# Differential Viewpoints for Ground Terrain Material Recognition

Jia Xue, Hang Zhang, Ko Nishino, and  Kristin J. Dana

**Abstract**—Computational surface modeling that underlies material recognition has transitioned from reflectance modeling using in-lab controlled radiometric measurements to image-based representations based on internet-mined single-view images captured in the scene. We take a middle-ground approach for material recognition that takes advantage of both rich radiometric cues and flexible image capture. A key concept is differential angular imaging, where small angular variations in image capture enables angular-gradient features for an enhanced appearance representation that improves recognition. We build a large-scale material database, Ground Terrain in Outdoor Scenes (GTOS) database, to support ground terrain recognition for applications such as autonomous driving and robot navigation. The database consists of over 30,000 images covering 40 classes of outdoor ground terrain under varying weather and lighting conditions. We develop a novel approach for material recognition called texture-encoded angular network (TEAN) that combines deep encoding pooling of RGB information and differential angular images for angular-gradient features to fully leverage this large dataset. With this novel network architecture, we extract characteristics of materials encoded in the angular and spatial gradients of their appearance. Our results show that TEAN achieves recognition performance that surpasses single view performance and standard (non-differential/large-angle sampling) multiview performance.

**Index Terms**—Material recognition, deep convolutional neural networks, texture reflectance, robot navigation.

---◆---

## 1 INTRODUCTION

REAL world scenes consist of surfaces made of numerous materials, such as wood, marble, dirt, metal, ceramic and fabric, which contribute to the rich visual variation we find in images. Material recognition has become an active area of research in recent years, with the goal of providing detailed material information for applications such as autonomous agents and human-machine systems. Real world surfaces and material characteristics are both apparent (the visual appearance) and latent (physical material properties of the surface such as friction, micro-geometry and roughness). Material properties affect both the spatial variation of surface appearance and the angular variation of reflectance with respect to both view and illumination.

Early studies of material appearance modeling largely concentrated on comprehensive lab-based measurements using dome systems, robots, or gonioreflectometers collecting measurements that are dense in angular space (such as BRDF, BTF) [1]. These reflectance-based studies have the advantage of capturing intrinsic invariant properties of the surface, which enables fine-grained material recognition [2], [3], [4], [5], [6]. The inflexibility of lab-based image capture, however, prevents widespread use in real world scenes, especially in the important class of outdoor scenes. A fundamentally different approach to reflectance modeling is image-based appearance modeling where surfaces are captured
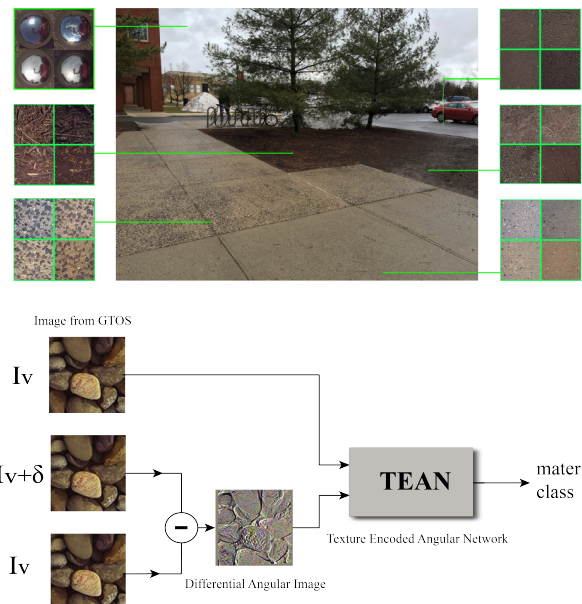


Fig. 1: (Top) Example from GTOS dataset comprising outdoor measurements with multiple viewpoints, illumination conditions and angular differential imaging. The example shows scene-surfaces imaged at different illumination/weather conditions. (Bottom) Texture Encoded Angular Network (TEAN) for ground terrain material recognition.

with a single-view image in-scene or "in-the-wild." Recent studies of image-based material recognition use single-view internet-mined images to train classifiers [7], [8], [9], [10], [11], [12] and can be applied to arbitrary images casually taken without the need

- *Jia Xue is with the Department of Electrical and Computer Engineering, Rutgers University–New Brunswick, New Brunswick, NJ 08901, USA. E-mail: jia.xue@rutgers.edu*
- *Hang Zhang is with the Amazon AI, Amazon Web Services Inc, East Palto Alto, CA 94025, USA E-mail: hzaws@amazon.com*
- *Ko Nishino is with the Department of Intelligence Science and Technology, Graduate School of Informatics, Kyoto University Yoshida Honmachi, Sakyo-ku Kyoto, Kyoto 606-8501. E-mail: kon@i.kyoto-u.ac.jp*
- *Kristin J. Dana is with the Department of Electrical and Computer Engineering, Rutgers University–New Brunswick, New Brunswick, NJ 08901,USA. E-mail: kristin.dana@rutgers.edu*
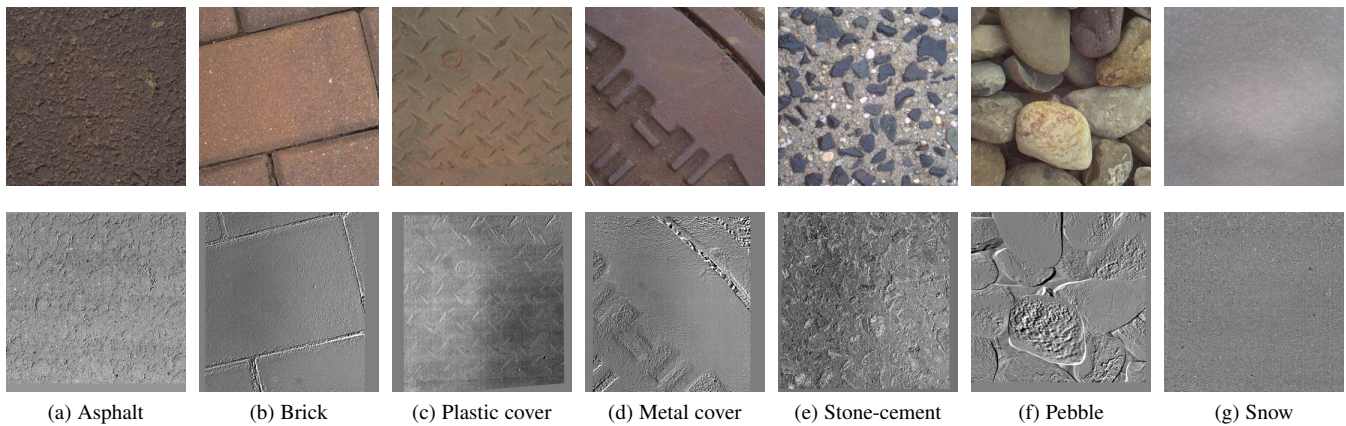
Fig. 2: Differential Angular Imaging. (Top) Examples of material surface images $I_v$. (Bottom) Corresponding differential images $I_\delta = I_v - I_{v+\delta}$ in our GTOS dataset. These sparse images encode angular gradients of reflection and 3D relief texture.

of multiview reflectance information. In these methods, however, recognition is typically based more on context including object and scene cues, than intrinsic material appearance properties except for a few purely local methods [13], [14].

Between the two approaches of reflectance-based and image-based material recognition, i.e. between comprehensive in-lab imaging and internet-mined images, we take an advantageous middle-ground. Specifically, we capture in-scene real-world surfaces but use multiple viewpoint angles for measurements that provide a partial reflectance sampling. This leads to a very basic question: how do multiple viewing angles help in material recognition? More interestingly, we consider a novel question: Do small changes in viewing angles, *differential changes*, result in significant increases in recognition performance? Prior work has shown the power of angular filtering to complement spatial filtering in material recognition. These methods, however, rely on a lightfield camera to achieve multiple differential viewpoint variations [15] or a mirror-based camera to capture a slice of the BRDF [5] which limits application in the wild due to a rigid imaging system setup and inadequacy for image capture at distance. We instead propose to capture surfaces with differential changes in viewing angles with an ordinary camera and compute discrete approximations of *angular gradients*. We present an approach called *angular differential imaging* that augments image capture for a particular viewing angle $v$ a differential viewpoint $v + \delta$. Contrast this method with lab-based reflectance measurements that often quantize the angular space measuring with domes or positioning devices with large angular spacing such as $22.5°$. These coarse-quantized measurements have limited use in approximating angular gradients. Angular differential imaging can be implemented with a small-baseline stereo camera or a moving camera (e.g. handheld). We demonstrate that differential angular imaging provides key information about material reflectance properties while maintaining the flexibility of convenient in-scene appearance capture.

To capture material appearance in a manner that preserves the convenience of image-based methods but important angular information of reflectance-based methods, we assemble a comprehensive, first-of-its-kind, *outdoor, in-place* material database that includes multiple viewpoints and multiple illumination directions (partial BRDF sampling), multiple weather conditions, a large

set of surface material classes surpassing existing comparable datasets, multiple physical instances per surface class (to capture intra-class variability) and differential viewpoints to support the framework of differential angular imaging. Specifically, the resulting database spans 40 surface classes that we find commonly in daily life, 4-14 examples or instances per class and for each surface measurement we collect an image set at 18 viewing angles using a mobile robot and at multiple illumination conditions (4) corresponding to times of day and weather conditions. Each image is collected with 3 exposures for high dynamic range imaging. The total number of surface images is 34,243. These surfaces are not measured in the lab, but rather in their typical state within a scene. The global scene image of the surface is also captured and indicates the scene context for the surface. We concentrate on outdoor scenes because of the limited availability of reflectance databases for outdoor surfaces. We also concentrate on materials from ground terrain in outdoor scenes (GTOS) for applicability in numerous application such as automated driving, robot navigation and scene semantics. The 40 surface classes include ground terrain such as grass, gravel, asphalt, concrete, black ice, snow, moss, mud and sand (see Figure 2).

For recognition, we build a recognition algorithm that leverages the strength of deep learning and differential angular imaging. We develop a two-branch network that combines deep encoding pooling for spatial (texture) information and a second branch for angular information as illustrated in Figure 1. We call this new architecture Texture Encoded Angular Network (TEAN). It combines two prior concepts DEP network [16] and DAIN [17] to account for spatial texture and angular information in a robust deep learning architecture. The original concepts have been improved with new architectures to incorporate better base networks with improved efficiency and accuracy.

## 2 RELATED WORK

**Material recognition:** Material recognition is a fundamental problem in computer vision. The classification of 3D material images and bidirectional texture functions, traditionally relies on handcrafted filter banks followed by grouping the outputs into texton histograms [22], [23] or bag-of-words [24], [25], [26], [27]. The success of deep learning methods in object recognition has also translated to the problem of material recognition, the

| Datasets | samples | classes | views | illumination | in scene | scene image | camera parameters | year |
|---|---|---|---|---|---|---|---|---|
| CUReT [18] | 61 | 61 | 205 | | N | N | N | 1999 |
| KTH-TIPS [19] | 11 | 11 | 27 | 3 | N | N | N | 2004 |
| UBO2014 [20] | 84 | 7 | 151 | 151 | N | N | N | 2014 |
| Reflectance disk [5] | 190 | 19 | 3 | 3 | N | N | Y | 2015 |
| 4D Light-field [15] | 1200 | 12 | 1 | 1 | Y | N | N | 2016 |
| NISAR [21] | 100 | 100 | 9 | 12 | N | N | N | 2016 |
| **GTOS(ours)** | **606** | **40** | **19** | **4** | **Y** | **Y** | **Y** | **2017** |

TABLE 1: Comparison between GTOS dataset and some publicly available BRDF material datasets. Note that the 4D Light-field dataset [15] is captured by the Lytro Illum light field camera.

classification and segmentation of material categories in arbitrary images. Bell *et al*. achieve per-pixel material category labeling by retraining the state-of-the-art object recognition network [28] on a large dataset of material appearance [7]. This method relies on large image patches that include object and scene context to recognize materials. In contrast, Schwartz and Nishino [13], [14] learn material appearance models from small image patches extracted inside object boundaries to decouple contextual information from material appearance. To achieve accurate local material recognition, they introduced intermediate material appearance representations based on their intrinsic properties (e.g., "smooth" and "metallic"). Zhang *et al* [9] 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 that learns an encoding for an orderless texture representation. These prior methods show the utility of spatial information within an image for material recognition.

While there has been recent emphasis of characterizing materials with apparent appearance in images, radiometric properties of materials such as the bidirectional reflectance distribution function (BRDF) [29] and the bidirectional texture function (BTF) [18] provide appearance as a function of viewing and illumination angle and therefore encode angular information. Materials have unique characteristics in the subtle variations of their reflectance functions (e.g., different types of metal [30] and paint [31]). However, reflectance measurements require elaborate image capture systems, such as a gonioreflectometer [29], [32], robotic arm [18], [33], or a dome with cameras and light sources [30], [31], [34]. (Numerous methods for capturing reflectance have been detailed in surveys [1], [35].) Recently, Zhang *et al* introduced the use of a one-shot reflectance field capture for material recognition [5]. They adapt the parabolic mirror-based camera developed by Dana and Wang [36] to capture the BRDF for a given light source direction in a single shot, called a reflectance disk. These results demonstrate that gradients of angular appearance encode rich cues for their recognition. Similarly, Wang *et al* [15] uses a light field camera and combines angular and spatial filtering for material recognition. The approach we present in this paper develops a novel material recognition framework that combines both spatial and angular filtering. Specifically, we combine state-of-the-art texture representations with reflectance cues from differential angular images that can easily be captured by a two-camera system or small motions of a single ordinary camera. The resulting method is instantiated in a two-branch network comprised of one branch for reflectance with an emphasis on angular gradients and another branch for texture with both orderless and ordered spatial cues.

**Datasets:** Datasets to measure reflectance of real world surfaces have a long history of lab-based measurements including:
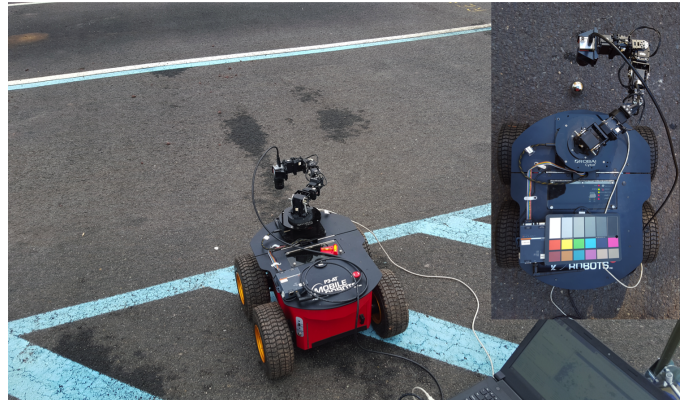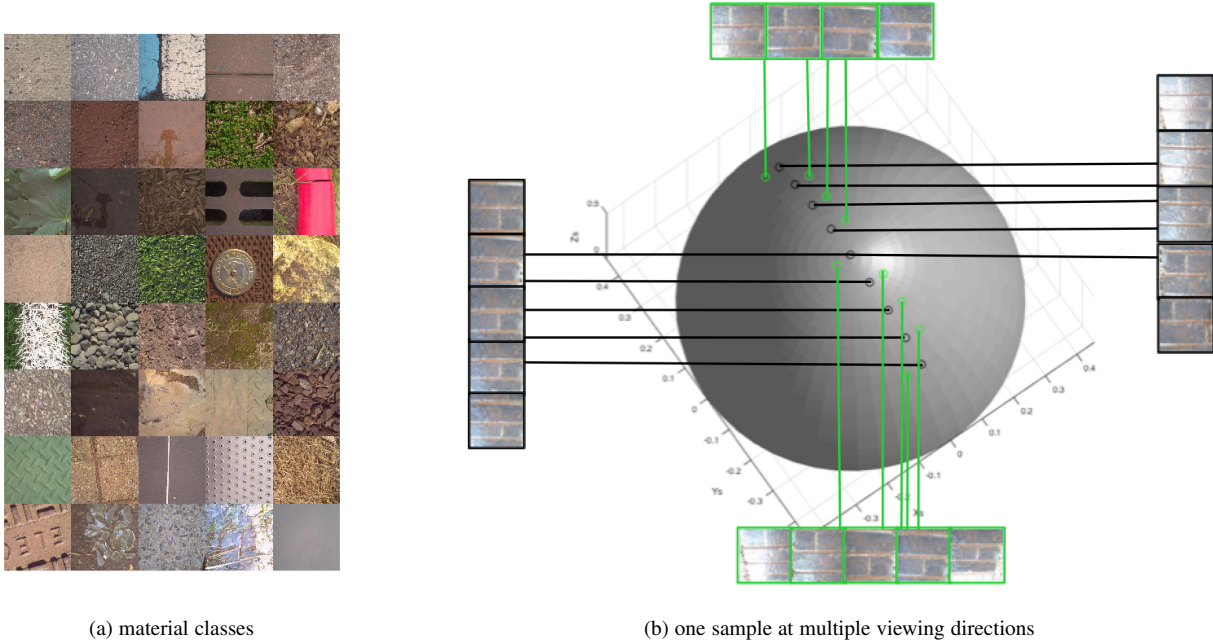


Fig. 3: The measurement equipment for the GTOS database: Mobile Robots P3-AT robot, Cyton gamma 300 robot arm, Basler aca2040-90uc camera with Edmund Optics 25mm/F1.8 lens, DGK 18% white balance and color reference card, and the 440C Stainless Steel Tight-Tolerance Sphere (McMaster-Carr).

CUReT database [18], KTH-TIPS database by Hayman *et al*. [19], MERL Reflectance Database [37], UBO2014 BTF Database [20], UTIA BRDF Database [38], Drexel Texture Database [39] and IC-CERTH Fabric Database [40]. In many of these datasets, dense reflectance angles are captured with special image capture equipment. Some of these datasets have limited instances/samples per surface category (different physical samples representing the same class for intraclass variability) or have few surface categories, and all are obtained from indoor measurements where the sample is removed from the scene. More recent datasets capture materials and texture in-scene, (a.k.a. in-situ, or in-the-wild). A motivation of moving to in-scene capture is to build algorithms and methods that are more relevant to real-world applications. These recent databases are from internet-mined databases and contain a single view of the scene under a single illumination direction. Examples include the the Flickr Materials Database by Sharan *et al*. [41] and the Material in Context Database by Bell *et al*. [7]. However, because photographs in both datasets are collected from internet, the reflectance properties from multiple views are lost. Recent work uses deep networks to estimate multiview reflectance models for novel-view material rendering [6], [42], [43], [44], [45], [46]. Training with multiview renderings may be a future direction for cases where real-world images are not available.

## 3 GTOS DATASET

In this section, we introduce the GTOS dataset and the measurement device. GTOS dataset is a first-of-its-kind in-scene material reflectance database, to investigate the use of spatial and angular

(a) material classes

(b) one sample at multiple viewing directions

Fig. 4: (a) The 40 material categories in the GTOS dataset introduced in this paper. (b) The material surface observation points. Nine viewpoint angles (black spots) separated along an arc spanning $80°$ are measured. For each viewpoint, a differential view (green spots) is captured $\pm 5°$ in azimuth from the original orientation (the sign is chosen based on robotic arm kinematics. )

reflectance information of outdoor ground terrain for material recognition.

### 3.1 Measurement Device

Our measurement device (depicted in Figure 3) is composed of a Mobile Robots P3-AT robot, Cyton gamma 300 robot arm, Basler aca2040-90uc camera with Edmund Optics 25mm/F1.8 lens, DGK 18% white balance and color reference card, and Hardened 440C Stainless Steel Tight-Tolerance Sphere (McMaster-Carr). The constraint that the Cyton arm can only hold 300g for full-range movement presents a practical obstacle in our choice for camera and lens, and we employed the Basler USB camera with Edmund Industrial optics 86572 fixed focus lens (the total weight of the camera and the lens is 203g). The aca2040-90uc camera can capture $2040 \times 2046$ pixels resolution photographs with 12 bits per pixel. As shown in Figure 4, a Hardened 440C Stainless Steel Tight-Tolerance Sphere is employed to reflect the sky and indicate the weather conditions. Camera parameter adjustment is challenging due to image capture in different lighting. We set the camera parameters to be adjusted automatically by simultaneously observing the DGK 18% white balance and color reference card. Camera parameters are adjusted only one time for each sample, ensuring that images at multiple viewing angles are captured under the same parameters. Sample appearance depends on sky/weather conditions and the time of day. We image the same region with four different weather conditions (cloudy dry, cloudy wet, sunny morning and sunny afternoon). As shown in Figure 4, we choose 9 points to form an approximate $80°$ arc as our viewing points. For each observation image, an additional image obtained by varying the viewing angle by a small angle $(3 - 5°)$ provides the pair needed to compute the angular gradient. To collect BRDF information, the observation points are fixed for the entire

database. The distance between observing points and sample is 4045 mm. The imaging region is 1510 mm $\times$ 1510 mm.

### 3.2 Dataset Overview

We collect the GTOS database, a first-of-its-kind in-scene material reflectance database, to investigate the use of spatial and angular reflectance information of outdoor ground terrain for material recognition. We capture reflectance systematically by imaging a set of viewing angles comprising a partial BRDF with a mobile exploration robot. The measurement device is depicted in Figure 3. Due to the joint limitation of the Cyton gamma 300 robot arm, we select $[-40°, 40°]$ as our measuring range. Differential angular images are obtained by measuring each of $N_v = 9$ base angles $v = (\theta_v, \phi_v)$, $\theta_v \in [-40°, -30°, \ldots, 40°]$, and a differential angle variation of $\delta = (0, 5°)$ resulting in 18 viewing directions per sample as shown in Figure 4 (b). Example surface classes are depicted in Figure 4 (a). The class names are (in order of top-left to bottom-right): cement, asphalt, painted asphalt , brick, soil, muddy stone, mud, mud-puddle, grass, dry leaves, leaves, asphalt-puddle, mulch, metal grating, plastic, sand, stone, artificial turf, aluminum, limestone, painted turf, pebbles, roots, moss, loose asphalt-stone, asphalt-stone, cloth, paper, plastic cover, shale, painted cover, stone-brick, sandpaper, steel, dry grass, rusty cover, glass, stone-cement, icy mud, and snow. The $N_c = 40$ surface classes mostly have between 4 and 14 instances (samples of intra-class variability) and each instance is imaged not only under $N_v$ viewing directions but also under multiple natural light illumination conditions. As illustrated in Figure 1, sample appearance depends on the weather condition and the time of day. To capture this variation, we image the same region with $N_i = 4$ different weather conditions (cloudy dry, cloudy wet, sunny morning, and sunny afternoon). We capture the samples

with 3 different exposure times to enable high dynamic range imaging. Additionally, we image a mirrored sphere to capture the environment lighting of the natural sky. In addition to surface images, we capture a scene image to show the global context. Although, the database measurements were obtained with robotic positioning for precise angular measurements, our recognition results are based on subsets of these measurements so that an articulated arm would not be required for an in-field system. The total number of surface images in the database is 34,243. As shown in Table 1, this is the most extensive outdoor in-scene multiview material database to date.

## 3.3 Differential Angular Imaging

Our GTOS dataset introduces a measurement method called differential angular imaging where a surface is imaged from a particular viewing angle $v$ and then from an additional viewpoint $v + \delta$. The motivation for this differential change in viewpoint is improved computation of the angular gradient of intensity $\partial I_v / \partial v$. Intensity gradients are the basic building block of image features and it is well known that discrete approximations to derivatives have limitations. In particular, spatial gradients of intensities for an image $I$ are approximated by $I(x + \Delta) - I(x)$ and this approximation is most reasonable at low spatial frequencies and when $\Delta$ is small. One implication is that the discrete approximation to the derivative is only valid at lower frequencies, as expected. The second implication is that increasing $\Delta$ decreases the range of frequencies over which the discrete approximation is valid. Therefore, small values of $\Delta$ provide better gradients. For angular gradients of reflectance, the discrete approximation to the derivative is a subtraction with respect to the viewing angle. Angular gradients are approximated by $I(v + \delta) - I(v)$ and this approximation requires a small $\delta$. Consequently, differential angular imaging provides more accurate angular gradients.

The differential angular images as shown in Figures 1 and 2 have several characteristics. First, the differential angular image reveals the gradients in BRDF/BTF at the particular viewpoint. Second, relief texture is also observable in the differential angular image due to non-planar surface structure. Finally, the differential angular images are sparse. This sparsity has the potential to provide a computational advantage, though we have not specifically utilized this advantage in our network design.

## 4 DEEP LEARNING ARCHITECTURES

To leverage differential angular imaging for material recognition, we build a two stream convolution network that takes two image streams as input, the original image and a differential image. We start our experiments with widely studied ImageNet [47] pre-trained networks [48], [49], [50] as the CNN streams, and we call the network Differential Angular Imaging Network (DAIN). Networks designed for object recognition take spatial order as critical for classification. However, texture recognition uses an orderless component to provide invariance to spatial layout [8], [51]. Through study the GTOS dataset, we find that for "images in the wild", homogeneous surfaces rarely fill the entire field-of-view, and many materials exhibit regular structure. We design a network to balance both orderless and ordered spatial information for the GTOS images, and we call this network Deep Encoding Pooling Network (DEP). Finally, by replacing the CNN architecture in DAIN with DEP to combine angular reflectance cues with orderless and ordered spatial infromation, we introduce Texture Encoded Angular Network (TEAN).

## 4.1 Differential Angular Imaging Network (DAIN)

We develop a two-stream convolutional neural network to fully leverage differential angular imaging for material recognition. The differential angular image $I_\delta$ sparsely encodes reflectance angular gradients as well as surface relief texture. The spatial variation of image intensity remains an important recognition cue and so our method integrates these two streams of information. A CNN is used on both streams of the network and then combined for the final prediction result. The combination method and the layer at which the combination takes place leads to variations of the architecture.

We employ the ImageNet [47] pre-trained VGG-M model [49] as the initial prediction unit (labeled CNN in Figure 5). The first input branch is the image $I_v$ at a specific viewing direction $v$. The second input branch is the differential image $I_\delta$. The first method of combination shown in Figure 5 (a) is a simple averaging of the output prediction vectors obtained by the two branches. The second method combines the two branches at the intermediate layers of the CNN, i.e. the feature maps output at layer $M$ are combined and passed forward to the higher layers of the CNN, as shown Figure 5 (b). We empirically find that combining feature maps generated by Conv5 layer after ReLU performs best. A third method (see Figure 5 (c)) is a hybrid of the two architectures that preserves the original CNN path for the original image $I_v$ by combining the layer $M$ feature maps for both streams *and* by combining the prediction outputs for both streams as shown in Figure 5 (c). This approach is the best performing architecture of the three methods and we call it the differential angular imaging network (DAIN).

For combining feature maps at layer $M$, consider features maps $x_a$ and $x_b$ from the two branches that have width $W$, height $H$, and feature channel depth $D$. The output feature map $y$ will be the same dimensions $W \times H \times D$. We can combine feature maps by: (1) *Sum:* pointwise sum of $x_a$ and $x_b$, and (2) *Max:* pointwise maximum of $x_a$ and $x_b$. As experimentally proved in section 5.1, the sum combination outperforms max combination. Without explicit declaration, we use sum combination as our default combination method. The CNN module of our DAIN network can be replaced by other state-of-the-art deep learning methods to further improve results. To demonstrate this, in section 5, we change the VGG-M model with ImageNet pre-trained MobileNet V2 [50] and provide several experiments for evaluation.

**Multiple Views** Our GTOS database has multiple viewing directions on an arc (a partial BRDF sampling) as well as differential images for each viewing direction. We evaluate our recognition network in two modes: (1) **Single view DAIN**, with inputs from $I_v$ and $I_\delta$, with $v$ representing a single viewing angle; (2) **Multiview DAIN**, with inputs $I_v$ and $I_\delta$, with $v \in [v1, v2, ..., vN]$. For our GTOS database, $v1, v2, ..., vN$ are viewing angles separated by $10°$ representing a $N \times 10°$ range of viewing angles. We empirically determine that $N = 4$ viewpoints are sufficient for recognition. For a baseline comparison we also consider non-differential versions: **Single View** with only $I_v$ for a single viewing direction and **Multiview** with inputs $I_v$, $v \in [v1, v2, ..., vN]$.

To incorporate multiview information in DAIN we use three methods: (1) voting (use the predictions from each view to vote), (2) pooling (pointwise maximum of the combined feature maps across viewpoints), (3) 3D filter + pooling (follow [52] to use a $3 \times 3 \times 3$ learned filter bank to convolve the multiview feature maps). See Figure 6. After 3D filtering, pooling is used (pointwise
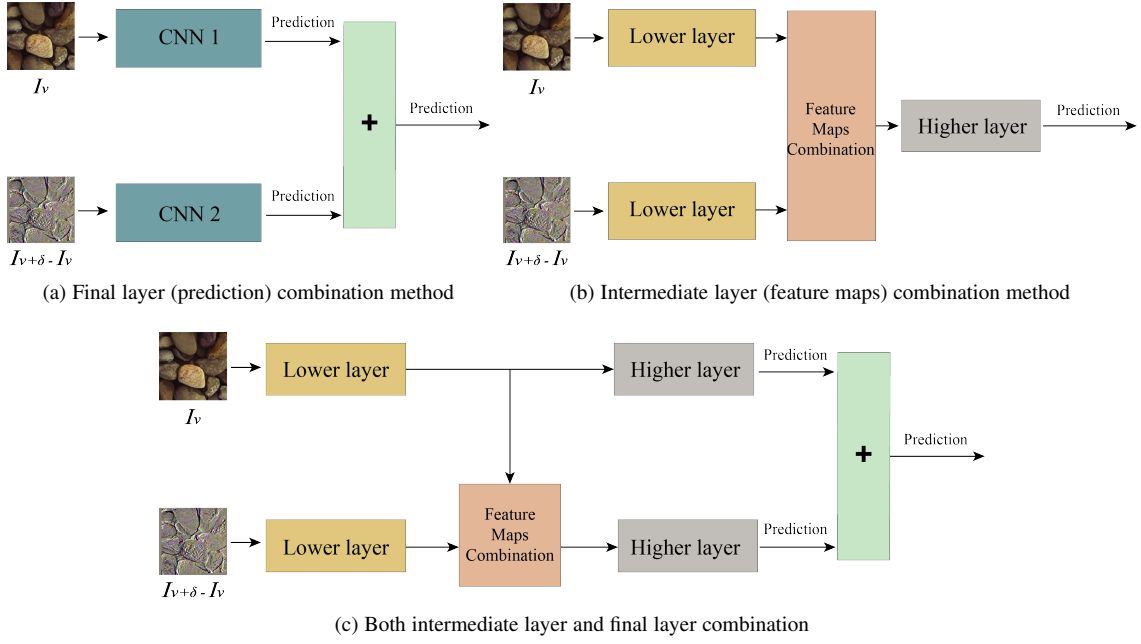
(a) Final layer (prediction) combination method

(b) Intermediate layer (feature maps) combination method

(c) Both intermediate layer and final layer combination

Fig. 5: Methods to combine two image streams, the original image $I_v$ and the differential angular image $I_\delta = I_{v+\delta} - I_v$. The architecture in (c) provides better performance than (a) and (b) and we call it the differential angular imaging network (DAIN).
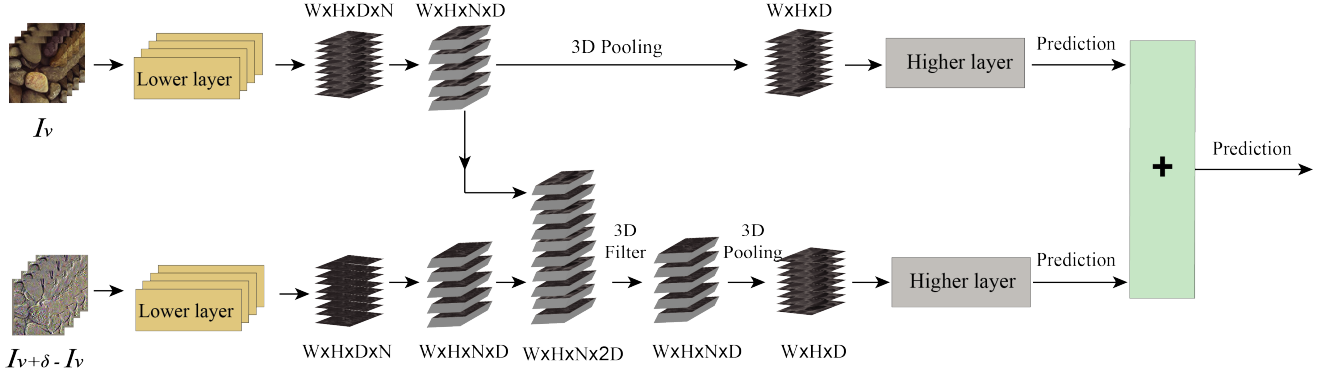


Fig. 6: Multiview DAIN. The 3D filter + pooling method to combine two streams (original and differential image) from multiple viewing angles $v \in [v1, v2, ..., vN]$. $W$, $H$, and $D$ are the width, height, and depth of corresponding feature maps, $N$ is the number of view points.

maximum across viewpoints). Due to learning the filter weights, the computational expense of this third method is significantly higher.

## 4.2 Deep Encoding Pooling Network (DEP)

We introduce a Deep Encoding Pooling Network (DEP) that leverages an orderless texture representation and high level spatial information for material recognition. The network is shown in Figure 7, with the material image as network input, outputs from convolutional layers are fed into two feature representation layers jointly; the texture encoding layer [9] and the global average pooling layer. The texture encoding layer captures 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 [53].

**Encoding Layer** The texture encoding layer [9] 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 $X = \{x_1, ...x_m\}$ be M visual descriptors, $C = \{c_1, ...c_n\}$ is the code book with N learned codewords. The residual vector $r_{ij}$ is calculated by $r_{ij} = x_i - c_j$, where $i = 1...m$ and $j = 1...n$. The residual encoding for codeword $c_j$ can be represented as

$$e_j = \sum_{i=1}^{M} w_{ij} r_{ij},$$

(1)

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

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

(2)

$s_1, ...s_m$ are learnable smoothing factors. With the texture encoding layer, the visual descriptors X are pooled into a set of

N residual encoding vectors $E = \{e_1, ...e_n\}$. 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 [54]. 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

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

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 Network** With aforementioned encoding layer and bilinear models, we introduce our Deep Encoding Pooling Network (DEP). Our Deep Encoding Pooling Network is shown in Figure 7. As in prior transfer learning algorithms [9], [55], we employ convolutional layers with non-linear layers from ImageNet [47] 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 high level spatial information. The dimension of outputs from the texture encoding layer is determined by the codewords N and the feature maps channel C (N×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 2 is an instantiation of DEP based on MobileNet V2 [50]. We set 8 codewords for the texture encoding layer. The size of input images are $224 \times 224$. 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 $8 \times 1280 = 10240$ and the dimension of outputs from global average pooling layer is 1280. We reduce the dimension of feature maps from the texture 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 [9], [56], 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.

### 4.3 Texture Encoded Angular Network (TEAN)

Adapting the DEP to the RGB image branch in DAIN, we introduce the Texture Encoded Angular Network (TEAN). The detailed network is shown in Figure 8. We develop a two-stream

| layer name | output size | encoding-pooling |
|---|---|---|
| conv2d | 112×112×32 | 3×3, stride 2 |
| bottleneck1_x | 112×112×16 | $\begin{bmatrix} 3 \times 3, 32 \\ 1 \times 1, 16 \end{bmatrix} \times 1$ |
| bottleneck2_x | 56×56×24 | $\begin{bmatrix} 1 \times 1, 96 \\ 3 \times 3, 96 \\ 1 \times 1, 24 \end{bmatrix} \times 2$ |
| bottleneck3_x | 28×28×32 | $\begin{bmatrix} 1 \times 1, 144 \\ 3 \times 3, 144 \\ 1 \times 1, 32 \end{bmatrix} \times 3$ |
| bottleneck4_x | 14×14×64 | $\begin{bmatrix} 1 \times 1, 192 \\ 3 \times 3, 192 \\ 1 \times 1, 64 \end{bmatrix} \times 4$ |
| bottleneck5_x | 14×14×96 | $\begin{bmatrix} 1 \times 1, 384 \\ 3 \times 3, 384 \\ 1 \times 1, 96 \end{bmatrix} \times 3$ |
| bottleneck6_x | 7×7×160 | $\begin{bmatrix} 1 \times 1, 576 \\ 3 \times 3, 576 \\ 1 \times 1, 160 \end{bmatrix} \times 3$ |
| bottleneck7_x | 7×7×320 | $\begin{bmatrix} 1 \times 1, 960 \\ 3 \times 3, 960 \\ 1 \times 1, 320 \end{bmatrix} \times 1$ |
| bottleneck8_x | 7×7×1280 | 1 × 1, 1280 |
| encoding / pooling | 8 x 1280 / 1280 | 8 codewords / ave pool |
| fc1_1 / fc1_2 | 64 / 64 | 10240×64 / 1280×64 |
| bilinear mapping | 4096 | - |
| fc2 | 128 | 4096×128 |
| classification | n classes | 128×n |

TABLE 2: The architecture of the Deep Encoding Pooling Network based on MobileNet V2 [50]. The input image size is $224 \times 224$.

convolutional neural network, one branch input is the differential angular image, representing the material reflectance information. The other branch input is the RGB image, representing the orderless texture details and ordered spatial information. For the color image branch, we utilize the Deep Encoding Pooling Network (DEP) to balance the orderless texture component and ordered spatial information. As in DAIN, we combine feature maps at both intermediate layer and final prediction layer. With the proposed Texture Encoded Angular Network (TEAN), we take advantage of material reflectance information, orderless texture details and ordered spatial information for ground terrain material recognition. This combination of angular cues, orderless spatial cues and ordered spatial cues leads to improved recognition results.

As shown in Figure 8, we employ ImageNet [47] pre-trained MobileNet V2 [50] as the initial prediction unit. As in single view DAIN (Sum), we combine feature maps from the bottlenect8_x with element-wise sum as intermediate layer combination. Feature maps from color images are fed into the texture encoding layer and the global average pooling layer jointly, followed by bilinear model and fully connected layer, the output from fully connected layer is a 128-D vector. The element-wise summed feature maps are fed into a fully connected layer for dimension reduction,
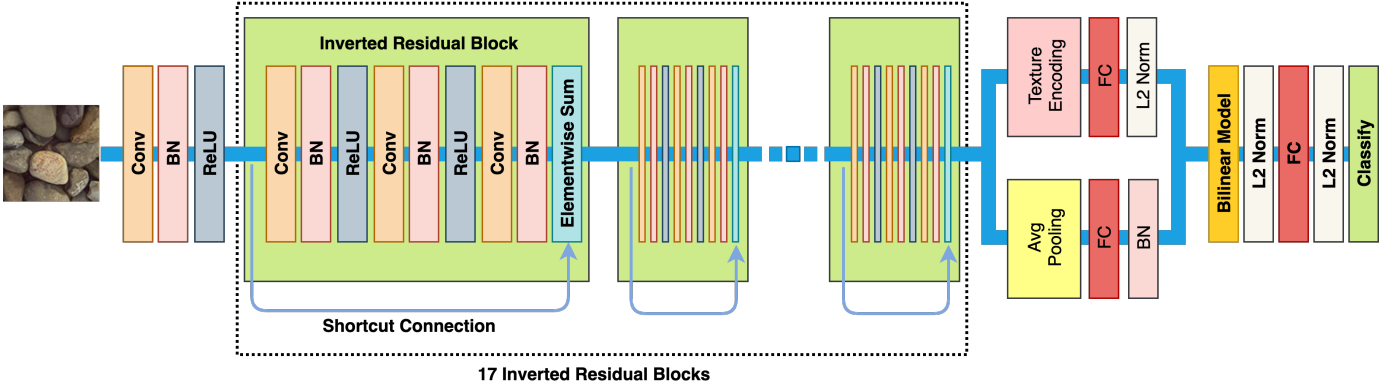
Fig. 7: A Deep Encoding Pooling Network (DEP) for material recognition. Outputs from convolutional layers are fed into the encoding layer and global average pooling layer jointly and their outputs are processed with bilinear model.

the output is also a 128-D vector. These two 128-D vectors are concatenated and fed into classify layer for material classification.

## 5 EXPERIMENTS

In this section, we evaluate the performance of DAIN, DEP and TEAN framework for material recognition. First, in section 5.1, we evaluate which structure of the two stream networks from Figure 5 works best on the GTOS dataset, leading to the choice in (c) as the DAIN architecture. Based on (c), we consider recognition performance with different DAIN variations for recognition and compare three other state-of-the-art approaches on our GTOS dataset, concluding that multiview DAIN works best. Second, in section 5.2, we compare the recognition performance of DEP with fine-tuning MobileNet, bilinear CNN and Deep-TEN. To prove the superior performance of the proposed DEP in material recognition, we experiment DEP on two other material/texture recognition datasets. Third, in section 5.3, we evaluate the performance of TEAN and verify the performance of multi-scale training. Finally, in section 5.4, to gain insight into the performance, we construct the confusion matrix and visualize the features before classification layers with BarnesHut t-SNE [57] for MobileNet, DEP, DAIN and TEAN.

**Training procedure** We design 5 training and testing splits by assigning about 70% of ground terrain surfaces of each class to training and the remaining 30% to testing. In order to ensure that there is no overlap between training and testing sets, if one sample is in the training set, all views and illumination conditions for that sample is in the training set. Each input image from our GTOS database is resized into $240 \times 240$. Since the snow class only has 2 samples in the dataset, we omit this class from experiments.

Comparing with recent mobile platform designed MobileNet V2 [50], the number of parameters for VGG-M [49] is tremendous. So for training a VGG-M based two branch network, we first fine-tune the VGG-M model separately with RGB and differential images with batch size 196, dropout rate 0.5, momentum 0.9. We employ the augmentation method that horizontally and vertically stretch training images within $\pm 10\%$, with an optional 50% horizontal and vertial mirror flips. The images are randomly cropped into $224 \times 224$ material patches. All images are pre-processed by subtracting a per color channel mean and normalizing for unit variance. The learning rate for the last fully connected layer is set to 10 times of other layers. We first fine-tune only the last fully connected layer with learning rate $5 \times 10^{-2}$ for 5 epochs; then, fine-tune all the fully connected layers with learning rate $10^{-2}$ for 5 epochs. Finally we fine-tune all the layers with leaning rate starting at $10^{-3}$, and decrease by a factor of 0.1 when the training accuracy saturates. For MobileNet V2 [50] based DAIN, we train the network end-to-end on GTOS dataset directly.

Following prior works [16], [17], for the fine-tuned two-branch VGG-M model and two-branch MobileNet V2 model, we experiment with batch size 64 and learning rate starting from 0.01 which is reduced by a factor of 0.1 when the training accuracy saturates. We augment training data with randomly stretch training images by $\pm 25\%$ horizontally and vertically, and also horizontal and vertical mirror flips with 50% chance. The images are randomly cropped into $224 \times 224$ material patches. We first backpropagate only to feature maps combination layer for 3 epochs, then fine tunes all layers. We employ the same augmentation method for the multiview images of each material surface. We randomly select the first viewpoint image, then subsequent $N = 4$ view point images are selected for experiments.

### 5.1 DAIN Recognition Results

**Recognition Benchmarks** To evaluate the recognition performance of DAIN, we employ the ImageNet [47] pre-trained VGG-M model [49] as the initial prediction unit. We compare DAIN with both single view and multiview CNNs. As in single view CNN, we follow the standard procedure to fine-tune pre-trained networks, by replacing the classification layer with a new 39-way classification layer. Since for our GTOS database, each sample is observed with multiple viewing angles, we set multiview CNN baseline to demonstrate the effectiveness of multiview observation. To incorporate multiview information, for baseline comparison, we use two different methods: (1) voting: use the predictions from each view to vote (2) 3D filter: follow [52] to use a $3 \times 3 \times 3$ learned filter bank to convolve the multiview feature maps.

**Evaluation for DAIN Architecture** Table 4 shows the mean classification accuracy of the different three branch combination methods depicted in Figure 5. Inputs are single view images ($I_v$) and single view differential images ($I_\delta$). Combining the two streams at the final prediction layer (77% accuracy) is compared with the intermediate layer combination (74.8%) or the hybrid approach in Figure 5 (c) (79.4%) which we choose as the network architecture for following experiments. The combination method
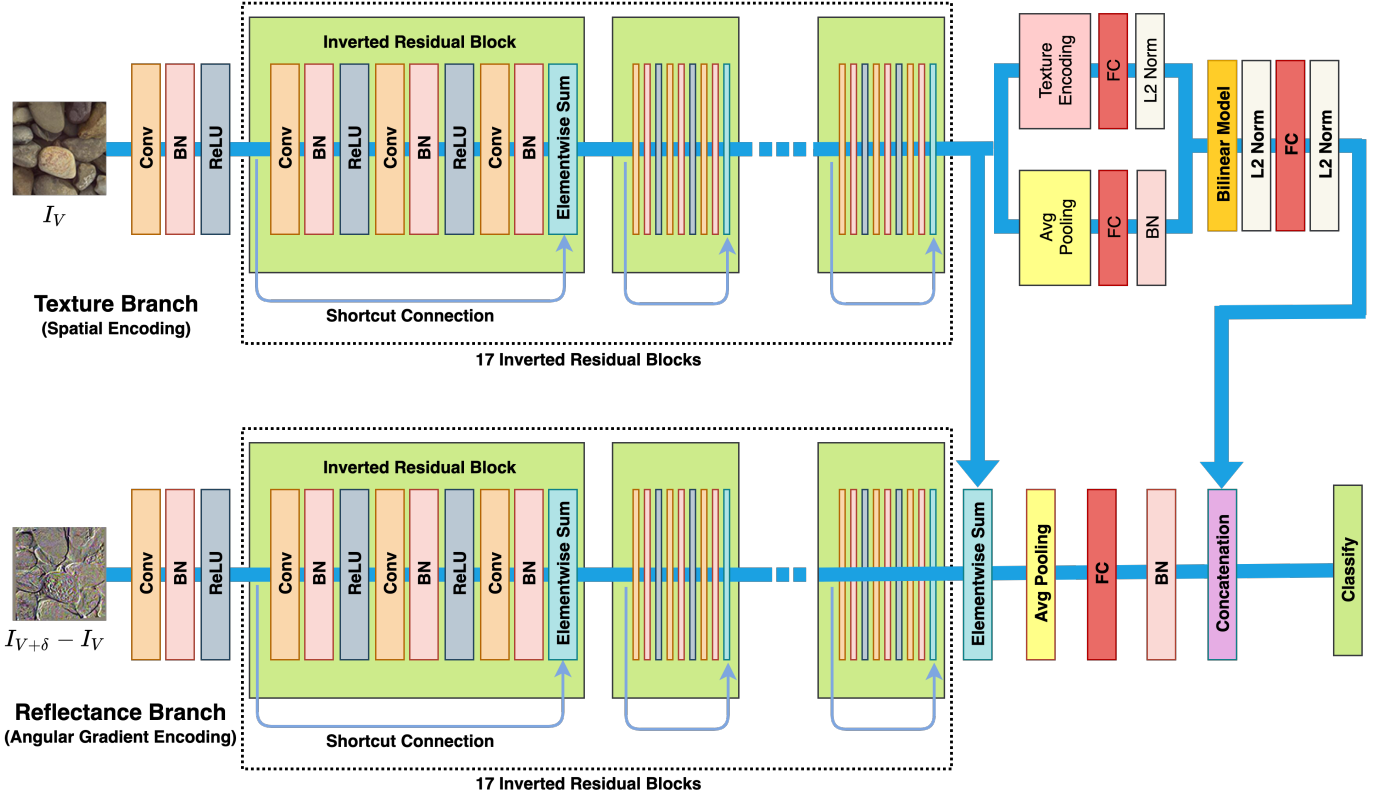
Fig. 8: A Texture Encoded Angular Network (TEAN) for material recognition. The input to the reflectance branch is the differential angular image, which captures material reflectance information via angular gradients. The input to the texture branch is the RGB color image, to provide the ordered and orderless spatial information. For the texture branch, we utilize DEP to balance the orderless texture component and ordered spatial information. The overall architecture of TEAN enables material classification using angular reflectance information, orderless texture and ordered spatial structure.

used is Sum and the feature maps are obtained from Conv5 layers after ReLU.

**DAIN Recognition Performance** We evaluate DAIN recognition performance for single view input (and differential image) and for multiview input from the GTOS database. Additionally, we compare the results to recognition using a standard CNN without a differential image stream. For all multiview experimental results we choose the number of viewpoints $N = 4$, separated by $10°$ with the starting viewpoint chosen at random (and the corresponding differential input). Table 3 shows the resulting recognition rates (with standard deviation over 5 splits shown as a subscript). The first three rows shows the accuracy *without* differential angular imaging, using both single view and multiview input. Notice the recognition performance for these non-DAIN results are generally lower than the DAIN recognition rates in the rest of the table. The middle three rows show the recognition results for single view DAIN. For combining feature maps we evaluate both Sum and Max which have comparable results. Notice that single view DAIN achieves better recognition accuracy than multiview CNN with voting (79.4% vs. 78.1%). This is an important result indicating the power of using the differential image. Instead of four viewpoints separated by $10°$, a single viewpoint and its differential image achieves a better recognition. These results provide design cues for building imaging systems tailored to material recognition. We also evaluate whether using inputs from the two viewpoints directly (i.e. $I_v$ and $I_{v+\delta}$) is comparable to using $I_v$ and the

differential image $I_\delta$. Interestingly, the differential image as input has an advantage (79.4% over 77.5%). The last three rows for VGG-M model of Table 3 show that recognition performance using multiview DAIN beats the performance of both single view DAIN and CNN methods with no differential image stream. We evaluate different ways to combine the multiview image set including voting, pooling, and the 3D filter+pooling illustrated in Figure 6.

Table 5 shows the recognition rates for multiview DAIN that outperforms three other multi-view classification method: FV+CNN [58], FV-N+CNN+N$_{3D}$ [59], and MVCNN [60]. The table shows recognition rates for a single split of the GTOS database with images resized to $240 \times 240$. All experiments are based on the same pre-trained VGG-M model. We use the same fine-tuning and training procedure as in the MVCNN [60] experiment. For FV-N+CNN+N$_{3D}$ applied to GTOS, 10 samples (out of 606) failed to get geometry information by the method provided in [59] and we removed these samples from the experiment. The patch size in [59] is $100 \times 100$, but the accuracy for this patch size for GTOS was only 43%, so we use $240 \times 240$. We implement FV-N+CNN+N$_{3D}$ with linear mapping instead of homogeneous kernel map [61] for SVM training to save memory with this larger patch size.

**DAIN with MobileNet V2** The CNN module of the two stream network can be replaced by other state-of-the-art deep learning methods to further improve results. To demonstrate this, we

| Model | Method | First input | Second input | Accuracy |
|---|---|---|---|---|
| VGG-M [49] | single view CNN | $I_v$ | - | $74.3_{\pm 2.8}$ |
| | multiview CNN, voting | $I_v$ | - | $78.1_{\pm 2.4}$ |
| | multiview CNN,3D filter | $I_v$ | - | $74.8_{\pm 3.2}$ |
| | single view DAIN (Sum) | $I_v$ | $I_{v+\delta}$ | $77.5_{\pm 2.7}$ |
| | single view DAIN (Sum) | $I_v$ | $I_\delta$ | $79.4_{\pm 3.4}$ |
| | single view DAIN (Max) | $I_v$ | $I_\delta$ | $79.0_{\pm 1.8}$ |
| | multiview (Sum/voting) | $I_v$ | $I_\delta$ | $80.0_{\pm 2.1}$ |
| | multiview DAIN (Sum/pooling) | $I_v$ | $I_\delta$ | $81.2_{\pm 1.7}$ |
| | multiview DAIN (3D filter/pooling) | $I_v$ | $I_\delta$ | $81.1_{\pm 1.5}$ |
| MobileNet V2 [50] | single view CNN | $I_v$ | - | $80.4_{\pm 3.2}$ |
| | multiview CNN, voting | $I_v$ | - | $82.5_{\pm 2.8}$ |
| | single view DAIN (Sum) | $I_v$ | $I_\delta$ | $82.5_{\pm 2.3}$ |
| | multiview DAIN (Sum/voting) | $I_v$ | $I_\delta$ | $85.8_{\pm 2.6}$ |
| | multiview DAIN (Sum/pooling) | $I_v$ | $I_\delta$ | $86.2_{\pm 2.5}$ |

TABLE 3: Results comparing performance of standard CNN recognition without angular differential imaging (first three rows) to our single-view DAIN (middle three rows) and our multi-view DAIN (bottom three rows). $I_v$ denotes the image from viewpoint $v$, $I_{v+\delta}$ is the image obtained from viewpoint $v + \delta$, and $I_\delta = I_v - I_{v+\delta}$ is the differential image. The differential angular imaging network (DAIN) has superior performance over CNN even when comparing single view DAIN to multiview CNN. Multiview DAIN provides the best recognition rates.

| Method | Final Layer Combination | Intermediate Layer Combination | Intermediate and Final layer Combination |
|---|---|---|---|
| Accuracy | $77.0_{\pm 2.5}$ | $74.8_{\pm 3.4}$ | $79.4_{\pm 3.4}$ |

TABLE 4: Comparison of accuracy from different two stream methods as shown in Figure 5. The feature-map combination method for (b) and (c) is Sum at Conv5 layers after ReLU. The reported result is the mean accuracy and the subscript shows the standard deviation over 5 splits of the data. Notice that the architecture in (c) gives the best performance and is chosen for the network architecture.

| Architecture | Accuracy |
|---|---|
| FV+CNN [58] | 75.4% |
| FV-N+CNN+N$_{3D}$ [59] | 58.3% |
| MVCNN [60] | 78.1% |
| **multiview DAIN (Sum/pooling)** | **81.2%** |

TABLE 5: Comparison with the state of art algorithms on GTOS dataset. Notice that our method, multiview DAIN, achieves the best recognition accuracy.
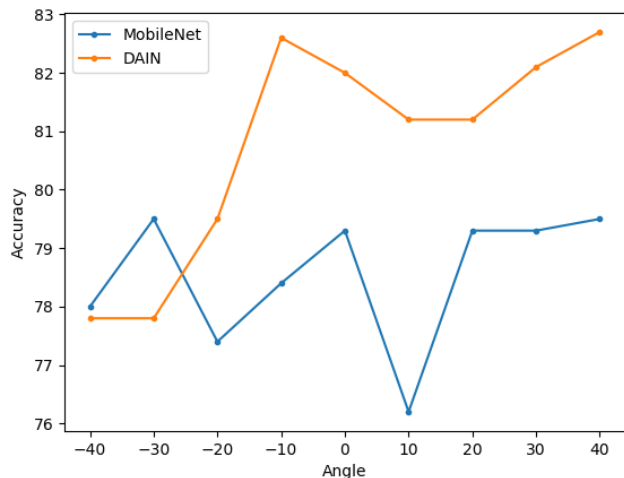


Fig. 9: The recognition accuracy of fine-tuned MobileNet V2 and DAIN on different observation angles. The recognition accuracy of DAIN outperforms fine-tuned MobileNet V2 in different observation angles.

change the CNN module to ImageNet pre-trained MobileNet V2 [50]. Combining feature maps generated from the bottleneck8_x (the eighth Bottleneck inverted residual block) with training batch size 64. The recognition results are shown in Table 3. After replacing VGG-M model with MobileNet V2, the recognition performance improves for all the methods. But still our DAIN network architecture performs better than the non-DAIN methods. Notice that sngle view DAIN achieves same recognition accuracy with smaller variance than multiview CNN with voting. The best recognition performance is the multiview DAIN (Sum/pooling) with single view images ($I_v$) and single view differential images ($I_\delta$) as inputs, which is 86.2%.

**Differential Angular Imaging Analysis** To further analyze the effect of differential angular images for recognition performance, based on GTOS split 1, we compare the recognition difference of fine-tuned MobileNet V2 and MobileNet V2 based DAIN for different observation angles and different material classes. For different observation angles, we compare the test set recognition accuracy. The result is shown in Figure 9; notice the recognition accuracy increases when observation angle moves to the center (0 degree). Also with the help of differential angular images, the recognition accuracy of DAIN outperforms fine-tuned MobileNet V2 in different observation angles. We select 10 classes from the GTOS dataset to compare the test set recognition accuracy for different material classes. The results are shown in Table 6. For materials with distinct shape or color information like painted asphalt, brick, painted turf and stone, the recognition performance is similar. But for material classes where recognition depends

| Material Class | painted asphalt | brick | **cement** | **dry grass** | **limestone** | moss | mud | painted turf | **sand** | stone |
|---|---|---|---|---|---|---|---|---|---|---|
| DAIN | 97.7 | 90.8 | 89.9 | 57.8 | 93.9 | 76.8 | 93.5 | 97.3 | 84.2 | 71.6 |
| MobileNet V2 | 98.1 | 97.3 | 72.3 | 18.5 | 87.5 | 91.2 | 98.1 | 98.9 | 56.4 | 68.1 |

TABLE 6: The recognition accuracy of MobileNet V2 based DAIN and fine-tuned MobileNet V2 on different material classes.



(a) painted asphalt    (b) brick    (c) painted turf    (d) stone

(e) dry grass    (f) limestone    (g) cement    (h) sand

Fig. 10: The sample images for different material classes. For materials with distinct shape or color information like painted asphalt, brick, painted turf and stone (top), the recognition performance is similar. But for material classes where recognition depends on material reflectance and fine-scale texture (like cement, dry grass, limestone and sand), angular gradients are an important cue and DAIN significantly outperforms MobileNet V2.

| MobileNet [50] | Bilinear CNN [55] | Deep-TEN [9] | DEP (ours) |
|---|---|---|---|
| $80.4_{\pm 3.2}$ | $80.8_{\pm 2.2}$ | $80.8_{\pm 1.5}$ | $\mathbf{83.3_{\pm 2.1}}$ |

TABLE 7: Comparison our Deep Encoding Pooling Network (DEP) with MobileNet V2 (left) [50], Bilinear CNN (mid) [55] and Deep-TEN (right) [9] on GTOS dataset. For MobileNet, we replace the 1000-way classification layer with a new classification layer, the output dimension of new classification layer is the number of classes, which is 39 for GTOS.

| Method | DTD [58] | Minc-2500 [62] |
|---|---|---|
| FV-CNN [8] | 72.3% | 63.1% |
| Deep-TEN [9] | 69.6% | 80.4% |
| DEP (ours) | **73.2%** | **82.0%** |

TABLE 8: Comparison with state-of-the-art algorithms on Describable Textures Dataset (DTD) and Materials in Context Database (MINC).

on material reflectance and fine-scale texture (like cement, dry grass, limestone and sand, see Figure 10), angular gradients are an important cue and DAIN significantly outperforms MobileNet V2.

## 5.2 DEP Recognition Results

**Recognition Benchmarks** We compare the DEP network with the following three baseline methods based on ImageNet [47] pre-trained MobileNet V2 [50]: (1) CNN with global average pooling (MobileNet), (2) CNN with texture encoding (Deep-Ten) and (3) CNN with bilinear models (Bilinear-CNN). All three methods support end-to-end training. For equal comparison, we employ the same training procedure aforementioned, and we use an identical training and evaluation procedure for each experiment.

CNN with global average pooling (MobileNet): As in single view CNN, we follow the standard procedure to fine-tune pre-trained MobileNet, by replacing the classification layer with a new 39-way classfication layer. The global average pooling works as feature pooling that encodes the $7 \times 7 \times 1280$ dimensional features from the pre-trained MobileNet V2 into a 1280 dimensional vector.

CNN with texture encoding (Deep-TEN): The Deep Texture Encoding Network (Deep-TEN) [9] embeds the texture encoding layer on top of the 50-layer pre-trained ResNet [48]. To make an equal comparison, we replace the 50-layer ResNet with MobileNet V2. As in [9], we reduce the number of CNN streams outputs channels from 1280 to 128 with a $1 \times 1$ convolutional layer. We replace the global average pooling layer in the MobileNet V2 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 [55] 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 pre-trained MobileNet V2 as CNN streams for feature extractor. Feature maps from the last convolutional layer are pooled with bilinear models. we reduce the number of CNN streams outputs channels from 1280 to 128 with a $1 \times 1$ convolutional layer before bilinear models. The dimension of feature maps for bilinear models is $7 \times 7 \times 128$ and the pooled bilinear feature is of size $128 \times 128$. The pooled bilinear feature is fed into classification layer for classification.

**DEP Recognition Performance** Table 7 is the classification accuracy of fine-tuning MobileNet V2 [50], Bilinear CNN [55], Deep-TEN [9] and the proposed DEP on the GTOS dataset. The recognition accuracy for combining spatial information and texture details (DEP) is 83.3%. That's 2.5% better than only focusing on spatial information (ResNet) and 2.5% better than only focusing on texture details (Deep-TEN).

**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) [58] and Materials in Context Database (MINC) [62]. For an equal comparison, we build DEP based on a 50-layer ResNet [48], the feature maps channels from CNN streams are reduced from 2048 to 512 with a $1 \times 1$ convolutional layer. The result is shown in Table 8, DEP outperforms the state-of-the-art on both datasets. Note that we only experiment with single scale training. As mentioned in [55], multi-scale training is likely to improve results for all methods.
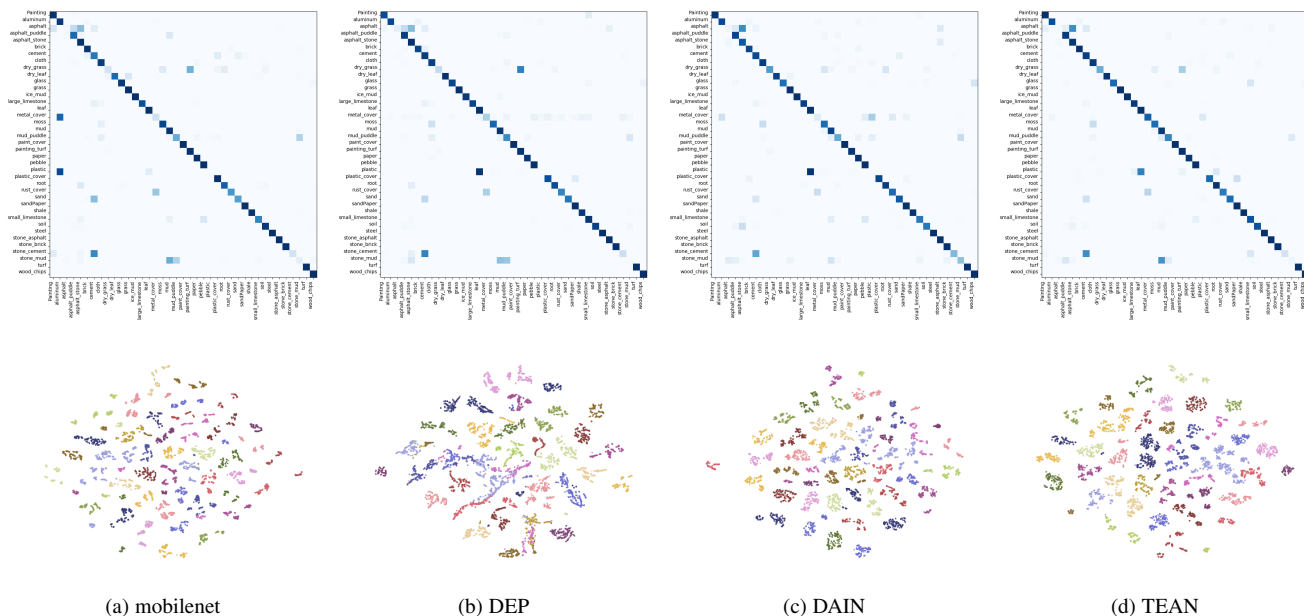
(a) mobilenet      (b) DEP      (c) DAIN      (d) TEAN

Fig. 11: The Barnes-Hut t-SNE [57] and confusion matrix of four material recognition models based on GTOS : MobileNet (left), DEP (mid left), DAIN (mid right) and TEAN (right). For Barnes-Hut t-SNE [57], we employ images from validation set and extract features before classification layers of four models for experiment. We see that TEAN separates and clusters the classes better. (Dark blue represents higher values and light blue represents lower values in the confusion matrix.)

| Method | Accuracy |
|---|---|
| single view CNN | $80.4_{\pm 3.2}$ |
| multiview CNN, voting | $82.5_{\pm 2.8}$ |
| single view DEP | $83.3_{\pm 2.1}$ |
| multiview DEP, voting | $85.8_{\pm 1.9}$ |
| single view DAIN (Sum) | $82.5_{\pm 2.3}$ |
| multiview DAIN (Sum/voting) | $85.8_{\pm 2.6}$ |
| multiview DAIN (Sum/pooling) | $86.2_{\pm 2.5}$ |
| single view TEAN (Sum) | $84.7_{\pm 2.1}$ |
| multiview TEAN (Sum/voting) | $87.4_{\pm 2.3}$ |
| multiview TEAN (Sum/pooling) | $87.6_{\pm 2.0}$ |

TABLE 9: Results comparing performance of CNN fine-tune, DEP, DAIN and TEAN based on MobileNet V2 [50].

### 5.3 TEAN Recognition Results

Table 9 is the mean classification accuracy comparison of MobileNet V2 based single view/multiview CNN fine-tune, DEP, DAIN and TEAN. As in DAIN, we experiment with voting and pooling to combine the multiview image set. From the result we can see that multiview TEAN (Sum/pooling) performs best, the recognition accuracy is 87.6%, which is 5.1% better than multiview CNN, voting baseline. Also the recognition performance for TEAN outperforms DAIN in both single view and multiview.

**Multi-scale Training** Multi-scale training is a common image augmentation trick to simulate observing materials at different distances [9], [16], [59]. We also experiment this with our GTOS dataset. We resize images into different resolutions, and randomly crop 224×224 patches for training. Following prior works [9], [16], We experiment TEAN with two groups of resolution settings: (256×256, 384×384, 512×512) and (224×224, 246×246,

268×268). For training/testing split 1, the recognition accuracy is 81.93% and 82.03% respectively, it is lower than the single view TEAN, in which the accuracy is 82.87%. Although the result is contrary with prior works [9], [16], [59], that simulating observing materials at different distances with multi-scale training is helpful for performance, we think the result is meaningful for GTOS. Since images in the GTOS dataset are captured with a fixed distance between the camera and ground terrain, the observing distance is constant for all the images. We conclude the multi-scale training is not helpful for our GTOS dataset.

### 5.4 Confusion Matrix and Feature Visualization

To gain insight into why TEAN performs best for material recognition, based on training/testing split 1, we compute the confusion matrix of MobileNet, DEP, DAIN and TEAN and visualize features before classification layers with Barnes-Hut t-SNE [57]. For features visualization, we employ images from validation set and extract features before classification layers of four models for experiment. The result is shown in Figure 11. Notice that TEAN separates and clusters the classes better.

## 6 CONCLUSION

In summary, there are three main contributions of this work: 1) The GTOS Dataset with ground terrain imaged by systematic in-scene measurement of partial reflectance instead of in-lab reflectance measurements. The database contains 34,243 images with 40 surface classes, 18 viewing directions, 4 illumination conditions, 3 exposure settings per sample and several instances/samples per class; 2) Differential Angular Imaging for a sparse representation of the spatial distribution of angular gradients that provides key cues for material recognition; 3) We develop and evaluate architectures for using differential angular imaging, texture details and spatial information for material recognition, showing

superior results for differential inputs as compared to original images. Our work in measuring and modeling outdoor surfaces has important implications for applications such as robot navigation (determining control parameters based on current ground terrain) and automatic driving (determining road conditions by partial real time reflectance measurements). The database and methods will provides a foundation for additional in-depth studies of material recognition in the wild.
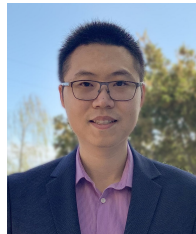
## ACKNOWLEDGMENTS

## REFERENCES

[1] K. J. Dana, "Capturing computational appearance: More than meets the eye," *IEEE Signal Processing Magazine*, vol. 33, no. 5, pp. 70–80, 2016.

[2] C. Liu and J. Gu, "Discriminative illumination: Per-pixel classification of raw materials based on optimal projections of spectral brdf," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 1, pp. 86–98, 2013.

[3] N. Salamati, C. Fredembach, and S. Süsstrunk, "Material classification using color and nir images," in *Color and Imaging Conference*, vol. 2009, no. 1. Society for Imaging Science and Technology, 2009, pp. 216–222.

[4] O. Wang, P. Gunawardane, S. Scher, and J. Davis, "Material classification using brdf slices," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 2805–2811.

[5] H. Zhang, K. Dana, and K. Nishino, "Reflectance hashing for material recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3071–3080.

[6] Z. Li, K. Sunkavalli, and M. Chandraker, "Materials for masses: Svbrdf acquisition with a single mobile phone image," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 72–87.

[7] 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*, 2015, pp. 3479–3487.

[8] 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*, 2015, pp. 3828–3836.

[9] H. Zhang, J. Xue, and K. Dana, "Deep ten: Texture encoding network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 708–717.

[10] X. Bu, Y. Wu, Z. Gao, and Y. Jia, "Deep convolutional network with locality and sparsity constraints for texture classification," *Pattern Recognition*, vol. 91, pp. 34–46, 2019.

[11] V. Andrearczyk and P. F. Whelan, "Using filter banks in convolutional neural networks for texture classification," *Pattern Recognition Letters*, vol. 84, pp. 63–69, 2016.

[12] S. Kong and C. Fowlkes, "Low-rank bilinear pooling for fine-grained classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 365–374.

[13] G. Schwartz and K. Nishino, "Visual material traits: Recognizing per-pixel material context," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2013, pp. 883–890.

[14] ——, "Automatically discovering local visual material attributes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3565–3573.

[15] 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*. Springer, 2016, pp. 121–138.

[16] J. Xue, H. Zhang, and K. Dana, "Deep texture manifold for ground terrain recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 558–567.

[17] J. Xue, H. Zhang, K. Dana, and K. Nishino, "Differential angular imaging for material recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 764–773.

[18] 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)*, vol. 18, no. 1, pp. 1–34, 1999.

[19] E. Hayman, B. Caputo, M. Fritz, and J.-O. Eklundh, "On the significance of real-world conditions for material classification," in *European conference on computer vision*. Springer, 2004, pp. 253–266.

[20] M. Weinmann, J. Gall, and R. Klein, "Material classification based on training data synthesized using a btf database," in *European Conference on Computer Vision*. Springer, 2014, pp. 156–171.

[21] G. Choe, S. G. Narasimhan, and I. So Kweon, "Simultaneous estimation of near ir brdf and fine-scale surface geometry," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2452–2460.

[22] O. G. Cula and K. J. Dana, "Compact representation of bidirectional texture functions," in *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, vol. 1. IEEE, 2001, pp. I–I.

[23] M. Varma and A. Zisserman, "A statistical approach to texture classification from single images," *International journal of computer vision*, vol. 62, no. 1-2, pp. 61–81, 2005.

[24] G. Csurka, C. Dance, L. Fan, J. Willamowski, and C. Bray, "Visual categorization with bags of keypoints," in *Workshop on statistical learning in computer vision, ECCV*, vol. 1, no. 1-22. Prague, 2004, pp. 1–2.

[25] 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)*, vol. 2. IEEE, 2006, pp. 2169–2178.

[26] L. Liu, J. Chen, P. Fieguth, G. Zhao, R. Chellappa, and M. Pietikainen, "A survey of recent advances in texture representation," *arXiv preprint arXiv:1801.10324*, vol. 3, 2018.

[27] L. Liu, J. Chen, P. Fieguth, G. Zhao, R. Chellappa, and M. Pietikäinen, "From bow to cnn: Two decades of texture representation for texture classification," *International Journal of Computer Vision*, vol. 127, no. 1, pp. 74–109, 2019.

[28] K. Simonyan and A. Zisserman, "Two-stream convolutional networks for action recognition in videos," in *Advances in Neural Information Processing Systems*, 2014, pp. 568–576.

[29] U. S. N. B. of Standards and F. E. Nicodemus, *Geometrical considerations and nomenclature for reflectance*. US Department of Commerce, National Bureau of Standards, 1977, vol. 160.

[30] C. Liu and J. Gu, "Discriminative illumination: Per-pixel classification of raw materials based on optimal projections of spectral brdf," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 1, pp. 86–98, 2014.

[31] O. Wang, P. Gunawardane, S. Scher, and J. Davis, "Material classification using brdf slices," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 2805–2811.

[32] G. J. Ward, "Measuring and modeling anisotropic reflection," in *ACM SIGGRAPH Computer Graphics*, vol. 26, no. 2. ACM, 1992, pp. 265–272.

[33] M. Levoy and P. Hanrahan, "Light field rendering," in *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996, pp. 31–42.

[34] P. Debevec, T. Hawkins, C. Tchou, H.-P. Duiker, W. Sarokin, and M. Sagar, "Acquiring the reflectance field of a human face," in *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*. ACM Press/Addison-Wesley Publishing Co., 2000, pp. 145–156.

[35] J. Filip and M. Haindl, "Bidirectional texture function modeling: A state of the art survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 11, pp. 1921–1940, Nov 2009.

[36] K. J. Dana and J. Wang, "Device for convenient measurement of spatially varying bidirectional reflectance," *JOSA A*, vol. 21, no. 1, pp. 1–12, 2004.

[37] W. Matusik, H. Pfister, M. Brand, and L. McMillan, "A data-driven reflectance model," *ACM Transactions on Graphics*, vol. 22, no. 3, pp. 759–769, Jul. 2003.

[38] J. Filip and R. Vávra, "Template-based sampling of anisotropic brdfs," in *Computer Graphics Forum*, vol. 33, no. 7. Wiley Online Library, 2014, pp. 91–99.

[39] G. Oxholm, P. Bariya, and K. Nishino, "The scale of geometric texture," in *European conference on computer vision*. Springer, 2012, pp. 58–71.

[40] C. Kampouris, S. Zafeiriou, A. Ghosh, and S. Malassiotis, "Fine-grained material classification using micro-geometry and reflectance," in *European Conference on Computer Vision*. Springer, 2016, pp. 778–792.

[41] L. Sharan, R. Rosenholtz, and E. Adelson, "Material perception: What can you see in a brief glance?" *Journal of Vision*, vol. 9, no. 8, pp. 784–784, 2009.

[42] V. Deschaintre, M. Aittala, F. Durand, G. Drettakis, and A. Bousseau, "Single-image svbrdf capture with a rendering-aware deep network," *ACM Transactions on Graphics (ToG)*, vol. 37, no. 4, pp. 1–15, 2018.
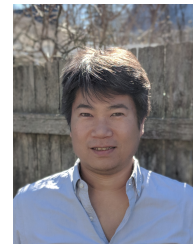
[43] Z. Li, Z. Xu, R. Ramamoorthi, K. Sunkavalli, and M. Chandraker, "Learning to reconstruct shape and spatially-varying reflectance from a single image," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 6, pp. 1–11, 2018.

[44] D. Gao, X. Li, Y. Dong, P. Peers, K. Xu, and X. Tong, "Deep inverse rendering for high-resolution svbrdf estimation from an arbitrary number of images," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 4, p. 134, 2019.

[45] V. Deschaintre, M. Aittala, F. Durand, G. Drettakis, and A. Bousseau, "Flexible svbrdf capture with a multi-image deep network," in *Computer Graphics Forum*, vol. 38, no. 4. Wiley Online Library, 2019, pp. 1–13.

[46] J. Riviere, I. Reshetouski, L. Filipi, and A. Ghosh, "Polarization imaging reflectometry in the wild," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 6, pp. 1–14, 2017.

[47] 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*. IEEE, 2009, pp. 248–255.

[48] 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*, 2016, pp. 770–778.

[49] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," *arXiv preprint arXiv:1405.3531*, 2014.

[50] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.

[51] M. Lin, Q. Chen, and S. Yan, "Network in network," *International Conference on Learning Representations*, 2014.

[52] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *2015 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2015, pp. 4489–4497.

[53] J. B. Tenenbaum and W. T. Freeman, "Separating style and content," in *Advances in neural information processing systems*, 1997, pp. 662–668.

[54] 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*. IEEE, 1997, pp. 554–560.

[55] 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*, 2015, pp. 1449–1457.

[56] F. Perronnin, J. Sánchez, and T. Mensink, "Improving the fisher kernel for large-scale image classification," *Computer Vision–ECCV 2010*, pp. 143–156, 2010.

[57] L. Van Der Maaten, "Accelerating t-sne using tree-based algorithms." *Journal of machine learning research*, vol. 15, no. 1, pp. 3221–3245, 2014.

[58] 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*, 2014, pp. 3606–3613.

[59] J. DeGol, M. Golparvar-Fard, and D. Hoiem, "Geometry-informed material recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1554–1562.

[60] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 945–953.

[61] A. Vedaldi and A. Zisserman, "Efficient additive kernels via explicit feature maps," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 3, pp. 480–492, 2012.

[62] 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*, 2015, pp. 3479–3487.

**Jia Xue** is a Ph.D. student of Electrical Computer Engineering at Rutgers University–New Brunswick, New Brunswick, NJ. He received the bachelor's degree from University of Electronic Science and Technology of China in 2015. His research interests lie in computer vision and machine learning, including material and texture recognition, segmentation, high-level and mid-level vision.



**Hang Zhang** Dr. Hang Zhang is an Applied Scientist working at Amazon Web Service Inc. Prior to joining Amazon, he received the PhD degree with Prof. Kristin Dana at Rutgers University in 2017. Before coming to Rutgers, he received the BS degree from Southeast University (China) in 2013. His research interests are material and texture recognition, neural style transfer, semantic segmentation and large scale image classification.



**Ko Nishino** is a Professor in the Department of Intelligence Science and Technology at Kyoto University. He received his B.E. and M.E. in Information and Communication Engineering in 1997 and 1999, respectively, and PhD in Computer Science in 2002, all from University of Tokyo. Before joining Kyoto University in 2018, he was a Professor in the Department of Computer Science at Drexel University. His primary research interests lie in computer vision and machine learning including appearance modeling and material recognition, human behavior analysis, and computational photography. He received the NSF CAREER award in 2008.



**Kristin J. Dana** Dr. Kristin J. Dana received the PhD from Columbia University (NY,NY) in 1999 and the MS degree from Massachusetts Institute of Technology in 1992, and a BS degree in 1990 from the Cooper Union (NY,NY). She is a Full Professor in the Department of Electrical and Computer Engineering at Rutgers University. Her research interests in computer vision include robot vision, socially cognizant robotics, deep learning, computational photography, illumination modeling, texture and reflectance. Dr. Dana is also a member of the Rutgers Center for Cognitive Science and a member of Graduate Faculty of the Computer Science Department. From 1992-1995 she was on the research staff at SRI-Sarnoff Corporation developing real-time motion estimation algorithms. She is the recipient of the National Science Foundation Career Award (2001) for a program investigating surface science for vision and graphics and a team member recipient of the Charles Pankow Innovation Award in 2014 from the ASCE. Dr. Dana currently leads an NSF National Research Traineeship (NRT) at Rutgers University entitled SOCRATES: Socially Cognizant Robotics for a Technology Enhanced Society.