



Describing Textures using Natural Language

Chenyun Wu , Mikayla Timm , and Subhransu Maji 

University of Massachusetts Amherst
{chenyun, mtimm, smaji}@cs.umass.edu

Abstract. Textures in natural images can be characterized by color, shape, periodicity of elements within them, and other attributes that can be described using natural language. In this paper, we study the problem of describing visual attributes of texture on a novel dataset containing rich descriptions of textures, and conduct a systematic study of current generative and discriminative models for grounding language to images on this dataset. We find that while these models capture some properties of texture, they fail to capture several compositional properties, such as the colors of dots. We provide critical analysis of existing models by generating synthetic but realistic textures with different descriptions. Our dataset also allows us to train interpretable models and generate language-based explanations of what discriminative features are learned by deep networks for fine-grained categorization where texture plays a key role. We present visualizations of several fine-grained domains and show that texture attributes learned on our dataset offer improvements over expert-designed attributes on the Caltech-UCSD Birds dataset.

1 Introduction

Texture is ubiquitous and provides useful cues for a wide range of visual recognition tasks. We rely on texture for estimating material properties of surfaces, for discriminating objects with a similar shape, for generating realistic imagery in computer graphics applications, and so on. Texture is localized and can be more easily modeled than shape that is affected by pose, viewpoint, or occlusion. The effectiveness of texture for perceptual tasks is also mimicked by deep networks trained on current computer vision datasets that have been shown to rely significantly on texture for discrimination (*e.g.*, [14, 16, 22, 26]).

While there has been significant work in the last few decades on visual representations of texture, limited work has been done on describing detailed properties of textures using natural language. The ability to describe texture in rich detail can enable applications on domains such as fashion and graphics, as well as to interpret discriminative attributes of visual categories within a fine-grained taxonomy (*e.g.*, species of birds and flowers) where texture cues play a key role. However, existing datasets of texture (*e.g.*, [11, 15]) are limited to a few binary attributes that describe patterns or materials, and do not describe detailed properties using the compositional nature of language (*e.g.*, descriptions of the color and shape of texture elements). At the same time, existing datasets

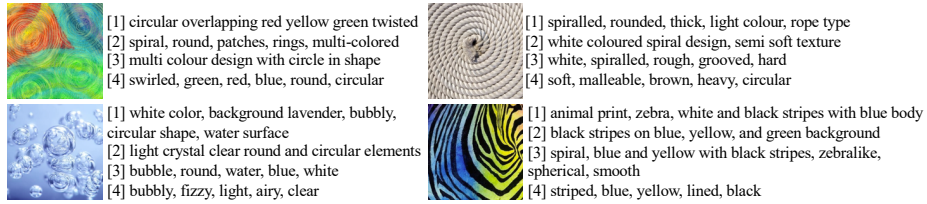


Fig. 1. We introduce the Describable Textures in Detail Dataset (DTD²) consisting of texture images from DTD [15] with natural language descriptions, which provide rich and fine-grained supervision for various aspects of texture such as color compositions, shapes, and materials. More examples are in the Supplementary material.

of language and vision [7, 27, 32, 38, 44, 48, 58] primarily focus on objects and their relations with a limited treatment of textures (Section 2). Addressing this gap in the literature, we introduce a new dataset containing natural language descriptions of textures called the Describable Textures in Detail Dataset (DTD²). It contains several manually annotated descriptions of each image from the Describable Texture Dataset (DTD) [15]. As seen in Figure 1, these contain descriptions of colors of the structural elements within the texture (*e.g.* “circles” and “stripes”), their shape, and other high-level perceptual properties of the texture (*e.g.* “soft” and “protruding”). The resulting vocabulary vastly extends the 47 attributes present in the original DTD dataset (Section 3).

We argue that the domain of texture is rich and poses many challenges for compositional language modeling that are present in existing language and vision datasets describing objects and scenes. For example, to estimate the color of dots in a dotted texture the model must learn to associate the color to the dots and not to the background. Yet the domain of texture is simple enough that it allows us to analyze the robustness and generalization of existing vision and language models by synthetically generating variations of a texture. We conduct a systematic study of existing visual representations of texture, models of language, and methods for matching the two domains on this dataset (Sections 4, 5 and 5.3). We find that adopting pre-trained language models significantly improve generalization. However, an analysis on synthetically generated variations of each texture by varying one attribute at a time (*e.g.*, foreground color and shape) shows that the representations fail to capture detailed properties.

We also present two novel applications of our dataset (Section 6). First, we visualize what discriminative texture properties are learned by existing deep networks for fine-grained classification on natural domains such as birds, flowers, and butterflies. To this end we generate “maximal images” for each category by “inverting” a texture-based classifier [35] and describe these images using captioning models trained on DTD². We find that the resulting explanations tend to be well aligned with the discriminative attributes of each category (*e.g.*, “Tiger Lily” flower is “black, red, white, and dotted” as seen in Figure 6-middle). We also show that models trained on DTD² offer improvements over expert-designed

binary attributes on the Caltech-UCSD Birds dataset [53]. This complements the capabilities of existing datasets for explainable AI on these domains that focus on shapes, parts, and their attributes such as color. Texture provides a domain-independent, albeit incomplete way of describing interpretable discriminative properties for several domains.

In summary, our contributions are:

- A novel dataset of texture descriptions (Section 3).
- An evaluation of existing models of grounding natural language to texture (Section 4 and 5).
- A critical analysis of these models using synthetic, but realistic variations of textures with their descriptions (Section 5.3).
- Application of our models for describing discriminative texture attributes and building interpretable models on fine-grained domains (Section 6).

Our dataset and code are at: <https://people.cs.umass.edu/~chenyun/texture>.

2 Related Work

Language and texture. Describing textures using language has a long history. Early works [5, 9, 49] showed that textures can be categorized along a few semantic axes such as “coarseness”, “contrast”, “complexity” and “stochasticity”. Bhusan *et al.* [12] systematically identified words in English that correspond to visual textures and analyzed their relationship to perceptual attributes of textures. This was the basis of the Describable Texture Dataset (DTD) [15] which consolidated a list of 47 texture attributes along with images downloaded from the Internet. The dataset captures attributes such as “dotted”, “chequered”, and “honeycombed”. However, it does not capture properties such as the color of the structural elements (“red and green dots”), or the attributes that describe the background color. Our goal is to model the rich space of texture attributes in a compositional manner beyond these attributes.

Datasets of images and text. The vision and language community has put significant efforts into building large-scale datasets. Image captioning datasets such as MS-COCO [32], Flickr30K [56] and Conceptual Captions [48] contain sentences describing the general content of images. The Visual Question Answering dataset [7] provides language question and answer pairs for each image, which requires more detailed understanding of the image content. In visual grounding datasets such as RefClef [27], RefCOCO [38, 58] and Flickr30K Entities [44], detailed descriptions of the target object instances are annotated to distinguish them from other objects. However, these tasks focus on recognizing object categories and descriptions of pose, viewpoint, and their relationships to other objects, and have a limited treatment of attributes related to texture.

Texture representations. Representations based on orderless aggregations of local features originally developed for texture has had an significant influence on early computer vision (*e.g.*, “Textons” [30], “Bag-of-Visual-Words” [17], higher-order statistics [45], and Fisher vector [41, 46]). Recent works (*e.g.*, [8, 16, 34]) have

shown that combining texture representations with deep networks lead to better generalization on scene understanding and fine-grained categorization tasks. Even without explicit modeling, deep networks are capable of modeling texture through convolution, pooling, and non-linear encoding layers [21]. Indeed, several works have shown that deep networks trained on existing datasets tend to rely more on texture than shape for classification [14, 22, 26, 33]. This motivates the need to develop techniques to describe texture properties using natural language as a way to explain the behavior of deep networks in an interpretable manner.

Methods for vision and language. There is a significant literature on techniques for various language and vision tasks. The Show-and-Tell [52] model was an early deep neural net based approach for captioning images that combined the convolutional image encoder followed by an LSTM [24] language decoder. Techniques for VQA are based on a joint encoding of the image and the question to retrieve or generate an answer [28, 50, 59]. For visual grounding, where the goal is to identify a region in the image given a “referring expression”, a common approach is to learn a metric over expressions and regions [36, 43, 57]. The basic architectures for these tasks have been improved in a number of ways such as by incorporating attention mechanisms [6, 20, 28, 37, 54, 59] and improved language models [19, 47]. To model the relation between texture images and their descriptions we investigate a discriminative approach, a metric-learning based approach, and a generative modeling based approach [55] on our dataset.

3 Dataset and Tasks

We begin by describing how we collected DTD² in Section 3.1, followed by the tasks and evaluation metrics in Section 3.2. DTD² contains multiple crowdsourced descriptions for each image in DTD. Each image I contains k descriptions $\mathbf{S} = \{S_1, S_2, \dots, S_k\}$ from k different annotators who are asked to describe the texture presented in the image. Instead of providing a grammatically coherent sentence, we found that it more effective for them to list a set of properties separated by commas. Thus each description S can be interpreted as a set of phrases $\{P_1, P_2, \dots, P_n\}$. Figure 1 shows some examples of the collected data. We found that the ordering of phrases in a description is somewhat arbitrary, which motivates this annotation structure. Figure 2 shows the overall dataset statistics. DTD² contains 5,369 images and 24,697 descriptions. We split the images into 60% training, 15% validation, and 25% test. Below we describe details of the dataset collection pipeline and tasks.

3.1 Dataset Collection

Annotation. We present each DTD image and its corresponding texture category to 5 different annotators on Amazon Mechanical Turk, asking them to describe the texture using natural language with at least 5 words. Describable aspects of each image include texture, color, shape, pattern, style, and material (we provided description examples of several texture categories in the guidelines).



Fig. 2. Statistics of DTD². The “overall” column in the table shows the statistics of all data, while the “frequent” column only considers the phrases (or words) that occur at least 10 (or 5) times in the training split which forms our evaluation benchmark. The cloud of phrases has the font sizes proportional to square-root of frequencies in the dataset. The vocabulary significantly expands the 47 attributes of DTD.

Verification. After collecting the raw annotations, we manually verified all of them and removed annotations that were irrelevant. For example, a breakfast waffle may have descriptions about the related food items such as strawberries instead of the texture which is our main goal. We also removed all images from “freckled” and “potholed” categories because they are primarily of human faces or scenes of roads with few texture-related terms in their descriptions. We also excluded images with fewer than 3 valid descriptions.

Post-processing. We found that the annotations (as seen in Figure 1) describing aspects of texture are often expressed as a set of phrases separated by commas, instead of a fully grammatical sentence. We did find some users who provided long unbroken sentences, but these were few and far between. Therefore, we represent each description as a set of phrases indicated by commas (“,”) or semicolons (“;”). For example, the first description of the top-right image in Figure 1 is: “spiralled, rounded, thick, light colour, rope type”, and it’s split into 5 phrases: “spiralled”, “rounded”, “thick”, “light colour”, “rope type”. For the purpose of evaluation, we consider words that appear at least 5 times and phrases that appear at least 10 times in the training split of the dataset, which results in 655 unique phrases. Although some long descriptions are lost in the process and most of the phrases are short (mostly within three words as seen in the lower histogram in Figure 2), the collection of phrases captures a rich set of describable attributes for each image. Modeling the space of phrases poses significant challenges to existing techniques for language and vision (Section 5.3).

3.2 Tasks and Evaluation Metrics

The annotation for each image is in the form of a set of descriptions, with each description in the form of a set of phrases. A phrase is an ordered list of words. We consider several tasks and evaluation metrics on this dataset described next.

Phrase retrieval. Given an image, the goal is to rank phrases $p \in \mathcal{P}$ that are relevant to the image. Here \mathcal{P} is the set of all possible phrases, restricted to 655 frequent ones. For each image, the set of “true” relevant phrases are obtained by taking the union of phrases from all descriptions of the image. We can evaluate the ranked list using various metrics described as follows:

- Mean Average Precision (MAP): area under the precision-recall curve;
- Mean Reciprocal Rank (MRR): One over the ranking of first correct phrase;
- Precision at K (P@K): precision of the top K ranked phrases ($K \in \{5, 20\}$);
- Recall at K (R@K): recall of the top K ranked phrases ($K \in \{5, 20\}$).

Image retrieval from a phrase. The task is to retrieve images given a query phrase. When taking phrases as the query, we consider all phrases $p \in \mathcal{P}$ as before and ask the retrieval model to rank all images in the test or validation set. The “true” list is all images that contain the phrase (in any of its descriptions). We consider the same metrics as the phrase retrieval task.

Image retrieval from a description. When using descriptions the query, we consider all description $s \in \mathcal{S}$ as the input. Here \mathcal{S} is the set of all descriptions in the test or validation set. We ask the retrieval model to rank all images in the corresponding set. We evaluate the rank of the image from which the description was collected (MRR metric). This metric allows us to evaluate the compositional properties of texture over phrases (*e.g.*, “red dots” + “white background”). While we only quantitatively evaluate phrases and descriptions in the dataset, the ranking models can potentially generalize to novel descriptions or phrases over the seen words. We present qualitative results and a detailed study of the models in Section 5 and 5.3.

Description generation. The task is to generate a description for an input image. Given each image I , we compare the generated description against the set of its collected descriptions $\{S_1, S_2, \dots, S_k\}$ using standard metrics for image captioning including BLEU-1,2,3,4 [40], METEOR [10], Rouge-L [31] and CIDEr [51]. However, we note that the task is open-ended and qualitative visualizations are just as important as these metrics.

4 Methods

We investigate three techniques to learn the mapping between visual texture and natural languages on our dataset — a discriminative approach, a metric learning approach, and a language generation approach. They are explained in detail in the next three sections.

4.1 A Discriminative Approach

A simple baseline is to treat each phrase $p \in \mathcal{P}$ as a binary attribute and train a multi-label classifier to map the images to phrase labels. Given a texture image I , let $\psi(I)$ be an embedding computed using a deep network. We use activations from layer 2 and layer 4 of ResNet101 [23] with mean-pooling over spatial locations as the image embedding. A comparison of features from different ResNet layers is included in the supplemental material. For the classification task, we attach a classifier head h to map the embeddings to a 655-dimensional space corresponding to each phrase in our frequent set \mathcal{P} . The function h is modeled as a two-layer network – the first is fully-connected layer with 512 units with

BatchNorm and ReLU activation; the second is a linear layer with 655 units followed by sigmoid activation. Given a training set of $\{(I_i, Y_i)\}_{i=1}^N$ where Y_i is the ground-truth binary labels across 655 classes for image I_i , the model is trained to minimize the binary cross-entropy loss: $L_{BCE} = \sum_i \ell_{bce}(h \circ \psi(I_i), Y_i)$, where $\ell_{bce}(y, z) = \sum_i (z_i \log(y_i) - (1 - z_i) \log(1 - y_i))$.

Training details. The ResNet101 is initialized with weights pre-trained on ImageNet [18] and fine-tuned on our training data for 75 epochs using the Adam optimizer [29] with an initial learning rate at 0.0001. We use image size 224×224 for all our experiments. The hyper-parameters are selected on the validation set.

Evaluation setup. The classification scores over each phrase for each image are directly used to rank images or phrases for phrase retrieval or image retrieval with phrase input. Retrieving images given a description is more challenging since we need to aggregate the scores corresponding to different phrases, and the phrases in input descriptions may not be in \mathcal{P} . We found the following strategy works well: Given a description $S = \{P_1, P_2, \dots, P_n\}$ and an image I , obtain the scores for each phrase $s(P_i) = \sigma(h \circ \psi(I))_k$ where k is the index of the phrase $P_i \in \mathcal{P}$. If the phrase is not in the set, we consider all its sub-sequences that are present in \mathcal{P} and average the scores of them instead. For example, if the phrase “red maroon dot” is not present in \mathcal{P} , we consider all sub-sequences {red maroon, maroon dot, red, maroon, dot}, score each that is present in \mathcal{P} separately and then average the scores. By concatenating the top 5 phrases for an image we can also use the classifier to generate a description for an image. The key disadvantage of the classification baseline is that it treats each phrase independently, and does not have a natural way to score novel phrases (our baseline using sub-sequences is an attempt to handle this).

4.2 A Metric Learning Approach

The metric learning approach aims to learn a common embedding over the images and phrases such that nearby image and phrase pairs in the embedding space are related. We adopt the standard metric learning approach based on triplet-loss [25]. Consider an embedding of an image $\psi(I)$ and of a phrase $\phi(P)$ in \mathbb{R}^d . Denote $\|\psi(I) - \phi(P)\|_2^2$ as the squared Euclidean distance between the two embeddings. Given an annotation (I, P) consisting of a positive (image, phrase) pair, we sample from the training set a negative image I' for P , and a negative phrase P' for I . We consider two losses; one from the negative phrase:

$$L_p(I, P, P') = \max(0, 1 + \|\psi(I) - \phi(P)\|_2^2 - \|\psi(I) - \phi(P')\|_2^2)$$

and another from the negative image:

$$L_i(P, I, I') = \max(0, 1 + \|\psi(I) - \phi(P)\|_2^2 - \|\psi(I') - \phi(P)\|_2^2)$$

The metric learning objective is to learn embeddings ψ and ϕ that minimize the loss $L = \mathbb{E}_{(I, P), (I', P')} (L_p + L_i)$ over the training set.

For embedding images, we use the same encoder as the classification approach with features from layer 2 and 4 from ResNet101. We add an additional linear

layer with 256 units resulting in the embedding dimension $\psi(I) \in \mathbb{R}^{256}$. One advantage of the metric learning approach is that it allows us to consider richer embedding models for phrases. Specially we consider the following encoders:

- **Mean-pooling**: $\phi_{mean}(P) = \frac{1}{N_w} \sum_{w \in \text{tokenize}(P)} \text{embed}(w)$, where `tokenize`(\cdot) splits the phrase into a list of words, `embed`(\cdot) encodes each token into \mathbb{R}^{300} .
- **LSTM** [47]: $\phi_{lstm}(P) = \text{biLSTM}[\text{embed}(w) \text{ for } w \text{ in } \text{tokenize}(P)]$, with the same `tokenize`(\cdot) and `embed`(\cdot) as above. `biLSTM`(\cdot) is a bi-directional LSTM with a single layer and hidden dimension 256 that returns the concatenation of the outputs on the last token from both directions.
- **ELMo** [42]: $\phi_{elmo}(P) = \text{ELMo}(P)$, where `ELMo`(\cdot) uses pre-trained ELMo model [4] with its own tokenizer, and outputs the average embedding of all tokens in the phrase P .
- **BERT** [19]: $\phi_{bert}(P) = \text{BERT}(P)$, where `BERT`(\cdot) uses pre-trained BERT model [3] with its own tokenizer, and outputs the average of last hidden states of all tokens in the phrase P .

To compute the final embedding of the phrase $\phi(P)$, we add a linear layer to map the embeddings to 256 dimensions compatible with the image embeddings.

Training details. We trained this model on our training split using the Adam optimizer [29] with an initial learning rate at 0.0001. We found this model to be more prone to over-fitting than the classifier. Stopping the training when the image retrieval and phrase retrieval MAP on the validation set stops improving was effective. Same as the classifier, ResNet101 is initialized with ImageNet [18] weights and fine-tuned on our data. `embed`(\cdot) in ϕ_{mean} and ϕ_{lstm} was initialized with FastText embeddings [1, 13] and tuned end-to-end. Pre-trained encoders ϕ_{elmo} and ϕ_{bert} were fixed in our training.

Evaluation setup. Given the joint embedding space, one can retrieve phrases for each image and images for each phrase based on the Euclidean distance. Similar to the classifier we concatenate the top 5 retrieved phrases as a baseline description generation model. We also investigate a metric learning approach over descriptions rather than phrases where the positive and negative triplets are computed over (image, description) pairs. The language embedding models are the same since they can handle descriptions of arbitrary length.

4.3 A Generative Language Approach

We adopt the Show-Attend-Tell model [55], a widely used model for image captioning. It combines a convolutional network to encode input images with an attention-based LSTM decoder to generate descriptions. Following the default setup, we encode images into the spatial features from the 4-th layer of ResNet101 (initialized with ImageNet [18] weights). The word embeddings are initialized from FastText [1, 13]. The entire model is then trained end-to-end on the training set, using the Adam optimizer [29] with initial learning rate 0.0001 for the image encoder and 0.0004 for the language decoder. We apply early stopping based on the BLEU-4 score of generated descriptions on the validation images.

Table 1. Phrase retrieval and image retrieval on DTD². Metric learning models are trained with phrase input. Among the language encoders BERT works the best.

Task:		Phrase Retrieval						Image Retrieval					
Data Split	Model	MAP	MRR	P@5	P@20	R@5	R@20	MAP	MRR	P@5	P@20	R@5	R@20
Validation	MetricLearning: MeanPool	18.80	48.66	23.13	16.20	11.52	31.54	7.19	16.18	7.60	6.56	3.36	11.44
	MetricLearning: biLSTM	23.53	58.78	31.85	18.73	15.83	36.31	8.31	17.46	8.15	7.06	4.21	13.40
	MetricLearning: ELMo	28.13	68.46	37.02	21.11	18.44	41.12	11.25	24.05	12.79	10.27	5.85	18.57
	MetricLearning: BERT	31.68	72.59	40.67	22.96	20.23	44.50	15.22	31.39	16.27	12.56	9.07	25.69
Test	Classifier: Feat 2,4	27.12	61.28	33.50	21.71	16.07	41.48	14.75	33.94	18.75	16.02	6.47	19.32
	MetricLearning: BERT	31.77	74.12	41.70	23.60	20.17	45.04	13.50	31.12	16.52	14.57	5.24	17.32

Table 2. Retrieving texture images with descriptions as input.

Model	MRR
Classifier	12.40
MetricLearning(phrase)	12.92
MetricLearning(description)	13.95

Table 3. Description generation on textures. Synthesizing descriptions from phrases retrieved by the metric-learning based approach outperforms other baselines.

Model	Bleu-1	Bleu-2	Bleu-3	Bleu-4	METEOR	Rouge-L	CIDEr
Classifier: top 5	68.07	46.17	28.39	14.44	19.89	48.13	44.73
MetricLearning: top 5	72.99	53.69	34.97	19.39	21.81	49.70	47.34
Show-Attend-Tell	59.90	40.41	26.52	16.35	19.92	46.64	37.47

This model is primarily for the description generation task. For evaluation, we apply beam search with a beam-size of 5 to compute the best description.

5 Experiments and Analysis

5.1 Phrase and Image Retrieval

Table 1 and 2 compare the classifier and the metric learning model on phrase and image retrieval tasks as described in Section 3.2. Figure 3 and 4 show examples of the top 5 retrieved images and phrases.

In Table 1 we first compare language encoders on the metric learning model. The performance of both phrase and image retrieval depends largely on the language encoder, and BERT performs the best. The metric learning model is better at phrase retrieval while the classifier is slightly better at image retrieval.

Table 2 shows results of image retrieval from descriptions and here too the metric learning model outperforms the other two models. As shown in Figure 3-right, although the models trained on phrases work reasonably well, the metric learning model trained on descriptions handles long queries better.

5.2 Description Generation

We compare the Show-Attend-Tell model [55] with a retrieval based approach. From the classifier and the metric learning model we retrieve the top 5 phrases and concatenate them in the order of their score to form a description. As shown in Table 3, the metric learning model reaches higher scores on the metrics. However, notice that in Figure 4 the generative model’s descriptions are more fluent and covers both the color and pattern of the images, while the retrieval baselines (especially the classifier) repeat phrases with similar meanings.

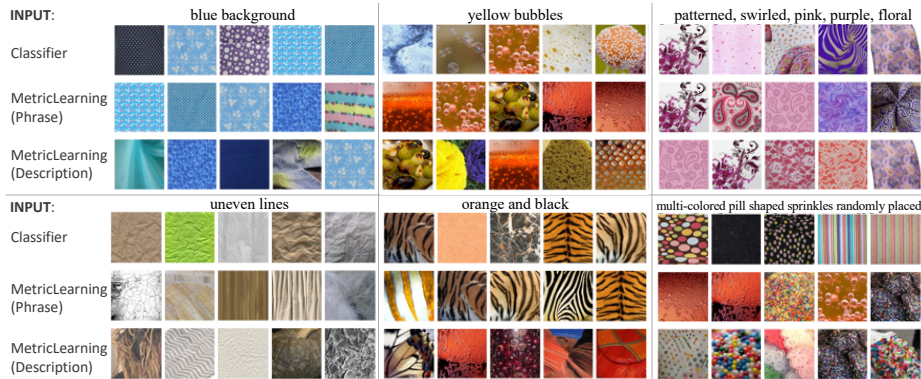


Fig. 3. Retrieve DTD² test images with language input. We show top 5 retrieved images from the classifier and the metric learning model (trained with phrase or description input). From left to right we show examples of (1) phrases the classifier has been trained on, (2) phrases beyond the frequent classes, and (3) full descriptions.

Table 4. Image retrieval performance of R-Precision on synthetic tasks.

Model	Foreground	Background	Color+Pattern	Two-colors
Classifier	45.45±20.34	59.82±9.63	35.95±21.48	26.82±14.17
MetricLearning - phrase	46.55±20.65	52.00±6.32	41.73±22.77	27.45±15.13
MetricLearning - description	47.64±18.97	53.64±4.66	35.77±21.12	21.59±13.77
Random guess	50.00	50.00	7.40	5.26

5.3 A Critical Analysis of Language Modeling

In this section, we evaluate the proposed models on tasks where we systematically vary the distribution of underlying texture attributes. This is relatively easy to do for textures than for natural images (*e.g.*, changing the color of dots) and allows us to understand the degree to which the models learn disentangled representations. We describe four tasks with varying degrees of difficulty to highlight the strengths and weaknesses of these models.

Automatically generating textures and their descriptions. To systematically generate textures with descriptions, we follow this procedure:

- Take the 11 most frequent colors in DTD² (white, black, brown, green, blue, red, yellow, pink, orange, gray, purple) and set their RGB values manually.
- Take 10 typical two-color images from ten different categories. We choose:
 - Type A: 5 images with “foreground on background”: [‘dots’, ‘polka-dots’, ‘swirls’, ‘web’, ‘lines’ (thin lines on piece of paper)], and
 - Type B: 5 images with no clear distinction between the foreground and background: [‘squares’ (checkered), ‘hexagon’, ‘stripes’ (zebra-like), ‘zigzagged’, ‘banded’ (bands with similar width)].
- For each of these 10 images, we manually extract masks for the foreground and background (Type A), or two foreground colors (Type B).



Fig. 4. Phrase retrieval and description generation on DTD² test images. For each input image, we list ground-truth descriptions beneath, and generated descriptions on the right. For the classifier and the metric learning model, we concatenate the top 5 retrieved phrases. Bold words are the ones included in ground-truth descriptions.

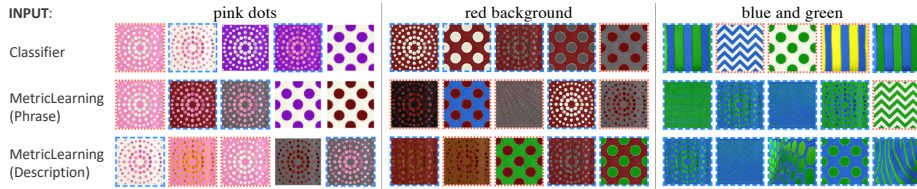


Fig. 5. Retrieval on synthetic images. Positive images are in dashed blue borders, hard negative ones are in dotted red borders.

- For each of the 10 images, generate a new image by picking 2 different colors from the 11 and modify pixel values of the two regions using the corresponding RGB value. This results in $10 \times 11 \times 10 = 1,100$ images.
- For each synthetic image, we construct the ground-truth description with the template as “[color1] [pattern], [color2] background” (such as “pink dots, white background”) for Type A, and “[color1] and [color2] [pattern]” (such as “yellow and gray squares”) for Type B.

Experiment 1: Foreground. On Type A set we construct:

- **Query:** A query of the form “[color=c] [pattern=p]” (e.g. “pink dots”).
- **Positive set:** [color=c] [pattern=p] on randomly colored background (e.g. “pink dots, white background”).
- **Negative set:** Randomly colored ($\neq c$) [pattern=p] on [color=c] background (e.g. “blue dots, pink background”).
- **Result:** Input the query description, we use the models to rank images from both the positive and negative set, and report R-Precision: the precision of

top R predictions, where R is the number of positive images. The results are listed in Table 4 first column. Since half the images have the right attribute the chance performance is 50% and the various models are nearly at the chance level. Figure 5 shows that the model is unable to distinguish between “pink dots” and “dots on a pink background”. This illustrates that the models are unable to associate color correctly with the foreground shapes.

Experiment 2: Background. This is similar to Experiment 1 but we focus on the background instead. On Type A set we construct: we know the name of its pattern (such as “dots”, “squares”, selected from the more frequent phrases that matches the category) and names of two colors (color1 and color2).

- **Query:** A query “[color=c] background” (*e.g.* “pink background”).
- **Positive set:** Randomly colored pattern on [color=c] background (*e.g.* “red dots on pink background”).
- **Negative set:** Random pattern of [color=c] on any [color≠c] background (*e.g.* “pink dots on white background”).
- **Result:** R-precision is shown in Table 4 second column. Once again the chance performance is 50% and the various models are nearly at the chance level. Figure 5-middle shows that the model is unable to distinguish between “red background” and “red dots on random background”.

Experiment 3: Color+Pattern. On both Type A and B images we construct:

- **Query:** A query “[color=c] [pattern=p]” (*e.g.* “pink dots”).
- **Positive set:** [color=c] [pattern=p] on random colored background, or with another color (*e.g.* “pink dots, white background”, “pink and blue squares”).
- **Negative set:** [color=c] [pattern≠p] or [color≠c] [pattern=p]. The negative set contains images with the correct pattern but wrong color or the wrong pattern with the right color (*e.g.*, “red dots” or “pink stripes”). Similar patterns (*e.g.*, “lines” *vs.* “banded”) are not considered negative.
- **Result:** The positive and negative set is unbalanced which results in a chance performance of 7.4%. The models presented in the earlier section are able to rank the correct color and pattern combinations ahead of the negative set and achieve a considerably higher performance.

Experiment 4: Two colors. On both Type A and B images we construct:

- **Query:** A query “[color=c1] and [color=c2]” (*e.g.* “pink and green”).
- **Positive set:** [color=c1] of random pattern on [color=c2] background (*e.g.* “pink dots on green background”), or [color=c1] and [color=c2] of random pattern (*e.g.* “pink and green squares”).
- **Negative set:** pattern with one color from {c1, c2} and another color ≠{c1,c2} (*e.g.*, “pink dots on yellow background”, “green and blue stripes”).
- **Result:** The positive and negative set are unbalanced which results in a chance performance of 5.26%. The models once again are able to rank the two color combinations ahead of the negative set and achieve a considerably higher performance. Figure 5-right shows an example.

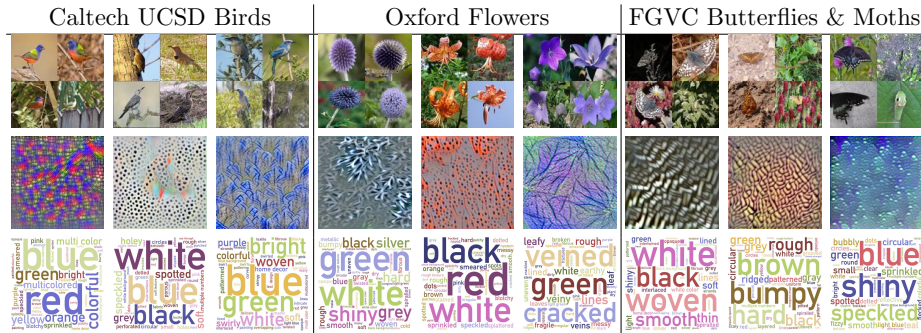


Fig. 6. Fine-grained categories visualized as their training images (top row), maximal texture images (middle row), and texture attributes (bottom row). The size of each phrase in the cloud is inversely decided by its Euclidean distance to the input maximal texture image calculated by the triplet model.

Summary. These experiments reveal that these models do indeed exhibit some high-level discriminative abilities (Exp. 3, 4), but they fail to disentangle properties such as the color of the foreground elements from background (Exp. 1, 2). This leaves much room for improvement, motivating future work, such as those that enforces spatial agreement between different attributes.

6 Applications

Describing textures of fine-grained categories. We analyze how the categories in fine-grained domains can be described by their texture. We consider categories from Caltech-UCSD Birds (CUB) [53], Oxford Flowers [39], and FGVC Butterflies and Moths [2] datasets. For each category, we follow the visualizing deep texture representations [33] to generate “maximal textures” – inputs that maximize the class probability using multi-layer bilinear CNN classifier [35]. These are provided as input to our metric learning model (with BERT encoder and phrase input) trained on DTD² to retrieve the top phrases. Figure 6 shows several categories with their maximal textures and a “phrase cloud” of the top retrieved phrases. These provide a qualitative description of each category.

Fine-grained classification with texture attributes. Here we apply models trained on our DTD² on the CUB dataset to show that embedding images into the space of texture attributes allows interpretable models for discriminative classification. Specifically, we input each image from the dataset to our phrase classifier (trained on DTD² and fixed) and obtain the log-likelihood over the 655 texture phrases as an embedding. We train a logistic regression model for the 200-way classification task. The dataset also comes with 312 binary attributes that describe the shape, pattern and color of specific parts of a bird, such as “has tail shape squared tail”, “has breast pattern spotted”, “has wing color yellow” (42 attributes for “shape”, 31 for “pattern” and 239 for “color”). We also train a logistic regression classifier on top of these attributes.

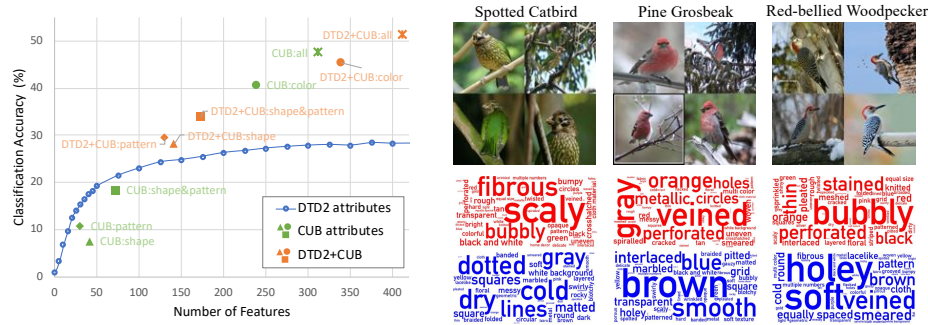


Fig. 7. Classification on CUB dataset with DTD² texture attributes. *Left:* classification accuracy vs. number of input features. Orange and green markers with the same shape are comparable with the same set of CUB attributes with or without the DTD² attributes. *Right:* The phrase clouds display important phrases for a few bird categories. Red phrases correspond to positive weights and blue are negative for a linear classifier for the category. Font sizes represent the absolute value of the coefficient.

Figure 7 shows the performance by varying the number of texture phrases ranked by their frequency on DTD² as the blue curve. It also shows a comparison of bird-specific attributes from CUB with generic texture attributes learned from DTD². Results using CUB attributes are shown in green, while those using combinations of CUB and texture attributes are shown in orange. Texture attributes are able to distinguish bird species with an accuracy of 28.5%, outperforming CUB shape and pattern attributes. However, they do not outperform the part-based color attributes that are highly effective. Yet, combining CUB attributes with texture attributes lead to consistent improvements. On the right is the visualization of discriminative texture attributes for some categories: we display phrases with the most positive weights in red, and those with the most negative weights in blue. They provide a basis for interpretable explanations of discriminative features without requiring a category-specific vocabulary.

7 Conclusion

We presented a novel dataset of textures with natural language descriptions and analyzed the performance of several language and vision models. The domain of texture is poses challenges to existing models which fail to learn a sufficiently disentangled representation leading to poor generalization on synthetic tasks. Yet, the models show some generalization to novel domains and enabling us to provide interpretable models for describing some fine-grained domains. In particular they are complementary to existing domain-specific attributes on the CUB dataset.

Acknowledgements. We would like to thank Mohit Iyyer for helpful discussions and feedback. The project is supported in part by NSF grants #1749833 and #1617917. Our experiments were performed in the UMass GPU cluster obtained under the Collaborative Fund managed by the Mass. Technology Collaborative.

References

1. FastText pretrained embeddings <https://dl.fbaipublicfiles.com/fasttext/vectors-english/wiki-news-300d-1M.vec.zip>
2. FGVC Butterflies and Moths Dataset, <https://sites.google.com/view/fgvc6/competitions/butterflies-moths-2019>
3. Pretrained BERT of version “bert-base-uncased” https://huggingface.co/transformers/pretrained_models.html
4. Pretrained ELMo https://allennlp.s3.amazonaws.com/models/elmo/2x4096_512_2048cnn_2xhighway/elmo_2x4096_512_2048cnn_2xhighway_weights.hdf5
5. Amadasun, M., King, R.: Textural features corresponding to textural properties. *IEEE Transactions on Systems, Man, and Cybernetics* **19**(5), 1264–1274 (1989)
6. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang, L.: Bottom-up and top-down attention for image captioning and visual question answering. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 6077–6086 (2018)
7. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: VQA: Visual Question Answering. In: *IEEE International Conference on Computer Vision (ICCV)* (2015)
8. Arandjelovic, R., Gronat, P., Torii, A., Pajdla, T., Sivic, J.: NetVLAD: CNN architecture for weakly supervised place recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 5297–5307 (2016)
9. Bajcsy, R.: Computer description of textured surfaces. In: *Proceedings of the 3rd international joint conference on Artificial intelligence*. pp. 572–579 (1973)
10. Banerjee, S., Lavie, A.: METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In: *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*. pp. 65–72. Association for Computational Linguistics, Ann Arbor, Michigan (Jun 2005), <https://www.aclweb.org/anthology/W05-0909>
11. Bell, S., Upchurch, P., Snavely, N., Bala, K.: Material recognition in the wild with the materials in context database. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 3479–3487 (2015)
12. Bhushan, N., Rao, A.R., Lohse, G.L.: The texture lexicon: Understanding the categorization of visual texture terms and their relationship to texture images. *Cognitive Science* **21**(2), 219–246 (1997)
13. Bojanowski, P., Grave, E., Joulin, A., Mikolov, T.: Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics (TACL)* **5**, 135–146 (2017)
14. Brendel, W., Bethge, M.: Approximating CNNs with Bag-of-Local-Features Models works surprisingly well on ImageNet. In: *International Conference on Learning Representations (ICLR)* (2019)
15. Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing textures in the wild. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014)
16. Cimpoi, M., Maji, S., Vedaldi, A.: Deep filter banks for texture recognition and segmentation. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015)
15. Csurka, G., Dance, C., Fan, L., Willamowski, J., Bray, C.: Visual categorization with bags of keypoints. In: *Workshop on statistical learning in computer vision, ECCV*. vol. 1, pp. 1–2. Prague (2004)

18. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2009)
19. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
20. Gao, P., Jiang, Z., You, H., Lu, P., Hoi, S.C., Wang, X., Li, H.: Dynamic fusion with intra-and inter-modality attention flow for visual question answering. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6639–6648 (2019)
21. Gatys, L.A., Ecker, A.S., Bethge, M.: Image style transfer using convolutional neural networks. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2414–2423 (2016)
22. Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F.A., Brendel, W.: ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In: International Conference on Learning Representations (ICLR) (2018)
23. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 770–778 (2016)
24. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Computation* **9**(8), 1735–1780 (1997)
25. Hoffer, E., Ailon, N.: Deep metric learning using triplet network. In: International Workshop on Similarity-Based Pattern Recognition. pp. 84–92. Springer (2015)
26. Hosseini, H., Xiao, B., Jaiswal, M., Poovendran, R.: Assessing shape bias property of convolutional neural networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops (2018)
27. Kazemzadeh, S., Ordonez, V., Matten, M., Berg, T.: Referitgame: Referring to objects in photographs of natural scenes. In: Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). pp. 787–798 (2014)
28. Kim, J.H., Jun, J., Zhang, B.T.: Bilinear attention networks. In: Advances in Neural Information Processing Systems. pp. 1564–1574 (2018)
29. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
30. Leung, T., Malik, J.: Representing and recognizing the visual appearance of materials using three-dimensional textons. *International Journal of Computer Vision (IJCV)* **43**(1), 29–44 (2001)
31. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. In: Proceedings of the ACL Workshop: Text Summarization Braches Out 2004. p. 10 (01 2004)
32. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European Conference on Computer Vision (ECCV). pp. 740–755. Springer (2014)
33. Lin, T.Y., Maji, S.: Visualizing and Understanding Deep Texture Representations. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2791–2799 (2016)
34. Lin, T.Y., RoyChowdhury, A., Maji, S.: Bilinear CNN models for fine-grained visual recognition. In: IEEE International Conference on Computer Vision (ICCV) (2015)
35. Lin, T.Y., RoyChowdhury, A., Maji, S.: Bilinear Convolutional Neural Networks for Fine-grained Visual Recognition. *IEEE transactions on pattern analysis and machine intelligence(PAMI)* **40**(6), 1309–1322 (2018)

36. Liu, J., Wang, L., Yang, M.H.: Referring expression generation and comprehension via attributes. In: IEEE International Conference on Computer Vision (ICCV) (2017)
37. Liu, X., Wang, Z., Shao, J., Wang, X., Li, H.: Improving referring expression grounding with cross-modal attention-guided erasing. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1950–1959 (2019)
38. Mao, J., Huang, J., Toshev, A., Camburu, O., Yuille, A., Murphy, K.: Generation and comprehension of unambiguous object descriptions. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
39. Nilsback, M.E., Zisserman, A.: Automated flower classification over a large number of classes. In: Indian Conference on Computer Vision, Graphics and Image Processing (ICVGIP) (Dec 2008)
40. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting on association for computational linguistics. pp. 311–318. Association for Computational Linguistics (2002)
41. Perronnin, F., Liu, Y., Sánchez, J., Poirier, H.: Large-scale image retrieval with compressed fisher vectors. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3384–3391. IEEE (2010)
42. Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L.: Deep contextualized word representations. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers). pp. 2227–2237 (Jun 2018)
43. Plummer, B.A., Kordas, P., Hadi Kiapour, M., Zheng, S., Piramuthu, R., Lazebnik, S.: Conditional image-text embedding networks. In: European Conference on Computer Vision (ECCV). pp. 249–264 (2018)
44. Plummer, B.A., Wang, L., Cervantes, C.M., Caicedo, J.C., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In: IEEE International Conference on Computer Vision (ICCV). pp. 2641–2649 (2015)
45. Portilla, J., Simoncelli, E.P.: A parametric texture model based on joint statistics of complex wavelet coefficients. *International Journal of Computer Vision (IJCV)* **40**(1), 49–70 (2000)
46. Sánchez, J., Perronnin, F., Mensink, T., Verbeek, J.: Image classification with the fisher vector: Theory and practice. *International Journal of Computer Vision (IJCV)* **105**(3), 222–245 (2013)
47. Schuster, M., Paliwal, K.K.: Bidirectional recurrent neural networks. *IEEE transactions on Signal Processing* **45**(11), 2673–2681 (1997)
48. Sharma, P., Ding, N., Goodman, S., Soricut, R.: Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 2556–2565 (2018)
49. Tamura, H., Mori, S., Yamawaki, T.: Textural features corresponding to visual perception. *IEEE Transactions on Systems, Man, and Cybernetics* **8**(6), 460–473 (1978)
50. Tan, H., Bansal, M.: Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490* (2019)
51. Vedantam, R., Lawrence Zitnick, C., Parikh, D.: Cider: Consensus-based image description evaluation. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4566–4575 (2015)

52. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3156–3164 (2015)
53. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The Caltech-UCSD Birds-200-2011 Dataset. Tech. Rep. CNS-TR-2011-001, California Institute of Technology (2011)
54. Wang, P., Wu, Q., Cao, J., Shen, C., Gao, L., Hengel, A.v.d.: Neighbourhood watch: Referring expression comprehension via language-guided graph attention networks. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1960–1968 (2019)
55. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: International Conference on Machine Learning (ICML). pp. 2048–2057 (2015)
56. Young, P., Lai, A., Hodosh, M., Hockenmaier, J.: From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics (TACL)* **2**, 67–78 (2014)
57. Yu, L., Lin, Z., Shen, X., Yang, J., Lu, X., Bansal, M., Berg, T.L.: Mattnet: Modular attention network for referring expression comprehension. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
58. Yu, L., Poirson, P., Yang, S., Berg, A.C., Berg, T.L.: Modeling context in referring expressions. In: European Conference on Computer Vision (ECCV). pp. 69–85. Springer (2016)
59. Yu, Z., Yu, J., Cui, Y., Tao, D., Tian, Q.: Deep modular co-attention networks for visual question answering. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6281–6290 (2019)