

Deep Texture Manifold for Ground Terrain Recognition

Jia Xue¹ Hang Zhang^{1,2} Kristin Dana¹

¹Department of Electrical and Computer Engineering, Rutgers University, ²Amazon AI

{jia.xue, zhang.hang}@rutgers.edu, kdana@ece.rutgers.edu

Abstract

We present a texture network called Deep Encoding Pooling Network (DEP) for the task of ground terrain recognition. Recognition of ground terrain is an important task in establishing robot or vehicular control parameters, as well as for localization within an outdoor environment. The architecture of DEP integrates orderless texture details and local spatial information and the performance of DEP surpasses state-of-the-art methods for this task. The GTOS database (comprised of over 30,000 images of 40 classes of ground terrain in outdoor scenes) enables supervised recognition. For evaluation under realistic conditions, we use test images that are not from the existing GTOS dataset, but are instead from hand-held mobile phone videos of similar terrain. This new evaluation dataset, GTOS-mobile, consists of 81 videos of 31 classes of ground terrain such as grass, gravel, asphalt and sand. The resultant network shows excellent performance not only for GTOS-mobile, but also for more general databases (MINC and DTD). Leveraging the discriminant features learned from this network, we build a new texture manifold called DEP-manifold. We learn a parametric distribution in feature space in a fully supervised manner, which gives the distance relationship among classes and provides a means to implicitly represent ambiguous class boundaries. The source code and database are publicly available¹.

1. Introduction

Ground terrain recognition is an important area of research in computer vision for potential applications in autonomous driving and robot navigation. Recognition with CNNs have achieved success in object recognition and the CNN architecture balances preservation of relative spatial information (with convolutional layers) and aggregation of spatial information (pooling layers). This structure is de-



Figure 1: Homogeneous textures (upper row) compared to more common real-world instances with local spatial structure that provides an important cue for recognition (lower row).

signed for object recognition, scene understanding, face recognition, and applications where spatial order is critical for classification. However, texture recognition uses an orderless component to provide invariance to spatial layout [5, 18, 38].

In classic approaches for texture modeling, images are filtered with a set of handcrafted filter banks followed by grouping the outputs into texton histograms [7, 16, 21, 34], or bag-of-words [6, 15]. Later, Cimpoi *et al.* [5] introduce FV-CNN that replace the handcrafted filter banks with pre-trained convolutional layers for the feature extractor, and achieve state-of-the-art results. Recently, Zhang *et al.* [38] introduce Deep Texture Encoding Network (Deep-TEN) that ports the dictionary learning and feature pooling approaches into the CNN pipeline for an end-to-end material/texture recognition network. Recognition algorithms that focus on texture details work well for images containing only a single material. But for “images in the wild”, homogeneous surfaces rarely fill the entire field-of-view, and

¹<http://ece.rutgers.edu/vision/>

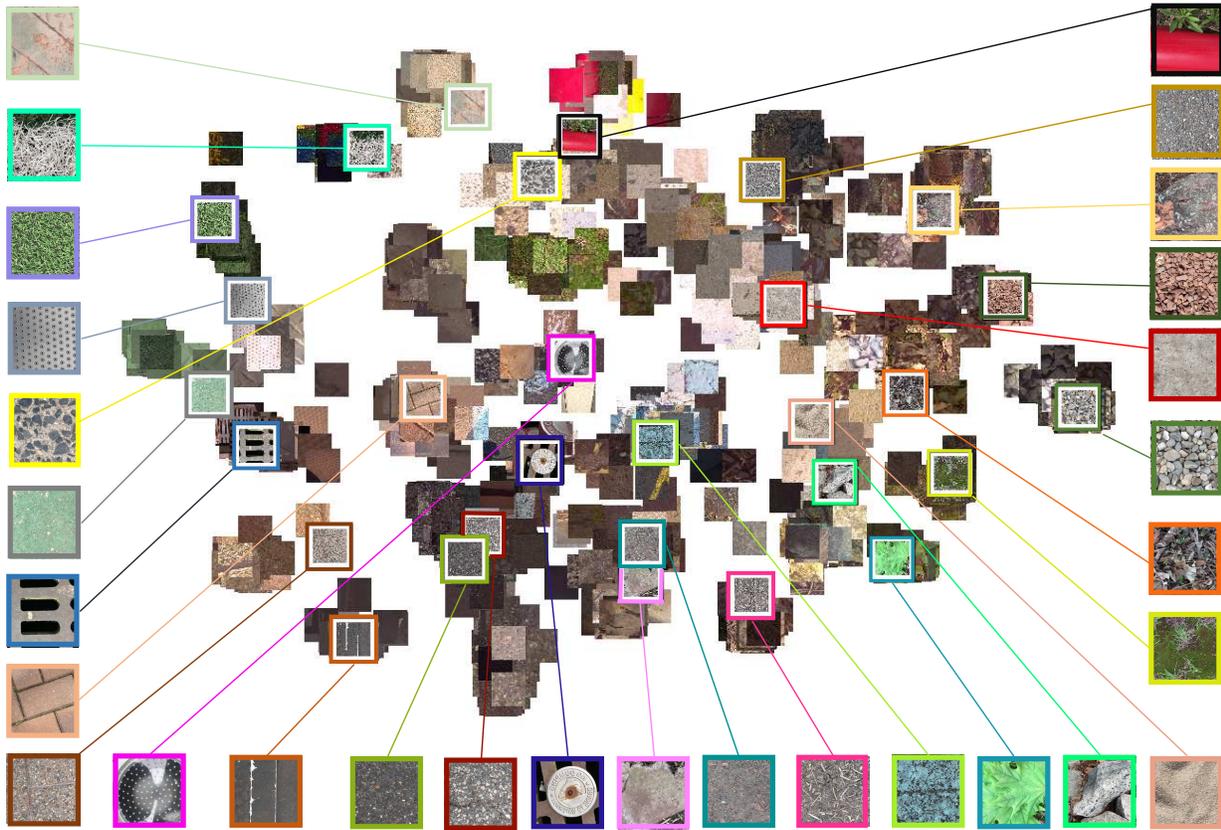


Figure 2: The result of texture manifold by DEP-manifold. Images with color frames are images in test set. The material classes are (from upper left to counter-clockwise): plastic cover, painted turf, turf, steel, stone-cement, painted cover, metal cover, brick, stone-brick, glass, sandpaper, asphalt, stone-asphalt, aluminum, paper, soil, mulch, painted asphalt, leaves, limestone, sand, moss, dry leaves, pebbles, cement, shale, roots, gravel and plastic. Not all classes are shown here for space limitations.

many materials exhibit regular structure.

For texture recognition, since surfaces are not completely orderless, *local spatial order* is an important cue for recognition as illustrated in Figure 1. Just as semantic segmentation balances local details and global scene context for pixelwise recognition [1, 17, 25, 29, 30, 40], we design a network to balance both an orderless component and ordered spatial information.

As illustrated in Figure 3, we introduce a Deep Encoding Pooling Network (DEP) that leverages an orderless representation and local spatial information for recognition. Outputs from convolutional layers are fed into two feature representation layers jointly; the encoding layer [38] and the global average pooling layer. The encoding layer is employed to capture texture appearance details and the global average pooling layer accumulates spatial information. Features from the encoding layer and the global average pooling layer are processed with bilinear models [31]. We apply

DEP to the problem of ground terrain recognition using an extended GTOS dataset [36]. The resultant network shows excellent performance not only for GTOS, but also for more general databases (MINC [1] and DTD [4]).

For ground terrain recognition, many class boundaries are ambiguous. For example, “asphalt” class is similar to “stone-asphalt” which is an aggregate mix of stone and asphalt. The class “leaves” is similar to “grass” because most of the example images for “leaves” in the GTOS database have grass in the background. Similarly, the grass images contain a few leaves. Therefore, it is of interest to find not only the class label but also the closest classes, or equivalently, the position in the manifold. We introduce a new texture manifold method, DEP-manifold, to find the relationship between newly captured images and images in dataset.

The t-Distributed Stochastic Neighbor Embedding (t-SNE) [20] provides a 2D embedding and Barnes-Hut t-SNE [33] accelerates the original t-SNE from $\mathcal{O}(n^2)$ to

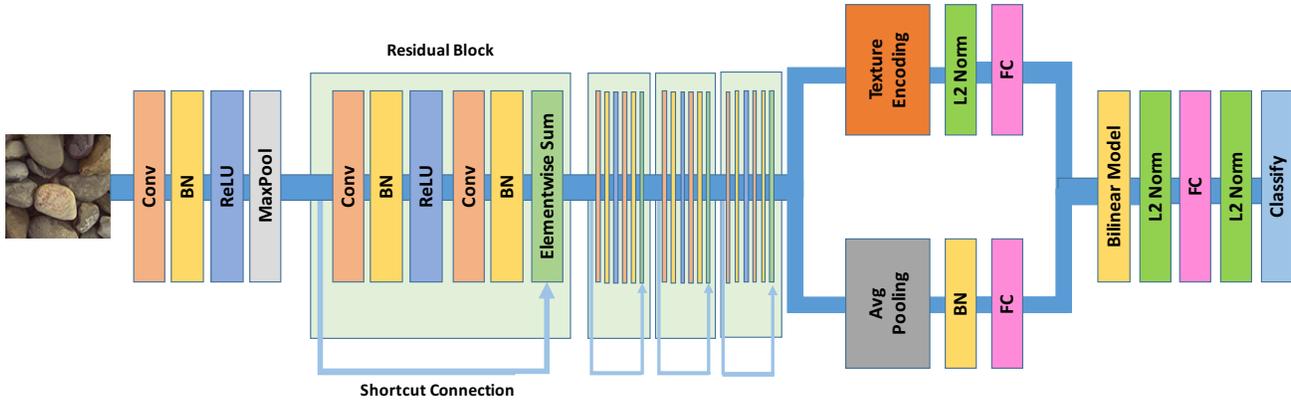


Figure 3: A Deep Encoding Pooling Network (DEP) for material recognition. Outputs from convolutional layers are fed into the encoding layer and global average pooling layer jointly and their outputs are processed with bilinear model.

$\mathcal{O}(n \log n)$. Both t-SNE and Barnes-Hut t-SNE are non-parametric embedding algorithms, so there is no natural way to perform out-of-sample extension. Parametric t-SNE [32] and supervised t-SNE [23, 24] introduce deep neural networks into data embedding and realize non-linear parametric embedding. Inspired by this work, we introduce a method for texture manifolds that treats the embedded distribution from non-parametric embedding algorithms as an output, and use a deep neural network to predict the manifold coordinates of a texture image directly. This texture manifold uses the features of the DEP network and is referred to as DEP-manifold.

The training set is a ground terrain database (GTOS) [36] with 31 classes of ground terrain images (over 30,000 images in the dataset). Instead of using images from the GTOS dataset for testing, we collect GTOS-mobile, 81 ground terrains videos of similar terrain classes captured with a handheld mobile phone and with arbitrary lighting/viewpoint. Our motivation is as follows: The training set (GTOS) is obtained in a comprehensive manner (known distance and viewpoints, high-res calibrated camera) and is used to obtain knowledge of the scene. The test set is obtained under very different and more realistic conditions (a mobile imaging device, handheld video, uncalibrated capture). Training with GTOS and testing with GTOS-mobile enables evaluation of knowledge transfer of the network.

2. Related Work

Tenenbaum and Freeman [31] introduce bilinear models to process two independent factors that underly a set of observations. Lin *et al.* [19] introduce the Bilinear CNN models that use outer product of feature maps from convolutional layers of two CNNs and reach state-of-the-art for fine grained visual recognition. However, this method has two drawbacks. First, bilinear models for feature maps from convolutional layers require that pairs of features maps have

compatible feature dimensions, i.e. the same height and width. The second drawback is computational complexity; this method computes the outer product at each location of the feature maps. To utilize the advantage of bilinear models and overcome these drawbacks, we employ bilinear models for outputs from fully connected layers. Then, outputs from fully connected layers can be treated as vectors, and there is no dimensionality restriction for the outer product of two vectors.

Material recognition is a fundamental problem in computer vision, the analysis of material recognition has varied from small sets collected in lab settings such as KTH-TIPS [3] and CuRET [8], to large image sets collected in the wild [1, 4, 35, 36]. The size of material datasets have also increased from roughly 100 images in each class [4, 35] to over 1000 images in each class [1, 36]. The Ground Terrain in Outdoor Scenes (GTOS) dataset has been used with angular differential imaging [36] for material recognition based on angular gradients. For our work, single images are used for recognition without variation in viewing direction, so reflectance gradients are not considered.

For many recognition problems, deep learning has achieved great success, such as face recognition [2, 22, 26, 37], action recognition [41, 42] and disease diagnosis [39]. The success of deep learning has also transferred to material recognition. We leverage a recent texture encoding layer [38] that ports dictionary learning and residual encoding into CNNs. We use this texture encoding layer as a component in our network to capture orderless texture details.

3. Deep Encoding Pooling Network

Encoding Layer The texture encoding layer [38] integrates the entire dictionary learning and visual encoding pipeline into a single CNN layer, which provides an orderless representation for texture modeling. The encoding

layer acts as a global feature pooling on top of convolutional layers. Here we briefly describe prior work for completeness. Let $X = \{x_1, \dots, x_m\}$ be M visual descriptors, $C = \{c_1, \dots, c_n\}$ is the code book with N learned codewords. The residual vector r_{ij} is calculated by $r_{ij} = x_i - c_j$, where $i = 1 \dots m$ and $j = 1 \dots n$. The residual encoding for codeword c_j can be represented as

$$e_j = \sum_{i=1}^m w_{ij} r_{ij}, \quad (1)$$

where w_{ij} is the assigning weight for residual vector r_{ij} and is given by

$$w_{ij} = \frac{\exp(-s_j \|r_{ij}\|^2)}{\sum_{k=1}^m \exp(-s_k \|r_{ik}\|^2)}, \quad (2)$$

s_1, \dots, s_m are learnable smoothing factors. With the texture encoding layer, the visual descriptors X are pooled into a set of N residual encoding vectors $E = \{e_1, \dots, e_n\}$. Similar to classic encoders, the encoding layer can capture more texture details by increasing the number of learnable codewords.

Bilinear Models Bilinear models are two-factor models such that their outputs are linear in one factor if the other factor is constant [10]. The factors in bilinear models balance the contributions of the two components. Let a^t and b^s represent the material texture information and spatial information with vectors of parameters and with dimensionality I and J . The bilinear function Y^{ts} is given by

$$Y^{ts} = \sum_{i=1}^I \sum_{j=1}^J w_{ij} a_i^t b_j^s, \quad (3)$$

where w_{ij} is a learnable weight to balance the interaction between material texture and spatial information. The outer product representation captures a pairwise correlation between the material texture encodings and spatial observation structures.

Deep Encoding Pooling (DEP) Network Our Deep Encoding Pooling Network (DEP) is shown in Figure 3. As in prior transfer learning algorithms [19, 38], we employ convolutional layers with non-linear layers from ImageNet [9] pre-trained CNNs as feature extractors. Outputs from convolutional layers are fed into the texture encoding layer and the global average pooling layer jointly. Outputs from the texture encoding layer preserve texture details, while outputs from the global average pooling layer preserve local spatial information. The dimension of outputs from the texture encoding layer is determined by the codewords N and

the feature maps channel C ($N \times C$). The dimension of outputs from the global average pooling layer is determined by the feature maps channel C . For computational efficiency and to robustly combine feature maps with bilinear models, we reduce feature maps dimension with fully connected layers for both branches. Feature maps from the texture encoding layer and the global average pooling layer are processed with a bilinear model and followed by a fully connected layer and a classification layer with non-linearities for classification. Table 1 is an instantiation of DEP based on 18-layer ResNet [12]. We set 8 codewords for the texture encoding layer. The size of input images are 224×224 . Outputs from CNNs are fed into the texture encoding layer and the global average pooling layer jointly. The dimension of outputs from the texture encoding layer is $8 \times 512 = 4096$ and the dimension of outputs from global average pooling layer is 512. We reduce the dimension of feature maps from the deep encoding layer and the global average pooling layer to 64 via fully connected layers. The dimension of outputs from bilinear model is $64 \times 64 = 4096$. Following prior works [27, 38], resulting vectors from the texture encoding layer and bilinear model are normalized with L2 normalization.

The texture encoding layer and bilinear models are both differentiable. The overall architecture is a directed acyclic graph and all the parameters can be trained by back propagation. Therefore, the Deep Encoding Pooling Network is trained end-to-end using stochastic gradient descent with back-propagation.

4. Recognition Experiments

We compare the DEP network with the following three methods based on ImageNet [28] pre-trained 18-layer ResNet [12]: (1) CNN with ResNet, (2) CNN with Deep-Ten and (3) CNN with bilinear models. All three methods support end-to-end training. For equal comparison, we use an identical training and evaluation procedure for each experiment.

CNN with global average pooling (ResNet) We follow the standard procedure to fine-tune pre-trained ResNet, by replacing the last 1000-way fully connected layer with the output dimension of 31. The global average pooling works as feature pooling that encodes the $7 \times 7 \times 512$ dimensional features from the 18-layer pre-trained ResNet into a 512 dimensional vector.

CNN with texture encoding (Deep-TEN) The Deep Texture Encoding Network (Deep-TEN) [38] embeds the texture encoding layer on top of the 50-layer pre-trained ResNet [12]. To make an equal comparison, we replace the 50-layer ResNet with 18-layer ResNet. Same as [38], we



Figure 4: Comparison of images from the GTOS dataset (left) and GTOS-mobile (right) video frames. The training set is the ground terrain database (GTOS) with 31 classes of ground terrain images (over 30,000 images in the dataset). GTOS is collected with calibrated viewpoints. GTOS-mobile, consists of 81 videos of similar terrain classes captured with a handheld mobile phone and with arbitrary lighting/viewpoint. A total of 6066 frames are extracted from the videos with a temporal sampling of approximately 1/10th seconds. The figure shows individual frames of 31 ground terrain classes.

layer name	output size	encoding-pooling
conv1	$112 \times 112 \times 64$	7×7 , stride 2
conv2_x	$56 \times 56 \times 64$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$
conv3_x	$28 \times 28 \times 128$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$
conv4_x	$14 \times 14 \times 256$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$
conv5_x	$7 \times 7 \times 512$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$
encoding / pooling	$8 \times 512 / 512$	8 codewords / ave pool
fc1_1 / fc1_2	$64 / 64$	$4096 \times 64 / 512 \times 64$
bilinear mapping	4096	-
fc2	128	4096×128
classification	n classes	$128 \times n$

Table 1: The architecture of Deep Encoding Pooling Network based on 18-layer ResNet [12]. The input image size is 224×224 .

reduce the number of CNN streams outputs channels from 512 to 128 with a 1×1 convolutional layer. We replace the global average pooling layer in the 18-layer ResNet with texture encoding layer, set the number of codewords to 32 for experiments. Outputs from the texture encoding layer are normalized with L2 normalization. A fully connected layer with soft max loss follows the texture encoding layer

for classification.

CNN with bilinear models (Bilinear-CNN) Bilinear-CNN [19] employs bilinear models with feature maps from convolutional layers. Outputs from convolutional layers of two CNN streams are multiplied using outer product at each location and pooled for recognition. To make an equal comparison, we employ the 18-layer pre-trained ResNet as CNN streams for feature extractor. Feature maps from the last convolutional layer are pooled with bilinear models. The dimension of feature maps for bilinear models is $7 \times 7 \times 512$ and the pooled bilinear feature is of size 512×512 . The pooled bilinear feature is fed into classification layer for classification.

4.1. Dataset and Evaluation

Dataset Extending the GTOS database [36], we collect *GTOS-mobile* consisting of 81 videos obtained with a mobile phone (Iphone SE) and extract 6066 frames as a test set. To simulate real world ground terrain collection, we walk through similar ground terrain regions in random order to collect the videos. Scale is changed arbitrarily by moving far or close and changes in viewing direction are obtained by motions in a small arc. The resolution of the videos is 1920×1080 , and we resize the short edge to 256 while keeping the aspect ratio for experiments. As a result, the resolution of the resized images are 455×256 . Some materials in GTOS were not accessible due to weather, therefore we removed the following classes: dry grass, ice mud, mud-puddle, black ice and snow from the GTOS dataset. Additionally, we merged very similar classes of asphalt and metal. The original GTOS set is 40 classes, as shown in Figure 4, there are 31 material classes in the modified dataset. The class names are (in the order of top-left to bottom-right): asphalt, steel, stone-cement, glass, leaves, grass,

	ResNet [12]	Bilinear CNN [19]	Deep-TEN [38]	DEP (ours)
Single scale	70.82	72.03	74.22	76.07
Multi scale	73.16	75.43	76.12	82.18

Table 2: Comparison our Deep Encoding Pooling Network (DEP) with ResNet (left) [12], Bilinear CNN (mid) [19] and Deep-TEN (right) [38] on GTOS-mobile dataset with single scale and multi scale training. For ResNet, we replace the 1000-way classification layer with a new classification layer, the output dimension of new classification layer is 31.

plastic cover, turf, paper, gravel, painted turf, moss, cloth, stone-asphalt, dry leaves, mulch, cement, pebbles, sand-paper, roots, plastic, stone-brick, painted cover, limestone, soil, sand, shale, aluminum, metal cover, brick, painted asphalt.

Multi-scale Training Images in the GTOS dataset were captured from a fixed distance between the camera and ground terrain, however the distance between the camera and ground terrain can be arbitrary in real world applications. We infer that extracting different resolution patches with different aspect ratio from images in GTOS simulate observing materials at different distance and viewing angle will be helpful for recognition. So for image pre-processing, instead of directly resizing the full resolution images into 256×256 as [36], we resize the full resolution images into different scales, and extract 256×256 center patches for experiment. Through empirical validation, we find that resizing the full resolution images into 256×256 , 384×384 and 512×512 works best.

Training procedure We employ an identical data augmentation and training procedure for experiments. For single scale training experiment, we resize the full resolution images into 384×384 and extract 256×256 center patches as training set. For multi scale training experiment, we resize the full resolution images into 256×256 , 384×384 and 512×512 , and extract 256×256 center patches as training set. For the training section data augmentation, following prior work [14, 38], we crop a random size (0.8 to 1.0) of the original size and a random aspect ratio (3/4 to 4/3) of the original aspect ratio, resize the cropped patches to 224×224 for experiment. All images are pre-processed by subtracting a per color channel mean value and normalized to unit variance with a 50% chance horizontal flip. The learning rate of newly added layers is 10 times of the pre-trained layers. The experiment starts with learning rate at 0.01, momentum 0.9, batch size 128; the learning rate decays by factor of 0.1 for every 10 epochs, and is finished after 30 epochs.

4.2. Recognition Results

Evaluation on GTOS-mobile Table 2 is the classification accuracy of fine-tuning ResNet [12], bilinear

Method	DTD [4]	Minc-2500 [1]
FV-CNN [5]	72.3%	63.1%
Deep-TEN [38]	69.6%	80.4%
DEP (ours)	73.2%	82.0%

Table 3: Comparison with state-of-the-art algorithms on Describable Textures Dataset (DTD) and Materials in Context Database (MINC).

CNN [19], Deep-TEN [38] and the proposed DEP on the GTOS-mobile dataset. When comparing the performance of single-scale and multi-scale training, multi-scale training outperforms single-scale training for all approaches. It proves our inference that extracting different resolution patches with different aspect ratio from images in GTOS to simulate observing materials at different distance and viewing angle will be helpful for recognition. The multi-scale training accuracy for combined spatial information and texture details (DEP) is 82.18%. That’s 9.02% better than only focusing on spatial information (ResNet) and 6% better than only focusing on texture details (Deep-TEN). To gain insight into why DEP outperforms ResNet and Deep-TEN for material recognition, we visualize the features before classification layers of ResNet, Deep-TEN and DEP with Barnes-Hut t-SNE [33]. We randomly choose 10000 images from training set for the experiment. The result is shown in Figure 5. Notice that DEP separates classes farther apart and each class is clustered more compactly.

Evaluation on MINC and DTD Dataset To show the generality of DEP for material recognition, we experiment on two other material/texture recognition datasets: Describable Textures Database (DTD) [4] and Materials in Context Database (MINC) [1]. For an equal comparison, we build DEP based on a 50-layer ResNet [12], the feature maps channels from CNN streams are reduced from 2048 to 512 with a 1×1 convolutional layer. The result is shown in Table 3, DEP outperforms the state-of-the-art on both datasets. Note that we only experiment with single scale training. As mentioned in [19], multi-scale training is likely to improve results for all methods.

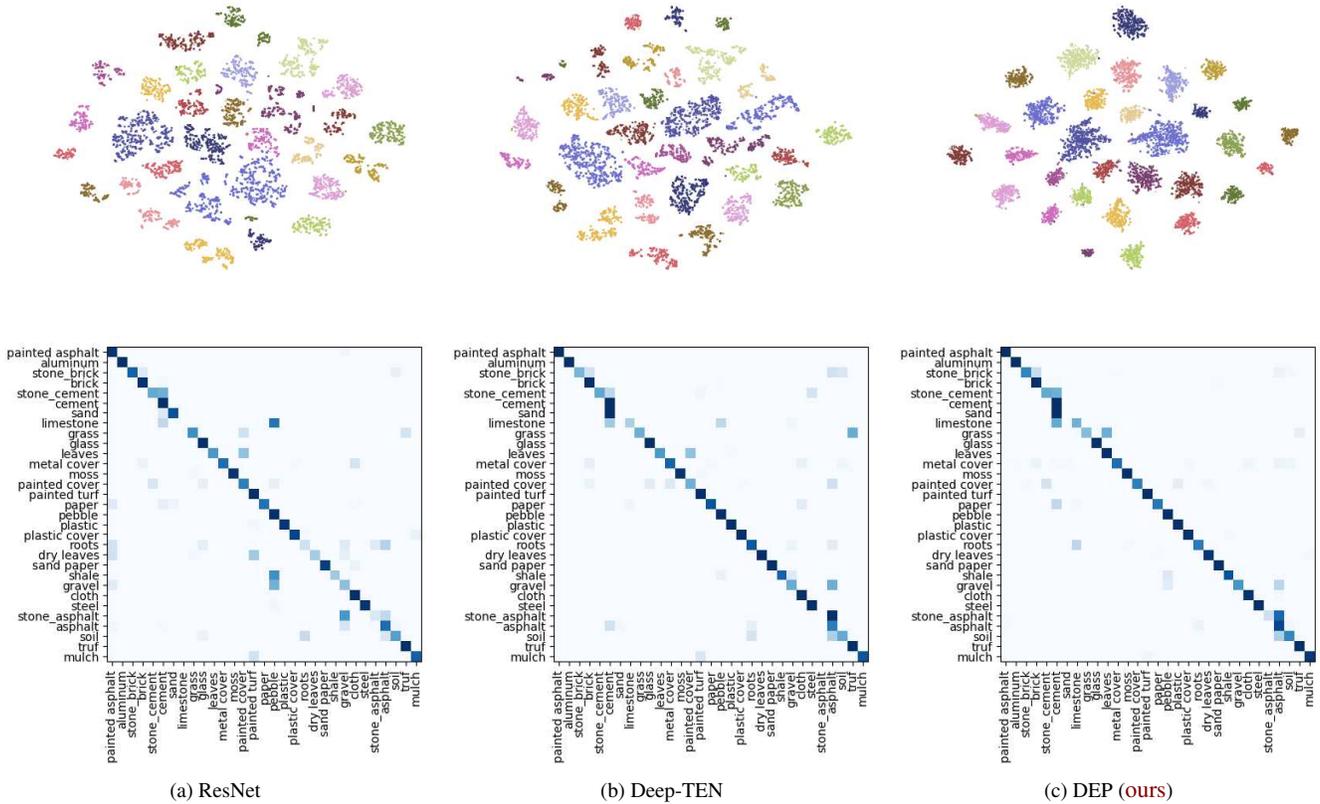


Figure 5: The Barnes-Hut t-SNE [33] and confusion matrix of three material recognition models: ResNet (left), Deep-TEN (mid) and DEP (right). For Barnes-Hut t-SNE, we randomly choose 10000 images from training set and extract features before classification layers of three models for experiment. We see that DEP separates and clusters the classes better. Some classes are misclassified, however, they are typically recognized as a nearby class. (Dark blue represents higher values and light blue represents lower values in the confusion matrix.)

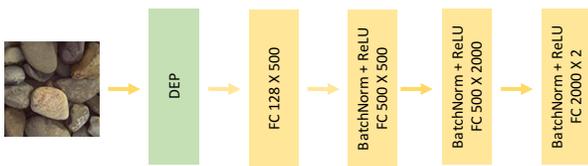


Figure 6: The deep network for texture manifold, we employ DEP as feature extractor, outputs from the last fully connected layer of DEP works as input for texture embedding.

5. Texture Manifold

Inspired by Parametric t-SNE [32] and supervised t-SNE [23, 24], we introduce a parametric texture manifold approach that learns to approximate the embedded distribution of non-parametric embedding algorithms [20, 33] us-

ing a deep neural network to directly predict the 2D manifold coordinates for the texture images. We refer to this manifold learning method using DEP feature embedding as DEP-manifold. Following prior work [24,32], the deep neural network structure is depicted in Figure 6. Input features are the feature maps before the classification layer of DEP, which means each image is represented by a 128 dimensional vector. Unlike the experiment in [24, 32], we add non-linear functions (Batch Normalization and ReLU) before fully connected layers, and we do not pre-train the network with a stack of Restricted Boltzmann Machines (RBMs) [13]. We train the embedding network from scratch instead of the three-stage training procedure (pre-training, construction and fine-tuning) in parametric t-SNE and supervised t-SNE. We randomly choose 60000 images from the multi-scale GTOS dataset for the experiment. We experiment with DEP-parametric t-SNE, and DEP-manifold based on outputs from the last fully connected layer of DEP.

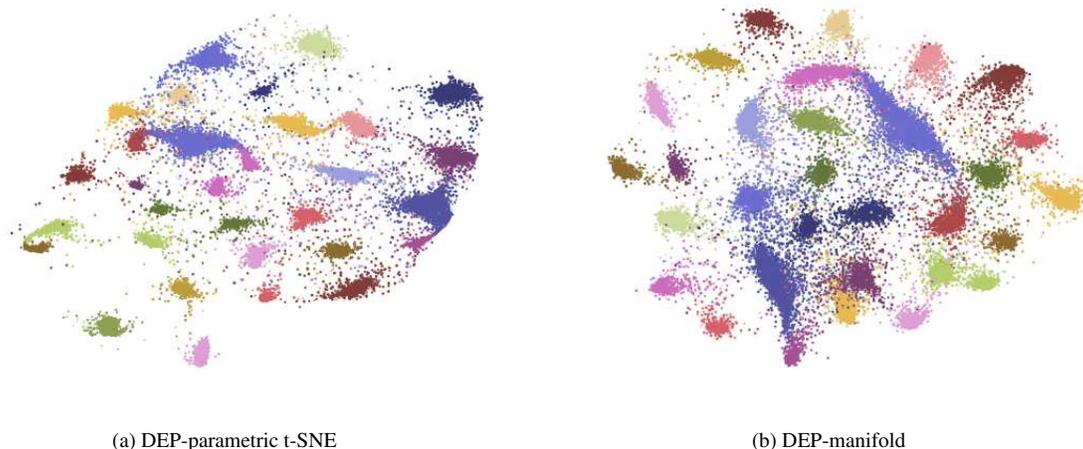


Figure 7: Comparison the performance between DEP-parametric t-SNE and DEP-manifold with 60000 images from multi-scale GTOS dataset. For the embedded distribution of DEP-Parametric t-SNE, the classes are distributed unevenly with crowding in some areas and sparseness in others. The DEP-manifold has a better distribution of classes within the 2D embedding.

Implementation For the DEP-manifold, we employ Barnes-Hut t-SNE [33] as a non-parametric embedding to build the embedded distribution. Following prior setting [33], we set perplexity to 30 and the output dimension of PCA to 50 for the experiment. For training the deep embedding network, we experiment with batch size 2048 and the parameters of the fully connected layers are initialized with the Xavier distribution [11]. We employ L2 loss as the objective function for the experiment. The initial learning rate is 0.01, and decays by a factor of 0.1 every 30 epochs. The experiment is finished after 80 epochs. On an NVIDIA Titan X card, the training takes less than 5 minutes.

Texture Manifold The texture manifold results are shown in Figure 7. For the embedded distribution of DEP-Parametric t-SNE, the classes are distributed unevenly with crowding in some areas and sparseness in others. The DEP-manifold has a better distribution of classes within the 2D embedding. We illustrate the texture manifold embedding by randomly choosing 2000 images from training set to get the embedded distribution; then we embed images from test set into the DEP-manifold. Note that the test set images are not used in the computation of the DEP-manifold. The result is shown in Figure 2. By observing the texture manifold, we find that for some classes, although the recognition accuracy is not perfect, the projected image is within the margin of the correct class, such as cement and stone-cement. Based on the similarity of material classes on the texture manifold, we build the confusion matrix for material recognition algorithms as shown in Figure 5. For visualization, the one dimensional ordering of the confusion matrix

axes are obtained from a one-dimensional embedding of the 2D manifold so that neighboring classes are close. Observe that for the DEP recognition (Figure 5 c), there are very few off-diagonal elements in the confusion matrix. And the off-diagonal elements are often near diagonal indicating find when these images are misclassified, they are recognized as closely-related classes.

6. Conclusion

We have developed methods for recognition of ground terrain for potential applications in robotics and automated vehicles. We make three significant contributions in this paper: 1) introduction of Deep Encoding Pooling network (DEP) that leverages an orderless representation and local spatial information for recognition; 2) Introduction of DEP-manifold that integrates DEP network on top of a deep neural network to predict the manifold coordinates of a texture directly; 3) Collection of the GTOS-mobile database comprised of 81 ground terrains videos of similar terrain classes as GTOS, captured with a handheld mobile phone to evaluate knowledge-transfer between different image capture methods but within the the same domain.

Acknowledgment

This work was supported by National Science Foundation award IIS-1421134. A TITAN X used for this research was donated by the NVIDIA Corporation.

References

- [1] S. Bell, P. Upchurch, N. Snavely, and K. Bala. Material recognition in the wild with the materials in context database. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3479–3487, 2015. [2](#), [3](#), [6](#)
- [2] J. Cai, Z. Meng, A. S. Khan, Z. Li, and Y. Tong. Island loss for learning discriminative features in facial expression recognition. *arXiv preprint arXiv:1710.03144*, 2017. [3](#)
- [3] B. Caputo, E. Hayman, and P. Mallikarjuna. Class-specific material categorisation. In *Computer Vision, 2005. ICCV 2005. Tenth IEEE International Conference on*, volume 2, pages 1597–1604. IEEE, 2005. [3](#)
- [4] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3606–3613, 2014. [2](#), [3](#), [6](#)
- [5] M. Cimpoi, S. Maji, and A. Vedaldi. Deep filter banks for texture recognition and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3828–3836, 2015. [1](#), [6](#)
- [6] G. Csurka, C. R. Dance, L. Fan, J. Willamowski, and C. Bray. Visual categorization with bags of keypoints. In *In Workshop on Statistical Learning in Computer Vision, ECCV*, pages 1–22, 2004. [1](#)
- [7] O. G. Cula and K. J. Dana. Compact representation of bidirectional texture functions. *IEEE Conference on Computer Vision and Pattern Recognition*, 1:1041–1067, December 2001. [1](#)
- [8] K. J. Dana, B. Van Ginneken, S. K. Nayar, and J. J. Koenderink. Reflectance and texture of real-world surfaces. *ACM Transactions On Graphics (TOG)*, 18(1):1–34, 1999. [3](#)
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009. [4](#)
- [10] W. T. Freeman and J. B. Tenenbaum. Learning bilinear models for two-factor problems in vision. In *Computer Vision and Pattern Recognition, 1997. Proceedings., 1997 IEEE Computer Society Conference on*, pages 554–560. IEEE, 1997. [4](#)
- [11] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 249–256, 2010. [8](#)
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [4](#), [5](#), [6](#)
- [13] G. E. Hinton, S. Osindero, and Y.-W. Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006. [7](#)
- [14] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten. Densely connected convolutional networks. *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. [6](#)
- [15] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 2169–2178, 2006. [1](#)
- [16] T. Leung and J. Malik. Representing and recognizing the visual appearance of materials using three-dimensional textures. *International journal of computer vision*, 43(1):29–44, 2001. [1](#)
- [17] G. Lin, C. Shen, A. van den Hengel, and I. Reid. Exploring context with deep structured models for semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP(99):1–1, 2017. [2](#)
- [18] M. Lin, Q. Chen, and S. Yan. Network in network. *International Conference on Learning Representations*, 2014. [1](#)
- [19] T.-Y. Lin, A. RoyChowdhury, and S. Maji. Bilinear cnn models for fine-grained visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1449–1457, 2015. [3](#), [4](#), [5](#), [6](#)
- [20] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008. [2](#), [7](#)
- [21] J. Malik, S. Belongie, T. Leung, and J. Shi. Contour and texture analysis for image segmentation. *International Journal of Computer Vision*, 43(1):7–27, Jun 2001. [1](#)
- [22] Z. Meng, P. Liu, J. Cai, S. Han, and Y. Tong. Identity-aware convolutional neural network for facial expression recognition. In *The Twelfth IEEE International Conference on Automatic Face and Gesture Recognition*, 2017. [3](#)
- [23] M. R. Min, H. Guo, and D. Song. Exemplar-centered supervised shallow parametric data embedding. *arXiv preprint arXiv:1702.06602*, 2017. [3](#), [7](#)
- [24] M. R. Min, L. Maaten, Z. Yuan, A. J. Bonner, and Z. Zhang. Deep supervised t-distributed embedding. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 791–798, 2010. [3](#), [7](#)
- [25] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014. [2](#)
- [26] X. Peng, R. S. Feris, X. Wang, and D. N. Metaxas. A recurrent encoder-decoder network for sequential face alignment. In *European Conference on Computer Vision*, pages 38–56. Springer, 2016. [3](#)
- [27] F. Perronnin, J. Sánchez, and T. Mensink. Improving the fisher kernel for large-scale image classification. *Computer Vision–ECCV 2010*, pages 143–156, 2010. [4](#)
- [28] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. [4](#)
- [29] G. Schwartz and K. Nishino. Material recognition from local appearance in global context. *arXiv preprint arXiv:1611.09394*, 2016. [2](#)
- [30] E. Shelhamer, J. Long, and T. Darrell. Fully convolutional networks for semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):640–651, April 2017. [2](#)

- [31] J. B. Tenenbaum and W. T. Freeman. Separating style and content. In *Advances in neural information processing systems*, pages 662–668, 1997. 2, 3
- [32] L. van der Maaten. Learning a parametric embedding by preserving local structure. *RBM*, 500(500):26, 2009. 3, 7
- [33] L. Van Der Maaten. Accelerating t-sne using tree-based algorithms. *Journal of machine learning research*, 15(1):3221–3245, 2014. 2, 6, 7, 8
- [34] M. Varma and A. Zisserman. A statistical approach to texture classification from single images. *International Journal of Computer Vision*, 62(1):61–81, Apr 2005. 1
- [35] T.-C. Wang, J.-Y. Zhu, E. Hiroaki, M. Chandraker, A. A. Efros, and R. Ramamoorthi. A 4d light-field dataset and cnn architectures for material recognition. In *European Conference on Computer Vision*, pages 121–138. Springer, 2016. 3
- [36] J. Xue, H. Zhang, K. Dana, and K. Nishino. Differential angular imaging for material recognition. *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2, 3, 5, 6
- [37] H. Zhang, V. M. Patel, B. S. Riggan, and S. Hu. Generative adversarial network-based synthesis of visible faces from polarimetric thermal faces. *International Joint Conference on Biometrics*, 2017. 3
- [38] H. Zhang, J. Xue, and K. Dana. Deep ten: Texture encoding network. *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 2, 3, 4, 6
- [39] Z. Zhang, Y. Xie, F. Xing, M. McGough, and L. Yang. Mdnet: A semantically and visually interpretable medical image diagnosis network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6428–6436, 2017. 3
- [40] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. S. Torr. Conditional random fields as recurrent neural networks. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015. 2
- [41] Y. Zhu, Z. Lan, S. Newsam, and A. G. Hauptmann. Hidden two-stream convolutional networks for action recognition. *arXiv preprint arXiv:1704.00389*, 2017. 3
- [42] Y. Zhu and S. Newsam. Depth2action: Exploring embedded depth for large-scale action recognition. In *European Conference on Computer Vision*, pages 668–684. Springer, 2016. 3