

# Knowledge Elicitation using Deep Metric Learning and Psychometric Testing

Lu Yin, Vlado Menkovski, and Mykola Pechenizkiy

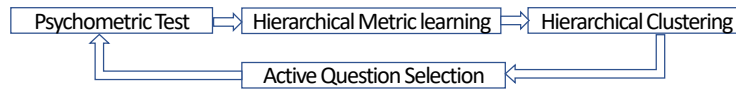
Eindhoven University of Technology, Eindhoven 5600 MB, Netherlands  
{l.yin, V.Menkovski, m.pechenizkiy}@tue.nl

**Abstract.** Knowledge present in a domain is well expressed as relationships between corresponding concepts. For example, in zoology, animal species form complex hierarchies; in genomics, the different (parts of) molecules are organized in groups and subgroups based on their functions; plants, molecules, and astronomical objects all form complex taxonomies. Nevertheless, when applying supervised machine learning (ML) in such domains, we commonly reduce the complex and rich knowledge to a fixed set of labels, and induce a model shows good generalization performance with respect to these labels. The main reason for such a reductionist approach is the difficulty in eliciting the domain knowledge from the experts. Developing a label structure with sufficient fidelity and providing comprehensive multi-label annotation can be exceedingly labor-intensive in many real-world applications. In this paper, we provide a method for efficient hierarchical knowledge elicitation (HKE) from experts working with high-dimensional data such as images or videos. Our method is based on psychometric testing and active deep metric learning. The developed models embed the high-dimensional data in a metric space where distances are semantically meaningful, and the data can be organized in a hierarchical structure. We provide empirical evidence with a series of experiments on a synthetically generated dataset of simple shapes, and Cifar 10 and Fashion-MNIST benchmarks that our method is indeed successful in uncovering hierarchical structures.

**Keywords:** Hierarchical knowledge elicitation · Psychometric testing · Deep metric learning · Active learning.

## 1 Introduction

Supervised learning models specified as a map from the data-space to a fixed set of labels is the panacea of machine learning (ML) applications. However, our goal is often to ‘solve’ a problem rather than to predict the labels. For example, consider a collection of texts or images – we may ask people to tag them, but our goal is not necessarily to predict tags, rather to understand the ‘latent’ taxonomy behind. Or one can imagine a scenario of a medical diagnosis involving images (e.g. diabetic retinopathy) and a set of diagnostic images and a set of labels denoting the severity of the disease (e.g. ‘no disease,’ ‘severity 1’, ... ‘severity 4’). From the very beginning, the domain expert is forced to project her



**Fig. 1.** Proposed hierarchy knowledge elicitation framework.

comprehensive knowledge on the five given discrete values. This limits our model both from its performance as it cannot learn from the full depth of knowledge as well from its interpretability as projecting to a fixed set of points further contributes to the ‘black-box-ness’ of our model. In contrast, adding a more rich hierarchical structure to the data adds a significantly larger understanding of the finer subtleties of the data stemming from the relationships between the concepts in the domain, which provide us a taxonomy of different concepts relevant to decision making.

The reason why so many solutions are formulated as mapping of datapoints (commonly high dimensional) to a fixed set of labels is that eliciting the knowledge from the domain and forming a training dataset is a difficult task. A data point can have a rich set of properties associated with it and eliciting an expert to provide an exhaustive annotation for each datapoint is prohibitively unscalable.

In this paper, we propose an approach that addresses this challenge by first developing a hierarchical representation of the data that facilitates more efficient annotation of data. Our method for knowledge elicitation consists of psychometric testing that captures the perceived relative distance of the datapoints for a given query (context) and a deep metric learning embedding algorithm. The embedded data are then structured in a hierarchical fashion that allows assigning properties and annotations to a large number of datapoints in a scalable way.

We evaluate the method empirically by showing we can capture the hierarchical knowledge structure in virtual responders for a given latent structure. Furthermore, using our method, we confirm the shape bias in the responses of a human participant.

## 2 The proposed approach

We present an approach that consists of four components: psychometric testing, hierarchical metric learning, active question selection, and hierarchical clustering. The overall framework is shown in Fig. 1.

First, in order to capture the user’s latent representation, we adopt a psychometric testing procedure [4] that relies on discriminative testing. The test collects information about the perceived relative differences to presented stimuli such as images or videos [17,2]. Measuring the perceived difference with a psychometric test is significantly more accurate [4] than directly quantifying a perceived value. On the other hand, we also do not need to present the expert with target labels at this stage, so they can express their knowledge without mapping it to a predefined set of concepts.

The captured responses from the psychometric testing are used to develop a distance metric using a deep metric learning method. The goal of the metric learning component is to develop an embedding of the data in a space where distances reflect the representation of the expert. In other words, images that lie closer together are perceived as more similar than images that are further apart in the embedding. To achieve this, we extend existing metric learning techniques and introduce dual-triplet loss with an adaptive margin. Unlike triplet loss in [15], in dual-triplet loss, we do not distinguish from anchor image and positive image and apply a symmetrical loss structure to align with psychometric testing and fully take advantage of every sample in a loss function. We explained it in detail in the proposed approach section.

Psychometric testing offers many advantages. However, the number of all possible discriminative tests is  $k$  combination of  $n$ , where  $k$  is 3, and  $n$  is the number of datapoints in the dataset. Asking all possible questions is typically not feasible. Nevertheless, to achieve a good embedding, we need only a small fraction of all possible questions. However, the quality of the embedding depends significantly on the selected questions as not all questions are equally informative. To address this, we develop a question selection approach using a Bayesian-based active learning method that selects questions with high uncertainty and high utility.

The embedding space is useful for many downstream tasks such as search and retrieval, but it also allows for efficient annotations. As now we have a metric to measure similarity, we can apply a label not only to a single datapoint but also to a region on the space covering multiple datapoints. Furthermore, if we can do this in a hierarchical fashion where labels at a different level of the hierarchy can be applied and then propagated to all lower levels. To achieve this, we combine the active question selection with a hierarchical clustering algorithm, such that we iteratively cluster and sub-clusters of the data to form a hierarchy and focus the question selection on sub-regions of the space. In this way, the approach starts first by forming the global distances and in turn organization of the data and then focuses on finer and finer differentiation as we go lower and lower in the hierarchy.

## 2.1 Psychometric Test

Psychometric testing is typically used to measure the subjectively perceived quantity of stimuli [4]. Different psychometric testing procedures are available, our method is a discriminative testing procedure, and a variation of the two-alternative forced-choice approaches. These methods typically present two alternatives to the participant that they are forced to choose from. One example of such set up would be to present two audio stimuli to a participant and ask which one is perceived as louder. To scale the perception of loudness, the experiment would consist of two pairs of stimuli, and the participant would need to answer which pair has a larger difference in loudness. Such experiments were also developed for estimating the quality of multimedia content [11]. Another option is to present three stimuli and ask to discriminate between the relative

difference between the two pairs formed by the three examples. However, many of these questions can be ambiguous, as the images may seem to have similar distances, or they may present different aspects that are not directly comparable in a pair-wise fashion [8].

To deal with this, we adapt the psychometric testing such that the participants are presented three objects  $a_1, a_2, a_3$ , and are forced to choose the most dissimilar one among them (Fig. 2). By carrying out this simple test, we can extract hierarchical knowledge from the data by simple ternary decision. The differences between these three objects are expressed as distances in our neural network model. The choice is based on the annotator’s personal perception. Therefore, different annotators may have different choices when facing the same question, and different hierarchical trees are created. We examine in detail how the method deals with these conditions in the experiment section.

Question: which object is the most dissimilar to the other two objects?

Option  $a_1$

Option  $a_2$

Option  $a_3$

**Fig. 2.** Three-alternative-forced choice interface.

## 2.2 Hierarchical Metric Learning

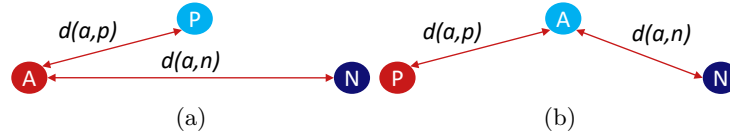
To successfully project the data in an embedded space that captures the latent representation of the expert, we train a model using a variation of the triplet loss function. The triplet loss given in Eq. 1 and its training process was introduced in [15], consists of training using triplet of datapoints  $(x_a^i, x_n^i, x_p^i)$ ,

$$L = \sum_{i=1}^N [d(x_a^i, x_p^i) - d(x_a^i, x_n^i) + m]_+ \quad (1)$$

where  $N$  is number of the possible triplets,  $d(x_a^i, x_p^i) = \|(f(x_a^i) - f(x_p^i))\|_2^2$  and  $d(x_a^i, x_n^i) = \|(f(x_a^i) - f(x_n^i))\|_2^2$ ,  $f(z)$  is the representation of the data point  $z$  in embedding space and  $f$  is our model or more precisely  $f_\theta$ , where  $\theta$  represent all the model parameters. hinge function  $[\cdot]_+$  indicate  $\max[0, \cdot]$  When the data is annotation with a fixed set, metric learning is typically implemented such that the triplets take the following roles:  $x_a^i$  is the anchor image, which has same label with positive image  $x_p^i$  and different label with negative image  $x_n^i$ . Therefore, the triplet loss produces a loss value when two images with the same label are further apart than two images with different labels.

Note that the form of triplet loss in Eq. 1 naturally fits with how we defined the 3AFC psychometric test. Based on the participant’s response, we can select the most different image as the negative sample, and the other two images as positive and anchor images.

However, the anchor and the positive image are not interchangeable in the triplet loss term as the distance of the anchor to the positive and negative image are compared. One can imagine a scenario where selecting one of the images as an anchor results in a gradient for the model and selecting the other one does not, or results limited gradient, as depicted in Fig. 3.



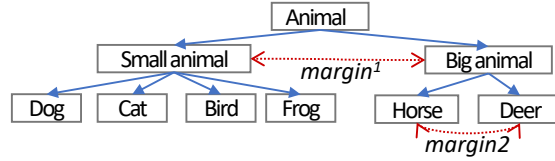
**Fig. 3.** (a) and (b) are different cases of choosing the anchor image. A represent the anchor image, N denotes negative and P for positive. Negative is chosen in 3AFC test and fixed, and we need to select the anchor image from the rest two. If  $d(a, n) - d(a, p)$  is too large there will be limited gradients, and if it is larger than a margin value there will be no gradient. Since in case (a) the anchor is farther away from the negative image than in case (b), it is more likely that case(a) results in limited or zero gradients while case (b) results in more gradients.

To deal with this, we propose a dual-triplet loss function in which there is no specific differentiation between the anchor and the positive image but rather define two positive and one negative image. The participant gives an answer about which image is the furthest from the other two, and set it to be negative. The dual-triplet loss is defined as,

$$L = \sum_{i=1}^N [d(x_{p1}^i, x_{p2}^i) - d(x_n^i, x_{p1}^i), +m_a^i]_+ + [d(x_{p1}^i, x_{p2}^i) - d(x_n^i, x_{p2}^i), +m_a^i]_+ \quad (2)$$

where,  $x_{p1}^i$  and  $x_{p2}^i$  are two positive images, and  $x_n^i$  are the negative image chosen by annotator during 3AFC test,  $N$  is the number of sampled triplets,  $m_a^i$  is an adaptive margin which will be explained later. Now negative image is compared with each of the rest two images, and both situations in Fig. 3 are considered. Every sampled is used twice in one triplet function to produce more gradients and to prevent zero gradients situation happen.

The triple loss margin has a significant impact on the performance of the metric learning model. Different methods have been developed to determine the value of the margin. Some of these approaches [3] propose using a margin value that adaptively changes during the training process. Since we collect the responses in a hierarchical manner, the margin needs to adapt accordingly, from large to small, as the questions become more focused. For example, if the responses are collected over the whole dataset, the margin needs to be appropriately large, and if we descent and become more focused in the testing, the margin needs to adapt and be smaller, as shown in Fig. 4. Specifically we compute the margin



**Fig. 4.** Questions about small animal and big animal are more general and have a larger diversity than questions about horse and deer. So  $margin^1$  should be larger than  $margin^2$ .

$m_a$  value based on the diversity between the three samples in a triplet  $i$  and is defined as,

$$m_a^i = m_h^i + \gamma d_H^i \quad (3)$$

where  $m_h^i$  is a constant value that guarantees  $m_a^i$  is not zero and  $\gamma d_H^i$  is the adaptive part.  $\gamma$  is a hyper-parameter that we set during training.  $d_H^i$  is a diversity factor of the triplet  $i$ . We calculate the diversity factor at a node  $i$ , by computing the average distance of the centroids of all sub-clusters of the cluster associated with the node  $i$ , as given in Eq. 4,

$$d_H^i = \frac{1}{n^{i^2} - n^i} \sum_{c_p, c_k \in L^i} \|c_p - c_k\|_2^2 \quad (4)$$

where  $L^i$  is a node where question  $i$  is sampled from in the hierarchy.  $n_i$  is the number of child nodes in  $L^i$ ,  $c_p$  is the center of cluster  $p$ , and  $c_k$  is the center of cluster  $k$ .

### 2.3 Active Questions Selection

Given a data set containing  $B$  images, there are  $\binom{B}{3}$  potential questions for the annotators to answer. It is impractical to answer all of them, and randomly selecting the questions is suboptimal with respect to the efficiency of the training process. To address this, we develop an active question selection scheme.

We start by randomly selecting  $m$  questions for the annotators to answer. With the answers, we train the model  $M$ , create a hierarchical tree  $H$ , and construct a set of knowledge pool containing the answered questions  $D$ . The method iteratively repeats the process of question selection and update of the model  $M$  as well as the hierarchical representation of the data  $H$ . To select the questions for the following iteration, our method takes two steps. In the first step, a set of questions is proposed uniformly sampled from each node in the hierarchy  $H$ , and in the second step, we reject some of the proposed questions for which we do not expect high utility.

By sampling from each level of the hierarchy in the first step, we maximize the probability that the model will receive information both about the global

distribution of the data as well as the specific differences at a finer level of detail. In contrast, randomly selecting questions is sufficiently less efficient in achieving the same goal.

Then, we use a pool-based active-learning question rejection scheme. We consider not only the uncertainty of the questions but also its utility, that is, select questions with high uncertainty and high variance. High uncertainty means annotators are not so sure which option to choose, thus the question could be informative, and its answer has a higher probability of introducing information to our model. High variance means similar questions have been answered many times, and this is to avoid ambiguous questions (i.e., similar questions have occurred too many times, but annotators are still not sure how to answer).

In order to evaluate the uncertainty and utility, we need to calculate the distribution of answering possibility by Bayes' theorem. For a given question  $q_i$  consists of 3 samples,  $\theta = (\theta_1, \theta_2, \theta_3)$  denote the possibility of choosing  $a_1, a_2, a_3$ . We assume it follow dirichlet distribution  $Dir(\alpha_1, \alpha_2, \alpha_3)$ , and it's density can be written as,

$$p(\theta|\alpha) = \frac{\Gamma(\alpha_1 + \alpha_2 + \alpha_3)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\Gamma(\alpha_3)} \prod_{i=1}^3 \theta_i^{\alpha_i-1} (\theta_i \geq 0; \sum_{i=1}^3 \theta_i = 1) \quad (5)$$

where  $\Gamma$  denotes the gamma function,  $\alpha = (\alpha_1, \alpha_2, \alpha_3)$ .  $\alpha_1 = \alpha_2 = \alpha_3 = 1$  at beginning, means the possibility of choosing  $a_1, a_2$  or  $a_3$  is equal. The distribution is updated based on the previous answered questions as,

$$p(\theta|\alpha, D) = p(\theta|\alpha + m) \quad (6)$$

where  $D$  is the set of answered questions. We define the questions sampled from neighbour questions of  $a_1, a_2, a_3$  as  $q_i$ 's similar questions, and count the times of choosing  $a_1, a_2, a_3$  when facing similar questions of  $q_i$  in  $D$ . The counting number are defined as  $m = (m_1, m_2, m_3)$ .

Now we evaluate the uncertainty by the maximum of the expected possibility of choosing any of  $a_1, a_2, a_3$ . It does not matter which one we choose, because if any of them has a high possibility to be chosen, we can conclude that the model is confident about this choice based on the previous answers, and this question has less uncertainty. So if the maximum expected possibility  $e(q_i)$  of choosing  $a_1, a_2$  or  $a_3$  is higher than  $s_e$  we will reject this question. Similar, we can evaluate the utility by the sum of possibility variance  $var(q_i)$ . If it is too high, we can interpret it as after facing similar questions many times, the annotator still has an ambiguous answer, so we also reject it.

## 2.4 Hierarchical Clustering

After answering enough actively selected questions, we embedded the data using the deep metric learning method such that its semantic hierarchical information represented by the Euclidean distance in that space. The hierarchical structure is extracted by a hierarchical clustering algorithm. Generally, there are two categories of methods for this task, a top-down approach called "Agglomerative"

and a top-down approach called “Divisive” [14]. Since we have to decide the threshold very carefully if applying the “Agglomerative” approach, in our experiment, we choose the “Divisive” approach K-means to perform the divisive cluster, and Silhouette Coefficient is applied to help us choose  $K$ .

### 3 Related work

Our work is related to three lines of research (1) Building hierarchical image structure (2) Deep metric learning (3) Informative samples mining.

**Building hierarchical image structure** There are two branches explored in building a hierarchical image structure. One is with the help of language information. Such as WordNet is used to build a semantic hierarchic classifier [10]. However, the language-based similarity is not always what we want. When comparing whitefish, goldfish, and cat, the former two are closer in word semantic because they are all fish, but this could not suit our needs if we want to build a hierarchical pet data set based on their cuteness. On the other hand, hierarchies could also be formed based on image features that are extracted by SIFT, HOG [18], or deep convolutional network [22]. While able to group images based on their visual similarities, the hierarchies are not build based on the user’s perception or knowledge, so it is also difficult to cater to a special user’s preference (e.g. expert domain).

**Deep metric learning** Deep metric learning is a data-driven approach to learn the measurement of similarity. It learns a nonlinear mapping for the original data to an embedding. Contrastive loss [5] is proposed to encourage the distance of positive paired samples be closer than negative samples of a margin of  $m$ . It is extended to triplet loss [15], quadruplet loss [1], N-paired loss [16] where three, four,  $N$  samples are used in loss function, and lifted structure loss [12] which take advantage all the positive and negative pairs in a mini-batch. Although yield promising results in distances based computer vision applications such as face recognition [15], person re-identification [1], and few-shot learning [13], their studies focus on the flat distance which has a weaker hierarchical distances representation ability.

**Informative samples mining** For pair-based deep metric learning tasks, there are  $O(N^n)$  potential tuples for training,  $N$  is the sample number in training set, and  $n$  is sample number in a tuple. It is impractical to train with all of them due to the GPU memory or train time concern. Besides, a lot of them are redundant or less informative and contribute little gradients during training. A semi-hard mining scheme is used in FaceNet [15] to select informative tuples online inside a mini-batch. Harwood *et al.* considers the global structure in their smart mining [6] for producing effective samples with a low computational cost. A distance weighted sampling strategy [19] is proposed to select samplings uniformly based on their relative distance. While these strategies boost the model performance by learning a more discriminative flat representation of each category object, our work aims to map data to an embedding space where a hierarchical structure is also considered. Therefore, a hierarchical sampling method



with Bayesian-based active learning is proposed in our paper to better cope with our needs.

## 4 Experiments

As a first in our empirical evaluation, we aim to confirm that our approach can elicitate the latent hierarchical structure of a single participant. To remove any other sources of variation, we design an experiment with a virtual participant that always responds according to its given latent hierarchical model. In this manner, we can objectively evaluate the precision of the elicited model against the one given to the virtual participant.

Next, we expand on this using multiple virtual participants that have partial agreements of their latent models. In large scale human experiments, we do not expect that the internal hierarchical representations of people will fully agree with each other. Therefore, in this analysis we aim to assess whether our approach can extract the common hierarchical model of a group of participants that have models with partial agreements. Towards this we simulate multiple individual participants and test whether the model can elicit common hierarchy accordingly.

Then, we evaluate with a human participant. The goal of this experiment is to evaluate how well our method can deal with the additional variability involved in using responses from a person. In this setting, the additional challenge is that we cannot be precisely sure of the human’s latent hierarchical model of the data. To deal with this, we specifically design an experiment that would confirm a well understood psychological bias that humans have, specifically the shape bias [9]. The shape bias states that shapes are more important than other properties of an object, such as colors when we aim to discriminate between them. We develop an image dataset that allows us to test if we can confirm this bias in hierarchical models that we elicit from the responses.

Finally, we develop an experiment on a natural image dataset, specifically the Fashion-Mnist [20]. Our aim here is to test in natural settings. In this case, our evaluation is limited to the flat set of labels that the dataset comes with.

In all experiments, two evaluation metrics are considered: (1) *The accuracy of each cluster*. Every cluster in the hierarchy should share a similar concept. We tag a cluster by the majority shared concept in it, and calculate the accuracy of that node by counting the number of correctly assigned images and dividing by overall image number in that node. (2) *The dendrogram purity between the extracted structure and the ground-truth*. Since we want the hierarchy able to cater for special users, we can only claim a hierarchy is good when structure formed according to the ground-truth, which is the user’s latent perception. That can be evaluated by dendrogram purity [7]. Rather than the flat cluster purity, it is a holistic way to evaluate the whole hierarchical structure tree. The value of it ranges from 0 to 1, 1 means every sample fits perfectly for the ground truth and on the contrary for 0. Note that dendrogram purity is only applicable when the ground-truth hierarchy is given in advance, so we only calculate it

in simulation experiments, and as we are applying hierarchy clustering in this paper, every node of the hierarchy is a cluster.

#### 4.1 Single Virtual Participant on Cifar 10

With our proposed hierarchical knowledge elicitation method, we experiment on a snipped tiny Cifar 10. The data set containing 1000 images sampled randomly from the original Cifar 10 with 100 images in each category.

Since we might not be able to know a user’s latent perception before the experiment, in order to evaluate how well our extracted structure matches the user’s latent perception, we pre-define a specific hierarchy (see Fig. 5) and run the experiment by simulation. A virtual participant gives perfect responses based on the given hierarchical structure. 1000 questions are simulated for model training with a fixed margin of 0.4 at the beginning. 600 more questions are sampled per iteration and used for training with adaptive margin in the following 5 iterations.

After collecting 3400 automatically respond triplets, we show our results in Fig. 5. We can see that our model can elicitate a clear hierarchical structure, which matches the given one well. Note that the small and big animals do not split down to individual animal nodes. That may because animal images are similar to each other, and it is not easy for the model to get a clear cluster with limited animal training triplets. We calculate the accuracy of the extracted nodes which are all above 90%.

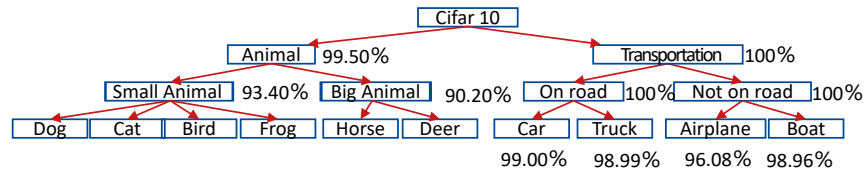


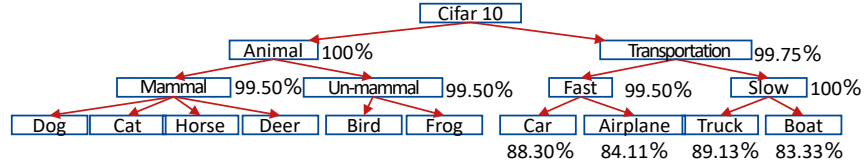
Fig. 5. Latent given hierarchy of a single virtual participant and the extracted result.

#### 4.2 Multiple Virtual Participants on Cifar 10

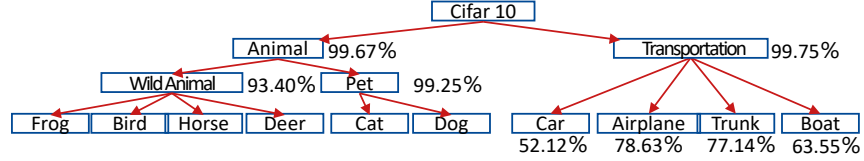
We further evaluate our method in multiply-users situations when they have partial disagreements. Two more participants are simulated with different pre-designed hierarchies, as shown in Fig. 6. We follow the same experiment setup with single-user situation in Section 4.1.

We illustrate our results in Fig. 6, where we can see our model is able to extract varying hierarchical structures according to different given latent hierarchies. Individual animal nodes still can not be clustered well due to their higher similarity and limited training samples.

We also mixed the collected triplets from all three users, and extract a general hierarchy from all the triplets (See Fig. 6). Due to enough training triplets,



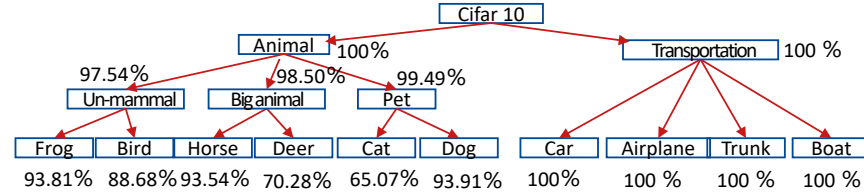
(a) Virtual participant 2 and the extracted result.



(b) Virtual participant 3 and the extracted result.

**Fig. 6.** Latent given hierarchy of multiple virtual participants and the extracted result.

individual animal nodes got extracted, and other nodes get better accuracies compared to single participant situations. We can also notice that both shared perception nodes (10 original Cifar 10 classes, animal and transportation) and individual perception nodes (big animal from participant 1, un-mammal from participant 2, pet from participant 3) can be seen in the structure. Therefore, our model can elicit a common structure that represents both shared and individual knowledge of all participants.



**Fig. 7.** Extracted common structure.

### 4.3 Real Participant on Synthetic Geometric Shape Data Set

The previous experiment has shown that our method is able to extract hierarchical structure by simulation. It is also crucial to test with real human responses. One challenge of real human involved experiments is how to evaluate our method since we don't know the latent hierarchical knowledge in advance. To deal with that, we test whether our method can uncover the well-studied shape bias in humans [9] on a synthetic geometric data set. The data set contains triangles, circles, rectangles 3 different shapes. Each shape contains 3 deformations, 5 colors, and 3 thicknesses (See the first layer in Fig. 8). Initially, we have the participant to answer 300 randomly chosen questions and train the model with a fixed

margin of 0.2. During each following iteration, 300 actively sampled questions are answered. After 5 iterations, 1500 more triplets are collected to train the model with an adaptive margin.

The extracted hierarchy are shown in Fig. 8. Because most of the nodes have 100% accuracy, we only mark the nodes which do not have perfect accuracies. We also notice there are few images located in the wrong clusters. Remember, all the questions are answered by real humans, so mistaken answers are inevitable, and that makes the dendrogram not perfect. The overall structure confirmed the human tend to believe shape is more important feature than other properties when they are asked to make choices.

#### 4.4 Real Participant on Fashion-Mnist

Another real participant involved experiment is conducted on real-life image data set Fashion-Mnist [20]. Accuracy can be calculated when the node of the hierarchy containing the original flat labels. The ground truth latent perception is unknown before experiments and is discovered by our proposed method. We sample questions on a tiny Fashion-Mnist containing 900 images with 90 images for each class. During the first iteration, we randomly choose 600 questions, and train the model with a fixed margin of 0.4. In the following 2 iterations, we actively select 600 questions every iteration and train the model with an adaptive margin.

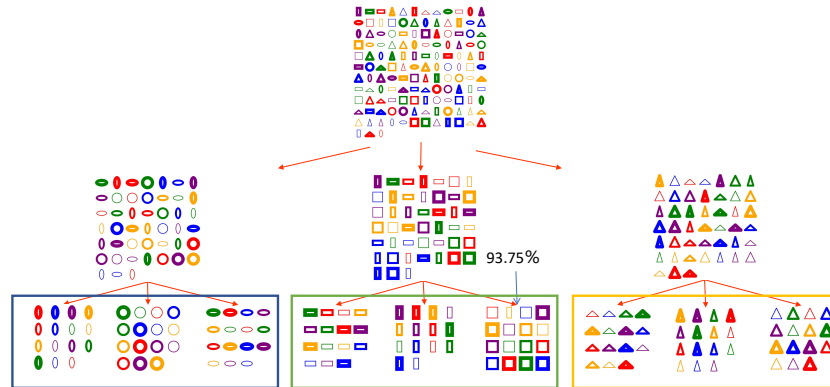
Fig. 9 illustrates the extracted hierarchy. Note that Due to page size limitation, it is just a diagram, not the real hierarchical images. The first splits is based on cloth(left) and accessories(right) which further split down into tops (T-shirt, pullover, dress, coat, shirt)/trouser and bags/shoes(sneaker, boot, sandal). Tops continue down into long sleeves/short sleeves, which are organized by the length of their sleeves regardless of the original cloth category. Further splits of shoes are wrap shoes(sneakers, boots) and sandals. For the nodes composed of original Fashion-Mnist labels, we can calculate the accuracy of each node.

Further splits of sandals and bags can be seen in Fig. 10, which presents us with a finer granularity than the original labels.

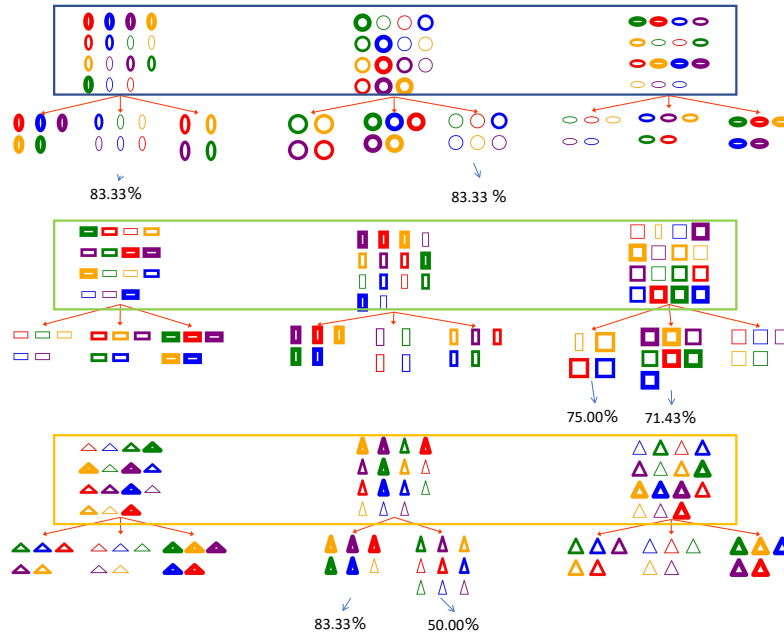
#### 4.5 Comparison and Ablation Study

We compare our method with three baselines and perform an ablation study to evaluate our proposed components on tiny Cifar 10. After training with 4000 triplets, dendrogram purities are calculated in different experiment settings and reported in Table 1. Trained curves are shown in Fig. 11.

Three benchmarks are chosen for comparison. First, pixel-level points are hierarchically clustered directly. Then, SIFT features are extracted and represented as BoVW [21] for hierarchical clustering. Besides, we train a CNN model with labels and hierarchically cluster on the last activation layer. For a fair comparison, we use the same CNN model as our HKE framework. From Table 1 and Fig. 11, we can see the dendrogram purities of pixel-level and SIFT BoVW are lower than 50%. For baseline CNN, the dendrogram purity is higher than 70%.



(a) The first three layers of the extracted hierarchy from the synthetic geometric shape data-set.



(b) The last layer of the extracted hierarchy from the synthetic geometric shape data-set.

**Fig. 8.** A four-layer hierarchy is extracted from the synthetic geometric shape data-set. The First 3 layers are shown in (a), and the last layer of the hierarchy is shown in (b) due to page size limitation. Clusters inside square-bound with the same color between (a) and (b) are the same. The first layer split into rectangles, circles, triangles three clusters, which is because the annotator thinks shape is the most discriminative property. In the second layer, shapes are again the discriminator; vertical stretched, horizontally stretched, and un-stretched images are cluster together. When the hierarchical tree goes down, we can notice that images begin to be organized based on different thicknesses, which are also shapes. Note only not perfectly pure nodes are marked with accuracy.

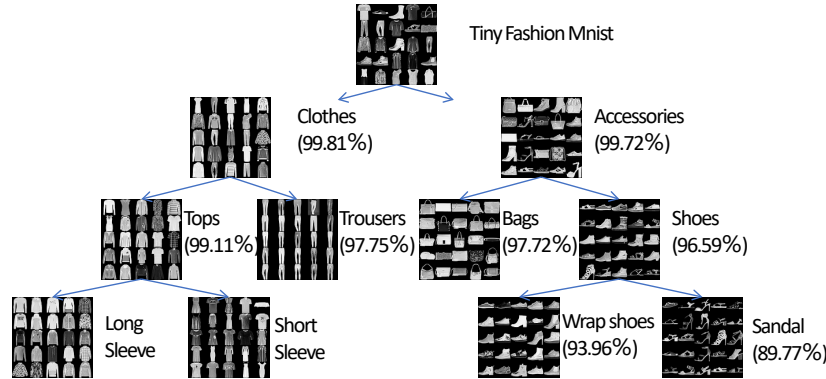


Fig. 9. Extracted hierarchy of tiny Fashion-Mnist.

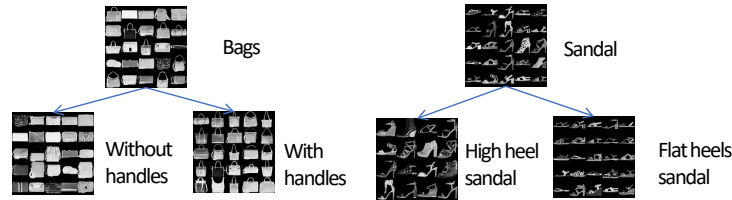


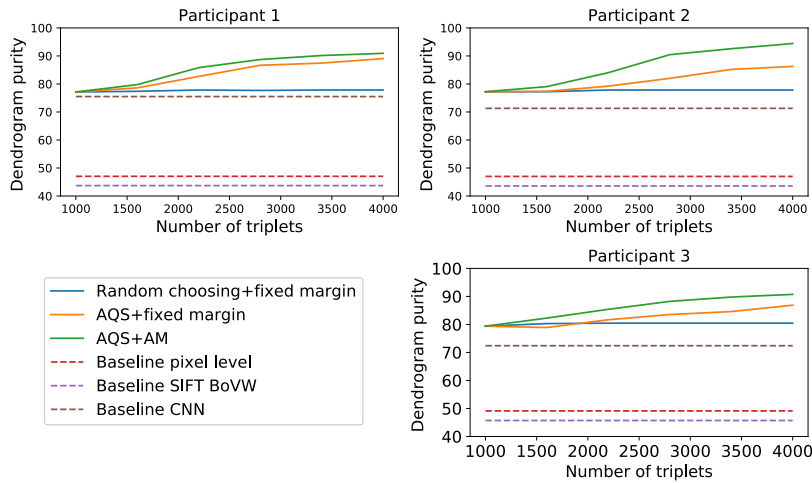
Fig. 10. Finer granularity of bags and sandals.

We also perform an ablation study to evaluate our proposed components. Experiment configuration are the same as in Section 4.2. In all the ablation experiments, dual-triplet is used to align with the 3AFC test. First, we sampling questions randomly and using a fixed margin of 0.4, the performance improves slightly in the first three iterations and keeps stuck even with more training data. Then, active question sampling(AQS) is applied to select questions with high-uncertainly and high utility, and performance gets a noticeable improvement. Finally, we integrate the adaptive margin(AM) and get the best performance.

When compare with HKE framework with baselines, we can notice that even with minimum components (Just 3AFC test, without AQS and AM), our proposed framework has better performance than all the baselines.

Table 1. Baselines and ablation study on Cifar 10 with 3 simulated participants

participant	Base line Pixel level	Base line SIFT BoVW	Base line CNN	Random choosing +fixed margin	AQS +fixed margin	AQS + AM
1	47.03	43.72	75.5	77.84	89.06	90.93
2	46.99	43.57	71.3	77.84	86.25	94.44
3	49.13	45.70	72.4	80.45	86.86	90.73



**Fig. 11.** Dendrogram purity training curves of different virtual responders. Our baselines are not trained with triplets, so they keep flat. Even with minimum components, our proposed HKE has better performance than baselines.

## 5 Conclusion

Developing a machine learning formulation for a given task usually involves compromises that are related to how much and what kind of data is available, but also very much around the quantity and kind of expert knowledge that can be collected. The latter tends to lead to oversimplifications, such as forcing a fixed set of labels and disregarding any more precise descriptions of the data. In this work, we present a method for knowledge elicitation from experts that can be applied to a broad set of machine learning formulations of problems. The results show that the proposed 3AFC testing with dual metric learning can extract latent hierarchical representations of data in a scalable way. This opens possibilities for applications in many domains where expert knowledge has a rich and complex underlying structure as in medical domain, biology as well as many industrial applications.

## References

1. Chen, W., Chen, X., Zhang, J., Huang, K.: Beyond triplet loss: a deep quadruplet network for person re-identification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 403–412 (2017)
2. Feng, H.C., Marcellin, M.W., Bilgin, A.: A methodology for visually lossless jpeg2000 compression of monochrome stereo images. *IEEE Transactions on Image Processing* **24**(2), 560–572 (2014)
3. Ge, W.: Deep metric learning with hierarchical triplet loss. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 269–285 (2018)

4. Gescheider, G.A.: *Psychophysics: the fundamentals*. Psychology Press (2013)
5. Hadsell, R., Chopra, S., LeCun, Y.: Dimensionality reduction by learning an invariant mapping. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06). vol. 2, pp. 1735–1742. IEEE (2006)
6. Harwood, B., Kumar, B., Carneiro, G., Reid, I., Drummond, T., et al.: Smart mining for deep metric learning. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2821–2829 (2017)
7. Heller, K.A., Ghahramani, Z.: Bayesian hierarchical clustering. In: Proceedings of the 22nd international conference on Machine learning. pp. 297–304. ACM (2005)
8. Hellinga, N., Menkovski, V.: Hierarchical annotation of images with two-alternative-forced-choice metric learning. arXiv preprint arXiv:1905.09523 (2019)
9. Landau, B., Smith, L.B., Jones, S.S.: The importance of shape in early lexical learning. *Cognitive development* **3**(3), 299–321 (1988)
10. Marszalek, M., Schmid, C.: Semantic hierarchies for visual object recognition. In: 2007 IEEE Conference on Computer Vision and Pattern Recognition. pp. 1–7. IEEE (2007)
11. Menkovski, V., Liotta, A.: Adaptive psychometric scaling for video quality assessment. *Signal Processing: Image Communication* **27**(8), 788–799 (2012)
12. Oh Song, H., Xiang, Y., Jegelka, S., Savarese, S.: Deep metric learning via lifted structured feature embedding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4004–4012 (2016)
13. Oreshkin, B., López, P.R., Lacoste, A.: Tadam: Task dependent adaptive metric for improved few-shot learning. In: Advances in Neural Information Processing Systems. pp. 721–731 (2018)
14. Rokach, L., Maimon, O.: Clustering methods. In: *Data mining and knowledge discovery handbook*, pp. 321–352. Springer (2005)
15. Schroff, F., Kalenichenko, D., Philbin, J.: Facenet: A unified embedding for face recognition and clustering. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 815–823 (2015)
16. Sohn, K.: Improved deep metric learning with multi-class n-pair loss objective. In: Advances in Neural Information Processing Systems. pp. 1857–1865 (2016)
17. Son, I., Winslow, M., Yazici, B., Xu, X.: X-ray imaging optimization using virtual phantoms and computerized observer modelling. *Physics in Medicine & Biology* **51**(17), 4289 (2006)
18. Wigness, M., Draper, B.A., Beveridge, J.R.: Efficient label collection for image datasets via hierarchical clustering. *International Journal of Computer Vision* **126**(1), 59–85 (2018)
19. Wu, C.Y., Manmatha, R., Smola, A.J., Krahenbuhl, P.: Sampling matters in deep embedding learning. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2840–2848 (2017)
20. Xiao, H., Rasul, K., Vollgraf, R.: Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. arXiv preprint arXiv:1708.07747 (2017)
21. Yang, J., Jiang, Y.G., Hauptmann, A.G., Ngo, C.W.: Evaluating bag-of-visual-words representations in scene classification. In: Proceedings of the international workshop on Workshop on multimedia information retrieval. pp. 197–206 (2007)
22. Zheng, Y., Fan, J., Zhang, J., Gao, X.: Hierarchical learning of multi-task sparse metrics for large-scale image classification. *Pattern Recognition* **67**, 97–109 (2017)