | Comparison | $d'$ | $BF_{10}$ | $c$ | $BF_{10}$ |
|---|---|---|---|---|---|
| ECII | FN vs. TN | 2.06 | 21.04 | -0.65 | 1.91 |
| | FP vs. TN | 1.74 | $4.82 \times 10^{15}$ | -0.75 | $1.73 \times 10^{12}$ |
| | TP vs. FP | 1.97 | $2.20 \times 10^{13}$ | -0.73 | 38.30 |
| | TP vs. TN | 1.68 | 169.21 | -0.78 | 55.97 |
| Human | FN vs. TN | 15.31 | $1.71 \times 10^{5}$ | -4.74 | $5.18 \times 10^{8}$ |
| | FP vs. TN | 3.05 | $4.37 \times 10^{24}$ | -1.46 | $3.60 \times 10^{109}$ |
| | TP vs. FP | 9.69 | $4.51 \times 10^{6}$ | -2.74 | $5.28 \times 10^{17}$ |
| | TP vs. TN | 3.83 | $4.16 \times 10^{17}$ | -1.66 | $2.20 \times 10^{103}$ |

Table 2: Estimates of Discriminability ($d'$) and Bias ($c$) for ECII- and Human-Generated Terms Across Four Comparison Types. Note that columns, from left to right, indicate the source of explanations, the comparison being evaluated, estimated discriminability, the corresponding Bayes factor for discriminability, estimated bias, and the corresponding Bayes factor for bias. Negative/positive values of bias indicate a propensity to choose Set A/B.

| Comparison | $d'_{\text{diff}}$ | $BF_{10}$ | $c_{\text{diff}}$ | $BF_{10}$ |
|---|---|---|---|---|
| FN vs. TN | 13.25 | 322.01 | -4.10 | 493.89 |
| FP vs. TN | 1.33 | 210.39 | -0.71 | 260.20 |
| TP vs. FP | 7.72 | 210.48 | -2.01 | 826.99 |
| TP vs. TN | 2.15 | 62.59 | -0.88 | 4.49 |

Table 3: Estimates of Differential Discriminability ($d'_{\text{diff}}$) and Bias ($c_{\text{diff}}$) between ECII- and Human-Generated Terms Across Four Comparison Types. Note that columns, from left to right, indicate the comparison being evaluated, the differential estimated discriminability (i.e., $d'_{\text{human}} - d'_{\text{ECII}}$), the corresponding Bayes factor, differential estimated bias (i.e., $c_{\text{human}} - c_{\text{ECII}}$), and the corresponding Bayes factor.

as a function of explanation type are reported in Table 2; differential values between humans and ECII are reported in Table 3. Human-generated explanations yielded higher discriminability for all four of our comparisons. For both human-generated and ECII-generated explanations, we found evidence for non-zero $d'$. Thus, while ECII is not as effective as humans at generating helpful descriptions of set differences, it is nonetheless still an effective intervention across all comparisons that we evaluated. Additionally, both sets of explanations yielded a bias toward responding with "Set A" (but note that correct responses were balanced across "Set A" and "Set B"), and this bias was consistently stronger for human-generated explanations.

### 4.4. Discussion

Experiment 2 suggests that ECII can generate sensible explanations that help humans to distinguish misclassified scenes from correctly classified ones. Explanations generated by humans were even more helpful, which is not entirely unexpected, but our findings nonetheless suggest that ECII may be deployed to effectively facilitate human analysts' error detection for AI classification. Another notable result from Experiment 2 is that human explanations generated more biased responding in addition to more discriminative responses; while we did not anticipate such an outcome, one plausible explanation is that participants were randomly responding more frequently when viewing
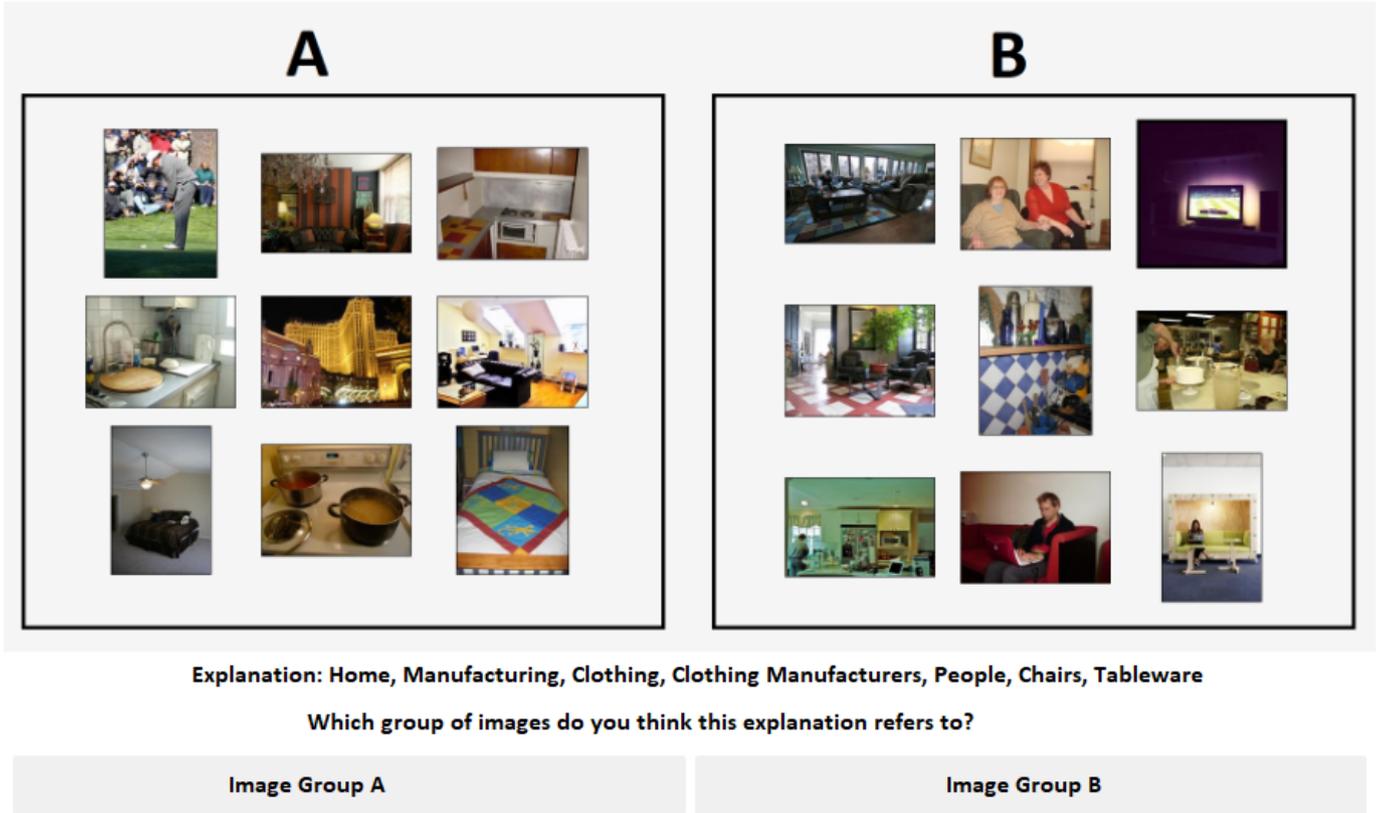
Figure 5: Experiment 2 task interface, with ECII explanation referring to image group B.
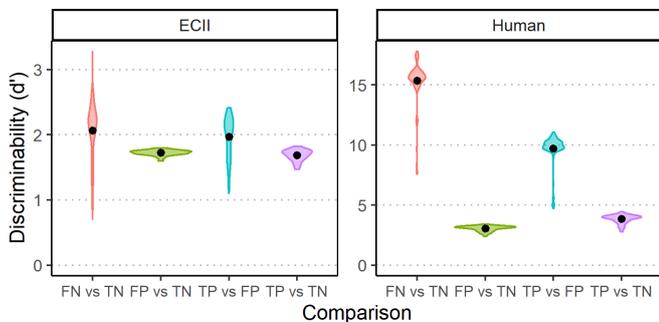


Figure 6: Discriminability as a function of comparison (noted on the abscissa) and explanation type (noted in the panel titles) for Experiment 2. Violins indicate the distribution of participants' $d'$ values. Black dots indicate group means. Note the different scales in the ordinates.

ECII explanations (i.e., because ECII was not as interpretable, despite our attempts to increase the concreteness of its terms), whereas human explanations encouraged more systematic (if also biased) responding. Overall, our findings suggest that ECII could serve as an effective teammate for facilitating human analysts' detection of misclassified entities.

One question that we have not addressed directly in Experiment 2 is how humans detect misclassified sets. For example,

did participants notice that the error images include/exclude certain objects that the correct images do/do not contain? Future research on this topic could employ eye-tracking methods to evaluate which particular scenes, or specific features thereof, people tend to assess prior to making their choices.

## 5. Related Work

Despite significant recent research efforts regarding approaches to XAI, this area of research is still in its infancy and in need of new ideas. Explanation methods so far can be divided into two broad categories, one focusing on local explanations, and one focusing on global explanations. Local explanation algorithms seek understanding of individual predictions, while global explanation algorithms seeks understanding of the overall behavior of the deep learning model.

Local explanation approaches can in turn be divided into 4 broad categories, a) those focusing on feature importance, b) those using saliency maps, c) those that are prototypes or examples based, and d) those based on counterfactuals. As the name suggests, feature importance algorithms attempt to identify which input features are more important for the prediction. Popular feature importance algorithms are LIME [43] and SHAP[34]. State-of-the-art saliency map based algorithms are [48, 38]. Prototype/example based algorithms seek to explain

individual decisions in terms of other examples (whether artificial or based on real data). Popular algorithms in this category make use of so-called influence functions [28], and a corresponding survey on this can be found in [37]. A popular counterfactuals based algorithm was provided in [20], and an overview of counterfactual explanations can be found in [51].

Global explanation approaches can be categorized into 3 broad categories, a) those that are based on collections of local explanations like SP-LIME [43], b) those that are representation based like TCAV [26], and c) those based on model distillation like [49, 16]). State of the art algorithms have issues regarding faithfulness [1], stability [2, 4] and fragility [27, 19].

The use of background or domain knowledge to produce or enhance explanations is gaining attention very recently. Domain knowledge comes in the form of simple concept tags, knowledge graphs (essentially data organized in graphs) or ontologies (knowledge bases using formal logic with formal semantics). For the use of concepts, one of the pioneering contributions was by Been Kim et. al. who showed the value of simple concepts to produce human understandable explanations [26]. Knowledge graph were used, for example, to explain individual decisions in the context of stock trend predictions [14]. Ontologies have been used in several lines of research for explanation. In the transfer learning domain, they were used to obtain an understanding which features are worth transferring or not [18]. The publication [39] introduced Doctor XAI as a model-agnostic explainer that focuses on explaining medical diagnosis predictions. They show how exploiting temporal dimensions in the data together with domain knowledge encoded as a medical ontology improves the quality of mined explanations. Trepan-reloaded [12] is an extension of the original Trepan [13] algorithm, and the authors demonstrate the use of ontologies for producing explanations that are more understandable by humans. Though, as we have just seen, there exists XAI research that makes use of domain knowledge, concept induction for making use of background knowledge is not yet a prevalent method. We are only aware of our own preliminary investigations [47] and [42] that make use of concept induction.

Among the many approaches to XAI, some explanations are generated for end users, and some explanations are meant for system developers, but in any case it remains an important criterion that an explanation makes sense to a human. To evaluate which explanations are better than others in this respect, different methods are being proposed [41, 17, 10, 24]. And while the different approaches are being evaluated, it is also important to keep in mind the variability of human understanding [36].

## 6. Conclusions and Future Work

The results reported herein clearly indicate that concept induction is an approach to explanation generation over background knowledge that produces explanations that are meaningful at a human level. But of course, the studies presented herein are only first steps on a longer journey towards making the approach useful in practice. We discuss several lines of investigation that can follow up on the presented results.

It appears a reasonable hypothesis that more expressive background knowledge – e.g. a highly axiomatized ontology – could yield even better and more fine-grained explanations. In fact, explanations could even include individuals in the form of nominal classes, i.e. reasoning would be over the schema (ontology) and a corresponding knowledge graph to derive explanations. A key technical hurdle remains for this, though: While ECII scales to the task, its heuristic relies on an ontology that is merely a class hierarchy, and other algorithms like those underlying DL-Learner do not scale sufficiently. Improved algorithms and systems – heuristic or not – for concept induction therefore need developing, before expressive background knowledge can be used at scale.

Another rather obvious avenue is in exploring concept induction for explanation generation in other settings, and in particular for deep learning systems that are inherently black boxes. One the one hand, an analysis of type 1 and type 2 errors by a trained network could lead to the identification of commonalities of inputs that produce errors, which could then be corrected, e.g., by including more training samples of the error-producing type. Another possible line of research is the use of concept induction to explain hidden layer activation patterns.

While concept induction is the key mechanism underlying our study, it appears to us that the choice of background knowledge is probably rather critical. Indeed, we also experimented with other hierarchies, e.g., SUMO [46], but the results were not as convincing. To advance, a systematic exploration of the effects of different types of background knowledge appears to be in order.

Besides the above mentioned immediate next steps, of course there also remains the broader vision for our work: The engineering of a human-machine data analysis support system that uses background information to provide explanations to a human analyst.

## Acknowledgement

## References

[1] J. Adebayo, J. Gilmer, M. Muelly, I. J. Goodfellow, M. Hardt, and B. Kim. Sanity checks for saliency maps. In S. Bengio, H. M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 9525–9536, 2018.

[2] D. Alvarez-Melis and T. S. Jaakkola. On the robustness of interpretability methods. *CoRR*, abs/1806.08049, 2018.

[3] A. B. Arrieta, N. Díaz-Rodríguez, J. D. Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera. Explainable Artifical Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58:82–115, 2020.

[4] N. Bansal, C. Agarwal, and A. Nguyen. SAM: the sensitivity of attribution methods to hyperparameters. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR Workshops 2020, Seattle, WA, USA, June 14-19, 2020*, pages 11–21. Computer Vision Foundation / IEEE, 2020.

[5] R. Bradley and M. Terry. Rank analysis of incomplete block designs: I. The method of paired comparisons. *Biometrika*, 39(3):324–345, 1952.

[6] M. Brysbaert, A. Warriner, and V. Kuperman. Concreteness ratings for 40 thousand generally known english word lemmas. *Behavior Research Methods*, 46(3):904–911, 2014.

[7] L. Bühmann, J. Lehmann, and P. Westphal. DL-Learner – A framework for inductive learning on the semantic web. *J. Web Sem.*, 39:15–24, 2016.

[8] P. Bürkner. Bayesian item response modeling in R with brms and Stan. *J. Stat. Softw.*, 100(5), 2021.

[9] B. CArptenter, A. Gelman, M. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, , P. Li, and A. Riddell. Stan: A probabilistic programming language. *Journal of Statistical Software*, 76(1):1–32, 2017.

[10] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso. Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(8):832, 2019.

[11] D. Castelvecchi. Can we open the black box of AI? *Nature News*, 538, 2016.

[12] R. Confalonieri, T. Weyde, T. R. Besold, and F. M. del Prado Martín. Using ontologies to enhance human understandability of global post-hoc explanations of black-box models. *Artif. Intell.*, 296:103471, 2021.

[13] M. W. Craven and J. W. Shavlik. Extracting tree-structured representations of trained networks. In D. S. Touretzky, M. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems 8, NIPS, Denver, CO, USA, November 27-30, 1995*, pages 24–30. MIT Press, 1995.

[14] S. Deng, N. Zhang, W. Zhang, J. Chen, J. Z. Pan, and H. Chen. Knowledge-driven stock trend prediction and explanation via temporal convolutional network. In S. Amer-Yahia, M. Mahdian, A. Goel, G. Houben, K. Lerman, J. J. McAuley, R. Baeza-Yates, and L. Zia, editors, *Companion of The 2019 World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019*, pages 678–685. ACM, 2019.

[15] F. K. Došilović, M. Brčić, and N. Hlupić. Explainable artificial intelligence: A survey. In *41st International convention on information and communication technology, electronics and microelectronics (MIPRO)*, pages 210–215. IEEE, 2018.

[16] N. Frosst and G. E. Hinton. Distilling a neural network into a soft decision tree. In T. R. Besold and O. Kutz, editors, *Proceedings of the First International Workshop on Comprehensibility and Explanation in AI and ML 2017 co-located with 16th International Conference of the Italian Association for Artificial Intelligence (AI*IA 2017), Bari, Italy, November 16th and 17th, 2017*, volume 2071 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2017.

[17] M. J. Gacto, R. Alcalá, and F. Herrera. Interpretability of linguistic fuzzy rule-based systems: An overview of interpretability measures. *Information Sciences*, 181(20):4340–4360, 2011.

[18] Y. Geng, J. Chen, E. Jiménez-Ruiz, and H. Chen. Human-centric transfer learning explanation via knowledge graph [extended abstract]. *CoRR*, abs/1901.08547, 2019.

[19] A. Ghorbani, A. Abid, and J. Y. Zou. Interpretation of neural networks is fragile. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 3681–3688. AAAI Press, 2019.

[20] Y. Goyal, Z. Wu, J. Ernst, D. Batra, D. Parikh, and S. Lee. Counterfactual visual explanations. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2376–2384. PMLR, 2019.

[21] D. Green and J. Swets. *Signal Detection Theory and Psychophysics*. Wiley, New York, 2014.

[22] P. Hitzler, M. Krötzsch, B. Parsia, P. Patel-Schneider, and S. Rudolph, editors. *OWL 2 Web Ontology Language: Primer (Second Edition)*. W3C Recommendation, 11 December 2012. Available at http://www.w3.org/TR/owl2-primer/.

[23] P. Hitzler, M. Krötzsch, and S. Rudolph. *Foundations of Semantic Web Technologies*. Chapman & Hall/CRC, 2010.

[24] R. R. Hoffman, S. T. Mueller, G. Klein, and J. Litman. Metrics for explainable AI: challenges and prospects. *CoRR*, abs/1812.04608, 2018.

[25] H. Jeffreys. *Theory Of Probability*. Oxford University Press, Oxford, 3rd edition, 1961.

[26] B. Kim, M. Wattenberg, J. Gilmer, C. J. Cai, J. Wexler, F. B. Viégas, and R. Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). In J. G. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 2673–2682. PMLR, 2018.

[27] P. Kindermans, S. Hooker, J. Adebayo, M. Alber, K. T. Schütt, S. Dähne, D. Erhan, and B. Kim. The (un)reliability of saliency methods. In W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K. Müller, editors, *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, volume 11700 of *Lecture Notes in Computer Science*, pages 267–280. Springer, 2019.

[28] P. W. Koh and P. Liang. Understanding black-box predictions via influence functions. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 1885–1894. PMLR, 2017.

[29] J. Lehmann and L. Bühmann. ORE – A tool for repairing and enriching knowledge bases. In P. F. Patel-Schneider, Y. Pan, P. Hitzler, P. Mika, L. Zhang, J. Z. Pan, I. Horrocks, and B. Glimm, editors, *The Semantic Web – ISWC 2010 – 9th International Semantic Web Conference, ISWC 2010, Shanghai, China, November 7-11, 2010, Revised Selected Papers, Part II*, volume 6497 of *Lecture Notes in Computer Science*, pages 177–193. Springer, 2010.

[30] J. Lehmann, N. Fanizzi, L. Bühmann, and C. d'Amato. Concept learning. In J. Lehmann and J. Völker, editors, *Perspectives on Ontology Learning*, volume 18 of *Studies on the Semantic Web*, pages 71–91. IOS Press, 2010.

[31] J. Lehmann and P. Hitzler. Concept learning in description logics using refinement operators. *Machine Learning*, 78(1-2):203–250, 2010.

[32] J. Lehmann, R. Isele, M. Jakob, A. Jentzsch, D. Kontokostas, P. N. Mendes, S. Hellmann, M. Morsey, P. van Kleef, S. Auer, and C. Bizer. DBpedia – A large-scale, multilingual knowledge base extracted from Wikipedia. *Semantic Web*, 6(2):167–195, 2015.

[33] Y. Lou, R. Caruana, and J. Gerke. Intelligble models for classification and regression. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 150–158. ACM, 2012.

[34] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017.

[35] S. Muggleton. Inductive logic programming. *New Gener. Comput.*, 8(4):295–318, 1991.

[36] M. Narayanan, E. Chen, J. He, B. Kim, S. Gershman, and F. Doshi-Velez. How do humans understand explanations from machine learning systems? an evaluation of the human-interpretability of explanation. *CoRR*, abs/1802.00682, 2018.

[37] A. Nguyen, J. Yosinski, and J. Clune. Understanding neural networks via feature visualization: A survey. In W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K. Müller, editors, *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, volume 11700 of *Lecture Notes in Computer Science*, pages 55–76. Springer, 2019.

[38] D. Omeiza, S. Speakman, C. Cintas, and K. Weldemariam. Smooth grad-cam++: An enhanced inference level visualization technique for deep convolutional neural network models. *CoRR*, abs/1908.01224, 2019.

[39] C. Panigutti, A. Perotti, and D. Pedreschi. Doctor XAI: an ontology-based approach to black-box sequential data classification explanations.

In M. Hildebrandt, C. Castillo, L. E. Celis, S. Ruggieri, L. Taylor, and G. Zanfir-Fortuna, editors, *FAT\* '20: Conference on Fairness, Accountability, and Transparency, Barcelona, Spain, January 27-30, 2020*, pages 629–639. ACM, 2020.

[40] H. Paulheim. Knowledge graph refinement: A survey of approaches and evaluation methods. *Semantic Web*, 8(3):489–508, 2017.

[41] F. Poursabzi-Sangdeh, D. G. Goldstein, J. M. Hofman, J. W. Vaughan, and H. M. Wallach. Manipulating and measuring model interpretability. In Y. Kitamura, A. Quigley, K. Isbister, T. Igarashi, P. Bjørn, and S. M. Drucker, editors, *CHI '21: CHI Conference on Human Factors in Computing Systems, Virtual Event / Yokohama, Japan, May 8-13, 2021*, pages 237:1–237:52. ACM, 2021.

[42] T. Procko, T. Elvira, O. Ochoa, and N. D. Rio. An exploration of explainable machine learning using semantic web technology. In *16th IEEE International Conference on Semantic Computing, ICSC 2022, Laguna Hills, CA, USA, January 26-28, 2022*, pages 143–146. IEEE, 2022.

[43] M. T. Ribeiro, S. Singh, and C. Guestrin. " why should i trust you?": Explaining the predictions of any classifier. *arXiv preprint arXiv:1602.04938*, 2016.

[44] J. Rouder. Optional stopping: No problem for Bayesians. *Psychonomic Bulletin & Review*, 21(2):301–308, 2014.

[45] M. K. Sarker and P. Hitzler. Efficient concept induction for description logics. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 – February 1, 2019*, pages 3036–3043. AAAI Press, 2019.

[46] M. K. Sarker, J. Schwartz, P. Hitzler, L. Zhou, S. Nadella, B. S. Minnery, I. Juvina, M. L. Raymer, and W. R. Aue. Wikipedia knowledge graph for explainable AI. In B. Villazón-Terrazas, F. Ortiz-Rodríguez, S. M. Tiwari, and S. K. Shandilya, editors, *Knowledge Graphs and Semantic Web - Second Iberoamerican Conference and First Indo-American Conference, KGSWC 2020, Mérida, Mexico, November 26-27, 2020, Proceedings*, volume 1232 of *Communications in Computer and Information Science*, pages 72–87. Springer, 2020.

[47] M. K. Sarker, N. Xie, D. Doran, M. Raymer, and P. Hitzler. Explaining trained neural networks with semantic web technologies: First steps. In T. R. Besold, A. S. d'Avila Garcez, and I. Noble, editors, *Proceedings of the Twelfth International Workshop on Neural-Symbolic Learning and Reasoning, NeSy 2017, London, UK, July 17-18, 2017.*, volume 2003 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2017.

[48] R. R. Selvaraju, A. Das, R. Vedantam, M. Cogswell, D. Parikh, and D. Batra. Grad-cam: Why did you say that? visual explanations from deep networks via gradient-based localization. *CoRR*, abs/1610.02391, 2016.

[49] S. Tan, R. Caruana, G. Hooker, P. Koch, and A. Gordo. Learning global additive explanations for neural nets using model distillation, 2019.

[50] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.

[51] S. Verma, J. P. Dickerson, and K. Hines. Counterfactual explanations for machine learning: A review. *CoRR*, abs/2010.10596, 2020.

[52] E.-J. Wagenmakers, T. Lodewyckx, H. Kuriyal, and R. Grasman. Bayesian hypothesis testing for psychologists: A tutorial on the Savage-Dickey method. *Cognitive Psychology*, 60(3):158–189, 2010.

[53] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Scene parsing through ADE20K dataset. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 5122–5130. IEEE Computer Society, 2017.

[54] B. Zhou, H. Zhao, X. Puig, T. Xiao, S. Fidler, A. Barriuso, and A. Torralba. Semantic understanding of scenes through the ADE20K dataset. *Int. J. Comput. Vis.*, 127(3):302–321, 2019.