

# Context Embedding Networks

Kun Ho Kim   Oisín Mac Aodha   Pietro Perona  
California Institute of Technology

## Abstract

Low dimensional embeddings that capture the main variations of interest in collections of data are important for many applications. One way to construct these embeddings is to acquire estimates of similarity from the crowd. Similarity is a multi-dimensional concept that varies from individual to individual. However, existing models for learning crowd embeddings typically make simplifying assumptions such as all individuals estimate similarity using the same criteria, the list of criteria is known in advance, or that the crowd workers are not influenced by the data that they see.

To overcome these limitations we introduce Context Embedding Networks (CENs). In addition to learning interpretable embeddings from images, CENs also model worker biases for different attributes along with the visual context i.e. the attributes highlighted by a set of images. Experiments on three noisy crowd annotated datasets show that modeling both worker bias and visual context results in more interpretable embeddings compared to existing approaches.

## 1. Introduction

Large annotated datasets are a vital ingredient for training automated classification and inference systems. Labeling these datasets has been made possible by crowdsourcing services, which enable the purchasing of annotations from crowd workers. Unfortunately fine-grained categorization is very challenging for untrained workers. The alternative, obtaining annotations from experts, is equally impractical due to the fact that for many domains experts are few [24]. Instead of obtaining semantic fine-grained category-level labels, one can ask workers to label images in terms of their similarities and differences. This is intuitively much easier for untrained workers because it requires the comparison of images, a task that humans are naturally good at. This approach, however, presents its own challenges: 1) different workers may use different criteria when estimating the similarity between pairs of images, and 2) workers may be influenced by the set of images that they see when making their decisions i.e. ‘context’.

In Fig. 1 we see an example of three different crowd

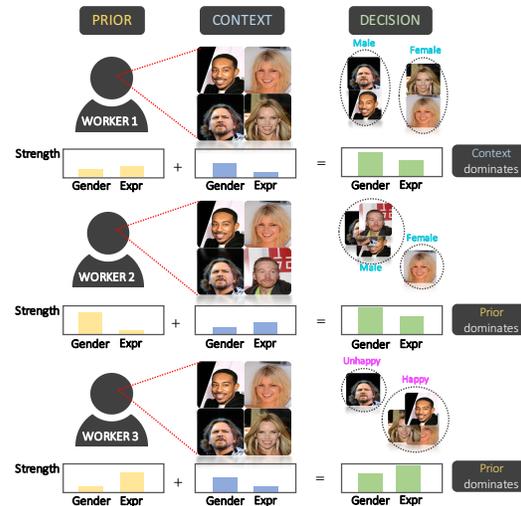


Figure 1. **Context influences similarity estimates.** We hypothesize that estimating similarity according to a particular visual attribute is influenced by a combination of innate biases and the context in which these decisions are made. Compared to worker 1, worker 2 has a strong prior bias towards using the gender attribute. Influenced by the context of the images worker 1 also groups based on gender. Worker 3 sees the same context as worker 1 but ultimately groups based on expression due to prior bias.

workers estimating similarity by clustering a collection of images. The workers’ decision for which visual attribute they use to compare the images can be explained by two factors: 1) The workers have an innate preference towards certain attributes based on their past experiences and 2) the set of related images that a worker observes biases them towards certain attributes. We call this first bias the *worker prior* and the second bias the *context*. Our hypothesis is that different sets of images highlight different visual attributes to the workers. The majority of existing work often assumes that all workers behave in the same way [23], the list of attributes are specified in advance [25], or in addition to similarity estimates, workers also indicate which attributes they used to make their decision [21].

We introduce Context Embedding Networks (CENs), an efficient end-to-end model that learns interpretable, low dimensional, image embeddings that respect the varied similarity estimates provided by different crowd workers. Our contributions are: 1) A flexible model that produces an em-

bedding for a set of input images. This is achieved by modeling worker bias and image context i.e. the degree to which each worker is influenced by the attributes present in a given set of images. 2) An empirical evaluation on annotations from real crowd workers showing that CENs outperform existing approaches, producing interpretable, disentangled, low-dimensional feature spaces.

## 2. Related Work

**Learning Embeddings** The goal of embedding algorithms is to learn a low dimensional representation of a collection of objects (e.g. images), such that objects that are “close” in the potentially high dimensional input space are also “close” in the embedding space. Embeddings are useful for a large number of tasks from face recognition [19] to estimating the clinical similarities between medical patients [34]. They can be learned from pre-defined feature vectors representing the input objects [22], from similarity estimates obtained from the crowd [20, 23], or a combination [33]. Crowdsourced annotations can come in the form of pairwise [7] or relative similarity estimates [20, 26]. Presenting workers with sets of images, as opposed to pairs or triplets, is an efficient way of acquiring estimates of similarity [7, 28, 32]. Another approach is to learn a function that can extract meaningful features from the raw input data by training on similarity labels e.g. [4, 25]. This has the advantage of being able to also embed objects not observed at training time.

**Different Notions of Similarity** A limitation of the above methods is that they typically assume that objects are compared using a single similarity criteria. Given a pair or triplet of images, one estimate of similarity may be valid for one visual attribute, or trait, but invalid for another. For example, in Fig. 1 comparing faces according to gender or expression will result in a different grouping. In practice, workers may use different criteria unless they are specifically told which attribute to use. To overcome this limitation there is a body of work that attempts to learn embeddings where alternative notions of similarity are represented in the embedding space. One common approach is to instruct the workers to provide additional information regarding the attribute they used when making their decision. This information can come in multiple forms such as category labels [25], user provided text descriptions [21], or part and correspondence annotations [16].

Similar to [28], [1] propose a model inspired by [23] that produces a separate embedding for each similarity criteria instead of learning a single embedding that tries to satisfy all constraints. In contrast, [25] learn a unified embedding where alternative notions of similarity are extracted by masking different dimensions in this space. However, the visual attribute used for each similarity estimate is assumed to be known. [29] also learn a weighted feature represen-

tation of the input examples but require category level labels in order to learn cross-category attributes. Their model learns a different weight vector for each triplet, resulting in a large number of parameters. [21] propose a generative model for learning attributes from the crowd where workers are instructed to specify an attribute of interest via a text box and then perform similarity estimates for a set of query images based on these pre-defined attributes. The majority of these methods assume that extra information, in addition to the pairwise or triplet labels, are available to the model. We instead make use of the context information that is present in the set of images that we show to our crowd workers.

**Modeling the Crowd** Crowdsourcing annotations is an effective way of gathering a large amount of labeled data [12]. One difficulty that arises when using such annotations is that they can be noisy, as workers behave differently. One solution to this problem is to model the ability and biases of each worker to resolve better quality annotations [31, 30, 3]. Specific to clustering, [7] propose a Bayesian model of how workers cluster data from noisy pairwise annotations. To efficiently gather a large number of labels, workers are presented with successive grids of images and are asked to cluster the images into multiple different groups. By modeling individual workers as linear classifiers in an embedding space they allow for different worker biases. However, they assume that workers are consistent in the criteria they use when making their decisions and that it does not change over time. Our approach also learns individual worker models while also making use of the strong context information provided by the image grid.

**Attribute Discovery** Low dimensional, attribute based, representations of images have the benefit of being more interpretable than raw pixel information [5, 6]. In addition to providing semantically understandable descriptions of images, they can also be used for applications such as zero shot learning [13]. Attributes can be discovered by various means, from mining noisy web images and their associated text descriptions [2] to crowdsourcing [18]. In this work, while we do not explicitly aim to produce ‘nameable’ attributes, we qualitatively observe that the embeddings that our model produces are often disentangled along the embedding dimensions.

## 3. Methods

We crowdsource the task of image similarity estimation for a dataset containing  $N$  images referenced by  $i, j = 1, \dots, N$ . Each crowd worker  $w = 1, \dots, W$ , is presented with an image grid  $g = 1, \dots, G$ , displaying a collection of images  $\{i_g\}$  which they group into as many categories as they wish [7]. A grid of  $S$  items results in  $(S^2 - S)/2$  pairwise labels, e.g. a single grid of 24 items produces the same number of annotations as 276 individual pairs. Across grids, real workers are often inconsistent with the attributes

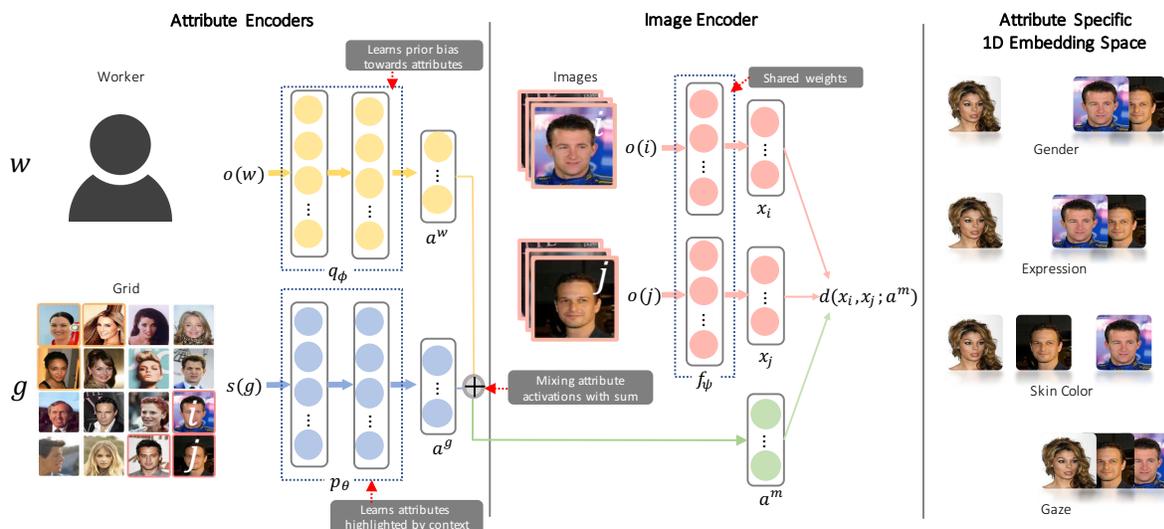


Figure 2. **Context Embedding Networks** are composed of three neural networks that are trained jointly. (Top-Left) A worker encoder network models workers’ annotation behavior and (Bottom-Left) a context encoder network models the attributes highlighted by a particular set of images. Jointly, these networks are referred to as the attribute encoders and are used to weight the embeddings produced by the image encoder network (Center). (Right) Our final embedding respects similarity estimates from each worker in the same low dimensional space where each dimension corresponds to a different visual attribute.

they use to cluster and the number of clusters they create. A pair of images  $(i_g, j_g)$  shown in the same grid,  $g$ , clustered by worker  $w$  is assigned a positive label  $l = 1$  if they are grouped together and  $l = 0$  otherwise. This results in a training set of pairwise similarity labels

$$\mathcal{D} = \{(w, g, i_g, j_g, l) | g = 1, \dots, G\}. \quad (1)$$

### 3.1. Context Embedding Network (CEN)

Here we present our CEN model and define the loss function used to train it. This involves joint training of three networks which model workers, grid context, and image embedding respectively, see Fig. 2. The first two networks are referred to as *attribute encoders* while the third is the *image encoder*.

#### 3.1.1 Worker Encoder

For the workers we define an attribute encoder network  $q_\phi$  which takes as input a one-hot encoding  $(o(\cdot))$  of worker  $w$  and outputs a  $K$  dimensional worker, *attribute activation*, vector  $a^w = q_\phi(o(w)) = [a_1^w, \dots, a_K^w]$ . Each  $a_k^w$ , for  $k = 1, \dots, K$ , represents the degree of prior bias towards attribute  $k$  for worker  $w$ . Once the network is trained, the output attribute activation vector models the worker’s prior preferences for each visual attribute. For example, a heavily biased worker that only attends to a single attribute  $k^*$  should have high activation for that particular attribute dimension  $a_{k^*}^w$ . On the other hand, a worker that does not have a strong preference for any particular attribute will have weak attribute activations in all  $K$  dimensions and may be more influenced by the grid context.

#### 3.1.2 Context Encoder

For an image grid containing  $S$  images, we define a context encoder network  $p_\theta$  that takes as input a  $S$ -hot encoding  $(s(\cdot))$  of the grid  $g$  and outputs a  $K$  dimensional grid *attribute activation* vector  $a^g = p_\theta(s(g)) = [a_1^g, \dots, a_K^g]$ . Each  $a_k^g$  for  $k = 1, \dots, K$  represents the degree of visual prominence of attribute  $k$  for grid  $g$ . Once the network is trained, the grid attribute activation dimensions with high values should correspond to the most salient visual attributes highlighted by the input grid. Intuitively, *attribute variance* in the collection of images should influence which attributes are more noticeable to workers. For instance, a collection of images that is similar along all other attributes except one  $k^*$  should have a peak activation at  $a_{k^*}^g$ . On the other hand, if the image set varies along many different attributes,  $a^g$  should be close to uniformly distributed. The attribute vectors  $a^w$  and  $a^g$  from the worker and context encoders are combined to produce the final attribute encoder output  $a^m$  (Fig. 2 Center).

#### 3.1.3 Image Encoder

We seek to learn a non-linear mapping from image  $i$  to a disentangled Euclidean coordinate  $x_i$  where each dimension embeds the image into a one dimensional attribute specific subspace. To achieve this we use a Siamese Network architecture for the image encoder network  $f_\psi$  with shared parameters  $\psi$  that take as input a one hot encoding of image  $i$  and outputs a  $K$  dimensional embedding vector  $x_i = f_\psi(o(i)) = [x_{i1}, \dots, x_{iK}]$ . Although our image embedding network learns an embedding for each input image

directly, with enough data it is possible to learn a feature extractor from the raw images [25]. Similarly, we present our model in terms of a pairwise loss, but it is also possible to use a triplet loss for the image encoder. For brevity, from this point forward we omit the one and S-hot encoding function notation  $o(\cdot)$ ,  $s(\cdot)$ .

### 3.2. Learning from the Crowd

By ignoring worker and context information, an embedding can be learned using Siamese networks [4], where the contrastive training loss  $L_c$  is defined as

$$L_c(x_i, x_j) = ld(x_i, x_j) + (1 - l) \max\{0, \xi_n - d(x_i, x_j)\}, \quad (2)$$

where  $d(x_i, x_j) = \|x_i - x_j\|_2$  is the  $L2$  distance between image  $i, j$  in embedding space.  $\xi_n$  is the negative margin which prevents over-expanding the embedding manifold, and  $l \in \{0, 1\}$  is the user provided label. This contrastive loss alone does not encourage the network to learn low dimensional attribute specific embeddings as it assumes that all crowd workers compare images using the same visual attributes. To overcome this, we weight the  $L2$  distance metric by the attribute activation vectors  $a^w$  and  $a^g$ . We hypothesize that a worker’s decision to cluster along a particular attribute depends on both their prior preferences for specific visual attributes and the context highlighted by the set of images in the grid. Based on this assumption, we define three variants of the distance metric weighted by the attribute activation vectors

$$\begin{aligned} d(x_i, x_j; a^w) &= \|a^w \cdot (x_i - x_j)\|_2 \\ &= \|q_\phi(w) \cdot (f_\psi(i) - f_\psi(j))\|_2 \end{aligned} \quad (3)$$

$$\begin{aligned} d(x_i, x_j; a^g) &= \|a^g \cdot (x_i - x_j)\|_2 \\ &= \|p_\theta(g) \cdot (f_\psi(i) - f_\psi(j))\|_2 \end{aligned} \quad (4)$$

$$\begin{aligned} d(x_i, x_j; a^m) &= \|a^m \cdot (x_i - x_j)\|_2 \\ &= \|(p_\theta(g) + q_\phi(w)) \cdot (f_\psi(i) - f_\psi(j))\|_2, \end{aligned} \quad (5)$$

where  $a^m = a^g + a^w$  is the mixed attribute activation vector.

After exploring different non-linear methods of mixing  $a^g$  and  $a^w$ , we found that a simple summation sufficiently captures the relationship between the two biases. In the experiments section, we compare the performance of the above three different models. For the model in Eq. 5, biased workers should have a concentrated worker attribute activation vector  $a^w$  which will dominate the mode of sum  $a^m = a^w + a^g$ . Alternatively, workers with weak prior preferences should have low worker attribute activations  $a^w$  and the grid attribute activations  $a^g$  will dictate the mode. Intuitively, the attribute activation vector serves as a mask which

indicates the embedding dimension that should be weighted heavily in the loss e.g. [25]. By encouraging sparsity in  $a^w$  and  $a^g$  along with ReLU non-linearities [17], we assume that grids that were clustered along one attribute will have a uni-modal  $a^m$  while grids that were clustered on a mixture of attributes will have a multi-modal  $a^m$  with peaks corresponding to the attribute dimensions used.

Inspired by the dual margin contrastive loss proposed in [29], we include a positive margin term  $\xi_p$  in the loss function to prevent two images from overlapping in the embedding space which could lead to over fitting. This ensures that images will be pushed closer only if their current embedding is separated by more than  $\xi_p$ . We use  $a$  to denote the general attribute activation vector which can be  $a^g$ ,  $a^w$ , or  $a^m$  depending on the model variant

$$\begin{aligned} L_c(x_i, x_j; a) &= l \max\{0, d(x_i, x_j; a) - \xi_p\} + \\ &\quad (1 - l) \max\{0, \xi_n - d(x_i, x_j; a)\}. \end{aligned} \quad (6)$$

A crowd worker’s decision to group two images is an active decision while choosing not to group images together can be seen as a more passive decision. This can become a problem when workers group images with different levels of detail. For example, a grid of shapes containing squares, triangles, circles, and stars might be clustered into two groups, squares and non-squares, by one worker. A second worker may group the images into the four different shape types. An embedding model might incorrectly assume that a different attribute was used to separate the images, when it is in fact just a different level of granularity of ‘shape’ that is being used by both workers. To overcome this problem, we introduce an additional positive similarity weight  $\gamma$ , that captures the relative importance of the positive similarity labels compared to the dissimilarity labels

$$\begin{aligned} L_c(x_i, x_j; a) &= \gamma l \max\{0, d(x_i, x_j; a) - \xi_p\} + \\ &\quad (1 - l) \max\{0, \xi_n - d(x_i, x_j; a)\}. \end{aligned} \quad (7)$$

This ensures that the model can learn the high level attributes when workers cluster with different levels of detail. In the example above, although cross category labels between circles, triangles, and stars are  $l = 0$ , the positive labels generated within each circle, triangle, and star groups agree with the positive labels generated within the non-square group thus allowing the network to learn that the high level attribute, i.e. shape, used by both workers are the same. We show the impact of  $\gamma$  on the performance of our CEN in the supplementary materials.

### 3.3. Regularization

We add  $L1$  penalties  $\lambda_1 \|a\|_1$  to the attribute encoders to encourage sparsity in the attribute activation vector. We also regularize the embedding network with a  $L2$  penalty

$\lambda_2 \|x\|_2$  to encourage regularity in the latent space. The final loss function for our CENs is

$$L_{CEN}(x_i, x_j; a) = \gamma l \max\{0, d(x_i, x_j; a) - \xi_p\} + (1 - l) \max\{0, \xi_n - d(x_i, x_j)\} + \lambda_1 \|a\|_1 + \lambda_2 \|x_i\|_2 + \lambda_2 \|x_j\|_2. \quad (8)$$

CENs require the number of dimensions  $K$  as a hyperparameter. However, we observe that by setting  $K$  to a large number and by  $L1$  regularizing  $a^w$  and  $a^g$ , our model tends to only use a subset of the available embedding dimensions.

## 4. Experiments

Here we show that CENs can recover meaningful low-dimensional embeddings from noisy data. Network architectures, training details, and hyperparameters tuning are described in the supplementary material. We perform experiments on the following three datasets:

**CELEBA** contains images of different celebrity faces from which we select a random subset of 300 images [15]. For this dataset we instruct workers in advance to cluster on one attribute per grid respecting four visual attributes: gender, expression, skin color, and gaze direction. Although we expect some workers to deviate from our instructions, having a definite ground truth set of attributes allow us to quantify the attribute retrieval accuracy. The CEN is unaware of the attribute selected for each grid. In total, 94 workers clustered 620 grids, yielding 170,000 similarity training pairs.

**RETINA** is a medical dataset comprising of fundus images of the retina belonging to patients with varying degrees of diabetic retinopathy [9]. The images contain a number of visual indicators for the disease such as hard exudates (yellow lesions dispersed throughout the retina). From 66 fundus images we crop out 300 image patches. These patches provide a localized view that may or may not contain indicator features of the disease. This dataset is more challenging to discover meaningful attributes as the disease indicator features are visually subtle and the images are unfamiliar to the crowd. We do not provide any instructions as to the attributes the workers should use for this dataset. 62 workers clustered 620 grids, yielding 170,000 similarity pairs.

**BIRDS** is a larger dataset composed of 1000 bird head images made up of 16 randomly selected species from [27]. We use this dataset to demonstrate the scalability of our CENs. 252 workers clustered 3,000 grids yielding 820,000 similarity labels.

### 4.1. Data Collection

We use Amazon Mechanical Turk’s crowdsourcing platform to request crowd workers to cluster grids of images using the GUI shown in Fig. 3. Workers were presented with a  $4 \times 6$  grid of images randomly sampled from the given dataset. Using up to ten possible groups, workers clustered



Figure 3. **Data collection GUI.** Workers group images they perceive to be visually similar by assigning them to different groups. They can create up to ten groups per grid of images.

images by first clicking on a group button on the right side of the page then clicking on the desired images. For each group they were asked to provide a short text description, used only for evaluation. The image, cluster, and worker ids were then converted into pairwise similarity labels (Eq. 1). Each worker clustered a minimum of ten grids in order to receive a reward, ensuring that the worker encoder network had sufficient data to learn from.

### 4.2. Baseline Comparisons

We compare results to four baseline methods and three variants of our model:

**Standard Siamese Network e.g. [4]:** Assumes that all pairwise similarity labels come from the same notion of similarity, as in Eq. 2.

**Standard Triplet Network e.g. [19]:** Learns embeddings given similarity labels of the form "A is more similar to B than C".

**Bayesian Crowd Clustering [7]:** Workers are modeled as linear classifiers in the embedding space where both an *entangled* image embedding and individual worker models are jointly learned with variational methods.

**CSN [25]:** Learns an entangled image embedding from similarity triplets which are disentangled by masks learned separately for each pre-known attribute. This baseline represents the situation where the similarity dimension used by the worker is *known*.

**CEN-worker encoder only:** This first variant of our model uses only worker modeling to learn attribute activations which weight the embeddings as in Eq. 3.

**CEN-context encoder only:** Here we only model context information to weight the embeddings as in Eq. 4.

**CEN-mixture:** Our full model, incorporates both worker and context information to learn a network that weights the worker bias  $a^w$  and grid context  $a^g$  as in Eq. 5.

### 4.3. Unsupervised Attribute Retrieval

First, we evaluate whether our CEN can accurately recover the four dominant attributes present in the CELEBA dataset. For each grid  $g$  clustered by worker  $w$ , we take the mode dimension of the attribute activation vector  $a$  to

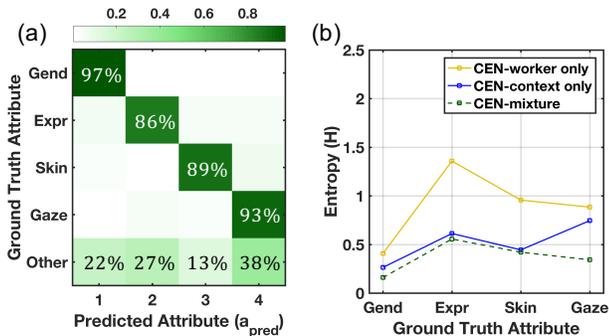


Figure 4. **Attribute retrieval accuracy.** On the left we see the predicted embedding dimensions from the CEN-mixture model compared to the ground truth visual attributes for the CELEBA dataset. On the right, we quantify how disentangled the learned embeddings are. Lower entropy indicates models that better capture the ground truth attributes along individual embedding dimensions.

be the model’s prediction,  $a_{pred} = \text{argmax}_k a_k$ . This is the attribute that we predict was used to cluster the set of images. Again  $a$  can be  $a^w$ ,  $a^g$  or  $a^m$  depending on the model variant used. We then examine the annotations provided by workers for each set of grids that map to a different  $a_{pred} \in \{1, 2, 3, 4\}$  and quantify the proportion of each attribute actually used. In Fig. 4(a) we show a confusion matrix illustrating that for each worker and grid pair, the CEN-mixture model is able to accurately predict which attribute was used. The row for *gender* and the first column denote the proportion of grid submissions that have  $a_{pred} = 1$  out of all the submissions that were clustered along gender. For all attributes we obtain over 85% attribute prediction accuracy. In Fig. 4(b) we plot the entropy  $H$  of the distribution  $p$  for each row of the confusion matrix where  $H_p = -\sum p \log p$ . High entropy indicates that the ground truth attributes are scattered throughout the attribute predictions and vice versa. The CEN-mixture model learns the most disentangled embeddings across the four ground truth attributes compared to its variants.

Although workers were encouraged to focus on four different attribute options for this experiment, in practice they did not abide by our instructions and the proportion of noise in the raw data is significant. For the CELEBA dataset approximately 19.1% of the HITs completed were either clustered on different attributes such as “wearing sun glasses” (see Fig. 5) or noisy submissions where images were not separated into different groups. We also observed workers using different levels of detail when clustering on the same attribute. For example, for the *gaze* attribute some workers labeled “looking left”, “looking right”, etc. To demonstrate our model’s robustness, we perform all of our experiments on this raw data without filtering out annotation noise. For evaluation of the worker model learned by the worker encoder, refer to the supplementary material.

<b>Gender</b> Male, Female Men, Women Guys, Girls All men All females	<b>Expression</b> Smiling, Not smiling Smiling, Frowning, Neutral Content, Undecided Thoughtful, Fearful Happy, Bored	<b>Other</b> Hat, No Hat Earing, No-Earing Hair up, Hair down Attractive, Ugly Humble, Arrogant I don't know Celebrity Artist Musician Beard, no beard Wearing tie Black hair, Blond
<b>Skin Color</b> Indian, African, Asian White, non-white Black, White Dark, Light, Tan White, Black Brown	<b>Gaze</b> Looking, Not-looking Front, Left, Right Looking Straight, Crooked Side Pose, Front Pose Facing, Not-facing	

Figure 5. **Cluster Names.** Keywords provided by workers for CELEBA. Colored labels indicate the manual grouping performed by us (only used for evaluation). Some workers use finer grained distinctions compared to others.

#### 4.4. Visualizing Disentangled Attributes

Fig. 6 shows the attribute specific embeddings of the four subspaces learned by the CEN for the CELEBA dataset. Fig. 6(a) shows that the embedding clearly separates the images according to *gender*. On the very left of the *expression* subspace (Fig. 6(b)) we can see that people are smiling with teeth showing while on the right they show serious or unhappy expressions. In the middle we see ambiguous expressions. Fig. 6(c) shows the subspace embedded along the *skin color* attribute. On the two ends we see darker skinned and lighter skinned people. Fig. 6(d) shows the subspace for *gaze direction* of people, showing people that are either looking at the camera or away from it. Again, in the middle we see people wearing sunglasses or looking in ambiguous directions making it difficult to assess their gaze direction.

In Fig. 7 we show attribute specific embeddings learned for the RETINA dataset in which no supervision was given to the workers for which attributes to pay attention to. Here we select four dimensions that are most highly activated from the learned ten dimensional embedding vector. Other attribute dimensions attain trivial activations. This shows that our CEN is robust to value of  $K$  (please see supplementary material for robustness analysis of  $K$ ). Fig. 7(a) shows the first dimension seemingly showing the presence or absence of the optic disc, a key feature of the retina. Fig. 7(b) shows the subspace which discriminates between patches with blood vessels present and those without. Blood vessels are mostly concentrated and visually prominent around the optic disc, meaning that the two attributes are highly correlated. Regardless, our CEN is capable of distinguishing between the two attributes, as we see that images displaying blood vessels without optic discs are correctly embedded in Fig. 7(a). Fig. 7(c) plots the attribute that groups laser scars (named after consulting with an ophthalmologist) and Fig. 7(d) groups hard exudates, a key indicator for diagnosing diabetic retinopathy [11]. A comparison of embedding qualities between baselines are presented in the supplementary material.

Fig. 8 shows a t-SNE plot of the four dimensional embedding space learned by the CEN for the BIRDS dataset.

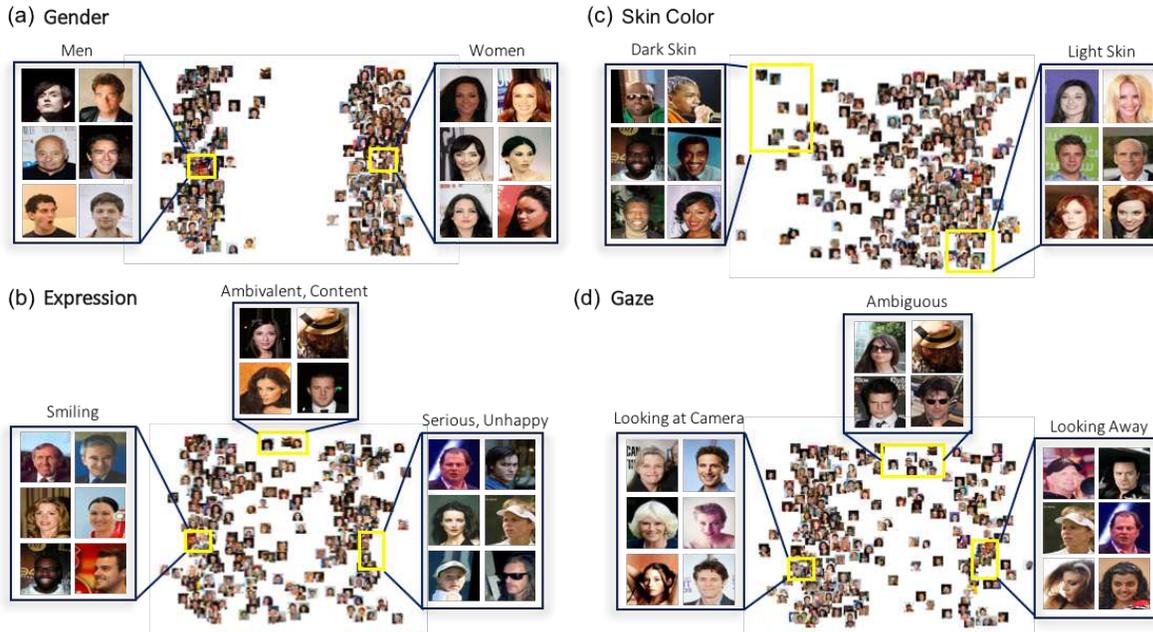


Figure 6. **CELEBA - Attribute specific embeddings.** Each plot shows one of the four different embedding dimensions produced by the CEN-mixture mode. The vertical axis in each subplot is randomly assigned for visualization purposes. We show representative images from the embeddings space in yellow boxes. We can see that the CEN learns to disentangle the attributes.

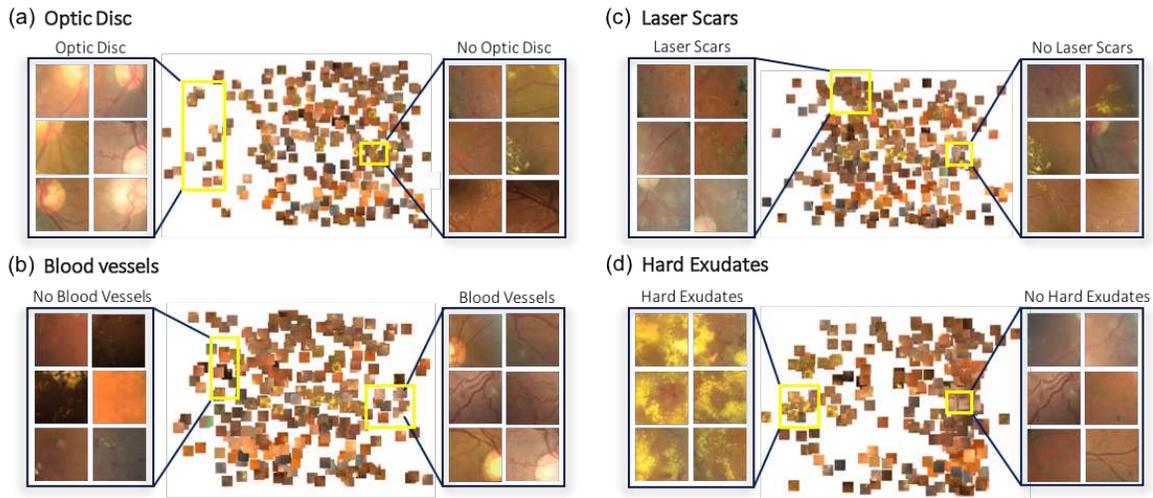


Figure 7. **RETINA - Attribute specific embeddings.** Here we show a subset of four of the ten embedding dimensions produced by the CEN-mixture model for the RETINA dataset. Dimensions correlated well with visual features of diabetic retinopathy.

Each ellipse center corresponds to the mean of a Gaussian distribution fit to the embedding coordinates for each ground truth species. We observe 16 compact clusters that directly correlate to the 16 ground truth species. Please refer to the supplementary materials for confusion plots of the ground truth species vs embedding clusters.

#### 4.5. Performance on Held-out Label Prediction

Here we quantify the generalization performance of the baseline methods on held-out pairwise label predictions while varying the amount of training data. We measure

the accuracy of the various model’s predictions on the similarity estimates for an unseen grid clustered by a known worker. For a grid input  $g$ , worker input  $w$ , and image pair  $i, j$ , the model predicts  $i$  and  $j$  to be in the same group if  $d < (\xi_n + \xi_p)/2$ . The test set is made up of 15% of the dataset and consists only of entire grids that were not present in the training set. This allows us to measure how well our CEN generalizes to new sets of images.

Fig. 9(a) shows results for the RETINA dataset. Standard Siamese Networks and Triplet Networks fail to capture the multiple attributes used to cluster the images and

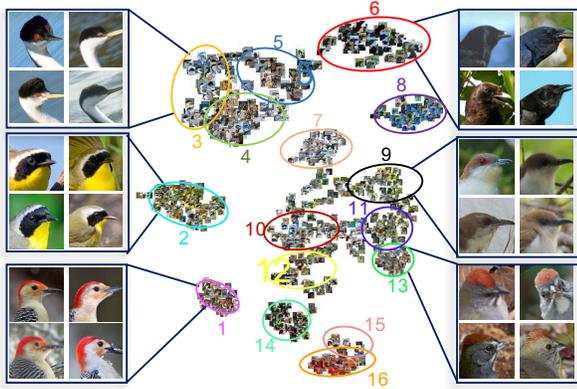


Figure 8. **BIRDS - t-SNE embedding.** Here we show a t-SNE [22] plot of the four dimensional embedding produced by the full CEN model for the BIRDS dataset. Indexed ellipses are centered at the Gaussian mean of different ground truth species. Clusters correlated well with ground truth species of birds.

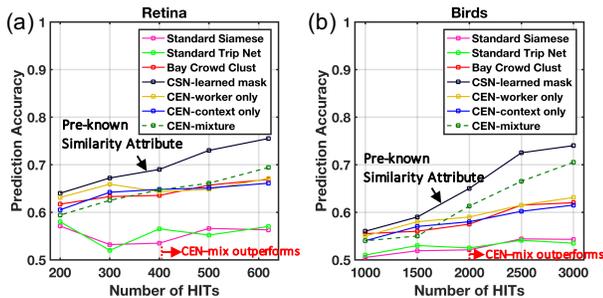


Figure 9. **Held-out label prediction.** Prediction accuracy on held out labels for the RETINA and BIRDS datasets plotted against the amount of available data during training.

have the lowest prediction accuracy of 58.1% and 58.5%. The Bayesian Crowd Clustering model, CEN worker, and CEN grid only models attain similar prediction accuracies of 67%. For the more challenging RETINA dataset workers found it difficult to discover various attributes to cluster on and thus often fixated on a single attribute on all their HITs. However, we still benefit from modeling the context as the CEN-mixture model achieves prediction accuracy of 69.4%. The CSN model with learned masks obtains the highest accuracy of 75.5%, but it is important to note that this model was trained on triplets pre-labeled with the true similarity attributes used to cluster them.

Fig. 9(b) shows the pairwise prediction accuracy for each model plotted against a varying number of training samples for the BIRDS dataset. The Bayesian Crowd Clustering model, CEN worker, and CEN grid only models attain similar prediction accuracies of 62%. The CEN-mixture substantially outperforms all baselines with a prediction accuracy of 70.5% which is only 3.5% below the accuracy of the CSN model which uses ground truth labels.

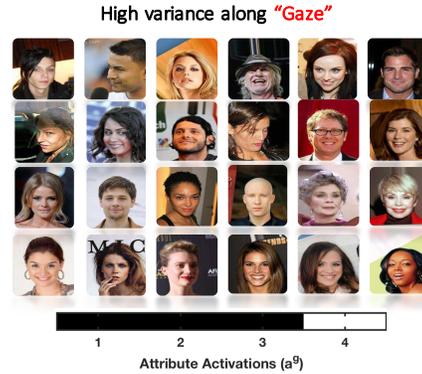


Figure 10. **Synthesized image grids.** Our context encoder can be used to generate collections of images that highlight specific attributes. The shown grid has high variance along the gaze direction attribute and low variance for the others.

#### 4.6. Image Grid Synthesis

Being able to synthesize image grids that highlight specific attributes may be useful in active learning where the data collector seeks to obtain similarity estimates along particular visual attributes. We randomly generate ten million image grids and individually pass them through the context encoder and extract the grid attribute activation vectors  $a^g$  for each grid. We take a softmax activation over the  $a^g$ s and select grids that have low entropy, thus choosing grids that are highly expressive for a particular attribute. Fig. 10 shows a generated grid with the lowest entropy for the gaze attribute. We see low variance among the images along other attributes such as gender and skin color, while there is high variance for ‘gaze’. This suggests that in order for a grid to emphasize a particular attribute, the contained items should be similar in all but one high variance attribute.

### 5. Conclusion

We proposed a novel deep neural network that jointly learns attribute specific embeddings, worker models, and grid context models from the crowd. By comparing to several baseline methods, we show that our model more accurately predicts the attributes used by individual workers and as a result produces better quality image embeddings.

In future we plan to incorporate relative similarity estimates and the learning of representations directly from images [25, 19]. Although currently we model each worker individually, in practice there may be similarity between different workers that could be discovered through clustering [10]. Finally, our grid context encoder enables us to generate sets of images that highlight specific attributes. By combining this with active learning we can potentially speed up the collection of annotations from the crowd [20, 14].

**Acknowledgements** We thank Google for supporting the Visipedia project and AWS Research Credits for their donation.

## References

- [1] E. Amid and A. Ukkonen. Multiview triplet embedding: Learning attributes in multiple maps. In *ICML*, 2015.
- [2] T. L. Berg, A. C. Berg, and J. Shih. Automatic attribute discovery and characterization from noisy web data. In *ECCV*, 2010.
- [3] S. Branson, G. Van Horn, and P. Perona. Lean crowdsourcing: Combining humans and machines in an online system. In *CVPR*, 2017.
- [4] J. Bromley, J. W. Bentz, L. Bottou, I. Guyon, Y. Lecun, C. Moore, E. Säckinger, and R. Shah. Signature verification using a "siamese" time delay neural network. *IJPRAI*, 1993.
- [5] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth. Describing objects by their attributes. In *CVPR*, 2009.
- [6] R. S. Feris, C. Lampert, and D. Parikh. *Visual Attributes*. 2017.
- [7] R. G. Gomes, P. Welinder, A. Krause, and P. Perona. Crowdclustering. In *NIPS*, 2011.
- [8] G. Jurman and C. Furlanello. A unifying view for performance measures in multi-class prediction. *arXiv:1008.2908*, 2010.
- [9] Kaggle. Kaggle Diabetic Retinopathy Detection. <https://www.kaggle.com/c/diabetic-retinopathy-detection>, 2015.
- [10] H. Kajino, Y. Tsuboi, and H. Kashima. Clustering crowds. In *AAAI*, 2013.
- [11] R. V. J. P. H. Kälviäinen and H. Uusitalo. Diaretdb1 diabetic retinopathy database and evaluation protocol. In *MIUA*, 2007.
- [12] A. Kovashka, O. Russakovsky, L. Fei-Fei, K. Grauman, et al. Crowdsourcing in computer vision. *Foundations and Trends in Computer Graphics and Vision*, 2016.
- [13] C. H. Lampert, H. Nickisch, and S. Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *CVPR*, 2009.
- [14] L. Liang and K. Grauman. Beyond comparing image pairs: Setwise active learning for relative attributes. In *CVPR*, 2014.
- [15] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *ICCV*, 2015.
- [16] S. Maji and G. Shakhnarovich. Part and attribute discovery from relative annotations. *IJCV*, 2014.
- [17] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.
- [18] D. Parikh and K. Grauman. Relative attributes. In *ICCV*, 2011.
- [19] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015.
- [20] O. Tamuz, C. Liu, S. Belongie, O. Shamir, and A. T. Kalai. Adaptively learning the crowd kernel. *ICML*, 2011.
- [21] T. Tian, N. Chen, and J. Zhu. Learning attributes from the crowdsourced relative labels. In *AAAI*, 2017.
- [22] L. van der Maaten and G. Hinton. Visualizing data using t-sne. *JMLR*, 2008.
- [23] L. van der Maaten and K. Weinberger. Stochastic triplet embedding. In *Machine Learning for Signal Processing*, 2012.
- [24] G. Van Horn, S. Branson, R. Farrell, S. Haber, J. Barry, P. Ipeirotis, P. Perona, and S. Belongie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *CVPR*, 2015.
- [25] A. Veit, S. Belongie, and T. Karalestos. Conditional similarity networks. In *CVPR*, 2017.
- [26] R. K. Vinayak and B. Hassibi. Crowdsourced clustering: Querying edges vs triangles. In *NIPS*, 2016.
- [27] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD birds-200-2011 dataset. 2011.
- [28] C. Wah, G. Van Horn, S. Branson, S. Maji, P. Perona, and S. Belongie. Similarity comparisons for interactive fine-grained categorization. In *CVPR*, 2014.
- [29] X. Wang, K. M. Kitani, and M. Hebert. Contextual visual similarity. *arXiv:1612.02534*, 2016.
- [30] P. Welinder, S. Branson, S. J. Belongie, and P. Perona. The multidimensional wisdom of crowds. In *NIPS*, 2010.
- [31] J. Whitehill, T.-f. Wu, J. Bergsma, J. R. Movellan, and P. L. Ruvolo. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In *NIPS*, 2009.
- [32] M. Wilber, I. S. Kwak, and S. Belongie. Cost-effective hits for relative similarity comparisons. In *HCOMP*, 2014.
- [33] M. Wilber, I. S. Kwak, D. Kriegman, and S. Belongie. Learning concept embeddings with combined human-machine expertise. In *ICCV*, 2015.
- [34] Z. Zhu, C. Yin, B. Qian, Y. Cheng, J. Wei, and F. Wang. Measuring patient similarities via a deep architecture with medical concept embedding. In *ICDM*, 2016.