

# Automatic Image Captioning

Jia-Yu Pan, Hyung-Jeong Yang, Pinar Duygulu and Christos Faloutsos  
Carnegie Mellon University  
{jypan, hjyang, pinar, christos}@cs.cmu.edu

## Abstract

In this paper, we examine the problem of automatic image captioning. Given a training set of captioned images, we want to discover correlations between image features and keywords, so that we can automatically find good keywords for a new image. We experiment thoroughly with multiple design alternatives on large datasets of various content styles, and our proposed methods achieve up to a 45% relative improvement on captioning accuracy over the state of the art.

## 1. Introduction and related work

“Given a large image database, find the images that have tigers. Given an unseen image, find the terms which best describe its content.” These are some of the problems that many image/video indexing and retrieval systems deal with (see [4][5][10] for recent surveys). Content based image retrieval systems, where images are matched based on visual similarities, have some limitations due to the missing semantic information. Manually annotated words could provide semantic information, however, it is a time consuming and error prone. Several automatic image annotation (captioning) methods have been proposed for better indexing and retrieval of large image databases [1][2][3][6][7].

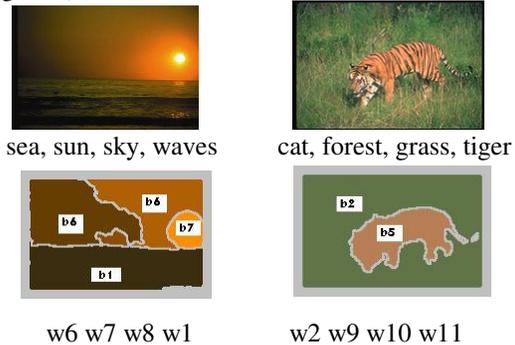
In this study, we are interested in the following problems: “Given a set of images, where each image is captioned with a set of terms describing the image content, find the association (model) between the image features and the terms”. Furthermore, “with the association (model) found, caption an unseen image”. It is possible to caption an image by captioning its constitute regions by a mapping from image regions to terms. Mori *et.al.* [10] used co-occurrence statistics of image grids and words for modeling the associations. Duygulu *et.al.* [3] view the association problem as translation of image regions to words. The correspondences between region groups and words are learned by an EM algorithm. Recently, cross-media relevance models by Jeon *et.al.* [6] and latent semantic analysis (LSA) based models by Monay *et.al.* [11] are also proposed for automatic annotation of images.

In this study, we experiment thoroughly with multiple design alternatives (better clustering decision; weighting on image features and keywords; dimensionality reduction for noise suppression), on large datasets of various content styles. The proposed methods achieve a 45% relative improvement on captioning accuracy over the result of [3].

The paper is organized as follows: Section 2 describes the data set used in the study and input representation. Section 3 and Section 4 describe new methods for obtaining region groups and term weighting respectively. The proposed *correlation-based image captioning* methods are given in Section 5. Section 6 presents the experimental results and Section 7 concludes the paper.

## 2. Input representation

In this study, the associations between image regions and words are learned from collections of manually annotated images (examples are shown in Figure 1).



**Figure 1.** *Top: annotated images with their captions, bottom: corresponding blob tokens and word tokens.*

To build a model given an annotated image set, we first segment the images and represent each segment as a vector of features including color, texture, shape, size and position information. Feature vectors are clustered into  $B$  clusters and each segmented region is assigned to the label of its closest cluster center as in [3]. These labels are called as **blob-tokens**, and they form the visual vocabulary for the image content.

Formally, let  $I=\{I_1, \dots, I_N\}$  be a set of annotated images where each image  $I_i$  is annotated with a set of terms  $W_i = \{w_{i,1}, \dots, w_{i,L_i}\}$  and a set of blob tokens  $B_i = \{b_{i,1}, \dots, b_{i,M_i}\}$ , where  $L_i$  is the number of words, and  $M_i$  is the number of regions in image  $I_i$ . The goal is to construct a model that captures the association between terms and the blob-tokens, given the set  $I$ .

### 3. Blob-token generation

The quality of blob-tokens affects the accuracy of image captioning. In [3], the blob-tokens are generated by applying K-means algorithm on feature vectors of all image regions in the image collection, with the number of blob-tokens,  $B$ , set at 500. However, the choice of  $B=500$  is by no means optimal.

In this study, we use the idea of G-means [12] to determine the number of blob-tokens  $B$  adaptively. G-means clusters the data set starting from small number of clusters,  $B$ , and increases  $B$  iteratively if some of the current clusters fail the Gaussianity test (e.g., Kolmogorov-Smirnov test). In our work, the blob-tokens are the labels of the clusters adaptively found by G-means. The numbers of blob-tokens generated for the 10 training set are all less than 500, ranging from 339 to 495, mostly around 400.

### 4. Weighting by uniqueness

If there are  $W$  possible terms and  $B$  possible blob-tokens, the entire annotated image set of  $N$  images can be represented by a data matrix  $\mathbf{D}_{[N\text{-by-}(W+B)]}$ . We now define two matrices: one is *unweighted*, the other is *uniqueness weighted* as initial data representation.

**Definition 1 (Unweighted data matrix)** Given an annotated image set  $I=\{I_1, \dots, I_N\}$  with a set of terms  $W$  and a set of blob-tokens  $B$ , the unweighted data matrix  $\mathbf{D}_0=[\mathbf{D}_{W0}|\mathbf{D}_{B0}]$  is a  $N$ -by- $(W+B)$  matrix, where the  $(i,j)$ -element of the  $N$ -by- $W$  matrix  $\mathbf{D}_{W0}$  is the count of term  $w_j$  in image  $I_i$ , and the  $(i,j)$ -element of the  $N$ -by- $B$  matrix  $\mathbf{D}_{B0}$  is the count of blob-token  $b_j$  in image  $I_i$ .

We weighted the counts in the data matrix  $\mathbf{D}$  according to the ‘‘uniqueness’’ of each term/blob-token. If a term appears only once in the image set, say with image  $I_1$ , then we will use that term for captioning only when we see the blob-tokens of  $I_1$  again, which is a small set of blob-tokens. The more common a term is, the more blob-tokens it has association with, and the uncertainty of finding the correct term-and-blob-token association goes up. The idea is to give higher weight to terms/blob-tokens which are more ‘‘unique’’ in the training set, and low weights to noisy, common terms/blob-tokens.

**Definition 2 (Uniqueness weighted data matrix)** Given an unweighted data matrix  $\mathbf{D}_0=[\mathbf{D}_{W0}|\mathbf{D}_{B0}]$ . Let  $z_j$

$(y_j)$  be the number of images which contain the term  $w_j$  (the blob-token  $b_j$ ). The weighted data matrix  $\mathbf{D}=[\mathbf{D}_W|\mathbf{D}_B]$  is constructed from  $\mathbf{D}_0$ , where the  $(i,j)$ -element of  $\mathbf{D}_W(\mathbf{D}_B)$ ,  $d_{W(i,j)}$  ( $d_{B(i,j)}$ ), is

$$d_{W(i,j)} = d_{W0(i,j)} \times \log\left(\frac{N}{z_j}\right), d_{B(i,j)} = d_{B0(i,j)} \times \log\left(\frac{N}{y_j}\right), \quad (3)$$

where  $N$  is the total number of images in the set.

In the following, whenever we mention the data matrix  $\mathbf{D}$ , it will be always the weighted data matrix.

### 5. Proposed methods for image captioning

**Definition 3 (Method Corr)** Let  $\mathbf{T}_{\text{corr},0}=\mathbf{D}_W^T\mathbf{D}_B$ . The correlation-based translation table  $\mathbf{T}_{\text{corr}}$  is defined by normalizing each column of  $\mathbf{T}_{\text{corr},0}$  such that each column sum up to 1. Note that the  $(i,j)$ -element of  $\mathbf{T}_{\text{corr}}(i,j)$  can be viewed as an estimate to  $p(w_i|b_j)$ , the conditional probability of term  $w_i$  given blob-token  $b_j$ .

$\mathbf{T}_{\text{corr}}$  measures the association between a term and a blob-token by their co-occurrences. Another possible measure could be to see how similar the overall occurrence patterns (over the training images) of a term and a blob-token are. Such occurrence patterns are in fact the columns of  $\mathbf{D}_W$  or  $\mathbf{D}_B$ , and the similarity can be taken as the cosine value between pairs of column vectors.

**Definition 4 (Method Cos)** Let the  $i$ -th column of the matrix  $\mathbf{D}_W$  ( $\mathbf{D}_B$ ) be  $d_{W_i}$  ( $d_{B_i}$ ). Let  $\text{cos}_{i,j}$  be the cosine value of the angle column vectors  $d_{W_i}$  and  $d_{B_j}$ , and let  $\mathbf{T}_{\text{cos},0}$  be a  $W$ -by- $B$  matrix whose  $(i,j)$ -element  $\mathbf{T}_{\text{cos},0}(i,j)=\text{cos}_{i,j}$ . Normalize the columns of  $\mathbf{T}_{\text{cos},0}$  such that each column sums up to 1, and we get the cosine-similarity translation table  $\mathbf{T}_{\text{cos}}$ .

*Singular Value Decomposition (SVD)* decomposes a given matrix  $\mathbf{X}_{[n \times m]}$  into a product of three matrices  $\mathbf{U}$ ,  $\mathbf{\Lambda}$ ,  $\mathbf{V}^T$ . That is,  $\mathbf{X}=\mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$ , where  $\mathbf{U}=[\mathbf{u}_1, \dots, \mathbf{u}_n]$ , and  $\mathbf{V}=[\mathbf{v}_1, \dots, \mathbf{v}_m]$  are orthonormal, and  $\mathbf{\Lambda}$  is a diagonal matrix. Note that  $\mathbf{u}_i(\mathbf{v}_j)$  are columns of the matrix  $\mathbf{U}(\mathbf{V})$ . Let  $\mathbf{\Lambda}=\text{diag}(\sigma_1, \dots, \sigma_{\min(n,m)})$ , then  $\sigma_j > 0$ , for  $j \leq \text{rank}(\mathbf{X})$ ,  $\sigma_j=0$ , for  $j > \text{rank}(\mathbf{X})$ .

Previous works [14] show that by setting small  $\sigma_j$  to zero, yielding an optimal low rank representation  $\hat{\mathbf{X}}$ , SVD could be used to clean up noise and reveal informative structure in the given matrix  $\mathbf{X}$ , and achieve better performance in information retrieval applications. We propose to use SVD to suppress the noise in the data matrix before learning the association. Following the general rule-of-thumb, we keep the first  $r$   $\sigma_j$ 's which preserve the 90% variance of the distribution, and set others to zero. In the following, we denote the data matrix after SVD as  $\mathbf{D}_{\text{svd}}=[\mathbf{D}_{W,\text{svd}}|\mathbf{D}_{B,\text{svd}}]$ .

**Definition 5** (*Method SvdCorr and SvdCos*) Method **SvdCorr** and **SvdCos** generates the correlation-based translation table  $\mathbf{T}_{\text{corr,svd}}$  and  $\mathbf{T}_{\text{cos,svd}}$  following the procedure outlined in Definition 3 and 4, but instead of starting with the weighted data matrix  $\mathbf{D}$ , here the matrix  $\mathbf{D}_{\text{svd}}$  is used.

**Algorithm 1** (*Captioning*) Given a translation table  $\mathbf{T}_{[W \times B]}$  ( $W$ : total number of terms;  $B$ : total number of blob-tokens), and the number of captioning terms  $k$  for an image. An image with  $l$  blob-tokens  $\mathbf{B}' = \{b'_1, \dots, b'_l\}$ , can be captioned by: First, form a query vector  $\mathbf{q} = [q_1, \dots, q_B]$ , where  $q_i$  is the count of the blob-token  $b_i$  in the set  $\mathbf{B}'$ . Then, compute the **term-likelihood vector**  $\mathbf{p} = \mathbf{T}\mathbf{q}$ , where  $\mathbf{p} = [p_1, \dots, p_W]^T$ , and  $p_i$  is the predicted likelihood of the term  $w_i$ . Finally, we select the  $k$  captioning terms corresponding to the highest  $k$   $p_i$ 's in the  $\mathbf{p}$  vector.

## 6. Experimental results

The experiments are performed on 10 sets of Corel images, each of them contains about 5200 training images and 1750 testing images. The sets cover a variety of themes ranging from urban scenes to natural scenes, and from artificial objects like jet/plane to animals. Each image has in average 3 captioning terms and 9 blobs.

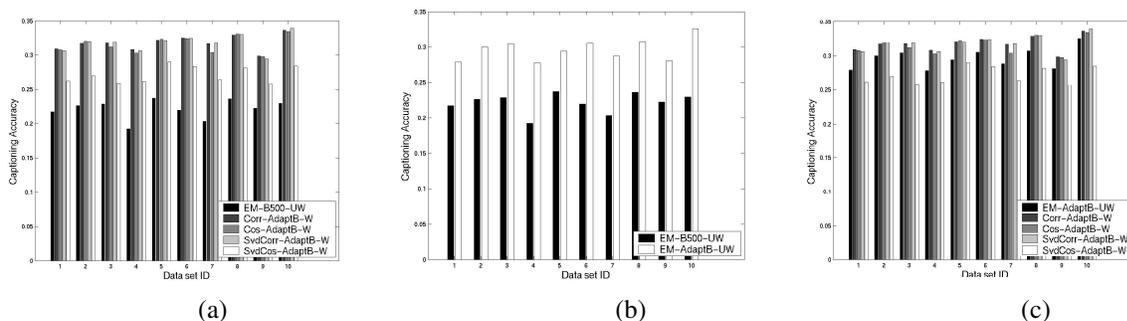
We evaluate the data set in which G-means and uniqueness weighting are applied to show the effects of clustering and weighting. We compare our proposed methods, namely **Corr**, **Cos**, **SvdCorr** and **SvdCos**, with the state-of-the-art machine translation approach [3] as the comparison baseline. For each method, a translation table, an estimate for the conditional probability of a term  $w_i$  given a blob-token  $b_j$  ( $p(w_i|b_j)$ ), is constructed. These translation tables are then used in **Algorithm 1**.

The captioning accuracy on a test image is measured as the percentage of correctly captioned words [1]. The captioning accuracy is defined as  $S = m_{\text{correct}} / m$ , where  $m_{\text{correct}}$  is the number of correctly captioned terms. The overall performance is expressed by the average accuracy over all images in a separate (test) image set.

Figure 2(a) compares the proposed methods with the baseline algorithm [3] which is denoted as **EM-B500-UW** (which means **EM** is applied to an unweighted matrix, denoted **UW**, in which the number of blob tokens is 500, denoted as **B500**). For the proposed methods, blob-tokens are generated adaptively (denoted **AdaptB**) and uniqueness weighting (denoted **W**) is applied. The proposed methods achieve an improvement around 12% absolute accuracy (45% relative improvement) over the baseline.

The proposed adaptive blob-token generation could also improve the baseline **EM** method. Figure 2(b) shows that the adaptively generated blob-tokens improve the captioning accuracy of the EM algorithm. The improvement is around 7.5% absolute accuracy (34.1% relative improvement) over the baseline method (whose accuracy is about 22%). In fact, we found that the improvement is not only on **EM** method, but also on our proposed methods. When the number of blob-tokens is set at 500, proposed methods are 9% less accurate. This suggests that the correct size of blob-token set is not 500, since all methods perform worse when the size is set at 500. Due to the lack of space, we do not show detail figures here.

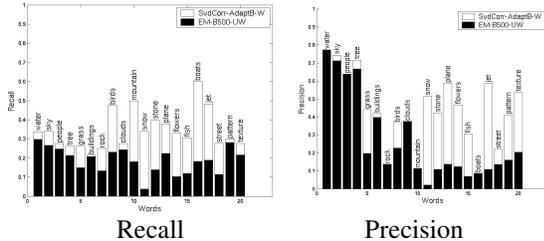
Before applying the ‘‘uniqueness’’ weighting, the 4 proposed methods perform similar to the baseline EM method (the difference is less than 3% of accuracy). The uniqueness weighting improves the performance of all proposed methods except **Cos** method, which stays put. We also observed that weighting does not affect the result of **EM**. Due to the lack of space, we do not show detail figures here.



**Figure 2.** (a) Captioning accuracy improvement over the baseline EM-B500-UW. (b) Score improvement on EM method for adaptively deciding the number of blob-tokens (c) The improvement of proposed methods over modified EM where tokenization is based on g-means.

Another measurement of the performance is the recall and precision values for each word (Figure 3). Given a word  $w$ , let there be  $r$  test images that are captioned with the word  $w$  by a captioning method (set  $R_w$ ). Let  $r^*$  be the actual number of test images that have the word  $w$  (set  $R_w^*$ ), and  $r'$  be size of the intersection of  $R_w$  and  $R_w^*$ . Then, the precision of word  $w$  is  $r'/r$ , and the recall is  $r'/r^*$ .

Note that some words could never be used in the automatic captioning, if the characteristics of the words are not shown in the training samples. We prefer a method which has fewer non-used words, since it could generalize better to unseen images. Table 1 shows that the proposed methods predict more words with non-zero recall and precision values (about three times more than the EM method on average). We observe that EM captions the frequent words with high precision and recall, but misses many words compare to SvdCorr/SvdCos.



**Figure 3.** Recall and precision values for the first 20 frequent words in real annotations. Black bars show the results for EM and white bars show the results for Svd-Corr.

As an example of how well the captioning is, for the image in Figure 1(a), EM-B500-UW and SvdCorr-AdaptiveB-W both give “sky”, “cloud”, “sun” and “water”. As for the image in Figure 1(b), EM-B500-UW gives “grass”, “rocks”, “sky” and “snow”, while SvdCorr-AdaptiveB-W gives “grass”, “cat”, “tiger”, and “water”.

## 6. Conclusion

In this paper, we studied the problem of automatic image captioning and proposed new methods (Corr, Cos, Svd and Svd-Corr) that consistently outperform the state of the art EM (45% relative improvement) in captioning accuracy. Specifically, in this paper,

- We do thorough experiments on large datasets of 10 different image content styles, and examine all possible combinations of the proposed techniques for improving captioning accuracy.
- The proposed uniqueness weighting scheme on terms and blob-tokens boosts the captioning accuracy.

**Table 1.** Comparison of the methods using average recall and precision values and the number of words that can be predicted with non-zero values.

	EM	Corr	Cos	SvdCorr	SvdCos
# predicted	36	57	72	56	132
Avg recall	0.0425	0.1718	0.1820	0.1567	0.2128
Avg prec.	0.0411	0.1131	0.1445	0.1197	0.2079

- Our improved, “adaptive” blob-tokens generation consistently leads to performance gains.
- The proposed methods are less biased to the training set and more generalized in terms of retrieval precision and recall.

The proposed methods can be applied to other areas, such as building an image glossary of different cell types from figures in medical journals [13].

## References

- [1] K. Barnard, P. Duygulu, N. de Freitas, D. Forsyth, D. Blei, M. Jordan, “Matching words and pictures”, Journal of Machine Learning Research, 3:1107:1135, 2003.
- [2] D. Blei, M. Jordan, “Modeling annotated data”, 26th Annual Int. ACM SIGIR Conf., Toronto, Canada, 2003.
- [3] P. Duygulu, K. Barnard, N. de Freitas, D. Forsyth, “Object recognition as machine translation: learning a lexicon for a fixed image vocabulary”, In Seventh European Conference on Computer Vision (ECCV), Vol. 4, pp. 97-112, 2002.
- [4] D. A. Forsyth and J. Ponce, “Computer Vision: a modern approach”, Prentice-Hall, 2001.
- [5] A. Goodrum, “Image information retrieval: An overview of current research”, Informing Science, 3(2), 2000.
- [6] J. Jeon, V. Lavrenko, R. Manmatha, “Automatic Image Annotation and Retrieval using Cross-Media Relevance Models”, 26th Annual Int. ACM SIGIR Conference, Toronto, Canada, 2003.
- [7] J. Li and J. Z. Wang, “Automatic linguistic indexing of pictures by a statistical modeling approach”, IEEE Trans. on Pattern Analysis and Machine Intelligence, 25(10), 2003.
- [8] M. Markkula, E. Sormunen, “End-user searching challenges indexing practices in the digital newspaper photo archive”, Information retrieval, vol.1, 2000.
- [9] Y. Rui, T. S. Huang, S.-F. Chang, “Image Retrieval: Past, Present, and Future”, Journal of Visual Communication and Image Representation, 10:1-23, 1999.
- [10] Y. Mori, H. Takahashi, R. Oka, “Image to word transformation based on dividing and vector quantizing images with words”, First Int. Workshop on Multimedia Intelligent Storage and Retrieval Management, 1999.
- [11] F. Monay, D. Gatica-Perez, “On Image Auto-Annotation with Latent Space Models”, Proc. ACM Int. Conf. on Multimedia (ACM MM), Berkeley, 2003.
- [12] Greg Hamerly and Charles Elkan, “Learning the k in k-means”, Proc. of the NIPS 2003.
- [13] Velliste, M. and R.F. Murphy, 2002. “Automated Determination of Protein Subcellular Locations from 3D Fluorescence Microscope Images,” Proc 2002 IEEE Intl Symp. Biomed Imaging (ISBI 2002), pp. 867-870.
- [14] G. W. Furnas, S. Deerwester, S. T. Dumais, T. K. Landauer, R. A. Harshman, L. A. Streeter and K. E. Lochbaum, “Information retrieval using a singular value decomposition model of latent semantic structure,” Proc. of the 11th ACM SIGIR conf., pp. 465-480, 1998.