

# Large Scale Online Learning of Image Similarity Through Ranking

**Gal Chechik**

GAL@GOOGLE.COM

*Google, 1600 Amphitheatre Parkway, Mountain View CA, 94043*

**Varun Sharma**

VASHARMA@GOOGLE.COM

*Google, RMZ Infinity, Old Madras Road, Bengalooru,  
Karnataka 560016, India*

**Uri Shalit**

URI.SHALIT@MAIL.HUJI.AC.IL

*The Gonda brain research center, Bar Ilan University, 52900, Israel,  
and ICNC, The Hebrew University of Jerusalem, 91904, Israel*

**Samy Bengio**

BENGIO@GOOGLE.COM

*Google, 1600 Amphitheatre Parkway, Mountain View CA, 94043*

**Editor:** Soeren Sonnenburg, Vojtech Franc, Elad Yom-Tov, Michele Sebag

## Abstract

Learning a measure of similarity between pairs of objects is an important generic problem in machine learning. It is particularly useful in large scale applications like searching for an image that is similar to a given image or finding videos that are relevant to a given video. In these tasks, users look for objects that are not only visually similar but also semantically related to a given object. Unfortunately, the approaches that exist today for learning such semantic similarity do not scale to large datasets. This is both because typically their CPU and storage requirements grow quadratically with the sample size, and because many methods impose complex positivity constraints on the space of learned similarity functions.

The current paper presents OASIS, an *Online Algorithm for Scalable Image Similarity* learning that learns a bilinear similarity measure over sparse representations. OASIS is an online dual approach using the passive-aggressive family of learning algorithms with a large margin criterion and an efficient hinge loss cost. Our experiments show that OASIS is both fast and accurate at a wide range of scales: for a dataset with thousands of images, it achieves better results than existing state-of-the-art methods, while being an order of magnitude faster. For large, web scale, datasets, OASIS can be trained on more than two million images from 150K text queries within 3 days on a single CPU. On this large scale dataset, human evaluations showed that 35% of the ten nearest neighbors of a given test image, as found by OASIS, were semantically relevant to that image. This suggests that query independent similarity could be accurately learned even for large scale datasets that could not be handled before.

## 1. Introduction

Large scale learning is sometimes defined as the regime where learning is limited by computational resources rather than by availability of data (Bottou, 2008). Learning a pairwise similarity measure is a particularly challenging large scale task: since pairs of samples have

to be considered, the large scale regime is reached even for fairly small data sets, and learning similarity for large datasets becomes exceptionally hard to handle.

At the same time, similarity learning is a well studied problem with multiple real world applications. It is particularly useful for applications that aim to discover new and relevant data for a user. For instance, a user browsing a photo in her album may ask to find similar or related images. Another user may search for additional data while viewing an online video or browsing text documents. In all these applications, similarity could have different flavors: a user may search for images that are similar visually, or semantically, or anywhere in between.

Many similarity learning algorithms assume that the available training data contains real-valued pairwise similarities or distances. However, in all the above examples, the precise numerical value of pairwise similarity between objects is usually not available. Fortunately, one can often obtain information about the *relative* similarity of different pairs (Frome et al., 2007), for instance, by presenting people with several object pairs and asking them to select the pair that is most similar. For large scale data, where man-in-the-loop experiments are prohibitively costly, relative similarities can be extracted from analyzing pairs of images that are returned in response to the same text query Schultz and Joachims (2004). For instance, the images that are ranked highly by one of the image search engines for the query “cute kitty” are likely to be semantically more similar than a random pair of images. The current paper focuses on this setting: similarity information is extracted from pairs of images that share a common label or are retrieved in response to a common text query.

Similarity learning has an interesting reciprocal relation with classification. On one hand, pairwise similarity can be used in classification algorithms like nearest neighbors or kernel methods. On the other hand, when objects can be classified into (possibly overlapping) classes, the inferred labels induce a notion of similarity across object pairs. Importantly however, similarity learning assumes a form of supervision that is weaker than in classification, since no labels are provided. OASIS is designed to learn a *class-independent* similarity measure with no need for class labels.

A large number of previous studies have focused on learning a similarity measure that is also a metric, like in the case of a positive semidefinite matrix that defines a Mahalanobis distance (Yang, 2006). However, similarity learning algorithms are often evaluated in a context of ranking. For instance, the learned metric is typically used together with a nearest-neighbor classifier (as done in the works of Weinberger et al. (2006), Globerson and Roweis (2006)). When the amount of training data available is very small, adding positivity constraints for enforcing metric properties is useful for reducing over fitting and improving generalization. However, when sufficient data is available, as in many modern applications, adding positive semi-definiteness constraints consumes considerable computation time, and its benefit in terms of generalization are limited. With this view, we take here an approach that avoids imposing positivity or symmetry constraints on the learned similarity measure.

The current paper presents an approach for learning semantic similarity that scales up to two to three orders of magnitude larger than current published approaches. Three components are combined to make this approach fast and scalable: First, our approach uses an unconstrained bilinear similarity. Given two images  $\mathbf{p}_1$  and  $\mathbf{p}_2$  we measure similarity through a bilinear form  $\mathbf{p}_1 \mathbf{W} \mathbf{p}_2$ , where the matrix  $\mathbf{W}$  is not required to be positive, or even

symmetric. Second we use a sparse representation of the images, which allows to compute similarities very fast. Finally, the training algorithm that we developed, OASIS, *Online Algorithm for Scalable Image Similarity learning*, is an online dual approach based on the passive-aggressive algorithm Crammer et al. (2006). It minimizes a large margin target function based on the hinge loss, and converges to high quality similarity measures already after being presented with a small fraction of the training pairs.

We find that OASIS is both fast and accurate at a wide range of scales: for a standard benchmark with thousands of images, it achieves better (but comparable) results than existing state-of-the-art methods, with computation times that are shorter by an order of magnitude. For web-scale datasets, OASIS can be trained on more than two million images within three days on a single CPU. On this large scale dataset, human evaluations of OASIS learned similarity show that 35% of the ten nearest neighbors of a given image are semantically relevant to that image.

The paper is organized as follows. We first present our online algorithm, OASIS, based on the Passive-aggressive family of algorithms. We then present the sparse feature extraction technique used in the experiments. We continue by describing experiments with OASIS on problems of image similarity, at two different scales: a large scale academic benchmark with tens of thousands of images, and a web-scale problem with millions of images. The paper ends with a discussion on properties of OASIS.

## 2. Learning Relative Similarity

We consider the problem of learning a pairwise similarity function  $S$ , given data on the relative similarity of pairs of images.

Formally, let  $\mathcal{P}$  be a set of images, and  $r_{ij} = r(p_i, p_j) \in \mathbb{R}$  be a pairwise relevance measure which states how strongly  $p_j \in \mathcal{P}$  is related to  $p_i \in \mathcal{P}$ . This relevance measure could encode the fact that two images belong to the same category or were appropriate for the same query. We do not assume that we have full access to all the values of  $r$ . Instead, we assume that we can compare some pairwise relevance scores (for instance  $r(p_i, p_j)$  and  $r(p_i, p_k)$ ) and decide which pair is more relevant. We also assume that when  $r(p_i, p_j)$  is not available, its value is zero (since the vast majority of images are not related to each other). Our goal is to learn a similarity function  $S(p_i, p_j)$  that assigns higher similarity scores to pairs of more relevant images,

$$S(p_i, p_i^+) > S(p_i, p_i^-), \quad \forall p_i, p_i^+, p_i^- \in \mathcal{P} \text{ such that } r(p_i, p_i^+) > r(p_i, p_i^-). \quad (1)$$

In this paper we overload notation by using  $p_i$  to denote both the image and its representation as a column vector  $p_i \in \mathbb{R}^d$ . We consider a parametric similarity function that has a bi-linear form,

$$S_{\mathbf{W}}(p_i, p_j) \equiv p_i^T \mathbf{W} p_j \quad (2)$$

with  $\mathbf{W} \in \mathbb{R}^{d \times d}$ . Importantly, if the images  $p_i$  are represented as sparse vectors, namely, only a number  $k_i \ll d$  of the  $d$  entries in the vector  $p_i$  are non-zeroes, then the value of Eq. (2) can be computed very efficiently even when  $d$  is large. Specifically,  $S_{\mathbf{W}}$  can be computed with complexity of  $O(k_i k_j)$  regardless of the dimensionality  $d$ .

## 2.1 An Online Algorithm

We propose an online algorithm based on the Passive-Aggressive (PA) family of learning algorithms introduced by Crammer et al. (2006). Here we consider an algorithm that uses triplets of images  $p_i, p_i^+, p_i^- \in \mathcal{P}$  such that  $r(p_i, p_i^+) > r(p_i, p_i^-)$ .

We aim to find a parametric similarity function  $S$  such that all triplets obey

$$S\mathbf{W}(p_i, p_i^+) > S\mathbf{W}(p_i, p_i^-) + 1 \quad (3)$$

which means that it fulfills Eq. (1) with a safety margin of 1. We define the following hinge loss function for the triplet:

$$l_{\mathbf{W}}(p_i, p_i^+, p_i^-) = \max \{0, 1 - S\mathbf{W}(p_i, p_i^+) + S\mathbf{W}(p_i, p_i^-)\}. \quad (4)$$

To minimize the loss, we apply the Passive-Aggressive algorithm iteratively to optimize  $\mathbf{W}$ . First,  $\mathbf{W}$  is initialized to some value  $\mathbf{W}^0$ . Then, at each training iteration  $i$ , we randomly select a triplet  $(p_i, p_i^+, p_i^-)$ , and solve the following convex problem with soft margin:

$$\begin{aligned} \mathbf{W}^i &= \underset{\mathbf{W}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{W} - \mathbf{W}^{i-1}\|_{Fro}^2 + C\xi \\ \text{s.t.} \quad & l_{\mathbf{W}}(p_i, p_i^+, p_i^-) \leq \xi \quad \text{and} \quad \xi \geq 0 \end{aligned} \quad (5)$$

where  $\|\cdot\|_{Fro}$  is the Frobenius norm (point-wise  $L_2$  norm). Therefore, at each iteration  $i$ ,  $\mathbf{W}^i$  is selected to optimize a trade-off between remaining close to the previous parameters  $\mathbf{W}^{i-1}$  and minimizing the loss on the current triplet  $l_{\mathbf{W}}(p_i, p_i^+, p_i^-)$ . The *aggressiveness* parameter  $C$  controls this trade-off.

We follow Crammer et al. (2006) to solve the problem in Eq. (5). When  $l_{\mathbf{W}}(p_i, p_i^+, p_i^-) = 0$ , it is clear that  $\mathbf{W}^i = \mathbf{W}^{i-1}$  satisfies Eq. (5) directly. Otherwise, we define the Lagrangian

$$\mathcal{L}(\mathbf{W}, \tau, \xi, \lambda) = \frac{1}{2} \|\mathbf{W} - \mathbf{W}^{i-1}\|^2 + C\xi + \tau(1 - \xi - p_i^T \mathbf{W}(p_i^+ - p_i^-)) - \lambda\xi \quad (6)$$

where  $\tau \geq 0$  and  $\lambda \geq 0$  are Lagrange multipliers. The optimal solution is such that the gradient vanishes  $\frac{\partial \mathcal{L}(\mathbf{W}, \tau, \xi, \lambda)}{\partial \mathbf{W}} = 0$ , hence

$$\frac{\partial \mathcal{L}(\mathbf{W}, \tau, \xi, \lambda)}{\partial \mathbf{W}} = \mathbf{W} - \mathbf{W}^{i-1} - \tau \mathbf{V}_i = 0$$

where the gradient matrix  $\mathbf{V}_i = \frac{\partial \mathbb{1}_{\mathbf{W}}}{\partial \mathbf{W}} = [p_i^1(p_i^+ - p_i^-), \dots, p_i^d(p_i^+ - p_i^-)]^T$ . The optimal new  $\mathbf{W}$  is therefore

$$\mathbf{W} = \mathbf{W}^{i-1} + \tau \mathbf{V}_i \quad (7)$$

where we still need to estimate  $\tau$ . Differentiating the Lagrangian with respect to  $\xi$  and setting it to zero also yields:

$$\frac{\partial \mathcal{L}(\mathbf{W}, \tau, \xi, \lambda)}{\partial \xi} = C - \tau - \lambda = 0 \quad (8)$$

which, knowing that  $\lambda \geq 0$ , means that  $\tau \leq C$ . Plugging Equations (7) and (8) back into the Lagrangian in Eq. (6), we obtain

$$\mathcal{L}(\tau) = \frac{1}{2}\tau^2\|\mathbf{V}_i\|^2 + \tau(1 - p_i^T(\mathbf{W}^{i-1} + \tau\mathbf{V}_i)(p_i^+ - p_i^-)) . \quad (9)$$

Regrouping the terms we obtain

$$\mathcal{L}(\tau) = -\frac{1}{2}\tau^2\|\mathbf{V}_i\|^2 + \tau(1 - p_i^T\mathbf{W}^{i-1}(p_i^+ - p_i^-)) .$$

Taking the derivative of this second Lagrangian with respect to  $\tau$  and setting it to 0, we have

$$\frac{\partial \mathcal{L}(\tau)}{\partial \tau} = -\tau\|\mathbf{V}_i\|^2 + (1 - p_i^T\mathbf{W}^{i-1}(p_i^+ - p_i^-)) = 0$$

which yields

$$\tau = \frac{1 - p_i^T\mathbf{W}^{i-1}(p_i^+ - p_i^-)}{\|\mathbf{V}_i\|^2} = \frac{l_{\mathbf{W}^{i-1}}(p_i, p_i^+, p_i^-)}{\|\mathbf{V}_i\|^2} .$$

Finally, Since  $\tau \leq C$ , we obtain

$$\tau = \min \left\{ C, \frac{l_{\mathbf{W}^{i-1}}(p_i, p_i^+, p_i^-)}{\|\mathbf{V}_i\|^2} \right\} . \quad (10)$$

Equations (7) and (10) summarize the update needed for every triplets  $(p_i, p_i^+, p_i^-)$ . It has been shown (Crammer et al., 2006) that applying such an iterative algorithm yields a cumulative online loss that is likely to be small. It was furthermore shown that selecting the best  $\mathbf{W}_i$  during training using a hold-out validation set achieves good generalization.

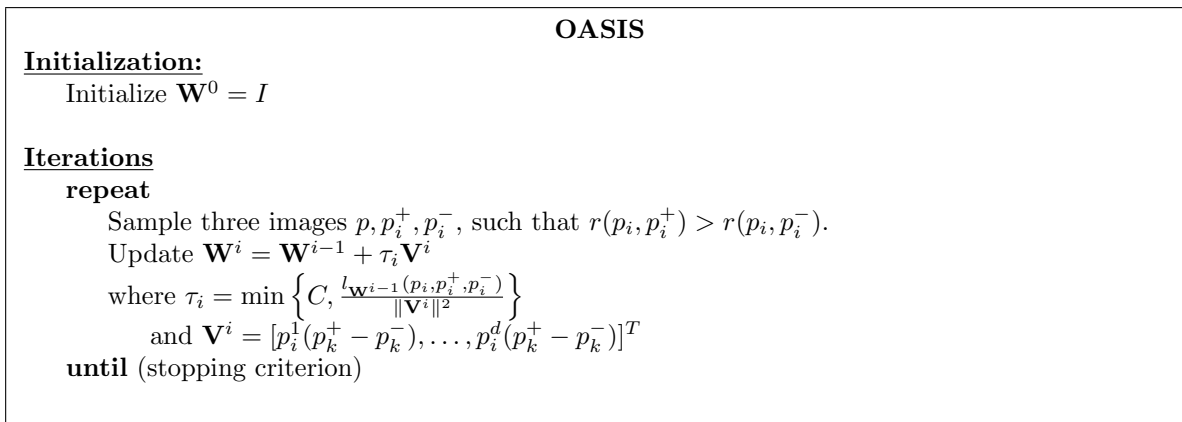


Figure 1: Pseudo-code of the OASIS algorithm.

## 2.2 Relation to Large Margin Nearest Neighbor Classification

The similarity matrix  $\mathbf{W}$  learned by OASIS is not guaranteed to be positive or even symmetric. This section discusses a variant of OASIS that yields symmetric solutions. We redefine the similarity function as

$$\hat{S}_{\mathbf{W}}(p_i, p_j) \equiv -(p_i - p_j)^T \mathbf{W} (p_i - p_j) . \quad (11)$$

Given the corresponding triplet hinge loss function:

$$\hat{l}_{\mathbf{W}}(p_i, p_i^+, p_i^-) = \max \left\{ 0, 1 - \hat{S}_{\mathbf{W}}(p_i, p_i^+) + \hat{S}_{\mathbf{W}}(p_i, p_i^-) \right\}$$

we can solve the following convex optimization problem, similar to Eq. (5),

$$\begin{aligned} \mathbf{W}^i &= \underset{\mathbf{W}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{W} - \mathbf{W}^{i-1}\|_{Fro}^2 + C\xi \\ \text{s.t.} \quad &\hat{l}_{\mathbf{W}}(p_i, p_i^+, p_i^-) \leq \xi \quad \text{and} \quad \xi \geq 0 \end{aligned} \quad (12)$$

and obtain

$$\begin{aligned} \mathbf{W}^i &= \mathbf{W}^{i-1} - \hat{\tau}_i \hat{\mathbf{V}}^i, \\ \text{where} \quad \hat{\tau}_i &= \min \left\{ C, \frac{\hat{l}_{\mathbf{W}^{i-1}}(p_i, p_i^+, p_i^-)}{\|\hat{\mathbf{V}}^i\|^2} \right\} \\ \text{and} \quad \hat{\mathbf{V}}^i &= (p_i - p_i^+)(p_i - p_i^+)^T - (p_i - p_i^-)(p_i - p_i^-)^T . \end{aligned}$$

Since the matrix  $\hat{\mathbf{V}}^i$  is symmetric, each update of  $\mathbf{W}$  preserves its symmetry. Hence, if initialized with a symmetric  $\mathbf{W}^0$ , we are guaranteed to obtain a symmetric solution  $\mathbf{W}_i$  at any step  $i$ . We name this variant of OASIS as DISSIM-OASIS, since  $(p_i - p_j)^T \mathbf{W} (p_i - p_j)$  is a dissimilarity measure.

DISSIM-OASIS is closely related to the *Large margin nearest neighbor* (LMNN) algorithm (Weinberger et al., 2006). In LMNN, samples are taken from multiple distinct classes, and a batch loss function is defined using a metric  $\mathbf{W} = L^T L$ :

$$\begin{aligned} \epsilon^{LMNN}(W) &= \omega \cdot \sum_{i,j \in N(i)} (p_i - p_j)^T \mathbf{W} (p_i - p_j) + \\ &\quad (1 - \omega) \cdot \sum_{p_i, p_j, p_l} \max [0, 1 + (p_i - p_j)^T \mathbf{W} (p_i - p_j) - (p_i - p_l)^T \mathbf{W} (p_i - p_l)] . \end{aligned} \quad (13)$$

Here the first sum is over target pairs  $p_i$  and  $p_j$  such that that  $p_j$  is one of the  $k$  nearest neighbors of  $p_i$  and is also in the same class as  $p_i$ .  $k$  is usually set to three. The second sum is also over an image  $p_l$  that is in a different class than  $p_i$  and  $p_j$ .

To study the relation between LMNN and OASIS, we cast LMNN into an online format, and assume that  $p_i$  and  $p_i^+$  share a class label while  $p_i^-$  has a different label. Using  $\hat{S}_{\mathbf{W}}(p_i, p_j)$  and  $\hat{l}_{\mathbf{W}}(p_i, p_i^+, p_i^-)$  defined above, we have

$$\epsilon^{online}(W) = -\omega \cdot \hat{S}_{\mathbf{W}}(p_i, p_i^+) + (1 - \omega) \cdot \hat{l}_{\mathbf{W}}(p_i, p_i^+, p_i^-) . \quad (14)$$

For  $\omega > 0$ , the first term  $-\hat{S}\mathbf{w}(p_i, p_i^+)$  is always positive (except the trivial case of  $p_i = p_i^+$ ), and the second term  $\hat{l}\mathbf{w}(p_i, p_i^+, p_i^-)$  is always non-negative. As a result the loss is always non-zero and an update will be performed on every step.

However, when,  $\omega = 0$ , this online version of LMNN becomes equivalent to the DISSIM-OASIS problem.

### 2.3 Sampling Strategy

For real world data sets, the actual number of triplets  $(p_i, p_i^+, p_i^-)$  is typically very large and cannot be stored in memory. Instead, we use the fact that the number of relevant images for a category or a query is typically small, and keep a list of relevant images for each query or category. For the case of single-labeled images, we can efficiently retrieve an image that is relevant to a given image, by first finding its class, and then finding another image from that class. The case of multi-labeled images is described in Sec. 5.3.

Specifically, to sample a triplet  $(p_i, p_i^+, p_i^-)$  during training, we first uniformly sample an image  $p_i$  from  $\mathcal{P}$ . Then we uniformly sample an image  $p_i^+$  from the images sharing the same categories or queries as  $p_i$ . Finally, we uniformly sample an image  $p_i^-$  from the images that share no category or query with  $p_i$ . When the set  $\mathcal{P}$  is very large and the number of categories or queries is also very large, one does not need to maintain the set of non-relevant images for each image: sampling directly from  $\mathcal{P}$  instead only adds a small amount of noise to the training procedure and is not really harmful.

When relevance feedbacks  $r(p_i, p_j)$  are provided as real numbers and not just  $\in \{0, 1\}$ , one could use these number to bias training towards those pairs that have a higher relevance feedback value. This can be done by considering  $r(p_i, p_j)$  as frequencies of apparition, and sampling pairs according to the distribution of these frequencies.

## 3. Image Representation

The problem of selecting an informative representation of images is still an unsolved computer vision challenge, and an ongoing research topic. Different approaches for image representation have been proposed including by Feng et al. (2004), Takala et al. (2005), Tieu and Viola (2004). In the information retrieval community there is wide agreement that a bag-of-words representation is a very useful representation for handling text documents in a wide range of applications. For image representation, there is still no such approach that would be adequate for a wide variety of image processing problems. However, among the proposed representations, a consensus is emerging on using *local descriptors* for various tasks, e.g. (Lowe, 2004, Quelhas et al., 2005). This type of representation segments the image into *regions of interest*, and extracts visual features from each region. The segmentation algorithm as well as the region features vary among approaches, but, in all cases, the image is then represented as a set of feature vectors describing the regions of interest. Such a set is often called a *bag-of-local-descriptors*.

In this paper we take the approach of creating a sparse representation based on the framework of local descriptors. Our features are extracted by dividing each image into overlapping square blocks, and each block is then described with edge and color histograms. For edge histograms, we rely on *uniform Local Binary Patterns* (uLBPs) proposed by Ojala

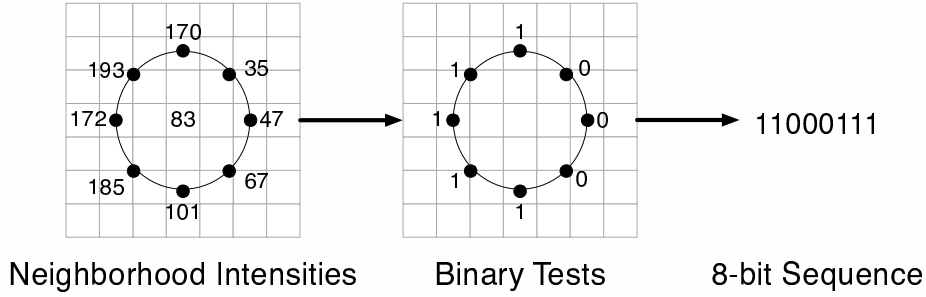


Figure 2: An example of Local Binary Pattern ( $LBP_{8,2}$ ). For a given pixel, the Local Binary Pattern is an 8-bit code obtained by verifying whether the intensity of the pixel is greater or lower than its 8 neighbors.

et al. (2002). These texture descriptors have shown to be effective on various tasks in the computer vision literature (Ojala et al., 2002, Takala et al., 2005), certainly due to their robustness with respect to changes in illumination and other photometric transformations (Ojala et al., 2002). Local Binary Patterns estimate a texture histogram of a block by considering differences in intensity at circular neighborhoods centered on each pixel. Precisely, we use  $LBP_{8,2}$  patterns, which means that a circle of radius 2 is considered centered on each block. For each circle, the intensity of the center pixel is compared to the interpolated intensities located at 8 equally-spaced locations on the circle, as shown on Figure 2, left. These eight binary tests (lower or greater intensity) result in an 8-bit sequence, see Figure 2, right. Hence, each block pixel is mapped to a sequence among  $2^8 = 256$  possible sequences and each block can therefore be represented as a 256-bin histogram. In fact, it has been observed that the bins corresponding to non-uniform sequences (sequences with more than 2 transitions  $1 \rightarrow 0$  or  $0 \rightarrow 1$ ) can be merged, yielding more compact 59-bin histograms without performance loss (Ojala et al., 2002).

Color histograms are obtained by K-means clustering. We first select a palette or typical colors by training a color codebook from the Red-Green-Blue pixels of a large training set of images using K-means. The color histogram of a block is then obtained by mapping each block pixel to the closest color in the codebook palette.

Finally, the histograms describing color and edge statistics of each block are concatenated, which yields a single vector descriptor per block. Our local descriptor representation is therefore simple, relying on both a basic segmentation approach and simple features. Naturally, alternative representations could also be used with OASIS, (Feng et al., 2004, Grangier et al., 2006, Tieu and Viola, 2004) However, this paper focuses on the learning model, and a benchmark of image representations is beyond the scope of the current paper.

As a final step, we use the representation of blocks to obtain a representation for an image. For computation efficiency we aim at a high dimensional and sparse vector space. For this purpose, each local descriptor of an image  $p$  is represented as a discrete index, called *visual term* or *vistterm*, and, like for text data, the image is represented as a *bag-of-vistterms* vector, in which each component  $p_i$  is related to the presence or absence of vistterm  $i$  in  $p$ .



The mapping of the descriptors to discrete indexes is performed according to a codebook  $C$ , which is typically learned from the local descriptors of the training images through k-means clustering (Duygulu et al., 2002, Jeon and Manmatha, 2004, Quelhas et al., 2005). The assignment of the weight  $p_i$  of visterm  $i$  in image  $p$  is as follows:

$$p_i = \frac{f_i d_i}{\sqrt{\sum_{j=1}^{|C|} (f_j d_j)^2}}, \quad (15)$$

where  $f_i$  is the term frequency of  $i$  in  $p$ , which refers to the number of occurrences of  $i$  in  $p$ , while  $d_j$  is the inverse document frequency of  $j$ , which is defined as  $-\log(r_j)$ ,  $r_j$  being the fraction of training images containing at least one occurrence of visterm  $j$ . This approach has been found successful for the task of content based image ranking described by Grangier and Bengio (2008).

In the experiments described below, we used a large set of images collected from the web to train the features. This set is described in more detail in Sec. 5.3. We used a set of 20 typical RGB colors (hence the number of clusters used in the k-means for colors was 20), the block vocabulary size  $|C| = 10000$  and our image blocks were of size 64x64 pixels, overlapping every 32 pixels. Furthermore, in order to be robust to scale, we extracted blocks at various scales by successively down scaling images by a factor of 1.25 and extracting the features at each level, until there were less than 10 blocks in the resulting image. There was on average around 70 non-zero values (out of 10000) describing a single image.

## 4. Related Work

Similarity learning can be considered in two main setups, depending on the type of available training labels. First, a regression setup, where the training set consists of pairs of objects  $x_i^1, x_i^2$  and their pairwise similarity  $y_i \in \mathbf{R}$ . In many cases however, precise similarities are not available, but rather a weaker notion of similarity order. In one such setup, the training set consists of triplets of objects  $x_i^1, x_i^2, x_i^3$  and a ranking similarity function, that can tell which of the two pairs  $(x^1, x^2)$  or  $(x^1, x^3)$  is more similar. Finally, multiple similarity learning studies assume that a binary measure of similarity is available  $y_i \in \{+1, -1\}$ , a pair of objects is either similar or not.

For small-scale data, there are two main groups of similarity learning approaches. The first approach, learning Mahalanobis distances, can be viewed as learning a linear projection of the data into another space (often of lower dimensionality), where a Euclidean distance is defined among pairs of objects. Such approaches include Fisher’s Linear Discriminant Analysis (LDA), relevant component analysis (RCA) (Bar-Hillel et al., 2003), supervised global metric learning (Xing et al., 2003), large margin nearest neighbor (LMNN) (Weinberger et al., 2006) and Metric Learning by Collapsing Classes (Globerson and Roweis, 2006). See also a review by Yang (2006) for more details.

The second family of approaches, learning kernels, is used to improve performance of kernel based classifiers. Learning a full kernel matrix in a non parametric way is prohibitive except for very small data. As an alternative, several studies suggested to learn a weighted sum of pre-defined kernels (Lanckriet et al., 2004) where the weights are being learned from data. In some applications this was shown to be inferior to uniform weighting of the

kernels (Noble, 2008). The work of Frome et al. (2007) further learns a weighting over local distance function for every image in the training set.

Finally, Jain et al. (2008) (based on work by Davis et al. (2007)) aim to learn metrics in an online setting. This work is one of the closest work with respect to OASIS: it learns online a linear model of a [dis-]similarity function between documents; the main difference is that Jain et al. (2008) try to learn a true distance, imposing positive definiteness constraints, which makes the algorithm more complex and more constrained. We argue in this paper that in the large scale regime, such a constraint is not necessary given the amount of available training examples.

Another work closely related to OASIS is that of Rasiwasia and Vasconcelos (2008), which also tries to learn a semantic similarity function between images. In their case, however, semantic similarity is learned by representing each image by the posterior probability distribution over a predefined set of semantic tags, and then computing the distance between two images as the distance between the two underlying posterior distributions. The representation size of images in this approach is therefore equal to the number of semantic classes, hence it will not scale when the number of semantic classes is very large as in free text search.

## 5. Experiments

Evaluating large scale learning algorithms poses special challenges. First, the benchmark datasets that are currently available for academic research are limited either in their scale (like 30K images in Caltech256, as described by Griffin et al. (2007)) or in their resolution (such as the tiny images dataset of Torralba et al. (2007)). Large scale methods are not expected to perform particularly well on small datasets, since they are designed to extract limited information from each sample. Second, many images on the web cannot be used without explicit permission, hence they cannot be collected and packed into a single database. Large, proprietary collections of images do exist, but are not available freely for academic research. Finally, except for very few cases, similarity learning approaches in current literature do not scale to handle large datasets effectively, which makes it hard to compare a new large scale method with the existing methods.

To address these issues, this paper takes the approach of conducting experiments at two different scales. First, to compare OASIS with small-scale methods we used subsets of the standard Caltech256 benchmark. This dataset is one of the largest academic datasets, and we found that OASIS performs well in such a setting. Second, we applied OASIS to a web-scale data with more than 2 million images. This data cannot be handled by current metric learning approaches, hence we report our results in terms of runtime and performance.

### 5.1 Evaluation Measures

We evaluated the performance of all algorithms using standard ranking precision measures based on nearest neighbors. For each query image in the test set, all other test images were ranked according to their similarity to the query image. The number of same-class images among the top  $k$  images (the  $k$  nearest neighbors) was computed. When averaged across test images (either within or across classes), this yields a measure known as precision-at-top- $k$ , providing a precision curve as a function of the rank  $k$ .

We also calculated the *mean average precision* (mAP), a measure that is widely used in the information retrieval community. To compute average precision, the precision-at-top- $k$  is first calculated for each test image. Then, it is averaged over all positions  $k$  that have a positive sample. For example, if all positives are ranked highest, the average-precision is 1. The average-precision measure is then further averaged across all test image queries, yielding the *mean average precision* (mAP).

## 5.2 Caltech256 Dataset

To compare OASIS with small-scale methods we used the *Caltech256* dataset (Griffin et al., 2007). This dataset consists of 30607 images that were obtained from Google image search and from *PicSearch.com*. Images were assigned to 257 categories and evaluated by humans in order to ensure image quality and relevance. After we have pre-processed the images as described in Sec. 3 and filtered images that were too small, we were left with 29461 images in 256 categories. To allow comparisons with other methods in the literature that were not optimized for sparse representation, we also reduced the block vocabulary size  $|C|$  from 10000 to 1000.

Using the Caltech256 dataset allows us to compare OASIS with existing similarity learning methods. For OASIS, we treated images that have the same labels as similar. The same labels were used for comparing with methods that learn a metric for classification, as described below.

### 5.2.1 COMPARED METHODS

We compared the following approaches:

1. **OASIS**. - The algorithm described above in Sec. 2.1.
2. **Euclidean**. - The standard Euclidean distance in feature space. The initialization of OASIS using the identity matrix is equivalent to this distance measure.
3. **MCML** - Metric Learning by Collapsing Classes (Globerson and Roweis, 2006). This approach learns a Mahalanobis distance such that samples from the same class are mapped to the same point. The problem is written as a convex optimization problem, and we have used the gradient-descent implementation provided by the authors.
4. **LMNN** - Large Margin Nearest Neighbor Classification (Weinberger et al., 2006). This approach learns a Mahalanobis distance for  $k$ -nearest neighbor classification, aiming to have the  $k$ -nearest neighbors of a given sample belong to the same class while examples from different classes are separated by a large margin. As a preprocessing phase, images were projected to a basis of the principal components (PCA) of the data, with no dimensionality reduction, since this improved the precision results.
5. **LEGO** - Online metric learning (Jain et al., 2008). LEGO learns a Mahalanobis distance in an online fashion using a regularized per instance loss, yielding a positive semidefinite matrix. The main variant of LEGO aims to fit a given set of pairwise distances. We used another variant of LEGO that, like OASIS, learns from relative distances. In our experimental setting, the loss is incurred for same-class examples

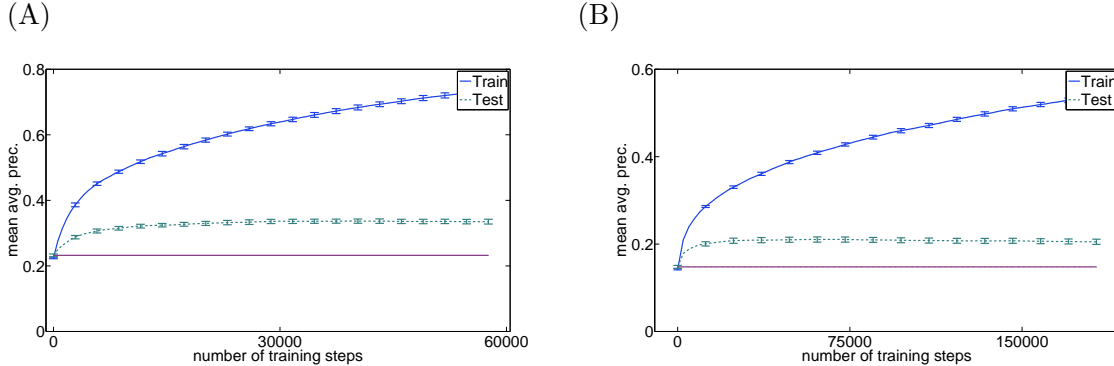


Figure 3: Mean average precision of OASIS as a function of the number of training steps. Error bars represent standard error of the mean over 5 selections of training (40 images) and test (25 images) sets. Performance is compared with a baseline obtained using the naïve Euclidean metric on the feature vector.  $C=0.1$  (A) **Var10**. Test performance saturates around 30K training steps, while going over all triplets would require 2.8 million steps. (B) **Var20**.

being more than a certain distance away, and different class examples being less than a certain distance away. LEGO uses the LogDet divergence for regularization, as opposed to the Frobenius norm used in OASIS.

For all these approaches, we used an implementation provided by the authors. For all approaches except OASIS and LEGO, algorithms were implemented in Matlab, with runtime bottlenecks implemented in C for speedup. For OASIS, we used a Matlab implementation (with no  $C$  components) for the Caltech256 experiments and a  $C^{++}$  implementation for the web-scale experiments described below.

We have also experimented with the methods of Xing et al. (2003) and RCA (Bar-Hillel et al., 2003). We found the method of Xing et al. (2003) to be too slow for the sets in our experiments. RCA is based on a per-class eigen decomposition that is not well defined when the number of samples is smaller than the feature dimensionality. We therefore experimented with a preprocessing phase of dimensionality reduction followed by RCA, but results were inferior to other methods and were not included in the evaluations below.

### 5.2.2 EXPERIMENTAL PROTOCOL

We tested all methods on subsets of classes taken from the Caltech256 repository. Each subset was built such that it included semantically diverse categories, controlled for classification difficulty.

The first set, *Easy10* consists of 10 classes taken amongst those classes that are easiest to classify as characterized by the Caltech256 benchmark. The second set *Var10* consisted of 10 classes that span the full range of classification difficulty. The third set *Var20* consisted of 20 classes, which again span the range of difficulties. The full lists of categories in each set are given in Appendix B.

Table 1: Average precision and precision at top 1, 10, and 50 of all compared methods

<b>Var10</b>	OASIS	MCML	LEGO	LMNN	Euclidean
Mean avg prec	<b>0.33</b> $\pm$ 0.016	0.29 $\pm$ 0.017	0.27 $\pm$ 0.008	0.24 $\pm$ 0.016	0.23 $\pm$ 0.009
Top 1 prec.	<b>0.43</b> $\pm$ 0.040	0.39 $\pm$ 0.051	0.39 $\pm$ 0.048	0.38 $\pm$ 0.054	0.37 $\pm$ 0.041
Top 10 prec.	<b>0.38</b> $\pm$ 0.013	0.33 $\pm$ 0.018	0.32 $\pm$ 0.012	0.29 $\pm$ 0.021	0.27 $\pm$ 0.015
Top 50 prec.	<b>0.23</b> $\pm$ 0.015	0.22 $\pm$ 0.013	0.20 $\pm$ 0.005	0.18 $\pm$ 0.015	0.18 $\pm$ 0.007
<b>Easy10</b>	OASIS	MCML	LEGO	LMNN	Euclidean
Mean avg prec	<b>0.57</b> $\pm$ 0.009	0.55 $\pm$ 0.013	0.51 $\pm$ 0.016	0.46 $\pm$ 0.015	0.42 $\pm$ 0.005
Top 1 prec.	<b>0.66</b> $\pm$ 0.022	0.65 $\pm$ 0.010	0.63 $\pm$ 0.028	0.62 $\pm$ 0.025	0.59 $\pm$ 0.029
Top 10 prec.	<b>0.61</b> $\pm$ 0.016	0.59 $\pm$ 0.014	0.57 $\pm$ 0.024	0.53 $\pm$ 0.016	0.50 $\pm$ 0.011
Top 50 prec.	<b>0.33</b> $\pm$ 0.006	<b>0.33</b> $\pm$ 0.007	0.31 $\pm$ 0.006	0.29 $\pm$ 0.010	0.27 $\pm$ 0.002
<b>Var20</b>	OASIS	MCML	LEGO	LMNN	Euclidean
Mean avg prec	<b>0.21</b> $\pm$ 0.014	0.17 $\pm$ 0.012	0.16 $\pm$ 0.012	0.14 $\pm$ 0.006	0.14 $\pm$ 0.007
Top 1 prec.	<b>0.29</b> $\pm$ 0.026	0.26 $\pm$ 0.023	0.26 $\pm$ 0.027	0.26 $\pm$ 0.030	0.25 $\pm$ 0.026
Top 10 prec.	<b>0.24</b> $\pm$ 0.019	0.21 $\pm$ 0.015	0.20 $\pm$ 0.014	0.19 $\pm$ 0.010	0.18 $\pm$ 0.010
Top 50 prec.	<b>0.15</b> $\pm$ 0.004	0.14 $\pm$ 0.005	0.13 $\pm$ 0.006	0.11 $\pm$ 0.002	0.12 $\pm$ 0.002

For each set, images from each class were split into a training set of 40 images and a test set of 25 images, as proposed by Griffin et al. (2007).

We used cross-validation to select the values of hyper parameters for all algorithms except MCML. Models were learned on 80% of the training set (32 images), and evaluated on the remaining 20%. Cross validation was used for setting the following hyper parameters: the early stopping time for OASIS; the  $\omega$  parameter for LMNN ( $\omega \in \{0.125, 0.25, 0.5\}$ ), and the regularization parameter  $\eta$  for LEGO ( $\eta \in \{0.02, 0.08, 0.32\}$ ). We found that LEGO was usually not sensitive to the choice of  $\eta$ , yielding a variance that was smaller than the variance over different cross-validation splits. Results reported below were obtained by selecting the best value of the hyper parameter and then training again on the full training set (40 images). For MCML, we used the default parameters supplied with the code from the authors, since its very long run time and multiple parameters made it non-feasible to tune hyper parameters on this data.

Table 2: Runtime (minutes) of all compared methods

	OASIS(Matlab)	MCML	LEGO	LMNN	DISSIM-OASIS(Matlab)
Var10	42 $\pm$ 15	1835 $\pm$ 210	143 $\pm$ 44	337 $\pm$ 169	58 $\pm$ 14
Easy10	18 $\pm$ 11	2554 $\pm$ 178	125 $\pm$ 32	207 $\pm$ 205	55 $\pm$ 3
Var20	45 $\pm$ 8	7425 $\pm$ 106	533 $\pm$ 49	631 $\pm$ 40	111 $\pm$ 34

### 5.2.3 RESULTS

Figure 3 traces the mean average precision over the training and the test sets as it progresses during learning. For the Easy10 and Var10 tasks, precision on the test set saturates early

(around 35K training steps), and then decreases very slowly. Training on the Var20 task was performed using a smaller aggressiveness parameter (determined by cross-validation) and thus test precision does not saturate as early.

Figure 4 and Table 1 compare the precision obtained with OASIS, with four competing approaches, as described above (Sec. 5.2.1). OASIS achieved consistently superior results throughout the full range of  $k$  (number of neighbors) tested, and on all four sets studied. Interestingly, we found that LMNN performance on the training set was often high, suggesting that it over fits the training set. This behavior was also noted by (Weinberger et al., 2006) in some of their experiments.

OASIS achieves superior or equal performance, with a runtime that is faster by about two orders of magnitudes than MCML, and about one order of magnitude faster than LMNN. The run time of OASIS and LEGO was measured until the point of early stopping.

Table 2 shows the total CPU time in minutes for training each of the algorithms compared. For the purpose of a fair comparison with competing approaches, we tested a Matlab implementation of OASIS. The versions of LMNN and MCML tested were supplied by the authors and implemented in Matlab, with core parts implemented in C. LEGO is fully implemented in Matlab as OASIS. All code was compiled to C, and was run on a standard CPU. Importantly, we found that Matlab does not make full use of the speedup that can be gained by sparse image representation. As a result, a C/C++ implementation of OASIS that was tested in the next section was found to be significantly faster.

#### 5.2.4 THE EFFECT OF SYMMETRY

We further looked in more details into the effect of enforcing symmetry. As discussed in Section 2.2 the OASIS similarity function based on  $\mathbf{W}$  may not be symmetric. We tested several variants of OASIS that do preserve symmetry. First, we tested the symmetric version discussed in Sec. 2.2, named here DISSIM-OASIS. Second, we tested an approach that updates the two off-diagonal parts of the matrix  $\mathbf{W}$  at each step,  $\mathbf{W}^{new} = \mathbf{W} + \tau \frac{1}{2} (\mathbf{V} + \mathbf{V}^T)$ , named here ONLINE-PROJ OASIS. Finally we modified the final  $\mathbf{W}$  obtained by OASIS in order to make it symmetric,  $\mathbf{W}^{new} = \frac{\mathbf{W} + \mathbf{W}^T}{2}$ , and called it PROJ OASIS.

Figure 5 compares the precision of the different symmetric methods with the original OASIS. In general, the symmetric variants perform slightly worse, or equal to the asymmetric OASIS. Asymmetric OASIS is also twice faster than DISSIM-OASIS, as shown in Table 2. It is interesting to note that the performance of PROJ OASIS was equivalent to that of OASIS, hinting that the final  $\mathbf{W}$  matrix obtained by OASIS was almost symmetric, without ever enforcing it during training.

To quantify the extent to which  $\mathbf{W}$  is symmetric, we separated  $\mathbf{W}$  into

$$\mathbf{W} \equiv \text{sym}(\mathbf{W}) + \text{skew}(\mathbf{W}) \quad (16)$$

where  $\text{sym}(\mathbf{W}) = \frac{1}{2}(\mathbf{W} + \mathbf{W}^T)$  and  $\text{skew}(\mathbf{W}) = \frac{1}{2}(\mathbf{W} - \mathbf{W}^T)$ . The Frobenius norm obeys  $\|\mathbf{W}\|_{Fro} = \|\text{sym}(\mathbf{W})\|_{Fro} + \|\text{skew}(\mathbf{W})\|_{Fro}$ , which allows us to define a symmetry index

$$\rho(\mathbf{W}) = \frac{\|\text{sym}(\mathbf{W})\|}{\|\mathbf{W}\|}, \quad (17)$$

whose values range between 0 for an anti symmetric matrix and 1 for a symmetric one. At the beginning of training,  $\rho(\mathbf{W}) = 1$ ; It then decreased slowly up until convergence, with a

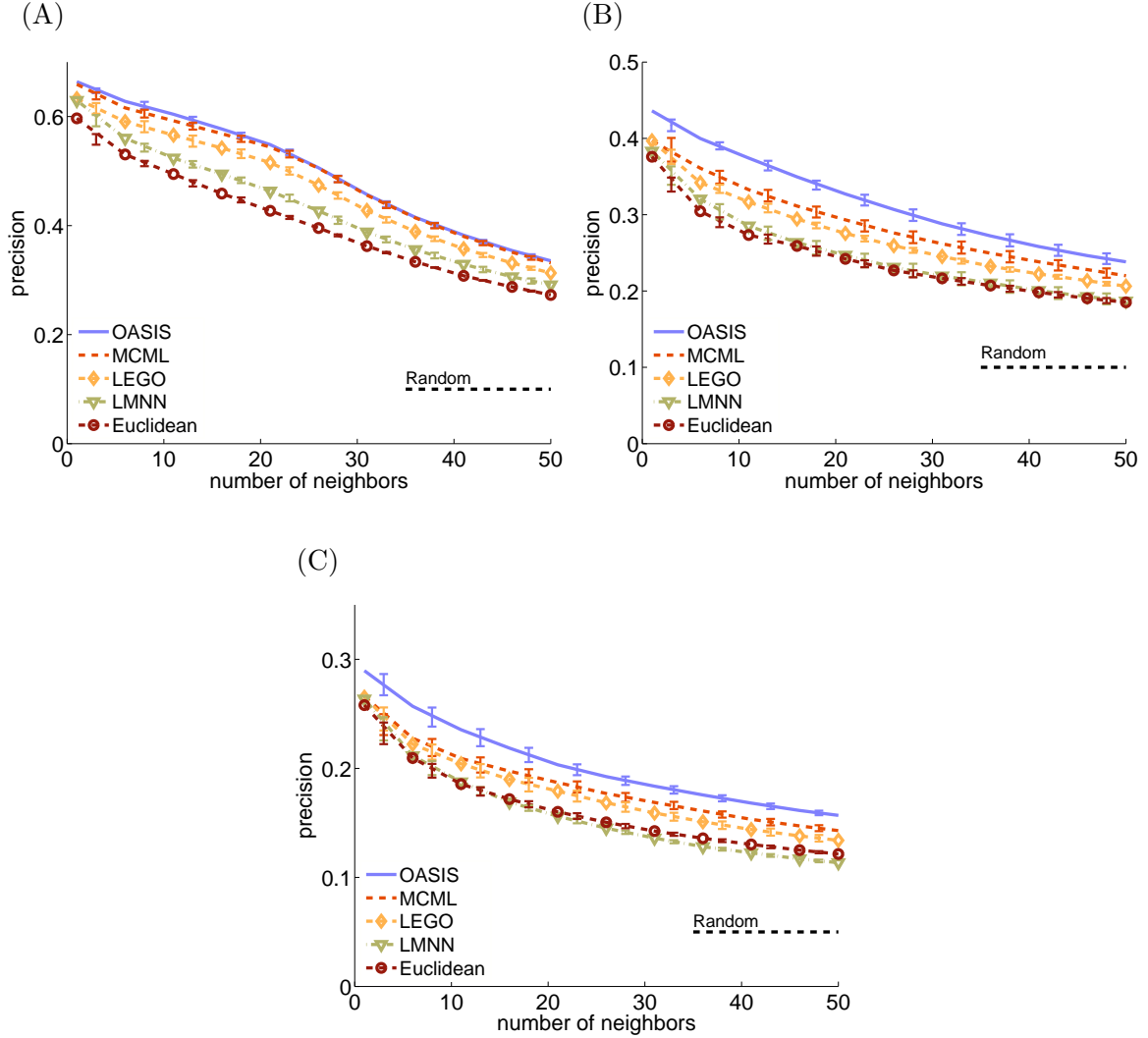


Figure 4: Comparison of the performance of OASIS, LMNN, MCML, LEGO and the Euclidean metric in feature space. Each curve shows the precision at top  $k$  as a function of  $k$  neighbors. The results are averaged across 5 train/test partitions (40 training images, 25 test images), error bars are standard error of the means (s.e.m.), black dashed line denotes chance performance. **(A)** Easy10. **(B)** Var10. **(C)** Var20.

value of  $\rho(\mathbf{W}) = 0.94$ . This suggests that even though we do not constrain the similarity matrix  $\mathbf{W}$  to be symmetric, the data keeps it to be near symmetric.

### 5.3 Web-Scale Experiment

Our second set of experiments is based on Google proprietary data and is around two orders of magnitude larger than the previous experiments. We collected a set of  $\sim 150\text{K}$

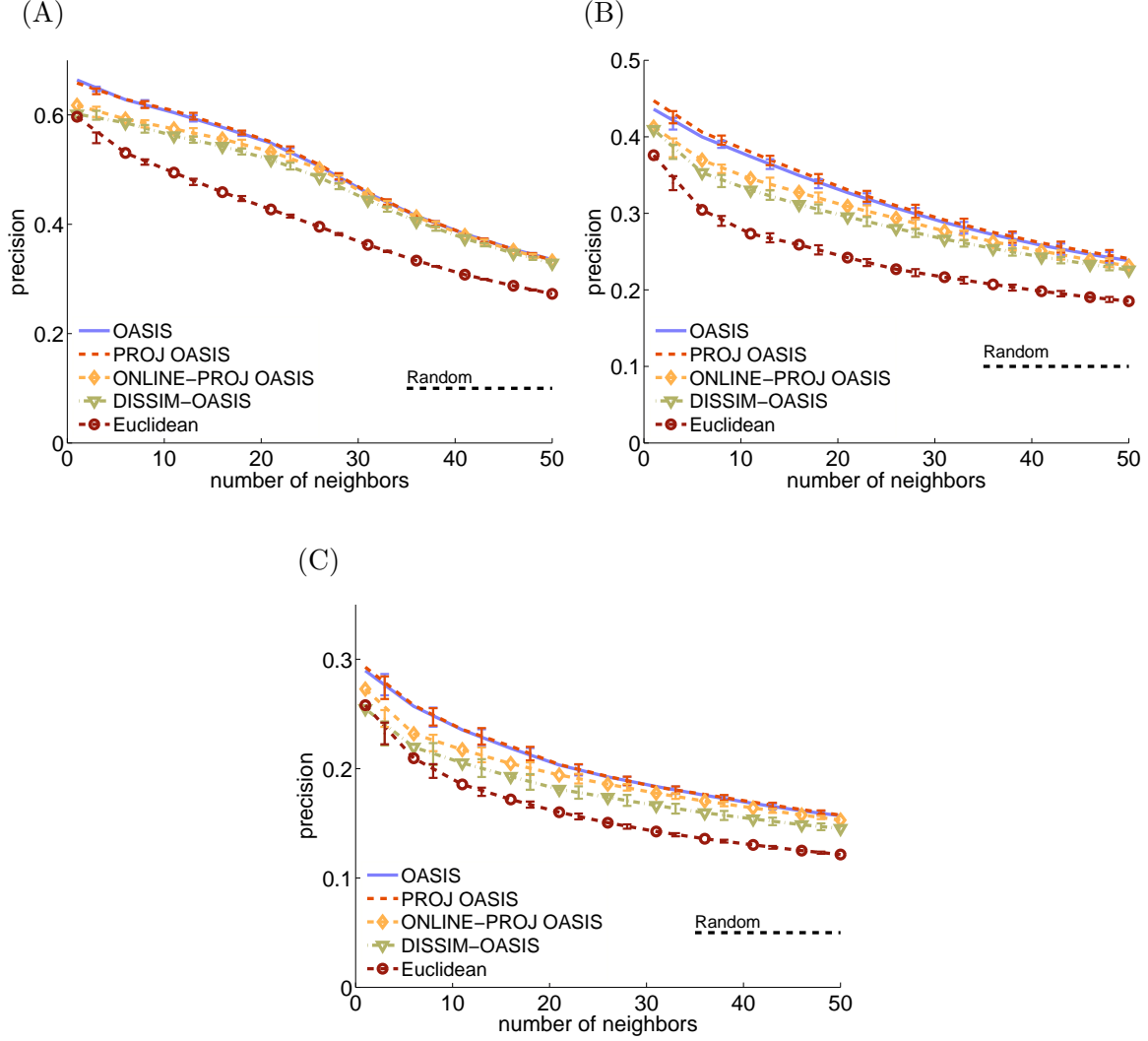


Figure 5: Comparison of the performance of Symmetric variants of OASIS. (A) Easy10. (B) Var10. (C) Var20.

text queries submitted to the Google Image Search system. For each of these queries, we had access to a set of relevant images, each of which associated with a numerical relevance score. This yielded a total of  $\sim 2.7$  million images, which we split into a training set of 2.3 million images and a test set of 0.4 million images (see Table 3).

### 5.3.1 EXPERIMENTAL SETUP

We used the query-image relevance information to create an image-image relevance as follows. Denote the set of text queries by  $\mathcal{Q}$  and the set of images by  $\mathcal{P}$ . For each  $q \in \mathcal{Q}$ , let  $\mathcal{P}_q^+$  denote the set of images that are relevant to the query  $q$ , and let  $\mathcal{P}_q^-$  denote the set of irrelevant images. The query-image relevance is defined by the matrix  $\mathbf{R}_{QI} : \mathcal{Q} \times \mathcal{P} \rightarrow \mathbb{R}^+$ ,



Table 3: Statistics of the Web dataset.

Set	Number of Queries	Number of Images
Training	139944	2292259
Test	41877	402164

and obeys  $\mathbf{R}_{QI}(q, p_q^+) > 0$  and  $\mathbf{R}_{QI}(q, p_q^-) = 0$  for all  $q \in \mathcal{Q}$ ,  $p_q^+ \in \mathcal{P}_q^+$ ,  $p_q^- \in \mathcal{P}_q^-$ . We also computed a normalized version of  $\mathbf{R}_{QI}$ , which can be interpreted as a joint distribution matrix, or the probability to observe a query  $q$  and an image  $p$  for that query,

$$Pr(q, p) = \frac{\mathbf{R}_{QI}(q, p)}{\sum_{q', p'} \mathbf{R}_{QI}(q', p')} \quad . \quad (18)$$

In order to compute the image-image relevance matrix  $\mathbf{R}_{II} : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}^+$ , we treated images as being conditionally independent given the queries,  $Pr(p_1, p_2 | q) = Pr(p_1 | q)Pr(p_2 | q)$ , and computed the joint image-image probability as a relevance measure

$$Pr(p_1, p_2) = \sum_{q \in \mathcal{Q}} Pr(p_1, p_2 | q) Pr(q) = \sum_{q \in \mathcal{Q}} Pr(p_1 | q) Pr(p_2 | q) Pr(q) \quad . \quad (19)$$

To improve scalability, we used a threshold over this joint distribution, and considered two images to be related only if their joint distribution exceeded a cutoff value  $\theta$

$$\mathbf{R}_{II}(p_1, p_2) = [Pr(p_1, p_2)]_\theta$$

where  $[x]_\theta = x$  for  $x > \theta$  and is zero otherwise. To set the value of  $\theta$  we have manually inspected a small subset of pairs of related images taken from the training set. We selected the largest  $\theta$  such that most of those related pairs had scores above the threshold, while minimizing noise in  $\mathbf{R}_{II}$ .

We trained OASIS over 2.3 million images in the training set using the sampling mechanism based on the relevance of each image, as described in Section 2.3.

To select the number of training iterations, we used a small subset of the training set to trace the precision of the model at some intervals as it changed throughout the training process, and stopped when its precision had saturated, which happened after 160 million iterations. Overall, training took a total of  $\sim 4000$  minutes on a single CPU of a standard modern machine. Finally, we evaluated the trained model on the 400 thousand images of the test set.

### 5.3.2 RESULTS

Table 4 shows the top five images as ranked by OASIS on four examples of query-images in the test set. The relevant text queries for each image are shown beneath the image. The first example (top row), shows a query-image that was originally retrieved in response to the text query “illusion”. All five images ranked highly by OASIS are semantically related, showing other types of visual illusions. Similar results can be observed for the three remaining examples on this table, where OASIS captures well the semantics of animal photos (cats and dogs), mountains and different food items.

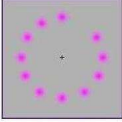


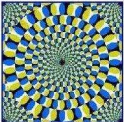
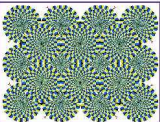
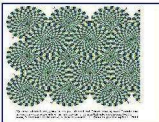


















Query image	Top 5 relevant images retrieved by OASIS				
 illusion, eye illusion, optical illusion	 illusion	 optical illusion	 trippy, trippy pictures, trip	 illusion, eye tricks	 circles, moving pictures
 scottish fold	 humor cat	 cubs tigers	 funny stuff, dog cartoon	 puppies	 agility
 swiss alps	 wedge, bodyboarding	 nighthawk	 china road silk	 winter landscape	 dogfight
 taco	 pizza	 bakery	 greek food	 panini, bread garlic, grill cheese	 food fish, fried fish

Table 4: OASIS: Successful cases from the Web dataset

In all these cases, OASIS captures similarity that is both semantic and visual, since the raw visual similarity of these images is not high.

On the other hand, Table 5 shows additional cases where OASIS was biased by visual similarity and provided high rankings to images that were semantically non relevant. In the first example, the assortment of flowers is confused with assortments of food items and a thigh section (5th nearest neighbor) which has visually similar shape. The second example presents a query image which in itself has no definite semantic element. The results retrieved are those that merely match texture of the query image and bear no semantic similarity. In the third example, OASIS fails to capture the butterfly in the query image.

To obtain a quantitative evaluation of OASIS we computed the precision at top  $k$ , in the same way as we did on the Caltech256 data. We used a threshold  $\theta = 0$ , which means that an image in the test set is considered relevant to a query image, if there exists at least one text query to which they were both relevant to.

Figure 6 shows the precision of top  $k$  as a function of  $k$  neighbors. The obtained precision values were drastically lower than those obtained for Caltech256. There are multiple possible reasons for this low precision. First, the number of unique textual queries in our data is very



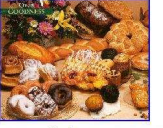


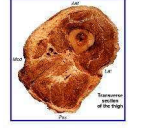












Query image	Top 5 relevant images retrieved by OASIS				
 roses bouquet	 dessert	 bakery	 panini, bread garlic, grill cheese	 newt	 thigh, muscle group
 garden vegetable	 schwitters	 canyon	 botswana	 canyon grand	 know
 insect	 flowers	 strawberry	 food japanese	 chinese food	 vegetable fruit, vitamin

Table 5: OASIS: Failure cases from the Web dataset

large (around 150K), hence the images in this dataset were significantly more heterogeneous than images in the Caltech256 data.

Second, and most importantly, our labels that measure pairwise relevance are very partial. This means that many pairs of images that are semantically related are not labeled as such. A clear demonstration of this effect is observed in Tables 4 and 5. The query images (like “*scottish fold*”) have labels that are usually very different from the labels of the retrieved images (as in “*humor cat*”, “*agility*”) even if their semantic content is very similar. This is a common problem in content-based analysis, since similar content can be described in many different ways. In the case discussed here, the partial data on the query-image relevance  $\mathbf{R}_{QI}$  is further propagated to the image-image relevance measure  $\mathbf{R}_{II}$ .

### 5.3.3 HUMAN EVALUATION EXPERIMENTS

In order to obtain a more accurate estimate of the real semantic precision, we performed a rating experiment with human evaluators. We chose the 25 most relevant images<sup>1</sup> from the test set and retrieved their 10 nearest neighbors as determined by OASIS. We excluded query-images which contained porn, racy or duplicates in their 10 nearest neighbors. We also selected randomly a set of 10 negative images  $p^-$  that were chosen for each of the query images  $p$  such that  $\mathbf{R}_{II}(p, p^-) = 0$ . These negatives were then randomly mixed with the 10 nearest neighbors.

1. The overall relevance of an image was estimated as the sum of relevances of the image with respect to all queries.

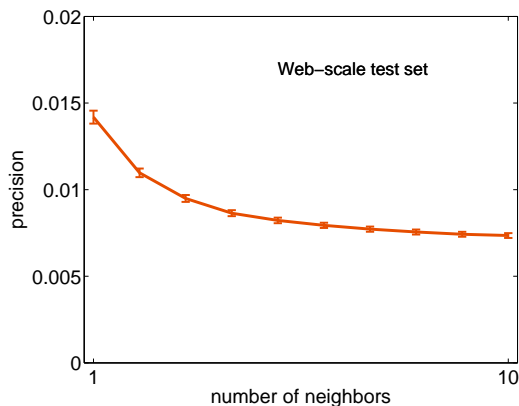


Figure 6: Precision at top  $k$  as a function of  $k$  neighbors computed against  $\mathbf{R}_{II}$  ( $\theta = 0$ ) for the web-scale test set.

All 25 query images were presented to twenty human evaluators, asking them to mark which of the 20 candidate images are *semantically relevant* to the query image<sup>2</sup>. Evaluators were volunteers selected from a pool of friends and colleagues, and many of which had experience with search or machine vision problems. We collected the ratings on the positive images and calculated the precision at top  $k$ .

Figure 7(A) shows the average precision across all queries and evaluators. Precision peaks at 42% and reaches 35% at the top 10 ranked image, being significantly higher than the values calculated automatically using  $\mathbf{R}_{II}$ .

We observed that the variability across different query images was also very high. Figure 7(B) shows the precision for 5 different queries, selected to span the range of average-precision values. The error bars at each curve show the variability in the responses of different evaluators. The precision of OASIS varies greatly across different queries. Some query images were “easy” for OASIS, yielding high scores from most evaluators. while other queries retrieved images that were consistently found to be irrelevant by most evaluators.

We also compared the magnitude of variability across human evaluators, with variability across queries. We first calculated the mAP from the precision curves of every query and evaluator, and then calculated the standard deviation in the mAP of every evaluator and of every query. The mean standard deviation over queries was 0.33, suggesting a large variability in the difficulty of image queries, as observed in Fig. 7(B). The mean standard deviation over evaluators was 0.25, suggesting that different evaluators had very different notions of what images should be regarded as “semantically similar” to a query image.

## 6. Discussion

We have presented OASIS, a scalable algorithm for learning image similarity that captures both semantic and visual aspects of image similarity. Three key factors contribute to the scalability of OASIS. First, using a large margin online approach allows training to converge

2. The description of the task as given to the evaluators is provided in Appendix A.

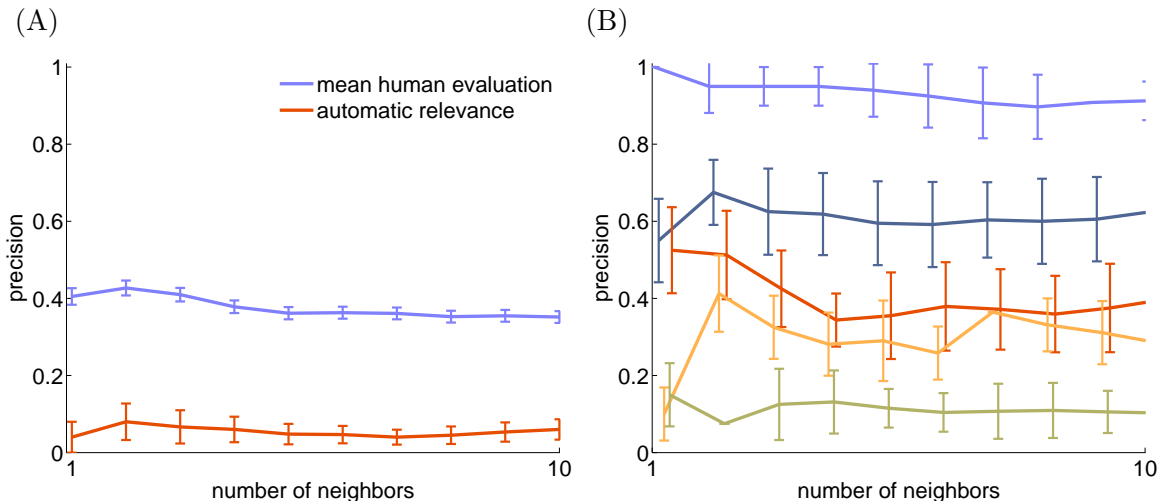


Figure 7: Precision at top  $k$  as a function of  $k$  neighbors for the human evaluation subset. **(A)** Mean precision across all 25 queries and 20 evaluators. Error bars denote the standard error of the mean. **(B)** Mean precision for 5 selected queries. Error bars denote the standard error of the mean. To select the queries for this plot, we first calculated the mean-average precision per query, sorted the queries by their mAP, and selected the queries ranked at position 1, 6, 11, 16, and 21.

even after seeing a small fraction of potential pairs. Second, the objective function of OASIS does not require the similarity measure to be necessarily a metric during training, although it appears to naturally converge to it. Finally, we use a sparse representation of low level features which allows to compute scores very efficiently.

We found that OASIS performs well in a wide range of scales: from problems with thousands of images, where it slightly outperforms existing metric-learning approaches, to large web-scale problems, where it achieves high accuracy, as estimated by human evaluators.

OASIS differs from previous methods in that the similarity measure that it learns is not forced to be a metric, or even symmetric. When the number of available samples is small, it is useful to add constraints that reflect prior knowledge on the type of similarity measure expected to be learned. However, we found that these constraints were not helpful even for problems with a few hundreds of samples. Interestingly, human judgments of pairwise similarity are known to be asymmetric, a property that can be easily captured by an OASIS model.

OASIS learns a class-independent model: it is not aware of which queries or categories were shared by two similar images. As such, it is more limited in its descriptive power and it is likely that class-dependent similarity models could improve precision. On the other hand, class-independent models could generalize to handle classes that were not observed during training, as in transfer learning. Large scale similarity learning, applied to images from a large variety of classes, could therefore be a useful tool to address real-world problems with a large number of classes.

## ACKNOWLEDGEMENTS

We thank Andrea Frome for very helpful discussions and comments on the manuscript. We thank Amir Globerson, Killian Weinberger and Prateek Jain, each providing an implementation of their method for our experiments.

## References

- A. Bar-Hillel, T. Hertz, N. Shental, and D. Weinshall. Learning Distance Functions using Equivalence Relations. In *Proc. of 20th International Conference on Machine Learning (ICML)*, page 11, 2003.
- Leon Bottou. Large-scale machine learning and stochastic algorithms. In *NIPS 2008 Workshop on Optimization for Machine Learning*, 2008.
- K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research (JMLR)*, 7:551–585, 2006.
- J.V. Davis, B. Kulis, P. Jain, S. Sra, and I.S. Dhillon. Information-theoretic metric learning. In *Proceedings of the 24th international conference on Machine learning*, pages 209–216. ACM Press New York, NY, USA, 2007.
- P. Duygulu, K. Barnard, N. de Freitas, and D. Forsyth. Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary. In *European Conference on Computer Vision (ECCV)*, pages 97–112, 2002.
- S.L. Feng, R. Manmatha, and V. Lavrenko. Multiple Bernoulli relevance models for image and video annotation. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2004.
- A. Frome, Y. Singer, F. Sha, and J. Malik. Learning globally-consistent local distance functions for shape-based image retrieval and classification. In *International Conference on Computer Vision*, pages 1–8, 2007.
- A. Globerson and S. Roweis. Metric Learning by Collapsing Classes. *Advances in Neural Information Processing Systems*, 18:451, 2006.
- D. Grangier and S. Bengio. A discriminative kernel-based model to rank images from text queries. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 30(8): 1371–1384, 2008.
- D. Grangier, Florent Monay, and S. Bengio. Learning to retrieve images from text queries with a discriminative model. In *International Conference on Adaptive Multimedia Retrieval (AMR)*, 2006.
- G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. Technical Report 7694, California Institute of Technology, 2007. URL <http://authors.library.caltech.edu/7694>.



- P. Jain, B. Kulis, I. Dhillon, and K. Grauman. Online metric learning and fast similarity search. In *Advances in Neural Information Processing Systems*, volume 22, 2008.
- J. Jeon and R. Manmatha. Using maximum entropy for automatic image annotation. In *International Conference on Image and Video Retrieval*, pages 24–32, 2004.
- G.R.G. Lanckriet, N. Cristianini, P. Bartlett, L. El Ghaoui, and M.I. Jordan. Learning the Kernel Matrix with Semidefinite Programming. *Journal of Machine Learning Research (JMLR)*, 5:27–72, 2004.
- D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision (IJCV)*, 60(2):91–110, 2004. URL <http://dx.doi.org/10.1023/B:VISI.0000029664.99615.94>.
- William Stafford Noble. Multi-kernel learning for biology. In *NIPS 2008 workshop on kernel learning*, 2008.
- T. Ojala, M. Pietikainen, and T. Maenpaa. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 24(7):971–987, 2002.
- P. Quelhas, F. Monay, J. M. Odobez, D. Gatica-Perez, T. Tuytelaars, and L. J. Van Gool. Modeling scenes with local descriptors and latent aspects. In *International Conference on Computer Vision*, pages 883–890, 2005.
- N. Rasiwasia and N. Vasconcelos. A study of query by semantic example. In *3rd International Workshop on Semantic Learning and Applications in Multimedia*, 2008.
- M. Schultz and T. Joachims. Learning a Distance Metric from Relative Comparisons. In *Advances in Neural Information Processing Systems 16: Proceedings of the 2003 Conference*. Bradford Book, 2004.
- V. Takala, T. Ahonen, and M. Pietikainen. Block-based methods for image retrieval using local binary patterns. In *Scandinavian Conference on Image Analysis (SCIA)*, 2005.
- K. Tieu and P. Viola. Boosting image retrieval. *International Journal of Computer Vision (IJCV)*, 56(1):17 – 36, 2004.
- A. Torralba, R. Fergus, and W. T. Freeman. Tiny images. Technical Report MIT-CSAIL-TR-2007-024, Computer Science and Artificial Intelligence Lab, Massachusetts Institute of Technology, 2007. URL <http://dspace.mit.edu/handle/1721.1/37291>.
- K. Weinberger, J. Blitzer, and L. Saul. Distance Metric Learning for Large Margin Nearest Neighbor Classification. *Advances in Neural Information Processing Systems*, 18:1473, 2006.
- E.P. Xing, A.Y. Ng, M.I. Jordan, and S. Russell. Distance Metric Learning with Application to Clustering with Side-Information. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15*, pages 521–528, Cambridge, MA, 2003. MIT Press.

Liu Yang. Distance metric learning: A comprehensive survey. Technical report, Michigan State University, 2006.



## Appendix A. Human Evaluation

The following text was provided to human evaluators when judging the relevance of images to a query image.

Scenario:

A user is searching images to use in a presentation he/she plans to give. The user runs a standard image search, and selects an image, the ‘‘query image’’. The user then wishes to refine the search and look for images that are SEMANTICALLY similar to the query image.

The difficulty lies, in the definition of ‘‘SEMANTICALLY’’. This can have many interpretations, and you should take that into account.

So for instance, if you see an image of a big red truck, you can interpret the user intent (the notion of semantically similar) in various ways:

- any big red truck
- any red truck
- any big truck
- any truck
- any vehicle

You should interpret ‘‘SEMANTICALLY’’ in a broad sense rather than in a strict sense but feel free to draw the line yourself (although be consistent).

Your task:

You will see a set of query images on the left side of the screen, and a set of potential candidate matches, 5 per row, on the right. Your job is to decide for each of the candidate images if it is a good semantic match to the query image or not. The default is that it is NOT a good match. Furthermore, if for some reason you cannot make-up your mind, then answer ‘‘can’t say’’.

## Appendix B. Caltech256 Class Sets

- **Easy10:** *car-side-101, faces-easy-101, zebra, tower-pisa, watch-101, sunflower-101, mars, desk-globe, sheet-music, trilobite-101.*
- **Var10:** *bear, skyscraper, billiards, yo-yo, minotaur, roulette-wheel, hamburger, laptop-101, hummingbird, blimp.*
- **Var20:** *airplanes-101, mars, homer-simpson, hourglass, waterfall, helicopter-101, mountain-bike starfish-101, teapot, pyramid, refrigerator, cowboy-hat, giraffe, joy-stick, crab-101, birdbath, fighter-jet tuning-fork, iguana, dog.*