

---

# Improving Calibration in Deep Metric Learning With Cross-Example Softmax

---

Andreas Veit\* Kimberly Wilber\*  
Google Research  
{aveit, kwilber}@google.com

## Abstract

Modern image retrieval systems increasingly rely on the use of deep neural networks to learn embedding spaces in which distance encodes the relevance between a given query and image. In this setting, existing approaches tend to emphasize one of two properties. Triplet-based methods capture *top-k relevancy*, where all top- $k$  scoring documents are assumed to be relevant to a given query. Pairwise contrastive models capture *threshold relevancy*, where all documents scoring higher than some threshold are assumed to be relevant. In this paper, we propose *Cross-Example Softmax* which combines the properties of top- $k$  and threshold relevancy. In each iteration, the proposed loss encourages all queries to be closer to their matching images than all queries are to all non-matching images. This leads to a globally more calibrated similarity metric and makes distance more interpretable as an absolute measure of relevance. We further introduce *Cross-Example Negative Mining*, in which each pair is compared to the hardest negative comparisons across the entire batch. Empirically, we show in a series of experiments on Conceptual Captions and Flickr30k, that the proposed method effectively improves global calibration and also retrieval performance.

## 1 Introduction

The goal of large-scale information retrieval is to efficiently find relevant documents for a given query among potentially billions of candidates. A canonical example is image search: finding relevant images given a text query. One common implementation is to learn a real-valued scoring function to rank all images. Since using large neural networks to compute the relevance of each query-image pair is prohibitively expensive, recent deep learning systems encode semantic relevance as distance in a vector space. These systems can model complex semantic relationships while still allowing for efficient retrieval using approximate nearest neighbor search through a large database of images.

In this embedding based retrieval setting, two approaches are commonly used to determine the relevance of retrieved documents. In *top-k relevancy*, all top- $k$  scoring documents are assumed to be relevant to a given query, while in *threshold relevancy*, all documents scoring higher than some threshold are assumed to be relevant. Both properties are commonly used; *top-k relevancy* is useful when searching for a document known to exist, *threshold relevancy* is useful when it is not known whether any relevant document exists. Although both properties are desirable, common learning approaches tend to explicitly optimize for only one of these properties. To address this issue, we propose a novel learning approach that simultaneously optimizes both properties.

Early approaches for embedding learning used pairwise *contrastive losses* [7], where the distance between matching query/document pairs is minimized and the distance between non-matching pairs is maximized until they reach a specified margin. Using explicit definitions for what constitutes matching and non-matching pairs, these approaches have been an effective means to achieve threshold

---

\*These two authors contributed equally

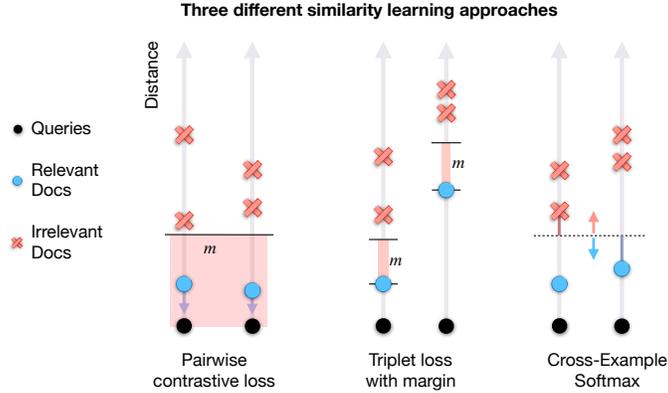


Figure 1: Left: Pairwise contrastive loss minimizes the distance of matching query/document pairs while pushing nonmatching pairs up to some pre-selected margin. Center: Triplet loss constrains non-matching documents to be at least some margin further away than matching documents for a given reference query. Right: The proposed Cross-Example Softmax compares distances across queries, encouraging all matching query/document pairs to be closer than all non-matching pairs.

relevancy. However, the binary nature of these pairwise approaches limits their ability to model the relative ordering among the matching pairs, limiting their effectiveness for top- $k$  relevancy.

To address this challenge, relative comparison approaches such as triplet loss [28] have been proposed that explicitly model the relative ordering of pairwise query/document distances with respect to a single reference query. While effectively modeling top- $k$  relevancy, the query-specific conditioning of pairwise distances implies that absolute distances are not necessarily meaningful across different queries, limiting threshold relevancy. Recent approaches commonly go beyond individual triplets with the usage of *Sampled Softmax* [30]. There, a batch of matching query/document pairs are embedded into their vector representations. Then, for each query, the documents from the other pairs in the batch can be re-used as negative comparisons. Since most of these random documents are unlikely to be informative to each query, it is also common to use stochastic negative mining [24] to only use the most informative negative documents, i.e., with the highest similarity score.

In this paper, we propose *Cross-Example Softmax*, an extension to *Sampled Softmax* which combines the properties of top- $k$  relevancy and threshold relevancy by introducing cross-example comparisons. Specifically, in addition to ranking a positive query/document pair over all other within-batch pairs containing the same query, *Cross-Example Softmax* also ranks each positive pair over all non-matching pairs containing different queries. Figure 1 illustrates the effect of *Cross-Example Softmax* and how it leads to calibrated distance scores. The proposed method further allows us to extend the concept of stochastic negative mining to *Cross-Example Negative Mining*. There, instead of only mining the most informative negative documents for the given query, we select those non-matching pairs with the highest similarity scores across all queries.

In experiments on Conceptual Captions and Flickr30k, we demonstrate that *Cross-Example Softmax* effectively learns embedding spaces that are more globally calibrated and that have higher retrieval performance. Quantitatively, in terms of area under the precision-recall curve, this results in an improvement from 14.6 to 20.1 on the Conceptual Captions test set. In terms of Recall@1, we observe improvements from 25.87% to 26.91% on Conceptual Captions and from 29.22% to 30.49% on Flickr30k.

Overall, we make the following **main contributions**:

- With *Cross-Example Softmax*, we propose a novel loss function for embedding learning which combines the properties of top- $k$  and threshold relevancy.
- We further introduce *Cross-Example Negative Mining*, where the per-example loss is focused on the hardest incorrect query/document pairs across the entire batch.
- In experiments on Conceptual Captions and Flickr30k, we show that the proposed methods effectively improve global score calibration and retrieval performance.

## 2 Related Work

**Vision-Language Models.** Approaches that combine the visual and language modalities have been of great interest in recent years, addressing tasks such as image captioning [6, 35, 39], visual question answering [2], visually-grounded dialog [8] and more. One common characteristic of many of these approaches is to combine features from both modalities early in the model [13]. Other work trains models that merge visual and language features in Transformer-based [33] architectures [1, 16, 17, 21, 31, 32]. These approaches are powerful because the model can consider nonlinear interactions between vision and language features. However, they are not suited to our large-scale image retrieval task, because inferring a pairwise similarity score between a query and millions of documents becomes intractable. Our work instead focuses on deep metric learning using a factorized model with two separate image and text encoders, limiting the interaction between modalities to an inner product.

**Deep Metric Learning.** Metric learning using deep models is a well-studied problem with many applications [3, 25, 28, 34], especially where the output space is very large. Early approaches are based upon Siamese networks [7] with contrastive loss on pairwise data or relative triplet similarity comparisons [12, 28]. Inspired by the success of large-scale classification tasks on ImageNet [9], more recent models are mainly trained using sampled softmax loss with cross entropy [4, 14]. Recently, several works have proposed modifications to the sampled softmax loss, by normalization [19], adding margins [18, 36, 37], and tuning the scaling temperature [42, 43]. Differing from our work, all of these works focus on optimizing the relative ordering of labels given an anchor query. In contrast, the method proposed in our work optimizes for *each* input query the score of the correct label against the *entire distribution of all possible negative query/label pairs in the entire batch*, even across queries.

**Stochastic Negative Mining.** In the setting of large output spaces, for any given query, most documents are not relevant and thus including them in the loss function is not informative for the optimization. To address this challenge, several works have proposed to mine for the hardest and most informative negative labels [22, 24]. However, as an approximation of per-query softmax, these methods perform negative mining only with respect to one single query at a time. In our work, we propose negative example mining across examples, wherein we mine for the globally hardest negative query/document comparison in the batch.

**Score Calibration.** There has been a sustained interest in score calibration to ensure scores are consistently normalized or interpretable [23, 26, 27]. One common approach is to interpret the output of the softmax function applied to model logits as probabilities. While the output of a softmax is technically a probability distribution in the sense that it is normalized, computing the probability for any label also requires the comparison to all other labels. In the setting of large output spaces, this is generally not possible, because the probability space is too large to calculate. To address this challenge, in this paper we propose a new loss function that explicitly encourages the underlying logits to be calibrated. This is done during training, not as a post-recognition calibration step as in [23, 27]. This allows the comparison of label scores across queries without needing to compute scores for all other labels.

## 3 Method

Consider the multiclass classification setting with a sample of instances  $x \in \mathcal{X}$  and their associated labels  $y \in \mathcal{Y}$  with  $|\mathcal{Y}| = K$ . There, the goal is to learn a scoring function  $f : \mathcal{X} \rightarrow \mathbb{R}^K$  that is able to score the labels for each instance according to their relevance. The information retrieval setting at hand can be defined analogous with a set of queries  $\mathcal{X}$  and associated relevant documents  $\mathcal{Y}$ . Our goal is to learn a scoring function which can sort all documents according to their relevance for a given query.

In our text-to-image retrieval setting,  $x_i$  is a text query and  $y_i$  is its corresponding relevant image. We aim to learn a text encoder  $f_{text} : (x_i) \rightarrow \mathbf{x}_i \in \mathbb{R}^d$  and an image encoder  $f_{image} : (y_i) \rightarrow \mathbf{y}_i \in \mathbb{R}^d$  that project the text and image into a shared  $d$ -dimensional embedding space. The relevance score between query  $x_i$  and image  $y_j$  is the dot product between their vector representations  $s_{i,j} = \langle \mathbf{x}_i, \mathbf{y}_j \rangle$ .

**Sampled Softmax.** In the standard multiclass classification setting, a Softmax Cross-Entropy loss over the whole label space is used to optimize the model. Ideally, in the retrieval setting we would also compare the score of a matching document to all other documents in the database. However,

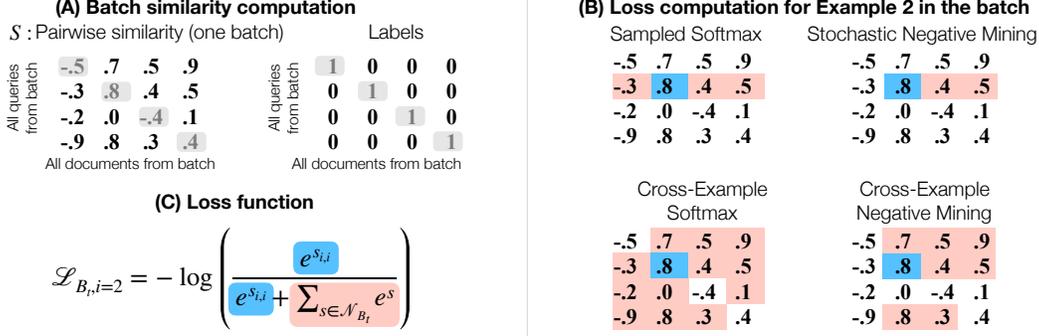


Figure 2: **Cross-Example Softmax** operates on a pairwise distance matrix (A) computed from a batch of paired queries and documents. Comparisons to other methods are shown in (B). In Sampled Softmax, for a given query the similarity score for the positive document is compared to all other documents within the batch. Stochastic Negative Mining focuses the loss on the hardest document for each query. For Cross-Example Softmax, each positive query/document pair is compared to all non-matching pairs across the entire batch. Cross-Example Negative Mining focuses the loss on the hardest non-matching query/document pairs across the entire batch. In all of these approaches, negatives appear in the partition function of the softmax activation (C).

since the number of documents may be in the billions, the Softmax Cross-Entropy loss is commonly only computed over a random subset of labels (Sampled Softmax).

Specifically, consider a mini-batch of  $N$  corresponding query/document pairs  $B_t = \{(x_1, y_1), \dots, (x_N, y_N)\}$ . Given the vector representations of all text queries and images from the mini-batch, we can compute the pairwise similarity matrix between all possible pairs  $S \in \mathbb{R}^{N \times N} = s_{i,j}, \forall i, j \in [1, \dots, N]$ . An example of such a matrix for  $N = 4$  is illustrated on the left in Fig. 2 (A). It is commonly assumed that for a given query  $x_i$ , within the batch only its corresponding document  $y_i$  is relevant. All other documents  $y_j, j \neq i$  within the same batch are assumed to be irrelevant to that query. The matching relationship between queries and documents within a batch is illustrated on the right in Fig. 2 (A) with 1 indicating a matching relationship and 0 indicating a non-match. Formally, let  $\mathcal{N}_{i,B_t}$  be the set of similarity scores between query  $x_i$  and all non-matching documents in the batch, i.e., except for the query’s corresponding document  $y_i$ . In Fig. 2 this would be all scores in row  $i$  of  $S$  except the relevant document.

$$\mathcal{N}_{i,B_t} = \{s_{i,j} : j \neq i\} \quad (1)$$

With this we can define *Sampled Softmax* Cross-Entropy as a relative ranking loss between the relevance score of a query and its matching document  $s_{i,i}$  and its relevance scores to all non-matching documents  $s \in \mathcal{N}_{i,B_t}$ . Formally,

$$\mathcal{L}_{B_t} = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{e^{s_{i,i}}}{e^{s_{i,i}} + \sum_{s \in \mathcal{N}_{i,B_t}} e^s} \right) \quad (2)$$

Fig. 2 (B) in the top left illustrates the per-example loss for the second query. The score of the matching pair  $s_{2,2}$  is emphasized in blue and  $\mathcal{N}_{2,B_t}$  is marked in red. The figure highlights that the loss only considers pairs from the same query.

**Stochastic Negative Mining.** For Sampled Softmax a relevant document is only compared to a small subset of random documents. Most of these documents can be easily recognized as irrelevant to the query, limiting their informativeness to the optimization.

As a consequence, Stochastic Negative Mining has been proposed [24], wherein one only selects the most difficult negative documents for each query within the batch. Formally, let  $\text{top}k(\mathcal{N}_{i,B_t})$  be the set of the top  $k$  largest scores within the set of negative scores for query  $x_i$ . With this, we can write the modified loss defined only over the most difficult documents for each query as

$$\mathcal{L}_{B_t} = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{e^{s_{i,i}}}{e^{s_{i,i}} + \sum_{s \in \text{top}k(\mathcal{N}_{i,B_t})} e^s} \right) \quad (3)$$

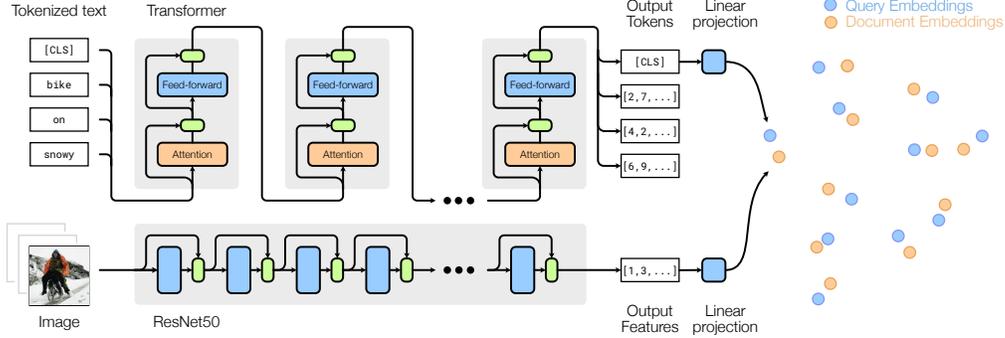


Figure 3: **The Text-to-Image model** comprises two encoders. Text captions are encoded using a 12-layer Transformer. The output of the CLS token is projected into the embedding space and L2 normalized. Images are encoded using a 50-layer Residual Network followed by a linear projection and L2 normalization. The relevance is then measured via dot product in the embedding space.

Fig. 2 (B) top right shows this scenario. The set of negatives now only comprise the hardest comparisons for the given query. In the diagram, the red shaded negative scores are top  $k(\mathcal{N}_{i,B_t})$ .

**Cross-Example Softmax.** From Equations 2 and 3 and the illustrations in Fig. 2 (B) it becomes clear that Sampled Softmax captures the distance of documents only relative with respect to a given query. Thus, distances in the learned vector space are not comparable across queries. To encourage global calibration such that distance can be used as an absolute measure of relevance, we propose *Cross-Example Softmax* which extends Softmax by introducing cross-example negatives. The proposed loss encourages that all queries are closer to their matching documents than all queries are to all irrelevant documents. Specifically, let  $\mathcal{N}_{B_t}$  be the pairwise comparisons between all queries in batch  $B_t$  and the documents of the same batch which they are not related to. Since queries are assumed to be only related to their respective document, this corresponds to all off-diagonal entries in  $S$ . Formally,

$$\mathcal{N}_{B_t} = \bigcup_{i \in [1, \dots, N]} \mathcal{N}_{i, B_t} \quad (4)$$

With this we can formally define Cross-Example Softmax Cross-Entropy as

$$\mathcal{L}_{B_t} = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{e^{s_{i,i}}}{e^{s_{i,i}} + \sum_{s \in \mathcal{N}_{B_t}} e^s} \right) \quad (5)$$

Fig. 2 (B) bottom left illustrates that for Cross-Example Softmax the loss for a single query includes all negative scores from  $\mathcal{N}_{B_t}$ , even query/document pairs from different queries.

**Cross-Example Negative Mining.** We can now extend Stochastic Negative Mining to mine for the hardest negative comparisons across the entire batch. Akin to the formulation above, let  $\text{top}k(\mathcal{N}_{i,B_t})$  be the set of the top  $k$  largest scores within the set of negative scores of the entire batch. With this we can define the Cross-Example Negative Mining loss as

$$\mathcal{L}_{B_t} = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{e^{s_{i,i}}}{e^{s_{i,i}} + \sum_{s \in \text{top}k(\mathcal{N}_{B_t})} e^s} \right) \quad (6)$$

This loss is illustrated in Fig. 2 (B) bottom right. Note that negative scores for each query are mined from the entire batch. This means that the set of mined scores could contain all negative scores from some queries, like row 1 in the figure, and no negative scores from others, like row 3 in the figure.

## 4 Experiments

We perform a series of experiments to evaluate the text-to-image retrieval performance of the proposed Cross-Example Softmax. We use two datasets, i.e., Conceptual Captions [29] and Flickr30k [40]. As baselines, we compare to sampled softmax [30] and sampled softmax with in-batch negative mining [24], as well as two angular margin-based similarity learning baselines [18, 19].

| Model                     | PR AUC       | Recall on Test Set |              |              | +3.3M distractors |             |             |              |
|---------------------------|--------------|--------------------|--------------|--------------|-------------------|-------------|-------------|--------------|
|                           |              | @1                 | @5           | @10          | @1                | @5          | @10         | @100         |
| Sampled Softmax [30]      | 14.61        | 25.87              | 50.67        | 61.12        | 1.38              | 3.99        | 6.16        | 20.44        |
| L-Softmax [18]            | 14.03        | 26.18              | 50.78        | 61.25        | 1.38              | 4.20        | 6.33        | 20.72        |
| A-Softmax [20]            | 14.11        | 26.20              | <b>50.98</b> | <b>61.48</b> | 1.38              | 4.18        | 6.38        | 20.82        |
| SNM [24]                  | 14.80        | 26.08              | 50.71        | 61.41        | 1.28              | 3.97        | 6.08        | 20.63        |
| <b>CE-Softmax</b>         | <b>20.12</b> | <b>26.95</b>       | 50.65        | 61.18        | 1.55              | 4.51        | 6.83        | 21.54        |
| <b>CE Negative Mining</b> | 20.09        | 26.91              | 50.85        | 61.21        | <b>1.57</b>       | <b>4.58</b> | <b>6.94</b> | <b>21.62</b> |

Table 1: **Retrieval Results on Conceptual Captions.** Left: Cross-Example Softmax improves score calibration as shown by the increase in PR-AUC. Middle: When retrieving images from the test set only, Cross-Example Softmax outperforms Sampled Softmax and the margin-based baselines for Recall@1. Right: With a large number of distractor images, Cross-Example Softmax clearly outperforms vanilla Sampled Softmax across all levels of  $k$ . Cross-Example Negative Mining further improves upon these results. All numbers shown are the average of 5 runs.

#### 4.1 Text-To-Image Model

For our text-to-image retrieval application we learn two separate encoders, limiting the interaction between text and images to a dot product. This allows to efficiently score a query against a large pre-computed database of image representations using approximate nearest neighbors (ANN), e.g., [5, 15, 38]. This is in contrast to *cross-attention* such as ViLBERT [21], where the interaction between text and image relies on heavy computation and thus does not scale to the retrieval setting.

**Text Encoder.** As text encoder  $f_{text}$  we use a 12-layer Transformer [33]. Specifically, we use the publicly available pre-trained BERT\_Base [10]. The inputs to the model are word piece tokenized image captions padded to a maximum sequence length of 32. As final vector representations for the caption, we use the 768-dimensional vector representing the [CLS] token and add a dimensionality reduction to 128 dimensions without any bias or activation.

**Image Encoder.** As image encoder  $f_{image}$  we use a 50-layer ResNet [11] pretrained on the classic ILSVRC2012 classification task. As final vector representations for the image, we replace the last layer with a dimensionality reduction to 128 dimensions without any bias or activation.

**Relevance Score.** Representations are compared using cosine similarity. Since cosine similarities lie within the range from -1 to 1, the resulting dynamic range is too small when used with softmax. We follow [43] and add a temperature  $\lambda$ , such that the final comparison is  $s_{i,j} = \langle \lambda \frac{\mathbf{x}_i}{|\mathbf{x}_i|}, \lambda \frac{\mathbf{y}_j}{|\mathbf{y}_j|} \rangle$ .

#### 4.2 Training Details

We use the Conceptual Captions [29] dataset to train our text-to-image model. The dataset is set of images from the web along with alt text descriptions. The training split includes 3,318,333 and the test split 12,559 image/text pairs. Our models are trained with batch size 512 for 25,000 steps. For models with negative mining, we select the 50% largest scores from the respective negative sets. Due to the different architectures in the two encoders, we use two separate optimizers [41]. The Transformer uses AdamW with a learning rate of 1e-4 that is linearly decayed to 0 after a warmup period of 1,500 steps. The ResNet uses SGD with momentum of 0.9 and a linear learning rate decay schedule starting at 3e-3 and ending at 5e-4. The output logits are scaled with a factor of  $\lambda = 20$ .

#### 4.3 Retrieval Results on Conceptual Captions

We now evaluate retrieval performance and score calibration of the learned embeddings.

**Recall@k.** We evaluate the retrieval performance of the different models on two retrieval sets of varying size. First, we focus solely on the test set, i.e., for all test captions we retrieve the most relevant image from the set of all test images. Second, for a more realistic and challenging setting we add the 3.3M training images as distractors to the retrieval set. Table 1 shows Recall@k for the first setting in the middle column and for the second setting on the right side. From the results, we observe that for the first setting, Cross-Example Softmax outperforms vanilla Sampled Softmax and the margin-based baselines for Recall@1. For recall at larger  $k$  we do not see large differences between the models. For the second setting, we observe that Cross-Example Softmax clearly outperforms all baselines across all levels of  $k$ . Cross-Example Negative Mining further improves upon these results and consistently achieves the best performance across all models.

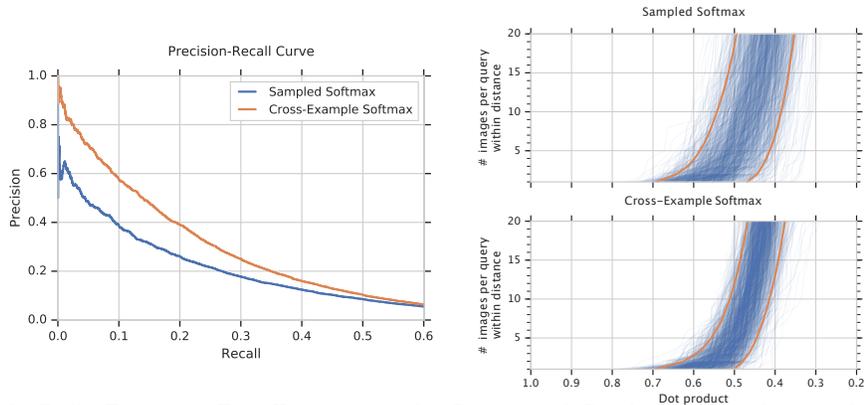


Figure 4: **Left: Precision-Recall curves** on the Conceptual Captions test set images for Sampled Softmax and Cross-Example Softmax. Incorporating negative comparisons across examples leads to much better global calibration and correspondingly higher AUC. **Right: Comparison of score distributions.** Each blue line is one of 1,000 randomly sampled queries from the test set, showing the number of neighbors ( $y$ -axis) at a given similarity to the query ( $x$ -axis). The orange lines indicate the 5th and 95th percentile, i.e., 90% of queries fall in between. The result demonstrates how Cross-Example Softmax leads to a more globally-calibrated similarity metric.

| Model                     | Recall@1     | Recall@5     | Recall@10    |
|---------------------------|--------------|--------------|--------------|
| ViLBERT [21]              | 31.86        | 61.12        | 72.80        |
| Sampled Softmax [30]      | 29.22        | 55.70        | 67.20        |
| L-Softmax [18]            | 29.12        | 55.79        | 66.95        |
| A-Softmax [20]            | 28.74        | 55.80        | 66.96        |
| SNM [24]                  | 29.07        | 55.95        | 67.27        |
| <b>CE-Softmax</b>         | 29.94        | 56.83        | 67.78        |
| <b>CE Negative Mining</b> | <b>30.49</b> | <b>57.04</b> | <b>68.06</b> |

Table 2: **Retrieval Results on Flickr30k.** Cross-Example softmax clearly outperforms Sampled Softmax across all recall levels. Moreover, Cross-Example Negative Mining improves further upon the results. Numbers shown are the average across 5 runs.

**Precision-Recall and AUC.** To measure global score calibration, we now evaluate the Precision-Recall curve, measuring the precision of all returned results at a specified global recall threshold. To generate PR curves, we compute the pairwise similarities between all captions and images in the test set and sort them by their similarity score. The Precision-Recall curves for Sampled Softmax and Cross-Example Softmax are shown in Figure 4 (left). Further, the PR AUC for all methods is shown on the left of Table 1. The results clearly indicate that Cross-Example Softmax improves the global score calibration with an increase in PR-AUC from 0.14 to 0.21. The PR curves in Figure 4 highlight that much of the improvement comes from a better separation between the positives and the hardest negative documents. This highlights the effectiveness of directly optimizing for a better separation between positives and the entire distribution of all negative scores.

#### 4.4 Zero-Shot Retrieval Results on Flickr30k

We now perform a set of zero-shot experiments on Flickr30k to evaluate the generalization of the learned embeddings. We use the zero-shot setting and test split from [21], which contains 1,000 images, each with 5 captions. Following the procedure of previous work, we compare all 5,000 captions to all 1,000 images in the Flickr30k test set in order to retrieve the most relevant images to each query. Note that we do not perform any dataset-specific fine-tuning on Flickr30k; all of our results are shown in the zero-shot setting of training on conceptual captions and testing on Flickr30k.

Table 2 shows the retrieval performance. Besides the methods presented above, the table further includes a comparison to ViLBERT [21]. However, as mentioned earlier, note that ViLBERT is a cross-attention model and thus does not scale to large scale retrieval. Similarly to the Conceptual Captions 3M results above, we find that Cross-Example Softmax outperforms the baselines across all levels of  $k$ . Focusing on the hardest examples with Cross-Example Negative Mining further improves

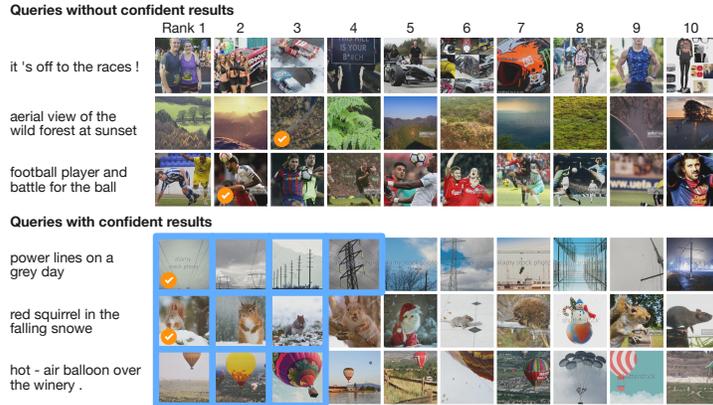


Figure 5: **Example queries and top retrieval results** from Cross-Example Softmax on the Conceptual Captions test set. Retrieved ground truth image have an orange check mark. Images with high relevance score, i.e.,  $\geq 0.6$ , have a blue outline. Top: queries without any high relevance score results. Bottom: queries that contain results with high relevance scores. Queries without high confidence result tend to be vague or not visually descriptive. Queries with high scoring results tend to be more descriptive and images with high scores are mostly plausibly correct. The last example shows a failure case, where the model ignores the ‘winery’ token and solely focuses on the hot air balloon.

recall, demonstrating the effectiveness of the proposed methods. It is interesting that our model performs roughly on-par with cross-attention models like ViLBERT that rely on combining visual and text features early in the model, ruling them infeasible for retrieval tasks over large-scale databases.

#### 4.5 Qualitative Analysis

**Score distributions.** First we look at the effect that Cross-Example Softmax has on the score distribution of retrieved documents. To analyze the distributions, Fig. 4 (right) shows the number of neighbors ( $y$ -axis) at a given similarity to the query ( $x$ -axis) for 1,000 randomly sampled queries from the Conceptual Captions test set. Each blue line represents one query. Blue lines on the leftmost side represent queries whose top retrieved results have higher similarity scores, and the lines on the rightmost side represent queries whose top results have comparatively low similarity. The orange lines indicate the 5th and 95th percentile, i.e., 90% of queries fall in between. The result demonstrates how Cross-Example Softmax leads to a more globally-calibrated similarity metric.

**Example results.** Fig. 5 shows example queries and their top retrieved results. Specifically, it shows three queries where all returned results have low confidence scores and three queries that return many high confidence results. Those queries would show on the rightmost and leftmost sides within the plot of Fig. 4 respectively. Generally, we observe that queries without any highly confident results are often vague or not visually descriptive. On the other hand, queries with high scoring results tend to be more descriptive and images with high scores are mostly plausible results.

## 5 Conclusion

In this work, we proposed Cross-Example Softmax, wherein given a batch of inputs and labels, the loss encourages that the score of a correct input/label pair is higher than all incorrect input/label pairs, even across multiple inputs. By extending the loss to all input/label pairs, we ensure that the relevance score is more globally calibrated and thus becomes more interpretable as an absolute measure of relevance. We further introduce Cross-Example Negative Mining, where the loss is focused on the hardest incorrect input/label pairs. Empirically, we showed that the proposed methods effectively improve global calibration as well as retrieval performance. This work opens up numerous paths for future work. From a practitioner’s point of view, it might be exciting to extend this work towards multiclass classification and object detection. From a robustness perspective, it would be intriguing to study the possible impact of the proposed method on susceptibility towards adversarial examples.

## Broader Impact

Interpretability is a key consideration of machine learning applications. In this work, we propose a novel method for making distances in embedding-based retrieval methods more interpretable. We emphasize that this only addresses one specific aspect of interpretability and only for retrieval.

Further, the general ethical concerns caused by improvements in image recognition do also apply to this work. Although the methods proposed in this paper are motivated to improve model interpretability, they could lead unforeseen applications.

## References

- [1] Chris Alberti, Jeffrey Ling, Michael Collins, and David Reitter. Fusion of detected objects in text for visual question answering. *arXiv preprint arXiv:1908.05054*, 2019.
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *International Conference on Computer Vision (ICCV)*, 2015.
- [3] Samy Bengio, Krzysztof Dembczynski, Thorsten Joachims, Marius Kloft, and Manik Varma. Extreme classification (dagstuhl seminar 18291). *Dagstuhl Reports*, 8:62–80, 2018.
- [4] Yoshua Bengio and Jean-Sébastien Senecal. Adaptive importance sampling to accelerate training of a neural probabilistic language model. *IEEE Transactions on Neural Networks*, 19:713–722, 2008.
- [5] Qi Chen, Haidong Wang, Mingqin Li, Gang Ren, Scarlett Li, Jeffery Zhu, Jason Li, Chuanjie Liu, Lintao Zhang, and Jingdong Wang. *SPTAG: A library for fast approximate nearest neighbor search*, 2018. URL <https://github.com/Microsoft/SPTAG>.
- [6] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [7] Sumit Chopra, Raia Hadsell, Yann LeCun, et al. Learning a similarity metric discriminatively, with application to face verification. In *Computer Vision and Pattern Recognition (CVPR)*, 2005.
- [8] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. Visual dialog. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition (CVPR)*. Ieee, 2009.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [12] Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *International Workshop on Similarity-Based Pattern Recognition*. Springer, 2015.
- [13] Allan Jabri, Armand Joulin, and Laurens Van Der Maaten. Revisiting visual question answering baselines. In *European Conference on Computer Vision (ECCV)*. Springer, 2016.
- [14] Sébastien Jean, Kyunghyun Cho, Roland Memisevic, and Yoshua Bengio. On using very large target vocabulary for neural machine translation. In *Association for Computational Linguistics (ACL)*, 2014.
- [15] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *arXiv preprint arXiv:1702.08734*, 2017.
- [16] Gen Li, Nan Duan, Yuejian Fang, Daxin Jiang, and Ming Zhou. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. *arXiv preprint arXiv:1908.06066*, 2019.

- [17] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [18] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-margin softmax loss for convolutional neural networks. In *International Conference on Computer Vision (ICML)*, 2016.
- [19] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6738–6746, 2017.
- [20] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. In *Conference on computer vision and pattern recognition (CVPR)*, 2017.
- [21] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *ArXiv*, abs/1908.02265, 2019.
- [22] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2013.
- [23] John Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. 1999.
- [24] Sashank J Reddi, Satyen Kale, Felix Yu, Daniel Holtmann-Rice, Jiecao Chen, and Sanjiv Kumar. Stochastic negative mining for learning with large output spaces. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019.
- [25] Antonio Rubio, LongLong Yu, Edgar Simo-Serra, and Francesc Moreno-Noguer. Multi-modal joint embedding for fashion product retrieval. *International Conference on Image Processing (ICIP)*, pages 400–404, 2017.
- [26] Walter J. Scheirer, Anderson Rocha, Ross Michaels, and Terrance E. Boult. Meta-recognition: The theory and practice of recognition score analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 33:1689–1695, 2011.
- [27] Walter J. Scheirer, Anderson Rocha, Jonathan Parris, and Terrance E. Boult. Learning for meta-recognition. *IEEE Transactions on Information Forensics and Security*, 7:1214–1224, 2012.
- [28] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [29] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Association for Computational Linguistics (ACL)*, 2018.
- [30] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Advances in neural information processing systems (NeurIPS)*, 2016.
- [31] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019.
- [32] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019.
- [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [34] Andreas Veit, Balazs Kovacs, Sean Bell, Julian McAuley, Kavita Bala, and Serge Belongie. Learning visual clothing style with heterogeneous dyadic co-occurrences. In *International Conference on Computer Vision (ICCV)*, 2015.
- [35] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [36] Feng Wang, Jian Cheng, Weiyang Liu, and Haijun Liu. Additive margin softmax for face verification. *IEEE Signal Processing Letters*, 25(7):926–930, 2018.

- [37] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [38] Xiang Wu, Ruiqi Guo, Ananda Theertha Suresh, Sanjiv Kumar, Daniel N Holtmann-Rice, David Simcha, and Felix Yu. Multiscale quantization for fast similarity search. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [39] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning (ICML)*, 2015.
- [40] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014.
- [41] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank J Reddi, Sanjiv Kumar, and Suvrit Sra. Why adam beats sgd for attention models. *arXiv preprint arXiv:1912.03194*, 2019.
- [42] Xiao Zhang, Rui Zhao, Yu Qiao, Xiaogang Wang, and Hongsheng Li. Adacos: Adaptively scaling cosine logits for effectively learning deep face representations. In *Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [43] Xu Zhang, Felix Xinnan Yu, Svebor Karaman, Wei Zhang, and Shih-Fu Chang. Heated-up softmax embedding. *arXiv preprint arXiv:1809.04157*, 2018.