

80 million tiny images: a large dataset for non-parametric object and scene recognition

Antonio Torralba, Rob Fergus and William T. Freeman

Abstract—With the advent of the Internet, billions of images are now freely available online and constitute a dense sampling of the visual world. Using a variety of non-parametric methods, we explore this world with the aid of a large dataset of 79,302,017 images collected from the Web. Motivated by psychophysical results showing the remarkable tolerance of the human visual system to degradations in image resolution, the images in the dataset are stored as 32×32 color images. Each image is loosely labeled with one of the 75,062 non-abstract nouns in English, as listed in the Wordnet lexical database. Hence the image database gives a comprehensive coverage of all object categories and scenes. The semantic information from Wordnet can be used in conjunction with nearest-neighbor methods to perform object classification over a range of semantic levels minimizing the effects of labeling noise. For certain classes that are particularly prevalent in the dataset, such as people, we are able to demonstrate a recognition performance comparable to class-specific Viola-Jones style detectors.

Index Terms—Object recognition, tiny images, large datasets, Internet images, nearest-neighbor methods.

I. INTRODUCTION

With overwhelming amounts of data, many problems can be solved without the need for sophisticated algorithms. One example in the textual domain is Google’s “Did you mean?” tool which corrects errors in search queries, not through a complex parsing of the query but by memorizing billions of query-answer pairs and suggesting the one closest to the users query. In this paper, we explore a visual analog to this tool by using a large dataset of 79 million images and nearest-neighbor matching schemes.

When very many images are available, simple image indexing techniques can be used to retrieve images with similar object arrangements to the query image. If we have a big enough database then we can find, with high probability, images visually close to a query image, containing similar scenes with similar objects arranged in similar spatial configurations. If the images in the retrieval set are partially labeled, then we can propagate the labels to the query image, so performing classification.

Nearest-neighbor methods have been used in a variety of computer vision problems, primarily for interest point matching [5], [19], [28]. They have also been used for global image matching (e.g. estimation of human pose [36]), character recognition [4], and object recognition [5], [34]. A number of recent papers have used large datasets of images in conjunction with purely non-parametric methods for computer vision and graphics applications [22], [39].

Finding images within large collections is the focus of the content based image retrieval (CBIR) community. Their emphasis

on really large datasets means that the chosen image representation is often relatively simple, e.g. color [17], wavelets [42] or crude segmentations [9]. This enables very fast retrieval of images similar to the query, for example the Cortina system [33] demonstrates real-time retrieval from a 10 million image collection, using a combination of texture and edge histogram features. See Datta et al. for a survey of such methods [12].

The key question that we address in this paper is: How big does the image dataset need to be to robustly perform recognition using simple nearest-neighbor schemes? In fact, it is unclear that the size of the dataset required is at all practical since there are an effectively infinite number of possible images the visual system can be confronted with. What gives us hope is that the visual world is very regular in that real world pictures occupy only a relatively small portion of the space of possible images.

Studying the space occupied by natural images is hard due to the high dimensionality of the images. One way of simplifying this task is by lowering the resolution of the images. When we look at the images in Fig. 6, we can recognize the scene and its constituent objects. Interestingly though, these pictures have only 32×32 color pixels (the entire image is just a vector of 3072 dimensions with 8 bits per dimension), yet at this resolution, the images already seem to contain most of the relevant information needed to support reliable recognition.

An important benefit of working with tiny images is that it becomes practical to store and manipulate datasets orders of magnitude bigger than those typically used in computer vision. Correspondingly, we introduce, and make available to researchers, a dataset of 79 million unique 32×32 color images gathered from the Internet. Each image is loosely labeled with one of 75,062 English nouns, so the dataset covers a very large number of visual object classes. This is in contrast to existing datasets which provide a sparse selection of object classes. In this paper we will study the impact on having very large datasets in combination with simple techniques for recognizing several common object and scene classes at different levels of categorization.

The paper is divided in three parts. In Section 2 we establish the minimal resolution required for scene and object recognition. In Sections 3 and 4 we introduce our dataset of 79 million images and explore some of its properties. In Section 5 we attempt scene and object recognition using a variety of nearest-neighbor methods. We measure performance at a number of semantic levels, obtaining impressive results for certain object classes.

II. LOW DIMENSIONAL IMAGE REPRESENTATIONS

A number of approaches exist for computing the *gist* of a image, a global low-dimensional representation that captures the scene and its constituent objects [18], [32], [24]. We show that very low-resolution 32×32 color images can be used in this role, containing enough information for scene recognition, object

The authors are with the Computer Science and Artificial Intelligence Lab (CSAIL) at the Massachusetts Institute of Technology.

Email: {torralba,fergus,billf}@csail.mit.edu

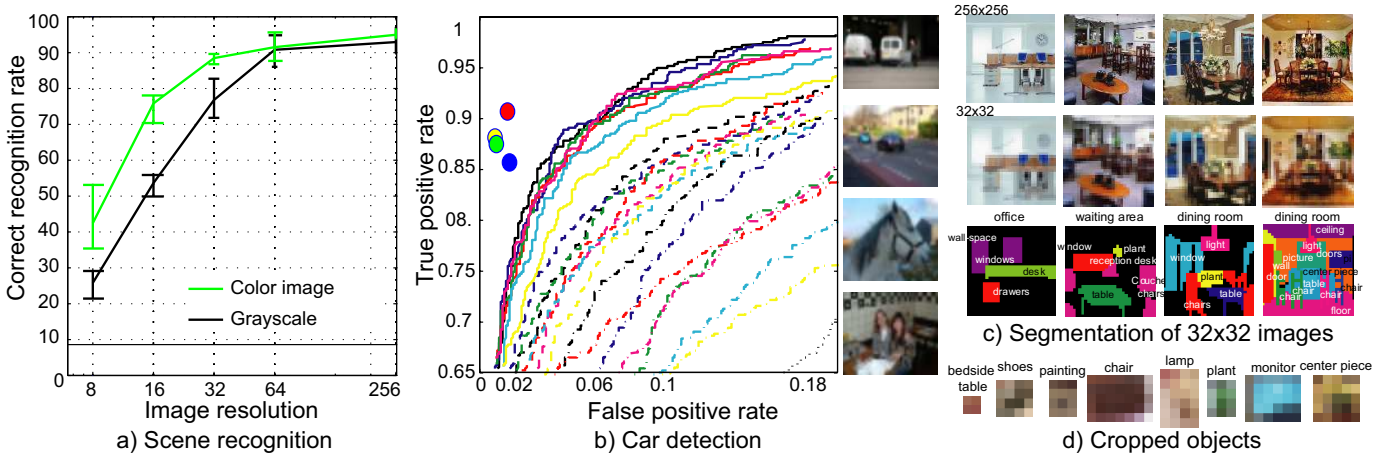


Fig. 1. a) Human performance on scene recognition as a function of resolution. The green and black curves show the performance on color and gray-scale images respectively. For color 32×32 images the performance only drops by 7% relative to full resolution, despite having 1/64th of the pixels. b) Car detection task on the PASCAL 2006 test dataset. The colored dots show the performance of four human subjects classifying tiny versions of the test data. The ROC curves of the best vision algorithms (running on full resolution images) are shown for comparison. All lie below the performance of humans on the tiny images, which rely on none of the high-resolution cues exploited by the computer vision algorithms. c) Humans can correctly recognize and segment objects at very low resolutions, even when the objects in isolation can not be recognized (d).

detection and segmentation (even when the objects occupy just a few pixels in the image).

A. Scene recognition

Studies on face perception [1], [21] have shown that only 16×16 pixels are needed for robust face recognition. This remarkable performance is also found in a scene recognition task [31].

We evaluate the scene recognition performance of humans as the image resolution is decreased. We used a dataset of 15 scenes that was taken from [14], [24], [32]. Each image was shown at one of 5 possible resolutions (8^2 , 16^2 , 32^2 , 64^2 and 256^2 pixels) and the participant task was to assign the low-resolution picture to one of the 15 different scene categories (bedroom, suburban, industrial, kitchen, living room, coast, forest, highway, inside city, mountain, open country, street, tall buildings, office, and store)¹. Fig. 1(a) shows human performance on this task when presented with grayscale and color images² of varying resolution. For grayscale images, humans need around 64×64 pixels. When the images are in color, humans need only 32×32 pixels to achieve more than 80% recognition rate. Below this resolution the performance rapidly decreases. Therefore, humans need around 3000 dimensions of either color or grayscale data to perform this task. In the next section we show that 32×32 color images also preserve a great amount of local information and that many objects can still be recognized even when they occupy just a few pixels.

¹Experimental details: 6 participants classified 585 color images as belonging to one of the 15 scene categories from [14], [24], [32]. Images were presented at 5 possible resolutions (8^2 , 16^2 , 32^2 , 64^2 and 256^2). Each image was shown at 5 possible sizes using bicubic interpolation to reduce pixelation effects which impair recognition. Interpolation was applied to the low-resolution image with 8 bits per pixel and color channel. Images were not repeated across conditions. 6 additional participants performed the same experiment but with gray scale images.

²100% recognition rate can not be achieved in this dataset as there is no perfect separation between the 15 categories.

B. Object recognition

Recently, the PASCAL object recognition challenge evaluated a large number of algorithms in a detection task for several object categories [13]. Fig. 1(b) shows the performances (ROC curves) of the best performing algorithms in the car classification task (i.e. is there a car present in the image?). These algorithms require access to relatively high resolution images. We studied the ability of human participants to perform the same detection task but using very low-resolution images. Human participants were shown color images from the test set scaled to have 32 pixels on the smallest axis, preserving their aspect ratio. Fig. 1(b) shows some examples of tiny PASCAL images. Each participant classified between 200 and 400 images selected randomly. Fig. 1(b) shows the performances of four human observers that participated in the experiment. Although around 10% of cars are missed, the performance is still very good, significantly outperforming the computer vision algorithms using full resolution images. This shows that even though the images are very small, they contain sufficient information for accurate recognition.

Fig. 1(c) shows some representative 32^2 images segmented by human subjects. Despite the low resolution, sufficient information remains for reliable segmentation (more than 80% of the segmented objects are correctly recognized), although any further decrease in resolution dramatically affects segmentation performance. Fig. 1(d) shows crops of some of the smallest objects correctly recognized when shown within the scene. Note that in isolation, the objects cannot be identified since the resolution is so low, hence the recognition of these objects within the scene is almost entirely based on context.

Clearly, not all visual tasks can be solved using such low resolution images. But the experiments in this section suggest that 32×32 color images are the minimum viable size for recognition tasks – the focus of the paper.

III. A LARGE DATASET OF 32×32 IMAGES

As discussed in the previous sections, 32×32 color images contain the information needed to perform a number of challenging

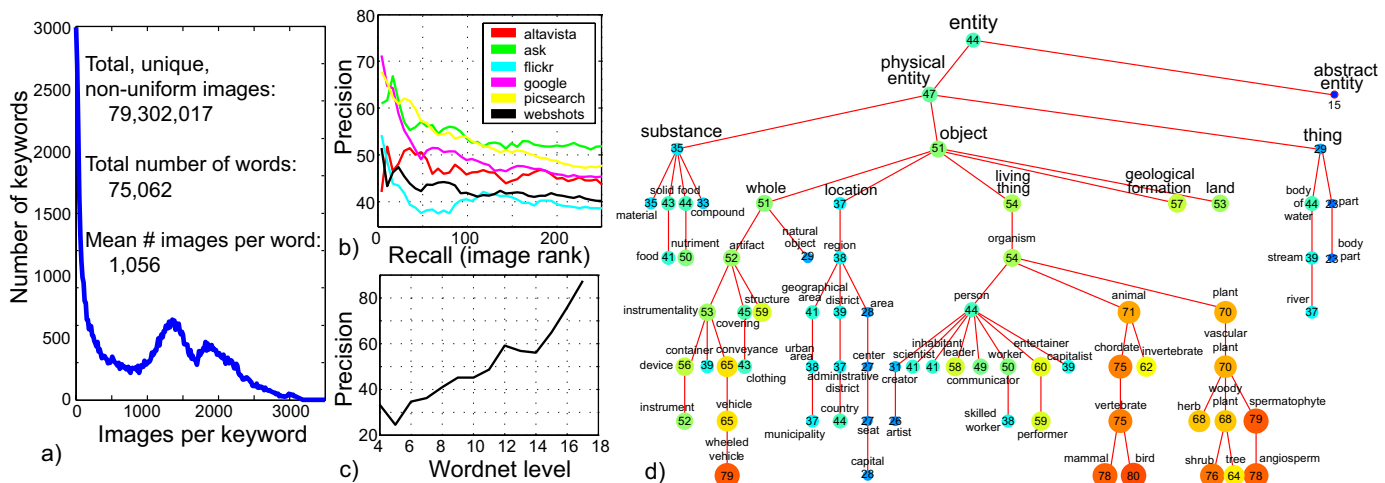


Fig. 2. Statistics of our database of tiny images. a) A histogram of images per keyword collected. Around 10% of keywords have very few images. b) Performance of the search various engines (evaluated on hand-labeled ground truth). c) Accuracy of the labels attached at each image as a function of the depth in the Wordnet tree (deeper corresponds to more specific words). d) Accuracy of labeling for different nodes of a portion of the Wordnet tree.

recognition tasks. One important advantage of very low resolution images is that it becomes practical to work with millions of images. In this section we will describe a dataset of 10^8 tiny images.

Current experiments in object recognition typically use 10^2 - 10^4 images spread over a few different classes; the largest available dataset being one with 256 classes[20]. Other fields, such as speech, routinely use 10^6 data points for training, since they have found that large training sets are vital for achieving low errors rates in testing [2]. As the visual world is far more complex than the aural one, it would seem natural to use very large set of training images. Motivated by this, and the ability of humans to recognize objects and scenes in 32×32 images, we have collected a database of nearly 10^8 such images.

A. Collection procedure

We use Wordnet [15] likely to have any kind of visual consistency. We do this by extracting all non-abstract nouns from the database, 75,062 of them in total. In contrast to existing object recognition datasets which use a sparse selection of classes, by collecting images for all nouns, we have a dense coverage of all visual forms.

We selected 7 independent image search engines: Altavista, Ask, Flickr, Cydral, Google, Picsearch and Webshots (others have outputs correlated with these). We automatically download all the images provided by each engine for all 75,846 non-abstract nouns. Running over 8 months, this method gathered 97,245,098 images in total. Once intra-word duplicates and uniform images (images with zero variance) are removed, this number is reduced to 79,302,017 images from 75,062 words (around 1% of the keywords had no images). Storing this number of images at full resolution is impractical on the standard hardware used in our experiments so we down-sampled the images to 32×32 as they

were gathered³. The dataset fits onto a single hard disk, occupying 760Gb in total. The dataset may be downloaded from <http://people.csail.mit.edu/torralba/tinyimages>.

Fig. 2(a) shows a histogram of the number of images per class. Around 10% of the query words are obscure so no images can be found on the Internet, but for the majority of words a reasonable number of images are found. We place an upper limit of 3000 images/word to keep the total collection time to a reasonable level. Although the gathered dataset is very large, it is not necessarily representative of all natural images. Images on the Internet have their own biases (e.g. objects tend to be centered and fairly large in the image). However, web images define an interesting visual world for developing computer vision applications [16], [37].

B. Characterization of labeling noise

Despite a number of recent efforts for image annotation [35], [43], collecting images from the web provides a powerful mechanism to build large image databases orders of magnitude larger than is possible with manual methods. However, the images gathered by the engines are loosely labeled in that the visual content is often unrelated to the query word (for example, see Fig. 10). In this section we characterize the noise present in the labels. Among other factors, the accuracy of the labels depend on the engine used, and the specificity of the term used for querying.

In Fig. 2(b) we quantify the labeling noise using 3526 hand-labeled images selected by randomly sampling images out of the first 250 images returned by each online search engine for each word. A recall-precision curve is plotted for each search engine in which the horizontal axis represents the rank in which the image was returned and the vertical axis is the percentage of images that corresponded to the query. Accuracy drops after the 100th image and then stabilizes at around 44% correct on average.

³Further comments: (i) Wordnet is a lexical dictionary, meaning that it gives the semantic relations between words in addition to the information usually given in a dictionary.; (ii) The tiny database is not just about objects. It is about everything that can be indexed with Wordnet and this includes scene-level classes such as streets, beaches, mountains, as well as category-level classes and more specific objects such as US Presidents, astronomical objects and Abyssinian cats.; (iii) At present we do not remove inter-word duplicates since identifying them in our dataset is non-trivial.

The accuracy of online searchers also varies depending on which terms were used for the query. Fig. 2(c) shows that the noise varies for different levels of the Wordnet tree, being more accurate when getting close to the leaves of the tree. Fig. 2(d) shows a subset of the Wordnet tree used to build our dataset (the full tree contains >40,000 leaves). The number and color at each node corresponds to the percentage of images correctly assigned to the leaves of each node. The more specific are the terms, the more likely are the images to correspond to the query.

Various methods exist for cleaning up the data by removing images visually unrelated to the query word. Berg and Forsyth [7] have shown a variety of effective methods for doing this with images of animals gathered from the web. Berg et al. [6] showed how text and visual cues could be used to cluster faces of people from cluttered news feeds. Fergus et al. [16] have shown the use of a variety of approaches for improving Internet image search engines. Li et al. [26] show further approaches to decreasing label noise. However, due to the extreme size of our dataset, it is not practical to employ these methods. In Section 5, we show that reasonable recognition performances can be achieved despite the high labeling noise.

IV. STATISTICS OF VERY LOW RESOLUTION IMAGES

Despite 32×32 being very low resolution, each image lives in a space of 3072 dimensions. This is a very large space — if each dimension has 8 bits, there are a total of 10^{7400} possible images. This is a huge number, especially if we consider that a human in a 100 years only gets to see 10^{11} frames (at 30 frames/second). However, natural images only correspond to a tiny fraction of this space (most of the images correspond to white noise), and it is natural to investigate the size of that fraction. A number of studies [10], [25] have been devoted to characterize the space of natural images by studying the statistics of small image patches. However, low-resolution scenes are quite different to patches extracted by randomly cropping small patches from images.

Given a similarity measure, the question that we want to answer is: *how many images are needed so that, for any given query image, we can always find a neighbor with the same class label?* Note that we are concerned solely with recognition performance, not with issues of intrinsic dimensionality or the like as explored in other studies of large collection of image patches [10], [25]. In this section, we explore how the probability of finding images with a similar label nearby increases with the size of the dataset. In turn, this tells us how big the dataset needs to be to give a robust recognition performance.

A. Distribution of neighbors as a function of dataset size

As a first step, we use the sum of squared differences (SSD) to compare two images. We will define later other similarity measures that incorporate invariances to translations and scaling. The SSD between two images I_1 and I_2 (normalized to have zero mean and unit norm)⁴ is:

$$D_{\text{ssd}}^2 = \sum_{x,y,c} (I_1(x,y,c) - I_2(x,y,c))^2 \quad (1)$$

Computing similarities among 7.9×10^7 images is computationally expensive. To improve speed, we index the images using

⁴Normalization of each image is performed by transforming the image into a vector concatenating the three color channels. The normalization does not change image color, only the overall luminance.

the first 19 principal components of the 7.9×10^7 images (19 is the maximum number of components per image such that the entire index structure can be held in memory). The $1/f^2$ property of the power spectrum of natural images means that the distance between two images can be approximated using few principal components (alternative representations using wavelets [42] could also be used in place of the PCA representation). We compute the approximate distance $\hat{D}_{\text{ssd}}^2 = 2 - 2 \sum_{n=1}^C v_1(n)v_2(n)$, where $v_i(n)$ is the n^{th} principal component coefficient for the i^{th} image (normalized so that $\sum_n v_i(n)^2 = 1$), and C is the number of components used to approximate the distance. We define S_N as the set of N exact nearest neighbors and \hat{S}_M as the set of M approximate nearest neighbors.

Fig. 3(a) shows the probability that an image, of index i , from the set S_N is also inside \hat{S}_M : $P(i \in \hat{S}_M | i \in S_N)$. The plot corresponds to $N = 50$. For the experiments in this section, we used 200 images randomly sampled from the datasets and for which we computed the exact distances to all the 7.9×10^7 images. Many images on the web appear multiple times. For the plots in these figures, we have removed manually all the image pairs that were duplicates.

Fig. 3(b) shows the number of approximate neighbors (M) that need to be considered as a function of the desired number of exact neighbors (N) in order to have a probability of 0.8 of finding N exact neighbors. As the dataset becomes larger, we need to collect more approximate nearest neighbors in order to have the same probability of including the first N exact neighbors.

For the experiments in this paper, we use the following procedure. First, we find the closest 16,000 images per image. From Fig. 3(a) we know that more than 80% of the exact neighbors will be part of this approximate neighbor set. Then, within the set of 16,000 images, we compute the exact distances to provide the final rankings of neighbors. Exhaustive search, used in all our experiments, takes 30 seconds per image using the principle components method. This can be dramatically improved through the use of a kd-tree to 0.3 seconds per query, if fast retrieval performance is needed. The memory overhead of the kd-tree means that only 17 of the 19 PCA components can be used. Devising efficient indexing methods for large image databases [30], [19], [40] is a very important topic of active research but it is not the focus of this paper.

Fig. 4 shows several plots measuring various properties as the size of the dataset is increased. The plots use the normalized correlation ρ between images (note that $D_{\text{ssd}}^2 = 2(1 - \rho)$). In Fig. 4(a), we show the probability that the nearest neighbor has a normalized correlation exceeding a certain value. Each curve corresponds to a different dataset size. Fig. 4(b) shows a vertical section through Fig. 4(a) at the correlations 0.8 and 0.9, plotting the probability of finding a neighbor as the number of images in the dataset grows. From Fig. 4(b) we see that a third of the images in the dataset are expected to have a neighbor with correlation > 0.8 .

In Fig. 4(c) we explore how the plots shown in Fig. 4(a) & (b) relate to recognition performance. Three human subjects labeled pairs of images as belonging to the same visual class or not (pairs of images that correspond to duplicate images are removed). The plot shows the probability that two images are labeled as belonging to the same class as a function of image similarity. As the normalized correlation exceeds 0.8, the probability of belonging to the same class grows rapidly. Hence a simple K-

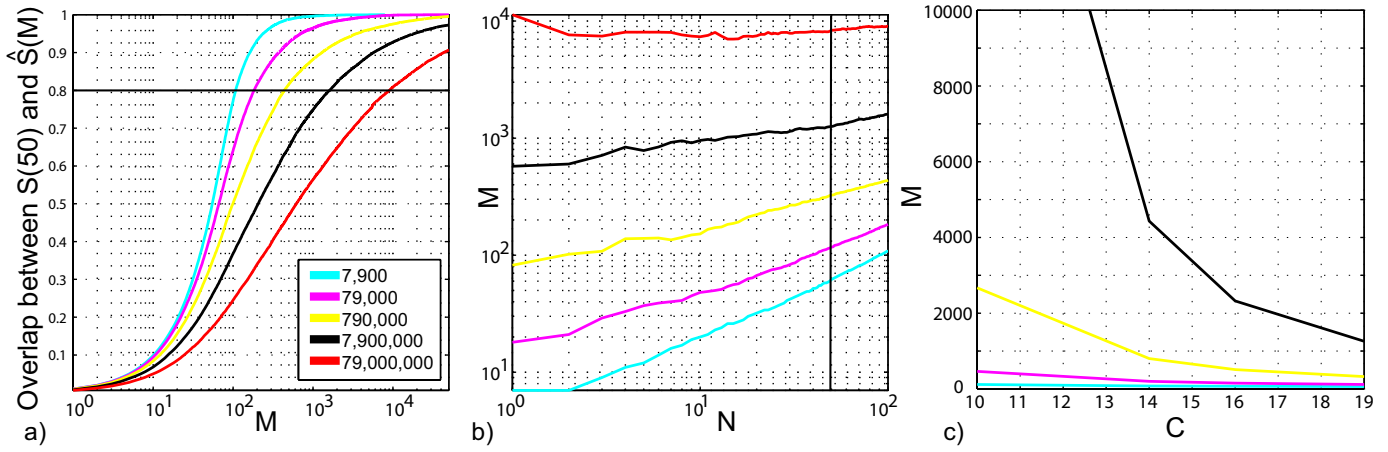


Fig. 3. Evaluation of the method for computing approximate nearest neighbors. These curves correspond to the similarity measure D_{ssd} . (a) Probability that an image from the set of exact nearest neighbors S_N , with $N = 50$, is inside the approximate set of nearest neighbors S_M as a function of M . (b) Number of approximate neighbors (M) that need to be considered as a function of the desired number of exact neighbors (N) in order to have a probability of 0.8 of finding N exact neighbors. Each graph corresponds to a different dataset size, indicated by the color code. (c) Number of approximate neighbors (M) that need to be considered as we reduce the number of principal components (C) used for the indexing (with $N = 50$).

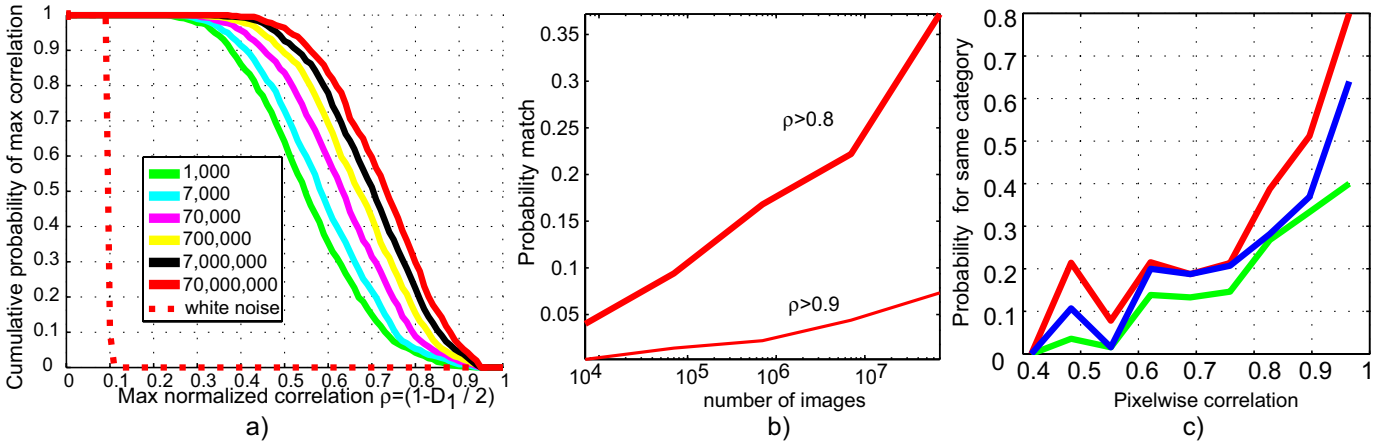


Fig. 4. Exploring the dataset using D_{ssd} . (a) Cumulative probability that the nearest neighbor has a correlation greater than ρ . Each of the colored curves shows the behavior for a different size of dataset. (b) Cross-section of figure (a) plots the probability of finding a neighbor with correlation > 0.9 as a function of dataset size. (c) Probability that two images belong to the same category as a function of pixel-wise correlation (duplicate images are removed). Each curve represents a different human labeler.

nearest-neighbor approach might be effective with our size of dataset. We will explore this further in Section V.

B. Image similarity metrics

We can improve recognition performance using better measures of image similarity. We now introduce two additional similarity measures between a pair of normalized images I_1 and I_2 , that incorporate invariances to simple spatial transformations.

- In order to incorporate invariance to small translations, scaling and image mirror, we define the similarity measure:

$$D_{\text{warp}}^2 = \min_{\theta} \sum_{x,y,c} (I_1(x,y,c) - T_{\theta}[I_2(x,y,c)])^2 \quad (2)$$

In this expression, we minimize the similarity by transforming I_2 (horizontal mirror; translations and scaling up to 10 pixel shifts) to give the minimum SSD. The transformation parameters θ are optimized by gradient descent [29].

- We allow for additional distortion in the images by shifting every pixel individually within a 5 by 5 window to give minimum SSD. This registration can be performed with more complex representations than pixels (e.g., Berg and Malik [5]). In our case, the minimum can be found by exhaustive evaluation of all shifts, only possible due to the low resolution of the images.

$$D_{\text{shift}}^2 = \min_{|D_{x,y}| \leq w} \sum_{x,y,c} (I_1(x,y,c) - \hat{I}_2(x + D_x, y + D_y, c))^2 \quad (3)$$

In order to get better matches, we initialize I_2 with the warping parameters obtained after optimization of D_{warp} , $\hat{I}_2 = T_{\theta}[I_2]$.

Fig. 5 shows a pair of images being matched using the 3 metrics and shows the resulting neighbor images transformed by the optimal parameters that minimize each similarity measure. The figure shows two candidate neighbors: one matching the target

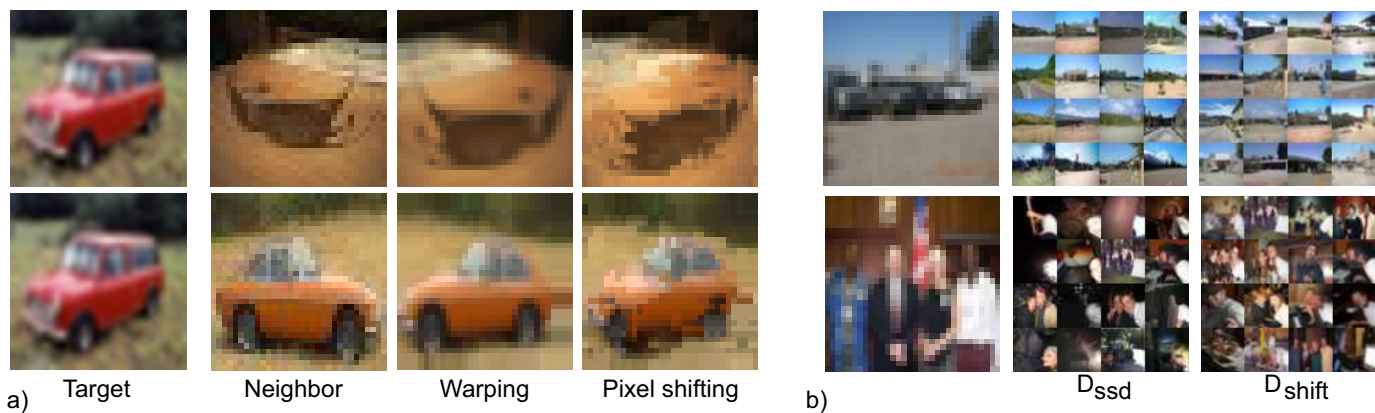


Fig. 5. a) Image matching using distance metrics D_{ssd} , D_{warp} and D_{shift} . Top row: after transforming each neighbor by the optimal transformation; the sunglasses always results in a poor match. However, for the car example on the bottom row, the matched image approximates the pose of the target car. b) Sibling sets from 79,302,017 images, found with distance metrics D_{ssd} , and D_{shift} . D_{shift} provides better matches than D_{ssd} .



Fig. 6. As we increase the size of the dataset from 10^5 to the 10^8 images, the quality of the retrieved set increases dramatically. However, note that we need to increase the size of the dataset logarithmically in order to have an effect. These results are obtained using D_{shift} as a similarity measure between images.

semantic category and another one that corresponds to a wrong match. For D_{warp} and D_{shift} we show the closest manipulated image to the target. D_{warp} looks for the best translation, scaling and horizontal mirror of the candidate neighbor in order to match the target. D_{shift} further optimizes the warping provided by D_{warp} by allowing pixels to move in order to minimize the distance with the target.

Fig. 5(b) shows two examples of query images and the retrieved neighbors (*sibling set*), out of 79,302,017 images, using D_{ssd} and D_{shift} . For speed we use the same low dimensional approximation as described in the previous section by evaluating D_{warp} and D_{shift} only on the first 16,000 candidates. This is a good indexing scheme for D_{warp} , but it results in slightly decrease of performance for D_{shift} which would require more neighbors to be considered. Despite this, both measures provide good matches, but D_{shift}

returns closer images at the semantic level. This observation will be quantified in Section V. Fig. 6 shows examples of query images and sets of neighboring images, from our dataset of 79,302,017 images, found using D_{shift} .

V. RECOGNITION

A. Wordnet voting scheme

We now attempt to use our dataset for object and scene recognition. While an existing computer vision algorithm could be adapted to work on 32×32 images, we prefer to use a simple nearest-neighbor scheme based on one of the distance metrics D_{ssd} , D_{warp} or D_{shift} . Instead of relying on the complexity of the matching scheme, we let the data to do the work for us: the hope is that there will always be images close to a given

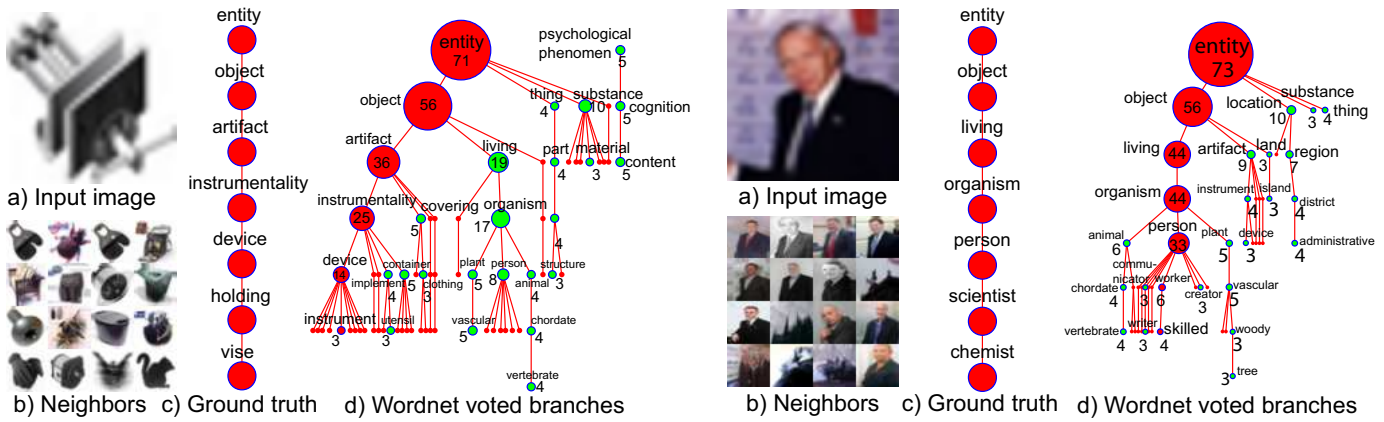


Fig. 7. This figure shows two examples. (a) Query image. (b) First 16 of 80 neighbors found using D_{shift} . (c) Ground truth Wordnet branch describing the content of the query image at multiple semantic levels. (d) Sub-tree formed by accumulating branches from all 80 neighbors. The number in each node denotes the accumulated votes. The red branch shows the nodes with the most votes. Note that this branch substantially agrees with the branch for vise and for person in the first and second examples respectively.

query image with some semantic connection to it. The goal of this section is to show that the performance achieved can match that of sophisticated algorithms which use much smaller training sets.

An additional factor in our dataset is the labeling noise. To cope with this we propose a voting scheme based around the Wordnet semantic hierarchy. Wordnet [15] provides semantic relationships between the 75,062 nouns for which we have collected images. For simplicity, we reduce the initial graph-structured relationships between words to a tree-structured one by taking the most common meaning of each word. The result is a large semantic tree whose nodes consist of the 75,062 nouns and their hypernyms, with all the leaves being nouns Fig. 7(c) shows the unique branch of this tree belonging to the nouns “vise” and “chemist”. Other work making use of Wordnet includes Hoogs and Collins [23] who use it to assist with image segmentation. While not using Wordnet explicitly, Barnard et al. [3] and Carbonetto et al. [8] learn models using both textual and visual cues.

Given the large number of classes in our dataset (75,062) and their highly specific nature, it is not practical or desirable to classify each of the classes separately. Instead, using the Wordnet hierarchy, we can perform classification at a variety of different semantic levels. So instead of just trying to recognize the noun “yellowfin tuna”, we may also perform recognition at the level of “tuna” or “fish” or “animal”. This is in contrast to current approaches to recognition that only consider a single, manually imposed, semantic meaning of an object or scene.

If classification is performed at some intermediate semantic level, for example using the noun “person”, we need not only consider images gathered from the Internet using “person”. Using the Wordnet hierarchy tree, we can also draw on all images belonging to nouns whose hypernyms include “person” (for example, “arithmetician”). Hence, we can massively increase the number of images in our training set at higher semantic levels. Near the top of the tree, however, the nouns are so generic (e.g. “object”) that the child images recruited in this manner have little visual consistency, so despite their extra numbers may be of

little use in classification⁵.

Our classification scheme uses the Wordnet tree in the following way. Given a query image, the neighbors are found using some similarity measure (typically D_{shift}). Each neighbor in turn votes for its branch within the Wordnet tree. Votes from the entire sibling set are accumulated across a range of semantic levels, with the effects of the labeling noise being averaged out over many neighbors. Classification may be performed by assigning the query image the label with the most votes at the desired height (i.e. semantic level) within the tree, the number of votes acting as a measure of confidence in the decision. In Fig. 7(a) we show two examples of this procedure, showing how precise classifications can be made despite significant labeling noise and spurious siblings. Using this scheme we now address the task of classifying images of people.

B. Person detection

In this experiment, our goal is to label an image as containing a person or not, a task with many applications on the web and elsewhere. A standard approach would be to use a face detector but this has the drawback that the face has to be large enough to be detected, and must generally be facing the camera. While these limitations could be overcome by running multiple detectors, each tuned to different view (e.g. profile faces, head and shoulders, torso), we adopt a different approach.

As many images on the web contain pictures of people, a large fraction (23%) of the 79 million images in our dataset have people in them. Thus for this class we are able to reliably find a highly consistent set of neighbors, as shown in Fig. 8. Note that most of the neighbors match not just the category but also the location and size of the body in the image, which varies considerably in the examples.

To classify an image as containing people or not, we use the scheme introduced in Section V-A, collecting votes from

⁵The use of Wordnet tree in this manner implicitly assumes that semantic and visual consistency are tightly correlated. While this might be the case for certain nouns (for example, “poodle” and “dachshund”), it is not clear how true this is in general. To explore this issue, we constructed an interactive poster that may be viewed at: <http://people.csail.mit.edu/torralba/tinyimages>.

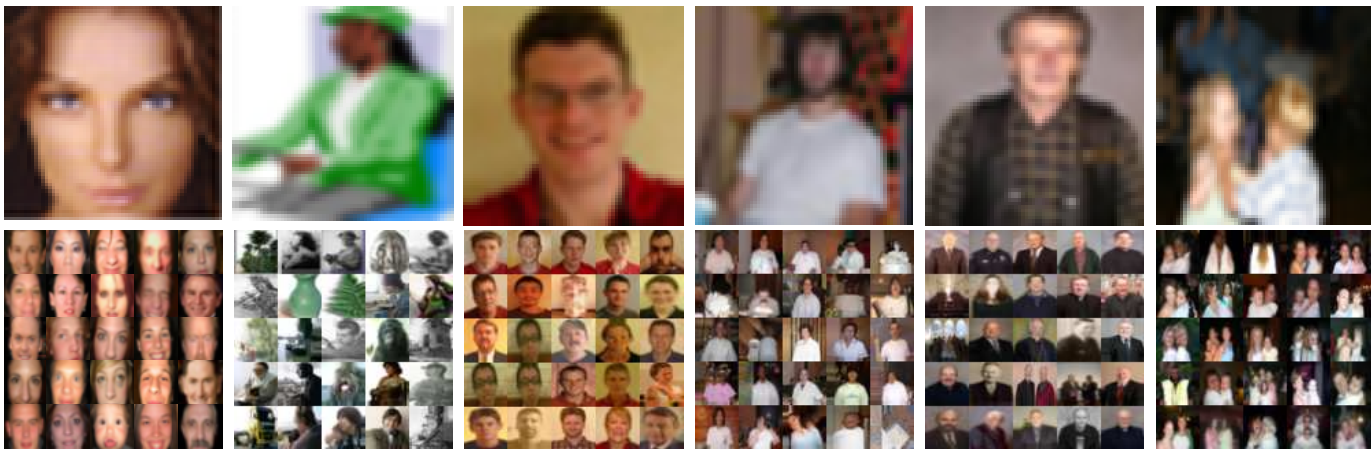


Fig. 8. Some examples of test images belonging to the “person” node of the Wordnet tree, organized according to body size. Each pair shows the query image and the 25 closest neighbors out of 79 million images using D_{shift} with 32×32 images. Note that the sibling sets contain people in similar poses, with similar clothing to the query images.

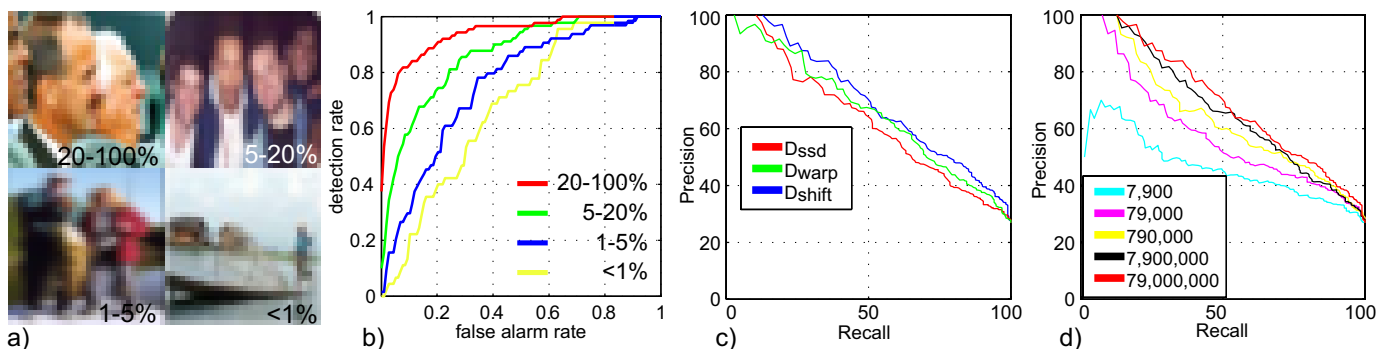


Fig. 9. (a) Examples showing the fraction of the image occupied by the head. (b)–(d): ROC curves for people detection (not localization) in images drawn randomly from the dataset of 79 million as a function of (b) head size; (c) similarity metrics and (d) dataset size using D_{shift} .

the 80 nearest neighbors. Note that the Wordnet tree allows us make use of hundreds of other words that are also related to “person” (e.g. artist, politician, kid, taxi driver, etc.). To evaluate performance, we used two different sets of test images. The first consisted of a random sampling of images from the dataset. The second consisted of images returned by Altavista using the query “person”.

1) *Evaluation using randomly drawn images:* 1125 images were randomly drawn from the dataset of 79 million (Fig. 8 shows 6 of them, along with some of their sibling set). For evaluation purposes, any people within the 1125 images were manually segmented⁶.

Fig. 9(b) shows the classification performance as the size of the person in the image varies. When the person is large in the image, the performance is significantly better than when it is small. This occurs for two reasons: first, when the person is large, the picture become more constrained, and hence finding good matches becomes easier. Second, the weak labels associated with each image in our dataset typically refer to the largest object in the image. Fig. 9(c)&(d) show precision-recall curves for different similarity measures and varying dataset size respectively, with the full 79 million images and D_{shift} yielding the best performance.

⁶The images and segmentations are available at: http://labelme.csail.mit.edu/browseLabelMe/static_web_tinyimages_testset.html

2) *Evaluation using Altavista images:* Our approach can also be used to improve the quality of Internet image search engines. We gathered 1018 images from Altavista image search using the keyword “person”. Each image was classified using the approach described in Section V-A. The set of 1018 images was then re-ordered according to the confidence of each classification. Fig. 10(a) shows the initial Altavista ranking while Fig. 10(b) shows the re-ordered set, showing a significant improvement in quality.

To quantify the improvement in performance, the Altavista images were manually annotated with bounding boxes around any people present. Out of the 1018 images, 544 contained people, and of these, 173 images contained people occupying more than 20% of the image.

Fig. 10 shows the precision-recall curves for the people detection task. Fig. 10(c) shows the performance for all Altavista images while Fig. 10(d) shows the performance on the subset where people occupy at least 20% of the image. Note that the raw Altavista performance is the same irrespective of the persons’ size (in both plots, by 5% recall the precision is at the level of chance). This illustrates the difference between indexing an image using non visual vs. visual cues. Fig. 10 also shows the results obtained when running a frontal face detector (an OpenCV implementation of Viola and Jones boosted cascade [27], [41]). We run the face detector on the original high-resolution images.

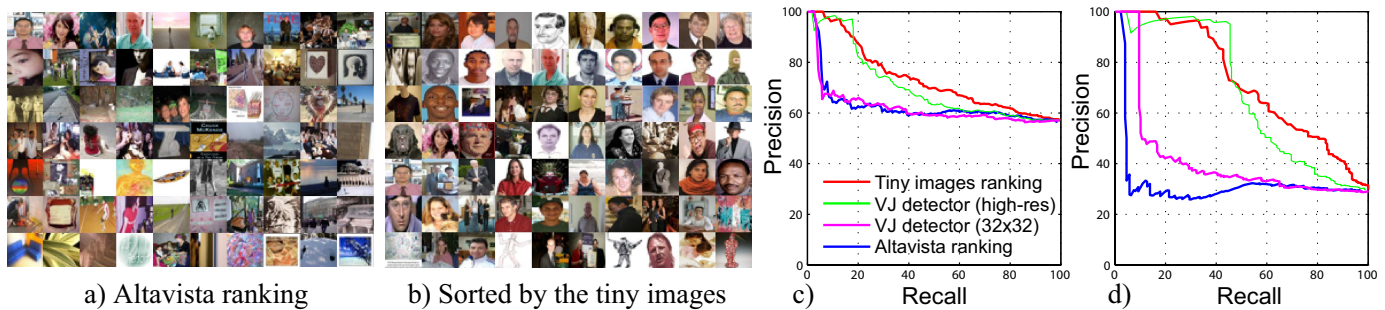


Fig. 10. (a) The first 70 images returned by Altavista when using the query “person” (out of 1018 total). (b) The first 70 images after re-ordering using our Wordnet voting scheme with the 79,000,000 tiny images. (c) Comparison of the performance of the initial Altavista ranking with the re-ordered images using the Wordnet voting scheme and also a Viola & Jones-style frontal face detector. (c) shows the recall-precision curves for all 1018 images gathered from Altavista, and (d) shows curves for the subset of 173 images where people occupy at least 20% of the image.

Note that the performance of our approach working on 32×32 images is comparable to that of the dedicated face detector on high resolution images. For comparison, Fig. 10 also shows the results obtained when running the face detector on low-resolution images (we downsampled each image so that the smallest axis has 32 pixels, we then upsampled the images again to the original resolution using bicubic interpolation. The upsampling operation was to allow the detector to have sufficient resolution to be able to scan the image.). The performance of the OpenCV detector drops dramatically with low-resolution images.

C. Person localization

While the previous section was concerned with an object detection task, we now address the more challenging problem of object localization. Even though the tiny image dataset has not been labeled with the location of objects in the images, we can use the weakly labeled (i.e. only a single label is provided for each image) dataset to localize objects. Much the recent work in object recognition uses explicit models that labels regions of images as being object/background. In contrast, we use the tiny image dataset to localize without learning an explicit object model. It is important to emphasize that this operation is performed without manual labeling of images: all the information comes from the loose text label associated with each image.

The idea is to extract multiple putative crops of the high resolution query image (Fig. 11(a)–(c)). For each crop, we resize it to 32×32 pixels and query the tiny image database to obtain its siblings set (Fig. 11(d)). When a crop contains a person, we expect the sibling set to also contain people. Hence, the most prototypical crops should get have a higher number of votes for the person class. To reduce the number of crops that need to be evaluated, we first segment the image using normalized cuts [11], producing around 10 segments (segmentation is performed on the high resolution image). Then, all possible combinations of contiguous segments are considered, giving a set of putative crops for evaluation. Fig. 11 shows an example of this procedure. Fig. 11(d) shows the best scoring bounding box for images from the Altavista test set.

D. Scene recognition

Many web images correspond to full scenes, not individual objects. In Fig. 12, we attempt to classify the 1125 randomly drawn images (containing objects as well as scenes) into “city”,

“river”, “field” and “mountain” by counting the votes at the corresponding node of the Wordnet tree. Scene classification for the 32×32 images performs surprisingly well, exploiting the large, weakly labeled database.

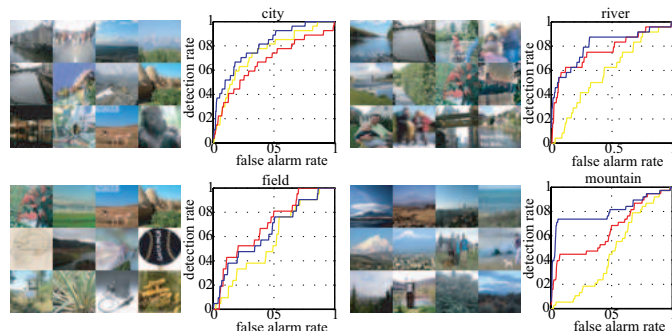


Fig. 12. Scene classification using the randomly drawn 1125 image test set. Note that the classification is “mountain” vs all classes present in the test set (which includes many kinds of objects), not “mountain” vs “field”, “city”, “river” only. Each quadrant shows some examples of high scoring images for that particular scene category, along with an ROC curve (yellow = 7,900,000 image training set; red = 79,000,000 images; blue = 79,000,000 images).

E. Automatic image annotation and dataset size

Here we examine the classification performance at a variety of semantic levels across many different classes as we increase the size of the database. For evaluation we use the test set of 1125 randomly drawn tiny images, with each image being fully segmented and annotated with the objects and regions that compose each image. To give a distinctive test set, we only use images for which the target object is absent or occupies at least 20% of the image pixels. Using the voting tree described in Section V-A, we classified them using $K = 80$ neighbors at a variety of semantic levels. To simplify the presentation of results, we collapsed the Wordnet tree by hand (which had 19 levels) down to 3 levels (see Fig. 13 for the list of categories at each level).

In Fig. 13 we show the average ROC curve area (across words at that level) at each of the three semantic levels for D_{ssd} and D_{shift} as the number of images in the dataset is varied. Note that (i) the classification performance increases as the number of images increases; (ii) D_{shift} outperforms D_{ssd} ; (iii) the performance drops off as the classes become more specific. A similar effect of dataset



Fig. 11. Localization of people in images. (a) input image, (b) Normalized-cuts segmentation, (c) three examples of candidate crops, (d) the 6 nearest neighbors of each crop in (c), accompanied by the number of votes for the person class obtained using 80 nearest neighbors under similarity measure D_{shift} . (e) Localization examples.

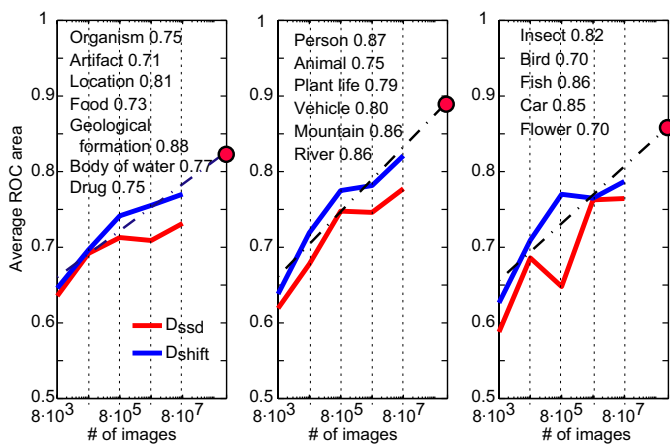


Fig. 13. Classification at multiple semantic levels using 1125 randomly drawn tiny images. Each plot shows a different manually defined semantic level, increasing in selectivity from left to right. The curves represent the average (across words at that level) ROC curve area as a function of number of images in the dataset (red= D_{ssd} , blue= D_{shift}). Words within each of the semantic levels are shown in each subplot, accompanied by the ROC curve area when using the full dataset. The red dot shows the expected performance if all images in Google image search were used (~ 2 billion), extrapolating linearly.

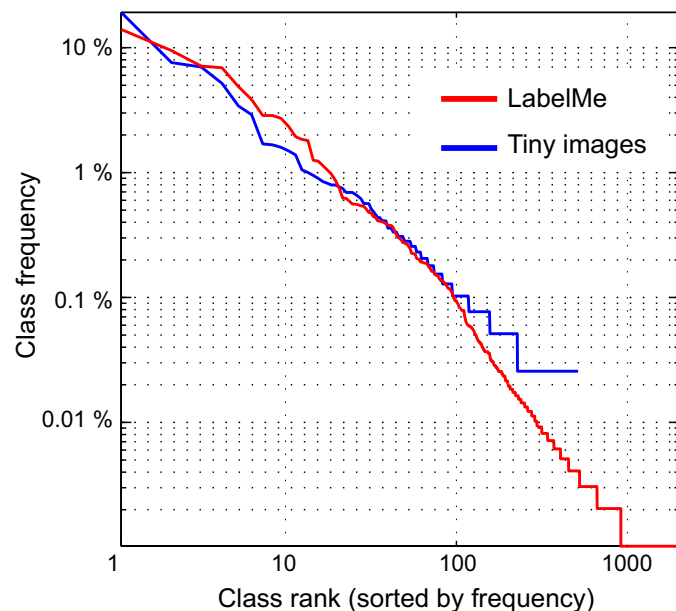


Fig. 15. Distribution of labels in image datasets. The vertical axis gives the percentage of polygons in the two datasets containing each object category (objects are sorted by frequency rank). The plot is in log-log axis.

size has already been shown by the language understanding community[2].

By way of illustrating the quality of the recognition achieved by using the 79 million weakly labeled images, we show in Fig. 14, for categories at three semantic levels, the images that were more confidently assigned to each class. Note that despite the simplicity of the matching procedure presented here, the recognition performance achieves reasonable levels even for relatively fine levels of categorization.

VI. THE IMPORTANCE OF SOPHISTICATED METHODS FOR RECOGNITION

The plot in Fig. 15 shows the frequency of objects in the tiny images database (this distribution is estimated using the hand labeled set of 1148 images). This distribution is similar to word frequencies in text (Zipf’s law). The vertical axis shows the percentage of annotated polygons for each object category. The horizontal axis is the object rank (objects are sorted by frequency). The four most frequent objects are people (29%), plant (16%), sky

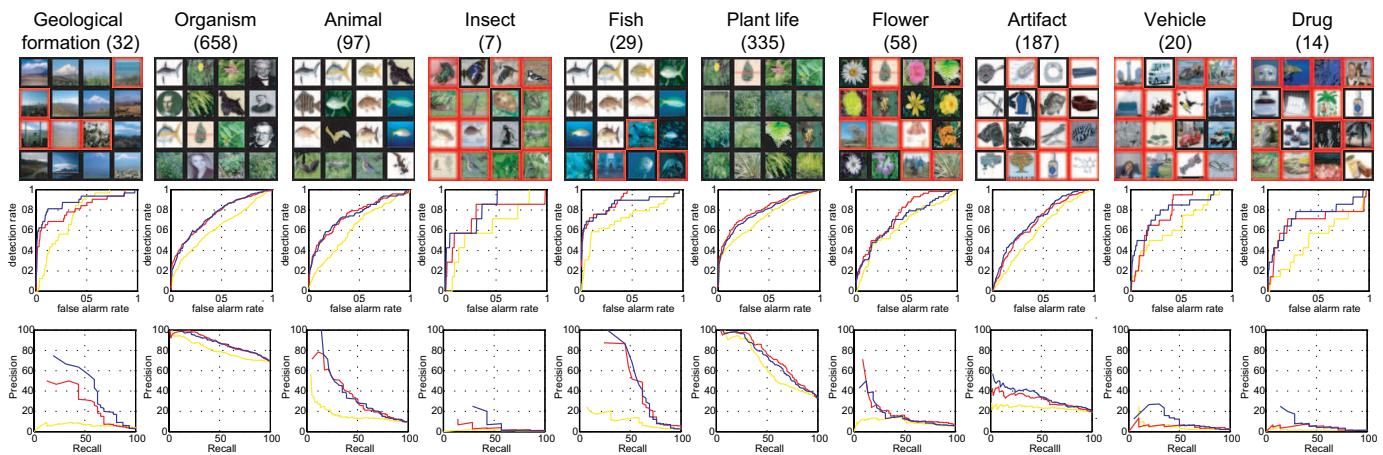


Fig. 14. Test images assigned to words, ordered by confidence. The number indicates the total number of positive examples in the test set out of the 1148 images. The color of the bounding box indicates if the image was correctly assigned (black) or not (red). The middle row shows the ROC curves for three dataset sizes (yellow = 7,900 image training set; red = 790,000 images; blue = 79,000,000 images). The bottom row shows the corresponding precision-recall graphs.

(9%) and building (5%). In the same plot we show the distribution of objects in the LabelMe dataset [35]. Similar distributions are also obtained from datasets collected by other groups [38]. As the distribution from Fig. 15 reveals, even when collecting extremely large databases, there will always be a large number of categories with very few training samples available. For some classes, a large amount of training data will be available and, as we discuss in this paper, nearest neighbor methods can be very effective. However, for many other classes learning will have to be performed with small datasets (for which we need to use sophisticated object models and transfer learning techniques).

VII. CONCLUSIONS

This paper makes the following important contributions: a) The compilation of a dataset of 79 million 32×32 color images, each with a weak text label and link to the original image. b) The characterization of the manifold of 32×32 images, showing that Internet sized datasets (10^8 – 10^9) yield a reasonable density over the manifold of natural images, at least for the purposes of object recognition. c) Showing that simple non-parametric methods, in conjunction with large datasets, can give reasonable performance on object recognition tasks. For richly represented classes, such as people, the performance is comparable to leading class-specific detectors.

Previous usage of non-parametric approaches in recognition have been confined to limited domains (e.g. pose recognition [36]) compared with the more general problems tackled in this paper, the limiting factor being the need for very large amounts of data. The results obtained using our tiny image dataset are an encouraging sign that the data requirements may not be insurmountable. Indeed, search engines such as Google index another 2–3 orders of magnitude more images, which could yield a significant improvement in performance.

In summary, all methods in object recognition have two components: the model and the data. The vast majority of the effort in recent years has gone into the modeling part – seeking to develop suitable parametric representations for recognition. In contrast, this paper moves into other direction, exploring how the data

itself can help to solve the problem. We feel the results in this paper warrant further exploration in this direction.

VIII. ACKNOWLEDGMENTS

Funding for this research was provided by NGA NEGI-1582-04-0004, Shell Research, Google, ONR MURI Grant N00014-06-1-0734, and NSF Career award (IIS0747120).

REFERENCES

- [1] T. Bachmann. Identification of spatially quantized tachistoscopic images of faces: How many pixels does it take to carry identity? *European Journal of Cognitive Psychology*, 3:85–103, 1991.
- [2] Michele Banko and Eric Brill. Scaling to very very large corpora for natural language disambiguation. In *ACL '01: Proceedings of the 39th Annual Meeting on Association for Computational Linguistics*, pages 26–33, Morristown, NJ, USA, 2001. Association for Computational Linguistics.
- [3] K. Barnard, P. Duygulu, N. de Freitas, D. Forsyth, D. Blei, and M. Jordan. Matching words and pictures. *JMLR*, 3:1107–1135, 2003.
- [4] S. Belongie, J. Malik, and J. Puzicha. Shape context: A new descriptor for shape matching and object recognition. In *Advances in Neural Info. Proc. Systems*, pages 831–837, 2000.
- [5] A. Berg, T. Berg, and J. Malik. Shape matching and object recognition using low distortion correspondence. In *CVPR*, volume 1, pages 26–33, June 2005.
- [6] T. Berg, A. Berg, J. Edwards, M. Maire, R. White, Y-W. Teh, E. Learned-Miller, and D. Forsyth. Names and faces in the news. In *CVPR*, volume 2, pages 848–854, 2004.
- [7] T. L. Berg and D. A. Forsyth. Animals on the web. In *CVPR*, volume 2, pages 1463–1470, 2006.
- [8] P. Carbonetto, N. de Freitas, and K. Barnard. A statistical model for general contextual object recognition. In *ECCV*, volume 1, pages 350–362, 2004.
- [9] C. Carson, S. Belongie, H. Greenspan, and J. Malik. Blobworld: Image segmentation using expectation-maximization and its application to image querying. *PAMI*, 24(8):1026–1038, 2002.
- [10] D. M. Chandler and D. J. Field. Estimates of the information content and dimensionality of natural scenes from proximity distributions. *JOSA*, 24:922–941, 2006.
- [11] T. Cour, F. Benezit, and J. Shi. Spectral segmentation with multiscale graph decomposition. In *CVPR*, volume 2, pages 1124–1131, 2005.
- [12] R. Datta, D. Joshi, J. Li, and J. Z. Wang. Image retrieval: Ideas, influences, and trends of the new age. *ACM Computing Surveys*, 40(2), 2008.
- [13] M. Everingham, A. Zisserman, C.K.I. Williams, and L. Van Gool. The PASCAL visual object classes challenge 2006 (voc 2006) results. Technical report, University of Oxford, September 2006.

- [14] L. Fei-Fei and P. Perona. A Bayesian hierarchical model for learning natural scene categories. In *CVPR*, pages 524–531, 2005.
- [15] C. Fellbaum. *Wordnet: An Electronic Lexical Database*. Bradford Books, 1998.
- [16] R. Fergus, P. Perona, and A. Zisserman. A visual category filter for Google images. In *ECCV*, pages 242–256, May 2004.
- [17] M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Petkovic, D. Steele, and P. Yanker. Query by image and video content: the QBIC system. *IEEE Computer*, 28(9):23–32, 1995.
- [18] M. M. Gorkani and R. W. Picard. Texture orientation for sorting photos at a glance. In *Proc. Intl. Conf. Pattern Recognition*, volume 1, pages 459–464, 1994.
- [19] K. Grauman and T. Darrell. Pyramid match hashing: Sub-linear time indexing over partial correspondences. In *CVPR*, pages 1–8, 2007.
- [20] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. Technical Report UCB/CSD-04-1366, 2007.
- [21] L. D. Harmon and B. Julesz. Masking in visual recognition: Effects of two-dimensional noise. *Science*, 180:1194–1197, 1973.
- [22] J. Hays and A. A. Efros. Scene completion using millions of photographs. *ACM SIGGRAPH*, 26, 2007.
- [23] A. Hoogs and R. Collins. Object boundary detection in images using a semantic ontology. In *AAAI*, 2006.
- [24] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*, pages 2169–2178, 2006.
- [25] Ann B. Lee, Kim S. Pedersen, and David Mumford. The nonlinear statistics of high-contrast patches in natural images. *Int. J. Comput. Vision*, 54(1-3):83–103, 2003.
- [26] J. Li, G. Wang, and L. Fei-Fei. Optimol: automatic object picture collection via incremental model learning. In *CVPR*, pages 1–8, 2007.
- [27] R. Lienhart, A. Kuranov, and V. Pisarevsky. Empirical analysis of detection cascades of boosted classifiers for rapid object detection. In *DAGM 25th Pattern Recognition Symposium*, pages 297–304, 2003.
- [28] David G. Lowe. Object recognition from local scale-invariant features. In *Proc. of the International Conference on Computer Vision ICCV, Corfu*, pages 1150–1157, 1999.
- [29] B. D. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *Proceedings of Imaging understanding workshop*, pages 121–130, 1981.
- [30] D. Nister and H. Stewenius. Scalable recognition with a vocabulary tree. In *CVPR*, pages 2161–2168, 2006.
- [31] A. Oliva and P.G. Schyns. Diagnostic colors mediate scene recognition. *Cognitive Psychology*, 41:176–210, 1976.
- [32] A. Oliva and A. Torralba. Modeling the shape of the scene: a holistic representation of the spatial envelope. *International Journal in Computer Vision*, 42:145–175, 2001.
- [33] T. Quack, U. Monich, L. Thiele, and B. Manjunath. Cortina: A system for large-scale, content-based web image retrieval. In *ACM Multimedia 2004*, Oct 2004.
- [34] B. Russell, A. Torralba, C. Liu, R. Fergus, and W. T. Freeman. Object recognition by scene alignment. In *Advances in Neural Info. Proc. Systems*, 2007.
- [35] B. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman. Labelme: a database and web-based tool for image annotation. *International Journal of Computer Vision*, 77(1-3):157–173, 2007.
- [36] G. Shakhnarovich, P. Viola, and T. Darrell. Fast pose estimation with parameter sensitive hashing. In *IEEE Intl. Conf. on Computer Vision*, volume 2, pages 750–757, 2003.
- [37] N. Snavely, S. M. Seitz, and R. Szeliski. Photo tourism: Exploring photo collections in 3D. *ACM Transactions on Graphics*, 25(3):137–154, 2006.
- [38] M. Spain and P. Perona. Measuring and predicting importance of objects in our visual world. Technical Report 9139, California Institute of Technology, 2007.
- [39] A. Torralba, R. Fergus, and W.T. Freeman. Tiny images. Technical Report MIT-CSAIL-TR-2007-024, CSAIL, MIT, 2007.
- [40] A. Torralba, R. Fergus, and Y. Weiss. Small codes and large databases for recognition. In *CVPR*, June 2008.
- [41] P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple classifiers. In *CVPR*, volume 1, pages 511–518, 2001.
- [42] J. Wang, G. Wiederhold, O. Firschein, and Wei. S. Content-based image indexing and searching using daubechies’ wavelets. *Int. J. Digital Libraries*, 1:311–328, 1998.
- [43] C. S. Zhu, N. Y. Wu, and D. Mumford. Minimax entropy principle and its application to texture modeling. *Neural Computation*, 9(8), November 1997.