

CENTRIST: A Visual Descriptor for Scene Categorization

Jianxin Wu, *Member, IEEE* and James M. Rehg, *Member, IEEE*

Abstract—CENTRIST (CENSus TRansform hISTogram), a new visual descriptor for recognizing topological places or scene categories, is introduced in this paper. We show that place and scene recognition, especially for indoor environments, require its visual descriptor to possess properties that are different from other vision domains (e.g. object recognition). CENTRIST satisfies these properties and suits the place and scene recognition task. It is a holistic representation and has strong generalizability for category recognition. CENTRIST mainly encodes the structural properties within an image and suppresses detailed textural information. Our experiments demonstrate that CENTRIST outperforms the current state-of-the-art in several place and scene recognition datasets, compared with other descriptors such as SIFT and Gist. Besides, it is easy to implement and evaluates extremely fast.

Index Terms—Place recognition, scene recognition, visual descriptor, Census Transform, SIFT, Gist

1 INTRODUCTION

KNOWING “Where am I” has always being an important research topic in robotics and computer vision. Various research problems have been studied in order to answer different facets of this question. For example, the following three research themes aim at revealing the location of a robot or determining where an image is taken.

- *Place recognition*, or global localization, which identifies the current position and orientation of a robot [1], [2], seeks to find the exact parameterization of a robot’s pose in a global reference frame. Place recognition is an inherent part of a Simultaneous Localization and Map Building (SLAM) system [3], [4].
- *Topological place recognition* answers the same question as place recognition, but at a coarser granularity [5]. In topological mapping, a robot is not required to determine its 3D location from the landmarks. It is enough to determine a rough location, e.g. corridor or office 113. A place in topological maps does not necessarily coincide with the human concept of rooms or regions [6]. Topological places are usually generated by a discretization of the robot’s environment based on certain distinctive features or events within it.
- *Scene recognition*, or scene categorization, is a term that is usually used to refer to the problem of recognizing the semantic label (e.g. bedroom, mountain, or coast) of a single image [7], [8], [9], [10], [11].

The input images in scene recognition are usually captured by a person, and are ensured to be representative or characteristic of the underlying scene category. It is usually easy for a person to look at an input image and determine its category label.

In this paper we are interested in recognizing places or scene categories using images taken by a usual rectilinear camera lens. Furthermore, since the exact robot pose estimation problem has been widely studied in SLAM systems, we focus on recognizing the topological location or semantic category of a place. Recognizing the semantic category of places from a robot platform is recently emerging as an interesting topic for both vision and robotics research, e.g. visual place categorization [12].

We believe that an appropriate representation (or, more precisely, visual descriptor) is the key to the success of a scene recognition task. In the literature, SIFT and Gist are probably the most popular descriptors in scene recognition [4], [7], [8], [9], [10], [11], [13], [14], [15], [16]. The SIFT descriptor is originally designed for recognizing the same object appearing under different conditions, and has strong discriminative power. Recognizing topological locations and scene categories, however, poses different requirements. Images taken from the same scene category may look very different, i.e. with huge intra-class variations. Similarly, images taken from different parts or view points of the same topological location (e.g. office 113) may also contain huge variations. Despite of the fact that densely sampled SIFT features plus the bag of visual words model have exhibited good performances in scene recognition, we would like to capture the stable spatial structure within images that reflects the functionality of the location, rather than capturing the detailed textural information of objects in the scene. Oliva and Torralba [10] proposed the Gist descriptor to represent such spatial structures.

- J. Wu is with the School of Computer Engineering, Nanyang Technological University, Singapore. Part of this research was finished when J. Wu was at Georgia Tech.
E-mail: jxwu@ntu.edu.sg
- J. M. Rehg is with the Center for Robotics and Intelligent Machines and College of Computing, Georgia Institute of Technology.
E-mail: rehg@cc.gatech.edu

Gist achieved high accuracy in recognizing natural scene categories, e.g. mountain and coast. However, when categories of indoor environments are added, its performance drops dramatically (c.f. Sec. 4.6).

The focus of this paper is CENTRIST, a visual descriptor that is suitable for recognizing topological places and scene categories. We will analyze the peculiarity of place images and list a few properties that are desirable for a place/scene recognition representation. We then focus on exhibiting how CENTRIST satisfies these properties better than competing visual descriptors, e.g. SIFT [17], HOG [18] or Gist [10]. We also show that CENTRIST has several important advantages in comparison to state-of-the-art descriptors for place/scene recognition and categorization:

- Superior recognition performance on multiple standard datasets;
- No parameter to tune;
- Extremely fast evaluation speed (> 50 fps);
- Very easy to implement, with source code publicly available.

The rest of the paper is organized as follows.¹ Related methods are discussed in Sec. 2. Sec. 3 introduces CENTRIST and focuses on how this visual descriptor suits the place/scene recognition domain. Experiments are shown in Sec. 4. Sec. 5 concludes this paper with discussions of drawbacks of the proposed method and future research.

2 RELATED WORK

2.1 Representation of scene images

Histograms of various image properties (e.g. color [5], [20], [21], or image derivatives [20]) have been widely used in scene recognition. However, after the SIFT [17] feature and descriptor are popularized in the vision community, it nearly dominates the visual descriptor choice in place and scene recognition systems [4], [7], [8], [9], [11], [14], [15], [16], [22]. SIFT features are invariant to scale and robust to orientation changes. The 128 dimensional SIFT descriptor has high discriminative power, while at the same time is robust to local variations [23]. It has been shown that the SIFT descriptor significantly outperforms edge points [9], pixel intensities [7], [8], and steerable pyramids [14] in recognizing places and scenes.

It is suggested in [10] that recognition of scenes could be accomplished by using *global configurations*, without detailed object information. Oliva and Torralba argued for the use of *Shape of the Scene*, an underlying similar and stable spatial structure that presumably exists within scene images coming from the same *function* category, to recognize scene categories. They proposed the Gist descriptor to represent such spatial structures. Gist computes the spectral information in an image through Discrete Fourier Transform (DFT). The spectral signals are then compressed by the Karhunen-Loeve

Transform (KLT). They showed that many scene signatures such as the degree of *naturalness* and *openness* were reliably estimated from the spectral signals, which in consequence resulted in satisfactory scene recognition results. Since spectral signals were computed from the global image, Oliva and Torralba suggested recognizing scenes without segmentation or recognizing local objects beforehand.

Gist achieved high accuracy in recognizing outdoor scene categories, e.g. mountain and coast. However, when categories of indoor environments are added, the Gist descriptor's performance drops dramatically. We will show in Sec. 4.6 that in a 15 class scene recognition dataset [9], which includes the categories used in [10] and several other categories (mainly indoor categories), accuracy of the Gist descriptor is much worse than its performance on outdoor images, and is significantly lower than the proposed CENTRIST descriptor.

The global configuration argument itself is accepted by many other researchers, whom used the SIFT descriptor to describe global configurations. Since the SIFT descriptor is designed to recognize the same object instance, statistical analyses of the distribution of SIFT descriptors are popular in scene recognition. Statistics of SIFT descriptors are more tolerant to the huge variations in scene images. In the bag of visual words model, SIFT descriptors are vector quantized to form the *visual codebook*, e.g. by the k-means clustering algorithm. The hope here is that the cluster centers will be representative common visual sub-structures, similar to the codebook in a communication system. We will compare SIFT and CENTRIST in Sec. 4.6.

A different representation was proposed by Vogel and Schiele [24]. They split each image into 10 by 10 cells. Each cell was given a semantic label from 9 categories (sky, water, grass, etc.). An SVM classifier ("concept classifier") is then trained to assign labels to cells. In other words, instead of generating intermediate concepts from data without supervision, they specify a small set of concepts and learn them in a supervised manner. Category of an image was determined from the concept labels of its 100 cells. Their experiments corroborated the observation that using intermediate concepts gave better performance than using crude image features.

2.2 Incorporating Spatial Information

Visual descriptors usually already encode some spatial information. For example, SIFT divides an image patch into 16 ($= 4 \times 4$) blocks. The SIFT descriptor is a concatenation of information extracted from these blocks. In the HOG visual descriptor, an image patch is divided into 105 ($= 7 \times 15$) overlapping local blocks. The concatenation of information from these blocks form the HOG feature vector.

It is long recognized that spatial arrangements in the image level are essential for recognizing scenes. For example, Szummer and Picard divided images into

¹. Preliminary version of portions of this work have been published in [19].

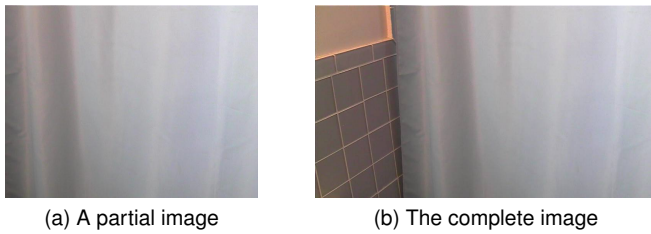


Fig. 1: A bathroom image is shown in Fig. 1b. Object in the scene (Fig. 1a) does not automatically reveal the room category. This image is one frame in the VPC dataset [12].

4×4 sub-blocks. The K-nearest neighbor classifier was applied to these sub-blocks. The final indoor-outdoor decision was then made based on classification results from the 16 sub-blocks [21]. Their experiments showed that a simple majority vote strategy for the second phase classification significantly improved recognition accuracy (approximately 10% higher compared to the sub-block accuracy).

Spatial arrangement information is completely ignored in the bag of visual words model. Lazebnik et al. proposed Spatial Pyramid Matching (SPM) to systematically incorporate spatial information in these models [9]. Features are quantized into M discrete types using k-means clustering with M centroids. They assume that only features of the same type can be matched. An image is divided in a hierarchical fashion (of level L). The image is divided into $2^l \times 2^l$ sub-blocks in level l , with each dimension (horizontal or vertical) being divided into 2^l evenly sized segments. For a feature type m , X_m and Y_m are sets of the coordinates of type m features. The histogram intersection kernel can be used to compute a matching score for feature type m . The final spatial pyramid matching kernel value is then the sum of all such scores.

3 CENTRIST: A VISUAL DESCRIPTOR FOR PLACE AND SCENE RECOGNITION

3.1 Desired properties

In this section, we first discuss some desired properties for a visual descriptor in place and scene recognition tasks. The CENTRIST descriptor is then proposed.

3.1.1 Holistic representation

Oliva and Torralba [10] showed that scene categories can be estimated without explicitly detecting and recognizing objects in the scene. As illustrated in Fig. 1, knowing the object in an image does not automatically tell us the place category. Instead, the curtain object and the tiles on the wall altogether clearly show that this is a bathroom image. Many useful information sources such as the tiles are usually contained in those regions that are not objects. Furthermore, detecting and recognizing objects in cluttered environments is probably more difficult than directly recognizing the scene category.

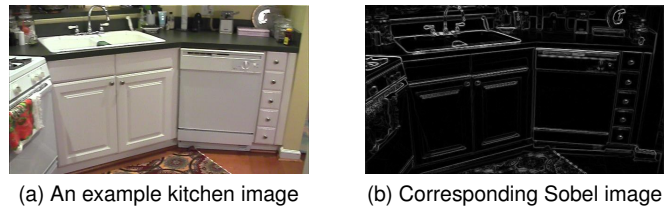


Fig. 2: Fig. 2a shows an example kitchen image. Fig. 2b shows its corresponding Sobel gradients. The Sobel gradients are normalized to $[0 \ 255]$. This image is one frame in the VPC dataset [12].

3.1.2 Capturing the structural properties

We want the descriptor to (implicitly or explicitly) capture general structural properties such as rectangular shapes, flat surfaces, tiles, etc., while suppressing detailed textural information. In recognizing place categories, fine-scale textures will distract the classifier. They can be noisy and harmful if the feature extraction method is not carefully designed. Fig. 2 illustrates this idea. For example, color and wooden texture of the cabinets and drawers do not provide useful hints for deciding the scene category. Other examples include patterns in the rug, or the detection of bottles on the counter-top.

On the other hand, spatial structures are very useful in suggesting the scene category. Fig. 2b shows the Sobel gradient image of Fig. 2a. In the Sobel image, many of the fine-scale textures are suppressed and spatial structures (e.g. the shapes that reflect the sink and dishwasher) become more prominent. It is possible that a human observer can recognize the category “kitchen” from the Sobel image alone.

In Fig. 2b we observe that many structural properties can be reflected by the distribution of local structures, for example, the percentages of local structures that are local horizontal edge, vertical edge, or junctions. Our CENTRIST descriptor models the distribution of local structures.

3.1.3 Rough geometry is useful

Strong geometrical constraints (e.g. the constellation model [25] or pictorial structures [26]) are very useful in object recognition. However, they are essentially not applicable in place categorization due to the large intra-class variations. Although object category recognition already deals with much larger variations than those in object instance recognition, the variabilities in scene/place categorization are even higher. In addition to these variations as those in object category recognition, objects can also appear in different spatial arrangements and some objects might be missing.

However, rough geometrical constraints are very helpful in recognizing place categories. For example, a reading lamp is usually placed close to the bed in a bedroom as in Fig. 3a. When a TV appears in a bedroom, it is often at the foot of the bed as shown in Fig. 3c. More general



Fig. 3: Place images do not exhibit strong geometrical constraints among objects.

constraints such as “sky is above the ground” will help reduce ambiguity, even when the images are taken from random viewpoints.

3.1.4 Generalizability

The learned category concepts will be applied to new images. An ideal situation is that the visual descriptors are compact within a category (even under large visual variations), and are far apart when they belong to different categories.

Fig. 4 shows three images from the corridor environment. These images were taken from approximately the same location in the same environment, but we already see large visual variations. We would expect even larger variations from pictures taken in different corridor environments. The visual descriptor must be able to capture the similar spatial structures: open spaces in the middle, stairs, strips on the wall, etc.

We propose to use CENTRIST (CENSus TRansform HISTogram) as our visual descriptor for the place category recognition task. CENTRIST is a holistic representation that captures structural properties by modeling distribution of local structures. We capture rough geometrical information by using a spatial CENTRIST representation. CENTRIST also has similar descriptor vectors for images in the same place category.

3.2 Census Transform (CT) and CENTRIST

Census Transform (CT) is a non-parametric local transform originally designed for establishing correspondence between local patches [28]. Census transform compares the intensity value of a pixel with its eight neighboring pixels, as illustrated in Eqn. 1. If the center pixel is bigger than (or equal to) one of its neighbors, a bit 1 is set in the corresponding location. Otherwise a bit 0 is set.

$$\begin{array}{c|c|c} 32 & 64 & 96 \\ \hline 32 & 64 & 96 \\ \hline 32 & 32 & 96 \end{array} \Rightarrow 1 \ 0 \Rightarrow (11010110)_2 \Rightarrow \text{CT} = 214 \quad (1)$$

The eight bits generated from intensity comparisons can be put together in any order (we collect bits from left to right, and from top to bottom), which is consequently converted to a base-10 number in $[0 \ 255]$. This is the Census Transform value (CT value) for this center pixel. Census Transform is robust to illumination changes,

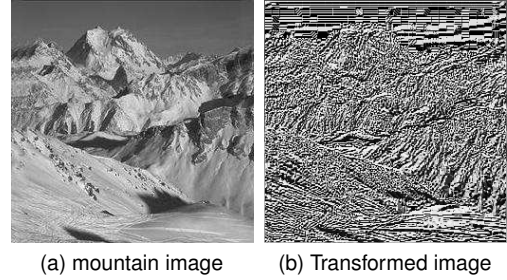


Fig. 5: An example Census Transformed image. This image is taken from the 15 class scene recognition dataset (c.f. Sec. 4).

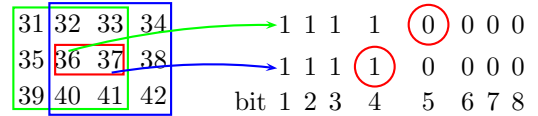


Fig. 6: Illustration of constraints between CT values of neighboring pixels.

gamma variations, etc. Note that the Census Transform is equivalent (modulo a difference in bit ordering) to the *local binary pattern* code $LBP_{8,1}$ [29].

As a visualization method, we create a *Census Transformed image* by replacing a pixel with its CT value. Shown by the example in Fig. 5, the Census Transform retains global structures of the picture (especially discontinuities) besides capturing the local structures.

Another important property of the transform is that CT values of neighboring pixels are highly correlated. In the example of Fig. 6, we examine the direct constraint posed by the two center pixels. The Census Transform for pixels valued 36 and 37 are depicted in right, and the two circled bits are both comparing the two center pixels (but in different orders). Thus the two bits must be strictly complementary to each other if the two pixels are not equal. More generally, bit 5 of $\text{CT}(x, y)$ and bit 4 of $\text{CT}(x + 1, y)$ must be complementary to each other, if the pixels at (x, y) and $(x + 1, y)$ are not equal. There are eight such constraints between one pixel and its eight neighboring pixels.

Besides these deterministic constraints, there also exist indirect constraints. For example, in Fig. 6, the pixel valued 32 compares with both center pixels when computing their CT values (bit 2 of $\text{CT}(x, y)$ and bit 1 of $\text{CT}(x + 1, y)$). Depending on the comparison results between the center pixels, there are probabilistic relationships between these bits.

The transitive property of such constraints also make them propagate to pixels that are far apart. For example, in Fig. 6, the pixels valued 31 and 42 can be compared using various paths of comparisons, e.g. $31 < 35 < 39 < 40 < 41 < 42$. Similarly, although no deterministic comparisons can be deduced between some pixels (e.g. 34 and 39), probabilistic relationships still can



Fig. 4: Example images from the KTH IDOL dataset [27]. Images showed approximately the same location under different weather conditions. Images were taken by a robot called Minnie.

be obtained. The propagated constraints make Census Transform values and Census Transform histograms (i.e. CENTRIST, CENSus TRansform hISTogram) implicitly contain information for describing global structures.

Finally, the Census Transform operation transforms any 3 by 3 image region into one of 256 cases, each corresponding to a special type of local structure of pixel intensities. The CT value acts as an index to these different local structures. No total ordering or partial ordering exists among the CT values. It is important to refrain from comparing two CT values as comparing two integers (like what we do when comparing two pixel intensity values). For example, in the homogeneous region of Fig. 5a, there are only a few distinct CT values which are close in the Hamming distance in the corresponding region in Fig. 5b.

A histogram of CT values for an image or image patch can be easily computed, and we use CENTRIST (CENSus TRansform hISTogram) as our visual descriptor. CENTRIST can be computed very efficiently. It only involves 16 operations to compute the CT value for a center pixel (8 comparisons and 8 additional operations to set bits to 0 or 1). The cost to compute CENTRIST is linear in the number of pixels of the region we are interested in. There is also potential for further acceleration to the computation of CENTRIST, by using special hardware (e.g. FPGA), because it mainly involves integer arithmetic that are highly parallel in nature.

3.3 Constraints among CENTRIST components

Usually there is no obvious constraint among the components of a histogram. For example, we would often treat the R, G, and B channels of a color histogram as independent to each other. CENTRIST, however, exhibits strong constraints or dependencies among its components.

Take as example the direct constraint shown in Fig. 6, bit 5 of $CT(x, y)$ and bit 4 of $CT(x, y + 1)$ must be complementary to each other if they are not equal. Both bits are 1 if they are equal. If we apply this constraint to all pixels in an image, we get to the conclusion that *the number of pixels whose CT value's bit 5 is 1 must be equal to or greater than the number of pixels whose CT value's bit 4 is 0*, if we ignore border pixels where such constraints break. Let \mathbf{h} be the CENTRIST descriptor of any image. The above statement is translated into the

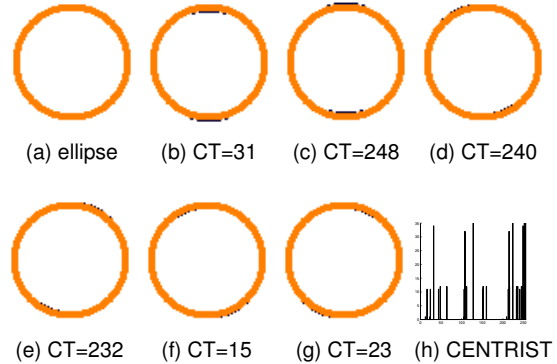


Fig. 7: Illustration of Census transforms. Fig. 7a is an example image of ellipse. Fig. 7b-7g show pixels having the 6 highest frequency CT values (shown in black). Fig. 7h is the CENTRIST feature vector of Fig. 7a.

following equation:

$$\sum_{i \& 0x08 = 0x08} \mathbf{h}(i) \geq \sum_{i \& 0x10 = 0} \mathbf{h}(i), \quad (2)$$

where $\&$ is *bitwise and*, $0x08$ is the number 8 in the hexadecimal format, and $0 \leq i \leq 255$. By switching 1 and 0, we get another equation:

$$\sum_{i \& 0x08 = 0} \mathbf{h}(i) \leq \sum_{i \& 0x10 = 0x10} \mathbf{h}(i). \quad (3)$$

Similarly, six other linear inequalities can be specified by comparing $CT(x, y)$ with $CT(x-1, y-1)$, $CT(x-1, y)$, and $CT(x-1, y+1)$. Any CENTRIST feature vector resides in a subspace that is defined by these linear inequalities.

We can not write down explicit equations for the indirect or transitive constraints in a CENTRIST feature vector. However, we expect these constraints to further reduce the intrinsic dimensionality of the CENTRIST feature vectors. The constraints among elements in a CENTRIST vector make PCA suitable for its dimension reduction task.

3.4 CENTRIST encodes image structures

In order to understand why CENTRIST efficiently captures the essence of a scene image, it is worthwhile to further examine the distribution of CT values and CENTRIST feature vectors. Using images from the 15

class scene dataset [9], we find that the 6 CT values with highest counts are $CT = 31, 248, 240, 232, 15, 23$ (excluding 0 and 255). As shown in Fig. 7b-7g, these CT values correspond to local 3×3 neighborhoods that have either horizontal or various close-to-diagonal edge structures.

The CENTRIST vector of the ellipse image in Fig. 7a is shown in Fig. 7h. It summarizes the distribution of various local structures in the image. If an image has a CENTRIST feature vector close to that of Fig. 7h, we would well expect the image to exhibit an ellipse shape with a high probability (c.f. Sec. 3.5 for more evidences.)

A simplification to the one dimensional world better explains the intuition behind our statement. In 1-d there are only 4 possible CT values, and the semantic interpretation of these values are obvious. As shown in Fig. 8a, the four CT values are $CT = 0$ (valley), $CT = 1$ (downhill), $CT = 2$ (uphill), and $CT = 3$ (peak). Downhill shapes and uphill shapes can only be connected by a valley, and uphill shapes require a peak to transit to downhill shapes. Because of these constraints, the only other shapes that have the same CENTRIST descriptor as that of Fig. 8a are those shapes that cut a small portion of the left part of Fig. 8a and move it to the right. Images that are different but keep the shapes (e.g. Fig. 8b) also are similar in their CENTRIST descriptors (Fig. 8d). On the contrary, a huge number of possible curves have the same intensity histogram as that of Fig. 8a. Even if we impose smoothness constraints between neighboring pixel intensities, the shape ambiguity is still large. Fig. 8c is smooth and has the same intensity histogram as that of Figs. 8a and 8b, but it has different shape and a very different CENTRIST descriptor.

3.5 Reconstructing image patches from CENTRIST descriptors

We also performed some reconstruction experiments to further illustrate how CENTRIST encodes image structures. When we randomly shuffle the pixels of an input image, the original structure of the image is completely lost. Using the shuffled image as an initial state, we repeatedly change two pixels at one time, until the current state has the same CENTRIST descriptor as the input image. This optimization is guided by the Simulated Annealing algorithm. If the structure of the original image is observed in the reconstruction result (i.e. the termination state), this is an evidence that structure of an image is encoded in its CENTRIST descriptor.

In the reconstruction results in Fig. 9, the left image in each subfigure is the input image. A pair of pixels in an input image are randomly chosen and exchanged. The exchange operation is repeated multiple times (equal to the number of pixels in the input image), which gives the initial state for the reconstruction. The cost function is set to the Euclidean distance between CENTRIST descriptors of the current state and the input. The terminating state is the output of the reconstruction (right image in each subfigure of Fig. 9).

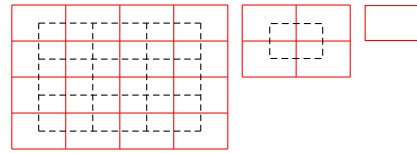


Fig. 10: Illustration of the level 2, 1, and 0 split of an image.

Although the initial states look like random collection of pixels, many of the reconstruction results perfectly match the input images (subfigure (a)-(g) in Fig. 9). More examples are reconstructed with minor discrepancies (subfigure (h)-(p)). Large scale structures of the input digits and characters are successfully reconstructed in these images, with small errors. In the rest examples, e.g. '2' and 'e', major structures of the original input images are still partially revealed.

Two points are worth pointing out about the reconstruction results. First, in larger images a CENTRIST descriptor is not enough to reconstruct the original image. However, as a visual descriptor, it has the ability to distinguish between images with different structural properties. Second, it is essentially impossible to reconstruct even a small image using other descriptors (e.g. SIFT, HOG, or Gist).

3.6 Spatial representations

Because CENTRIST can only encode global shape structure in a small image patch, in order to capture the global structure of an image in larger scales, we propose a spatial representation based on the Spatial Pyramid Matching scheme in [9]. A *spatial pyramid* (dividing an image into subregions and integrating correspondence results in these regions) encodes rough global structure of an image and usually improves recognition. It is straightforward to build a spatial pyramid for the proposed CENTRIST representation.

As shown in Fig. 10, the level 2 split in a spatial pyramid divides the image into $2^2 \times 2^2 = 16$ blocks. We also shift the division (dash line blocks) in order to avoid artifacts created by the non-overlapping division, which makes a total of 25 blocks in level 2. This is different from the spatial hierarchy in [9]. Similarly, level 1 and 0 have 5 and 1 blocks respectively. The image is resized between different levels so that all blocks contain the same number of pixels. CENTRIST in all blocks are then concatenated to form an overall feature vector. For example, if we use PCA to reduce the dimensionality of CENTRIST to 40, a level 2 pyramid will then result in a feature vector which has $40 \times (25 + 5 + 1) = 1240$ dimensions.

We want to emphasize that this spatial representation is independent of the descriptor used for each sub-window. In this paper, we use two different representations. In the first we use PCA to reduce the CENTRIST descriptor to 40 dimensions, which we call spatial PACT

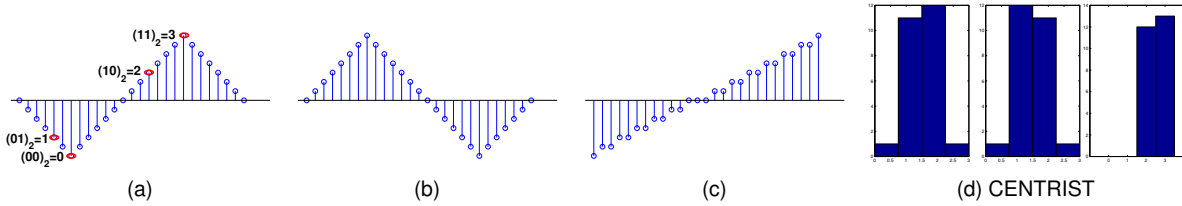


Fig. 8: Census Transform encodes shapes in 1-d. Subfigure (d) shows CENTRIST descriptors of figures (a)-(c), respectively.

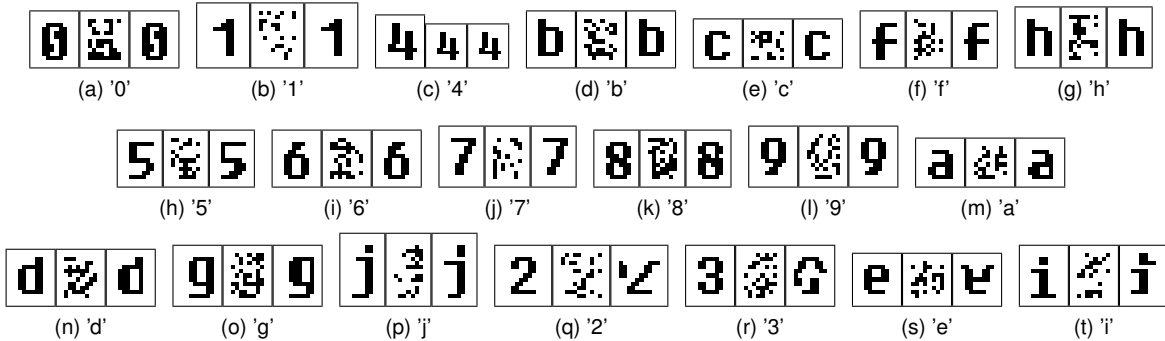


Fig. 9: Reconstruct images from CENTRIST descriptors. In each group of images, we show the input image, the initial state of optimization, and the terminating state (reconstruction result).

(spatial Principal component Analysis of Census Transform histograms), or sPACT. In the second approach we use a bag of visual words model. Details of both approaches are described in Sec. 4.

3.7 Limitations of CENTRIST

As we have stated from the very beginning, CENTRIST is designed to be a representation that suits place recognition and categorization problems. This design choice renders limitations that prevent it from being applied in some other applications.

- CENTRIST is not invariant to rotations or scale changes. In scene recognition, images are always taken in the upright view and we do not request rotation-invariance. Similarly, scale invariance is not critical for scene recognition either. However, these limitations indicate that CENTRIST may not be suitable for some other areas. For example, we will show in Fig. 16 that CENTRIST is inferior to SIFT in the Caltech 101 object recognition problem [30].
- CENTRIST is not a precise shape descriptor. It is designed to recognize shape categories, but not for exact shape registration applications, *e.g.* the shape retrieval task in [31], [32].
- CENTRIST ignores color information.

4 EXPERIMENTS

The CENTRIST visual descriptor is tested on 5 datasets: Swedish leaf [33], KTH IDOL [20], 15 class scene category [9], 8 class sports event [15], and 67 class indoor scene recognition [34]. In each dataset, the available data

are randomly split into a training set and a testing set following published protocols on these datasets. The random splitting is repeated 5 times, and the average accuracy is reported. Although color images are available in 4 datasets (leaf, IDOL, events, and indoor), we only use the intensity values and ignore color information.

In spatial PACT, we use 40 eigenvectors in the PCA operation. The largest 40 eigenvalues accounted for 90.6% to 94.1% of the sum of all eigenvalues in these datasets. We remove the two bins in CENTRIST with $CT = 0, 255$. We normalize the CENTRIST descriptors and PCA eigenvectors such that they have zero mean and unit norm. Our experiments empirically showed that instead of using the standard PCA, slightly faster speed and higher accuracies were obtained if we did not subtract the mean vector in the PCA operation across all datasets. Thus we do not subtract the mean vector in PACT.² CENTRIST will also be used in a bag of visual words framework in Sec. 4.7.

Since the CT values are based solely on pixel intensity comparisons, it might be helpful to include a few image statistics, *e.g.* average value and standard deviation of pixels in a block. The feature vector of a level 2 spatial PACT then becomes $(40 + 2) \times (25 + 5 + 1) = 1302$ dimensional.

SVM classifiers were widely used in our experiments. Whenever an SVM classifier was applied, we used the RBF kernel. Kernel parameters were chosen by cross validation on the training set of each dataset, in the grid $\log_2 C \in [-5, 15]$, $\log_2 \gamma \in [-11, 3]$ (with grid step size 2).

². PACT Code is available at <http://c2inet.sce.ntu.edu.sg/Jianxin/PACT/PACT.htm>.



Fig. 11: Example images from the Swedish Leaf dataset. The first 15 images are chosen from the 15 leaf species, one per species. The last image is the contour of the first leaf image.

TABLE 1: Results on the Swedish leaf dataset.

Method	Input	Rates
Shape-Tree [31]	Contour only	96.28%
IDSC+DP [32]	Contour only	94.13%
spatial PACT	Contour only	90.61%
SC+DP [32]	Contour only	88.12%
Söderkvist [33]	Contour only	82.40%
spatial PACT	Gray-scale image	97.87%
SPTC+DP [32]	Gray-scale image	95.33%

4.1 Swedish Leaf

The Swedish leaf dataset [33] collects pictures of 15 species of Swedish leaves (c.f. Fig. 11). There are 75 images in each class. Following [33], 25 images from each class are used for training and the remaining 50 for testing. This dataset has been used to evaluate shape matching methods [31], [32], in which the contour of leaves (instead of the gray-scale or color leaf picture) were used as input. In the contour image (e.g. the last picture in Fig. 11), no other information is available (e.g. color, texture) except shape or structure of the leaf. We use the contour input to further verify our statement that the CENTRIST descriptor encodes such information.

In each train/test split of images, the 25 training images from each class are used to compute the PCA eigenvectors. 10 and 40 eigenvectors are used when the inputs are contours and intensity images, respectively, in order to capture roughly 90% of the sum of eigenvalues. Results on this dataset are shown in Table 1. Although not specifically designed for matching shapes, spatial PACT can achieve 90.61% accuracy on leaf contours, better than Shape Context+Dynamic Programming (SC+DP). When pictures instead of contours are used as input, spatial PACT can recognize 97.87% leaves, which outperforms other methods.

4.2 KTH IDOL and INDECS

The KTH IDOL (Image Database for rObot Localization) dataset [27] was captured in a five-room office environment, including a one-person office, a two-person office, a kitchen, a corridor, and a printer area. Images were taken by two Robots: Minnie and Dumbo. The purpose of this dataset is to recognize which room the robot is in based on a single image, i.e. a topological place instance recognition problem.

Image resolution in IDOL is 320×240 . A complete image sequence contained all the images captured by

a robot when it was driven through all five rooms. Images were taken under 3 weather conditions: Cloudy, Night, and Sunny. For each robot and each weather condition, 4 runs were taken on different days. Thus, there are in total $2 \times 3 \times 4 = 24$ image sequences. Various changes during different robot runs (e.g. moving persons, changing weather and illumination conditions, relocated/added/removed furniture) make this dataset both realistic and challenging. Fig. 4 in page 5 shows images taken by the Minnie robot under 3 different weather conditions at approximately the same location, but with substantial visual changes.

In our experiments we use the run 1 and 2 in each robot and weather condition. We perform 3 types of experiments as those in [20]. First we train and test using the same robot, same weather condition. Run 1 is used for training and run 2 for testing, and vice versa. Second we use the same robot for training and testing, but with different weather conditions. These experiments test the ability of spatial PACT to generalize over variations caused by person, furniture, and illumination. The third type of experiment uses training and testing set under the same weather conditions, but captured by different robots. Cameras were mounted at different heights on the robots, which made the pictures taken by the two robots quite different

The KTH-INDECS dataset [35] was collected in the same environment as the IDOL dataset. Instead of using robots, cameras were mounted in several fixed locations inside each room. Pictures of multiple viewing angles were taken in each location. In the last type of experiment we use INDECS images as training examples, and test on both INDECS images under different weather conditions and on images taken by robots.

In [20], images were represented by the “High Dimensional Composed Receptive Field Histograms”, and were classified by χ^2 kernel SVMs. Results using level 2 spatial PACT and 1-NN are shown in Table 2, compared against results in [20]. The mean and standard deviation of an image block will vary greatly with illumination changes. Since the IDOL and INDECS datasets both contain dramatic illumination changes, they were not appended to the PACT vectors in this dataset. Note that we used the 15 class scene dataset [9] to compute eigenvectors for this problem.

In the first type of experiments, both spatial PACT and the method in [20] attain high accuracy ($> 95\%$), and the two methods are performing almost equally well. However, in the second type of experiments spatial PACT has significantly higher accuracies (18% higher in Minnie and 14% higher in Dumbo). The superior performance of our CENTRIST based representation shows that it is robust to illumination changes and other minor variations (e.g. moving persons, moved objects in an image, etc). The Dumbo robot achieves a 94.57% accuracy using a single input image without knowing any image histories (a “kidnapped robot” [36]). Thus, after walking a robot in an environment, spatial PACT enables the

TABLE 2: Average accuracies on recognizing topological place instances using the KTH-IDOL dataset and the KTH-INDECS dataset. Level 2 pyramids are used for spatial PACT. “Robots” means both Minnie and Dumbo.

Train	Test	Condition	sPACT+1-NN	sPACT+SVM	[20]
Minnie	Minnie	Same	95.35%	94.79%	95.51%
Dumbo	Dumbo	Same	97.62%	96.35%	97.26%
Minnie	Minnie	Different	90.17%	83.10%	71.90%
Dumbo	Dumbo	Different	94.98%	89.35%	80.55%
Minnie	Dumbo	Same	77.78%	70.15%	66.63%
Dumbo	Minnie	Same	72.44%	65.18%	62.20%
Camera	Camera	Different	90.01%	78.39%	75.67%
Camera	Robots	Same	64.39%	42.16%	50.56%

TABLE 3: Average accuracies on the KTH-IDOL dataset and the KTH-INDECS dataset using different levels of spatial pyramid. “Robots” means both Minnie and Dumbo.

Train	Test	Condition	$L = 3$	$L = 2$	$L = 1$	$L = 0$
Minnie	Minnie	Same	95.01%	95.35%	95.08%	86.08%
Dumbo	Dumbo	Same	95.51%	97.62%	96.87%	88.26%
Minnie	Minnie	Different	90.30%	90.17%	85.75%	60.51%
Dumbo	Dumbo	Different	94.67%	94.98%	91.75%	74.67%
Minnie	Dumbo	Same	74.96%	77.78%	75.56%	62.34%
Dumbo	Minnie	Same	68.59%	72.44%	71.36%	53.74%
Camera	Camera	Different	92.37%	90.01%	84.80%	71.45%
Camera	Robots	Same	60.73%	64.39%	57.87%	41.55%

robot to robustly answer the question “Where am I?” based on a single image, a capacity that is attractive to indoor robot applications. When the training and testing data come from different robots, the performance of both methods drop significantly. This is expected, since the camera heights are quite different. However, spatial PACT still outperforms [20] by about 10%. In the last type of experiment involving camera images, spatial PACT achieved about 14% higher accuracies than those reported in [20].

We observed that the 1-NN classifier achieved higher accuracies than SVM in Table 2. In this topological place recognition problem, different videos contained images taken from the same set of rooms under different conditions. Different views of one room could appear very differently, which added to the difficulty of SVM classification. However, for an image in the “cloudy” condition video, it is highly possible that its nearest neighbor in the “sunny” video was an image that was taken in the same room and approximately same view.

We also tested the effects of using different pyramid levels. As shown in Table 3, applying a spatial pyramid matching scheme greatly improves system performances ($L > 0$ vs. $L = 0$). However, the improvement after $L > 2$ is negligible. $L = 3$ performance is even worse than that of $L = 2$ in most cases. Our observation corroborates that of [9], which used a scene recognition dataset. In the remainder of this paper, we will use $L = 2$ in spatial PACT.



Fig. 12: One image from each of the 15 scene categories from [9]. The categories are bedroom, coast, forest, highway, industrial, inside city, kitchen, living room, mountain, office, open country, store, street, suburb, and tall building, respectively (from top to bottom, and from left to right).

Finally, CENTRIST can be computed and evaluated quickly, and so is spatial PACT. The IDOL dataset has around 1000 images in each image sequence, and spatial PACT processes at about 50 frames per second on an Intel Pentium IV 2GHz computer for computing the features, and finding the 1-NN match. The speed is 20 fps if we also consider hard drive I/O time.

4.3 The 15 class scene category dataset

The 15 class scene recognition dataset was built gradually by Oliva and Torralba ([10], 8 classes), Fei-Fei and Perona ([8], 13 classes), and Lazebnik et al. ([9], 15 classes). This is a scene category dataset (scene classes including office, store, coast, etc. Please refer to Fig. 12 for example images and category names.) Images are about 300×250 in resolution, with 210 to 410 images in each category. This dataset contains a wide range of scene categories in both indoor and outdoor environments. Unlike the KTH IDOL images which are taken by robots, images in this datasets are taken by people and representative of the scene category. We use SVM and spatial PACT in this dataset. Same as previous research on this dataset, 100 images in each category are used for training, and the remaining images constitute the testing set. The training images in each train/test split were used to perform PCA. The results are shown in Table 4, where our level 2 spatial PACT achieves the highest accuracy. In Table 4 the average accuracy in all categories are reported (i.e., average of diagonal entries in the confusion matrix).

In [9], low level features were divided into weak features (which were computed from local 3×3 neighborhoods) and strong features (which are computed from larger 16×16 image patches). Strong features were shown to have much higher accuracy than weak features (c.f. Table 4). The Census Transform is computed from 3×3 local neighborhoods, and falls into the weak feature category. However, when $L = 0$ (not using spatial pyramid), CENTRIST substantially outperforms the weak

TABLE 4: Recognition rates on the 15 class scene dataset.

L	Method	Feature type	Rates
0	SPM [9]	16 channel weak features	45.3 \pm 0.5
0	SPM [9]	SIFT, 200 cluster centers	72.2 \pm 0.6
0	SPM [9]	SIFT, 400 cluster centers	74.8 \pm 0.3
0	CENTRIST	CENTRIST, not using PCA	73.29 \pm 0.96
3	SPM [9]	16 channel weak features	66.8 \pm 0.6
2	SPM [9]	SIFT, 200 cluster centers	81.1 \pm 0.3
2	SPM [9]	SIFT, 400 cluster centers	81.4 \pm 0.5
3	SPM [16]	SIFT, 400 concepts	83.3
2	SP-pLSA [7]	SIFT, 1200 pLSA topics	83.7
2	spatial PACT	CENTRIST, 40 eigenvectors	83.88 \pm 0.76

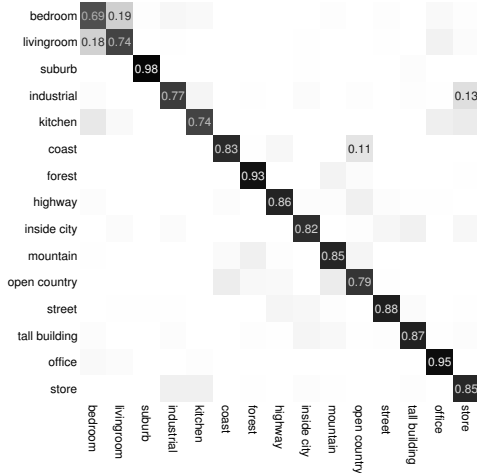


Fig. 13: Confusion matrix of the 15 class scene dataset. Only rates higher than 0.1 are shown in the figure.

features and the strong features with 200 codebook size in [9], and is only inferior to the strong features with 400 codebook size. When a spatial pyramid is used, spatial PACT outperforms other compared methods. Note that length of the spatial PACT feature vector is only about 5% of the SP-pLSA feature vector length in [7].

Confusion matrix from one run on this dataset ($L = 2$ spatial PACT) is shown in Fig. 13, where row and column names are true and predicted labels respectively. The biggest confusion happens between category pairs such as bedroom/living room, industrial/store, and coast/open country, which coincides well with the confusion distribution in [9].

More experiments were also carried out to compare our CENTRIST based descriptor with other descriptors, and to examine various aspects of the scene recognition problem.

Orientation Histogram. Orientation histogram [37] is a representation that uses histogram of quantities computed from 3×3 neighborhoods. We implemented this method with 40 bins. Combined with a level 2 spatial pyramid, Orientation Histogram achieves $76.78 \pm 0.90\%$ recognition rate, which is significantly worse than spatial PACT ($83.88 \pm 0.76\%$).

Linear classifiers. Linear SVM classifiers are also applied to the scene dataset. They achieve accuracy of 82.54%



Fig. 14: Images from 8 different sports event categories.

and 73.59%, using spatial PACT with $L = 2$ and $L = 0$, respectively. The implication of these results are two fold. First, the difference in performance of RBF kernels and linear kernels are quite small. In all the datasets we experimented with, the difference in recognition rates between these two kernel types are smaller than 2%. This observation suggests that images from the same category are compact in the spatial PACT descriptor space. Second, because of the fast testing speed of linear classifiers and small performance difference, linear SVM classifiers could be used to ensure real-time classification.

Speed and classifier analysis. The time to extract CENTRIST is proportional to the input image size. However, large images can be down-sampled to ensure high speed. Our experiments observed only slight (usually $< 1\%$) performance drop. Also, experiments show that spatial PACT is not sensitive to SVM parameters. $(C, \gamma) = (8, 2^{-7})$ is recommended for RBF kernels with probability output, and $C = 2^{-5}$ for linear SVM.

Effect of extra information. When we removed the extra information (mean and standard deviation of pixel blocks), the average accuracy became $83.46 \pm 0.74\%$, which is slightly lower when the extra information was utilized ($83.88 \pm 0.76\%$).

4.4 The 8 class event dataset

The event dataset [15] contains images of eight sports: badminton, bocce, croquet, polo, rock climbing, rowing, sailing, and snowboarding (see Fig. 14 for example images from each category). In [15], Li and Fei-Fei used this dataset in their attempt to classify these events by integrating scene and object categorizations (i.e. deduce *what* from *where* and *who*). We use this dataset for scene classification purpose only. That is, we classify events by classifying the scenes, and do not attempt to recognize objects or persons.

The images are high resolution ones (from 800x600 to thousands of pixels per dimension). The number of images in each category ranges from 137 to 250. Following [15], we use 70 images per class for training, and 60 for testing. The training images in each train/test split are used to compute the eigenvectors. We use RBF kernel SVM classifiers with level 2 pyramid spatial PACT features in this dataset.

Overall we achieve $78.25 \pm 1.27\%$ accuracy on this dataset. In [15], the scene only model achieved approx-

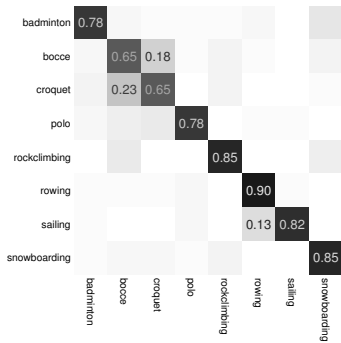


Fig. 15: Confusion matrix of the event dataset. Only rates higher than 0.1 are shown in the figure.

imately 60% accuracy, which is significant lower than the spatial PACT result. When both scene and object categorization were used, the method in [15] had an accuracy of 73.4%, still inferior to our result. Note that this scene+object categorization used manual segmentation and object labels as additional inputs.

The scene only model of spatial PACT exhibits different behaviors than the scene+object model in [15], as shown in the confusion matrix in Fig. 15. The most confusing pairs of our method are bocce/croquet, and rowing/sailing. These results are intuitive because these two pairs of events share very similar scene or background. In [15], the most confusing pairs are bocce/croquet, polo/bocce, and snowboarding/badminton. The object categorization helped in distinguishing rowing and sailing. However, it seems that it also confused events that have distinct backgrounds, such as snowboarding and badminton.

4.5 The 67 class indoor scene recognition dataset

A challenging 67 class indoor scene recognition dataset was proposed in [34]. There are 15620 images in this dataset. The indoor scenes range from specific categories (e.g. dental office) to generic concepts (e.g. mall). It was argued in [34] that both local and global information are needed to recognize complex indoor scenes.

In [34], the global Gist feature achieved about 21% average recognition accuracy on this challenging dataset. When it was supplemented by local information (in the form of local prototypes based on image segmentation), the accuracy was improved to 25%.

The proposed spatial PACT representation achieved higher recognition accuracies on this indoor scene recognition problem. Following [34], we use 80 images in each category for training, and 20 images for testing. The eigenvectors in PACT are computed using the training set in each train/test random split. We used RBF SVM classifier with level 2 spatial PACT. The average recognition accuracy in 5 random split of train/test images is $36.88 \pm 1.10\%$. In this challenging indoor scene recognition problem, spatial PACT achieved much higher accuracies than Gist.

TABLE 5: Comparing recognition accuracies of CENTRIST and Gist in scene recognition datasets.

Dataset	Environment	CENTRIST	Gist
8 class	outdoor	$86.22 \pm 1.02\%$	$82.60 \pm 0.86\%$
15 class	outdoor + indoor	$83.88 \pm 0.76\%$	$73.28 \pm 0.67\%$

4.6 Comparing CENTRIST, SIFT, and Gist

As discussed in Sec. 2 and 4.5, we observe that Gist usually had relatively lower accuracies for complex indoor scenes. Our experiments on the 8 outdoor scene categories [10] and the 15 scene categories (which is a super set of the 8 category dataset) further corroborated this observations. Using the Gist descriptor and SVM classifier³, the recognition accuracy was $82.60 \pm 0.86\%$ on the 8 outdoor categories, which is worse than $86.2 \pm 1.02\%$, the accuracy using CENTRIST on this dataset. However, on the 15 class dataset which adds several indoor categories, the accuracy using Gist dramatically dropped to $73.28 \pm 0.67\%$, which is significantly lower than CENTRIST’s accuracy, $83.88 \pm 0.76\%$. Our conjecture is that the frequency domain features in the Gist descriptor might not be discriminative enough to distinguish between the subtle differences between indoor categories, e.g. bedroom vs. living room. The same procedures and parameters are used in all experiments, except that CENTRIST and Gist are used in different experiments. Table 5 summarizes these results.

On the contrary, SIFT is originally designed to have high discriminative power. Thus it may not be able to cope with the huge intra-class variation in scene images. For any two feature vectors, we can compute their Histogram Intersection Kernel (HIK) value [38] as a simple measure for the similarity between them. By observing the similarity distribution between- and within- categories, which are shown in Fig. 16 for both SIFT and CENTRIST, we can have an estimate of their capability in place and scene recognition.

For any image, we can find its nearest neighbor in the same category and the nearest neighbor in a different category. If the out-of-category nearest neighbor has a higher similarity value than the in-category nearest neighbor, the simple nearest neighbor classifier will make a wrong decision for this image. In Fig. 16 the x-axis shows the difference of these two similarity values. In other words, a value in the left hand side of 0 (the dashed line) means an error. For any given curve, if we find area of the part that is at the left hand side of the black dashed line, and divide it by area of the entire curve, we get the leave one out estimation of the classification error of a nearest neighbor rule. Thus Fig. 16 is an indication of the discriminative power of the descriptors. CENTRIST has a clear advantage in recognizing place and scene images. It has 35.83% 1-NN leave-one-out error in the 15 class scene recognition

³. RBF kernel was used and parameters were set by cross-validation. Gist features were constructed using the code from <http://people.csail.mit.edu/torr/alba/code/spatialenvelope/>.

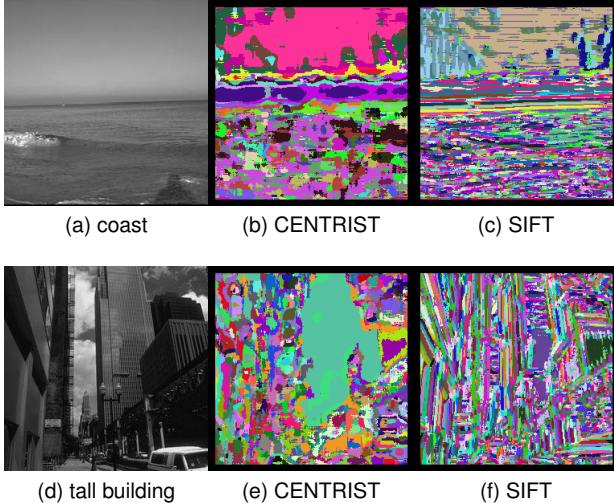


Fig. 17: Visualization of images mapped to code words. In each row, the first image is an input image, with the second and third being visualization for CENTRIST and SIFT codebooks, respectively. (This picture needs to be viewed in color.)

dataset, compared to 57.24% for SIFT. However, SIFT is suitable for object recognition (67.39% error in the Caltech 101 dataset, compared to 83.80% for CENTRIST).

4.7 Bag of Visual words with CENTRIST

Since CENTRIST can be extracted for any rectangle image patches, we can also apply the Bag of Visual words framework with CENTRIST being the base visual descriptor. This is the second approach to use CENTRIST in this paper.⁴ Following [9], we use image patches of size 16 by 16, and sample over a grid with a spacing of 8 pixels. In every train/test split, one fourth of the image patches sampled from the training set are used to generate a codebook which contain 200 code words. Since CENTRIST is only 256 dimensions, PCA operations are not performed (i.e. CENTRIST is directly used for each 16 by 16 image patch). The k-means algorithm is used to cluster CENTRIST vectors into 200 code words. For a level 2 spatial hierarchy, the final feature vector has a length of $200 \times (25 + 4 + 1) = 6200$. SVM classifiers with the histogram intersection kernel are used.

On the 15 class scene recognition dataset, codebook of CENTRIST correctly recognize $80.73 \pm 0.59\%$ of the testing images, which is similar to the result of codebook with 200 SIFT code words in [9] ($81.1 \pm 0.3\%$), but inferior to the spatial PACT result ($83.10 \pm 0.60\%$). Similarly, on the 8 sports event dataset, codebook of CENTRIST achieved an accuracy of $75.21 \pm 1.06\%$, which is lower than the spatial PACT accuracy, but higher than those reported in [15].

4. Source code is available at <http://c2inet.sce.ntu.edu.sg/Jianxin/projects/libHIK/libHIK.htm>

Although the CENTRIST visual codebook’s performance is not as good as spatial PACT, it provides a way to visualize the behavior of the CENTRIST descriptor, and consequently improve our understanding of CENTRIST. We build a visual codebook with 256 visual code words using the 15 class scene recognition dataset. Given an input image, an image patch with coordinates $[x - 8, x + 8] \times [y - 8, y + 8]$ can be mapped to a single integer by the following procedure. We first extract the CENTRIST descriptor from this window (whose size is 16 by 16). This CENTRIST vector is compared to all code words, and the index of the nearest neighbor is the mapping result for pixel position (x, y) . By choosing a random RGB tuple for each code word index, a gray scale image can be transformed into a visualization of corresponding code word indexes.

Fig. 17 are examples of the codebook visualization results for a coast and a tall building image. The SIFT code words tend to emphasize discontinuities in the images. Edges (especially straight lines) usually are mapped to the same code word (i.e. displayed in the same color in the visualization). The visualization also suggests that SIFT pays more attention to detailed textural information, because the visualization is fragmented (connected component of the same color is small). Image patches with similar visual structure and semantics are mapped to different visual code words, e.g. the tall building in the right half of Fig. 17d.

Instead, CENTRIST visualizations tend to group image regions with similar visual structure into the same code word. The connected component in CENTRIST visualizations are larger than those in the SIFT visualizations. For example, the sky in the coast image share similar semantics and visual structures. This region is mostly mapped to the same color (i.e. same code word) using CENTRIST, which is desirable for the scene category recognition task. Instead, the SIFT descriptor maps this region to different colors.

The different behaviors of CENTRIST and SIFT might be explained by the way local image measurements are accumulated. In CENTRIST, we only concern whether a center pixel’s intensity is higher or lower than its neighbors. The magnitude of difference of pixel intensities is ignored. On the contrary, visual descriptors like SIFT and HOG accumulate orientation gradients of pixel intensities. The magnitude of pixel intensity differences has strong effect on the histogram of orientation gradients. Thus, we conjecture that SIFT and HOG are more sensitive to smaller changes of visual contents than CENTRIST. In scene recognition we want our descriptors to be insensitive to small variations in images.

5 CONCLUSIONS

In this paper we propose CENTRIST, CENSus TRansform HISTogram, as a visual descriptor for recognizing places and scene categories. We first show that place and scene recognition pose different requirement for a

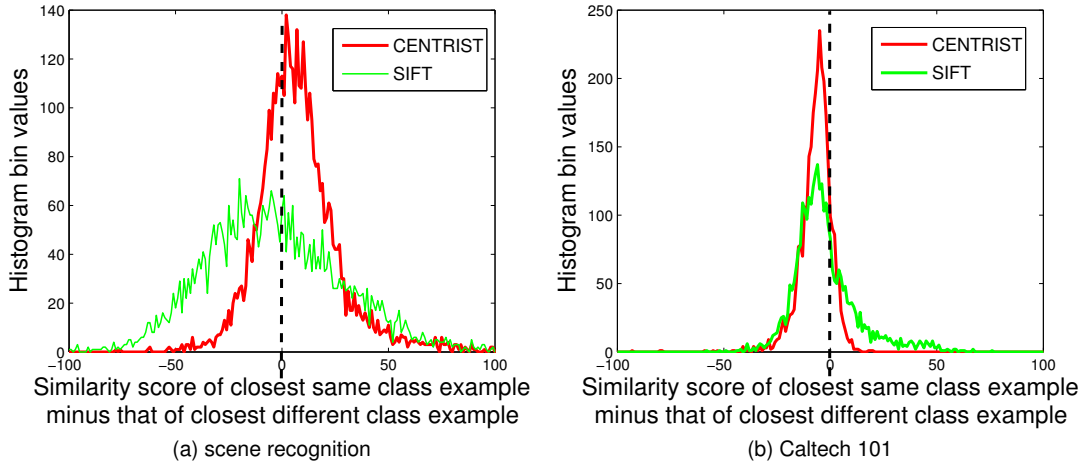


Fig. 16: Histogram comparing similarity values of best in-category nearest neighbor with best out-of-category nearest neighbor of an image.

visual descriptor, especially for such tasks in indoor environments. Thus we need a visual descriptor that is different from commonly used ones (e.g. SIFT in object recognition). We analyze these tasks and show that the descriptor needs to be holistic and generalizable. It also needs to acquire structural properties in the image while suppressing textural details, and contain rough geometrical information in the scene.

We then focus on understanding the properties of CENTRIST, and show how CENTRIST suits the place and scene recognition domain. CENTRIST is a holistic representation that captures the structural properties of an image. Through the strong constraints among neighboring Census Transform values, CENTRIST is able to capture the structural characteristic within a small image patch. In larger scales, spatial hierarchy of CENTRIST is used to catch rough geometrical information. CENTRIST also shows high generalizability, exhibiting similar visual descriptors for images with similar structures.

On five datasets including both place and scene category recognition tasks, CENTRIST achieves higher accuracies than previous state-of-the-art methods. Comparing with SIFT and Gist, CENTRIST not only exhibits superior performance. It is easy to implement, and evaluates extremely fast. Implementation of methods proposed in this paper is publicly available.

In this paper we also analyzed several limitations of CENTRIST and there are research directions that may improve it. First, CENTRIST is not invariant to rotations. Although robot acquired images and scene images are usually upright, making it rotational invariant will enlarge its application area. Second, we want to recognize place categories in more realistic settings, i.e. learning the category concepts using images acquired without human effort in acquiring canonic views. Third, CENTRIST now only utilize the gray scale information in images. As shown in [39], different channels in the color space con-

tain useful information for object and place recognition. The performance of CENTRIST should improve if color channels are incorporated appropriately.

ACKNOWLEDGMENTS

The authors would like to thank Henrik I. Christensen and Aaron Bobick for fruitful discussions, and to thank the anonymous reviewers for their comments and suggestions. J. Wu is supported by the NTU startup grant and the Singapore MoE AcRF Tier 1 project RG 34/09.

REFERENCES

- [1] B. Kuipers and P. Beeson, "Bootstrap learning for place recognition," in *AAAI Conference on Artificial Intelligence*, 2002, pp. 174–180.
- [2] S. Thrun, D. Fox, W. Burgard, and F. Dellaert, "Robust Monte Carlo localization for mobile robots," *Artificial Intelligence*, vol. 128, no. 1-2, pp. 99–141, 2001.
- [3] H. Durrant-Whyte and T. Bailey, "Simultaneous localization and mapping: part I," *IEEE Robotics & Automation Magazine*, vol. 13, no. 2, pp. 99–108, 2006.
- [4] S. Se, D. G. Lowe, and J. J. Little, "Vision-based mobile robot localization and mapping using scale-invariant features," in *Proc. IEEE Int'l Conf. Robotics and Automation*, 2001, pp. 2051–2058.
- [5] I. Ulrich and I. R. Nourbakhsh, "Appearance-based place recognition for topological localization," in *Proc. IEEE Int'l Conf. Robotics and Automation*, 2006, pp. 1023–1029.
- [6] H. Choset and K. Nagatani, "Topological simultaneous localization and mapping (SLAM): toward exact localization without explicit localization," *IEEE Trans. on Robotics and Automation*, vol. 17, no. 2, pp. 125–137, 2001.
- [7] A. Bosch, A. Zisserman, and X. Muñoz, "Scene classification using a hybrid generative/discriminative approach," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 30, no. 4, pp. 712–727, 2008.
- [8] L. Fei-Fei and P. Perona, "A bayesian hierarchical model for learning natural scene categories," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, vol. II, 2005, pp. 524–531.
- [9] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, vol. II, 2006, pp. 2169–2178.
- [10] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *International Journal of Computer Vision*, vol. 42, no. 3, pp. 145–175, 2001.

- [11] P. Quelhas, F. Monay, J.-M. Odobez, D. Gatica-Perez, and T. Tuytelaars, "A thousand words in a scene," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 29, no. 9, pp. 1575–1589, 2007.
- [12] J. Wu, H. I. Christensen, and J. M. Rehg, "Visual Place Categorization: Problem, Dataset, and Algorithm," in *Proc. IEEE/RSJ Int'l Conf. Intelligent Robots and Systems*, 2009.
- [13] J. Hays and A. A. Efros, "IM2GPS: estimating geographic information from a single image," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [14] J. J. Kivinen, E. B. Sudderth, and M. I. Jordan, "Learning multi-scale representaiton of natural scenes using dirichlet processes," in *The IEEE Conf. on Computer Vision*, 2007.
- [15] L.-J. Li and L. Fei-Fei, "What, where and who? Classifying events by scene and object recognition," in *The IEEE Conf. on Computer Vision*, 2007.
- [16] J. Liu and M. Shah, "Scene modeling using Co-Clustering," in *The IEEE Conf. on Computer Vision*, 2007.
- [17] D. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [18] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, vol. 1, 2005, pp. 886–893.
- [19] J. Wu and J. M. Rehg, "Where am I: Place instance and category recognition using spatial PACT," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [20] A. Pronobis, B. Caputo, P. Jensfelt, and H. I. Christensen, "A discriminative approach to robust visual place recognition," in *Proc. IEEE/RSJ Int'l Conf. Intelligent Robots and Systems*, 2006.
- [21] M. Szummer and R. W. Picard, "Indoor-outdoor image classification," in *CAIVD*, 1998, pp. 42–51.
- [22] J. C. van Gemert, J.-M. Geusebroek, C. J. Veenman, and A. W. Smeulders, "Kernel codebooks for scene categorization," in *European Conf. Computer Vision*, 2008.
- [23] K. Mikolajczyk and C. Schmid, "A performance evaluation of local descriptors," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 27, no. 10, pp. 1615–1630, 2005.
- [24] J. Vogel and B. Schiele, "Semantic modeling of natural scenes for content-based image retrieval," *International Journal of Computer Vision*, vol. 72, no. 2, pp. 133–157, 2007.
- [25] R. Fergus, P. Perona, and A. Zisserman, "Object class recognition by unsupervised scale-invariant learning," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, vol. II, 2003, pp. 264–271.
- [26] P. Felzenszwalb and D. Huttenlocher, "Pictorial structures for object recognition," *International Journal of Computer Vision*, vol. 61, no. 1, pp. 55–79, 2005.
- [27] J. Luo, A. Pronobis, B. Caputo, and P. Jensfelt, "The KTH-IDOL2 database," Kungliga Tekniska Hoegskolan, CVAP/CAS, Tech. Rep. CVAP304, October 2006.
- [28] R. Zabih and J. Woodfill, "Non-parametric local transforms for computing visual correspondence," in *European Conf. Computer Vision*, vol. 2, 1994, pp. 151–158.
- [29] T. Ojala, M. Pietikäinen, and T. Mäenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 971–987, 2002.
- [30] L. Fei-Fei, R. Fergus, and P. Perona, "Learning generative visual models from few training example: an incremental bayesian approach tested on 101 object categories," in *CVPR 2004, Workshop on Generative-Model Based Vision*, 2004.
- [31] P. F. Felzenszwalb and J. D. Schwartz, "Hierarchical matching of deformable shapes," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2007.
- [32] H. Ling and D. W. Jacobs, "Shape classification using the inner-distance," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 29, no. 2, pp. 286–299, 2007.
- [33] O. J. O. Söderkvist, "Computer vision classification of leaves from swedish trees," Master's thesis, Linköping University, 2001.
- [34] A. Quattoni and A. Torralba, "Recognizing indoor scenes," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2009.
- [35] A. Pronobis and B. Caputo, "The KTH-INDECS database," Kungliga Tekniska Hoegskolan, CVAP, Tech. Rep. CVAP297, September 2005.
- [36] J. Wolf, W. Burgard, and H. Burkhardt, "Robust vision-based localization for mobile robots using an image retrieval system based on invariant features," in *Proc. IEEE Int'l Conf. Robotics and Automation*, 2002, pp. 359–365.
- [37] W. T. Freeman and M. Roth, "Orientation histogram for hand gesture recognition," in *FG workshop*, 1995, pp. 296–301.
- [38] M. J. Swain and D. H. Ballard, "Color indexing," *International Journal of Computer Vision*, vol. 7, no. 1, pp. 11–32, 1991.
- [39] K. E. A. van de Sande, T. Gevers, and C. G. M. Snoek, "Evaluation of color descriptors for objects and scene recognition," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.



Jianxin Wu received the BS degree and MS degree in computer science from the Nanjing University, and his PhD degree in computer science from the Georgia Institute of Technology. He is currently an assistant professor in the School of Computer Engineering, Nanyang Technological University, Singapore. His research interests are computer vision, machine learning, and robotics. He is a member of the IEEE.



James M. Rehg received his PhD degree in Electrical and Computer Engineering from the Carnegie Mellon University. He is a Professor in the College of Computing at the Georgia Institute of Technology. He is a member of the Graphics, Visualization, and Usability Center and co-directs the Computational Perception Lab. His research interests are computer vision, robotics, machine learning, and computer graphics. He is a member of the IEEE.