

# Text2Model: Model Induction for Zero-shot Generalization Using Task Descriptions

Ohad Amosy  
Bar Ilan University  
Ramat Gan, Israel  
amosyoh@biu.ac.il

Tomer Volk  
Technion - IIT  
Haifa, Israel

Eyal Ben-David  
Technion - IIT  
Haifa, Israel

Roi Reichart  
Technion - IIT  
Haifa, Israel

Gal Chechik  
Bar Ilan University, Ramat Gan, Israel  
NVIDIA Research, Tel Aviv, Israel

## Abstract

We investigate the challenge of generating a task-specific visual classifier without visual training samples, only using textual descriptions of the output classes. Unlike approaches that learn a fixed representation of the output classes, we generate at inference time a model tailored to a query classification task - a learning setup which we call *text2model*. To generate task-based zero-shot classifiers, we train a hypernetwork that receives class descriptions and outputs a multi-class model. The hypernetwork architecture is designed to be equivariant with respect to the set of descriptions and the classification layer, thus obeying the symmetries of the problem and improving generalization. Our approach generates non-linear classifiers and can handle rich textual descriptions, e.g., descriptions that include negation. We name this approach *T2M-HN* and evaluate it in a wide series of zero-shot recognition tasks, for image, point-cloud, and action recognition, using a range of text descriptions: From single words to rich descriptions. Our results demonstrate strong improvements over previous approaches, showing that zero-shot learning can be applied with little multi-model training data.

## 1. Introduction

People can classify perceived objects by following language instructions, like “separate soft toys from hard ones” or “collect the furry toy animals” [34]. Developing models that have this capacity has many applications in open-world domains where the label space is not fully known during training or when labels are expensive or hard to obtain. Indeed, numerous attempts have been made to achieve zero-shot classification capacity (see related methods). Unfor-

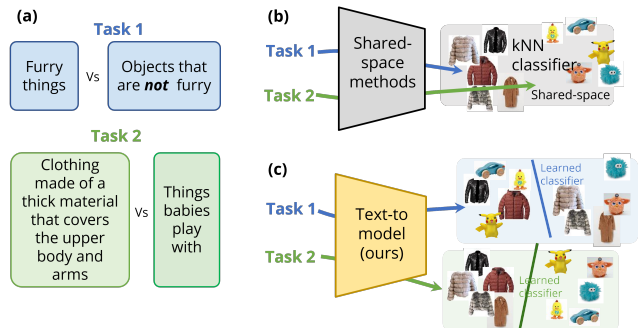


Figure 1: The text-to-model problem. (a) Classification tasks are described in rich language. (b) Traditional zero-shot methods produce static representations, shared for all tasks. (c) Text-to-model uses a hypernetwork to generate task-specific representations and classifiers. This allows T2M to extract task-specific discriminative features.

tunately, as we now explain, existing studies are limited in two major ways: (1) Query-dependence; and (2) Richness of Language description.

First, *Query-dependence*. To illustrate the issue, consider a popular family of zero-shot learning (ZSL) approaches, which maps text (like class labels) and images to a shared space [19, 72, 1, 70, 58, 67, 40, 66, 47]. To classify a new image from an unseen class, one finds the closest class label in the shared space. The problem with this family of shared-space approaches is that the learned representation (and the kNN classifier that it induces) remain “frozen” after training, and are not tuned to the classification task given at inference time. For instance, *furry toys* would be mapped to the same shared representation regardless of whether they are to be distinguished from other *toys*, or from other *furry things* (see Figure 1). The same limita-

tion also hinders another family of ZSL approaches, which synthesize samples from unseen classes at inference time using conditional generative models, and use these samples with kNN classification [14, 24]. Some approaches address the query-dependence limitation by assuming that test descriptions are known during training [21, 54], or by (costly) training a classifier or generator at inference time [65, 54]. Instead, here we learn a model that produces task-dependent classifiers and representations without test-time training.

The second limitation is *language richness*. Natural language can be used to describe classes in complex ways. Most notably, people use negative terms, like "dogs without fur", to distinguish class members from other items. Previous work could only handle limited richness of language descriptions. For instance, it cannot represent adequately textual descriptions with negative terms [19, 72, 1, 70, 58, 67, 40, 66, 14, 24]. In this paper, we wish to handle the inherent linguistic richness of natural language.

Here, we describe a novel deep network architecture and a learning workflow that addresses these two aspects: (1) generating a discriminative model tuned to requested classes at query time and (2) supporting rich language and negative terms.

To achieve these properties, we propose an approach based on hypernetworks (HNs) [20]. An HN is a deep network that emits the weights of another deep network (see Figure 2 for an illustration). Here, the HN receives a set of class descriptions and emits a multi-class model that can classify images according to these classes. Interestingly, this text-based ZSL setup has an important symmetric structure. Specifically, if the order of input descriptions is permuted, one expects that the same classifiers are emitted, but following the same permutation. This property is called *equivariance*, and it can be leveraged to design better architectures [16, 13, 25, 28, 17]. Taking invariance and equivariance into account has been shown to provide significant benefits for learning in spaces with symmetries like sets [69] [35] graphs [23, 63] and deep weight spaces [39]. In general however, HNs are not always permutation equivariant. We design invariant and equivariant layers and describe an HN architecture that respects the symmetries of the problem, and term it T2M-HN: *a text-to-model hypernetwork*.

We put versatility of T2M-HN to the test across an array of zero-shot classification tasks, spanning diverse data types including images, 3D point clouds, and 3D skeletal data for action recognition (see Table 1). Our framework exhibits a remarkable ability to incorporate various forms of class descriptions including long and short texts, as well as class names. Notably, T2M-HN surpasses the performance of previous state-of-the-art methods in all of these setups.

Our paper offers four key contributions: (1) identifying limitations in prior shared space methods for ZSL that rely on fixed representations and distance-based classifiers

for text and image data, and proposing task-dependent representations as an alternative; (2) introducing the Text-to-Model (T2M) approach for generating deep classification models from textual descriptions; (3) investigating the equivariance and invariance properties of T2M models and proposing T2M-HN, an architecture based on HNs that adheres to the setup’s symmetries; and (4) demonstrating the efficacy of T2M-HN on a range of zero-shot tasks, including image and point-cloud classification and action recognition, using diverse text descriptions. T2M-HN outperforms existing state-of-the-art approaches in all tasks.

## 2. Related work

**Zero-shot learning (ZSL).** The core challenge in ZSL lies in recognizing unseen classes based on their semantic associations with seen classes. This association is learned using human-annotated attributes [30, 57, 38, 3]. Another source of information for learning semantic associations is to use textual descriptions. Three main sources were used in the literature to obtain text descriptions of classes: (1) Using class names as descriptions [70, 18, 10, 11]; (2) using encyclopedia articles that describe the class [29, 15, 46, 9, 42, 73]; and (3) providing per-image descriptions manually annotated by domain experts [49, 41, 61]. These can then be aggregated into class-level descriptions.

**Shared space zero-shot learning.** One popular approach to ZSL is to learn a joint visual-semantic representation, using either attributes or natural text descriptions. Some studies project visual features onto the textual space [18, 27, 67], others learn a mapping from a textual to a visual space [70, 40], and some project both images and texts into a new shared space [1, 4, 58, 72, 5, 6, 53, 66, 47]. Once both image and text can be encoded in the same space, classifying an image from a new class can be achieved without further training by first encoding the image and then selecting the nearest class in the shared space. In comparison, instead of nearest-neighbour based classification, our approach is learned in a discriminative way, which may result in richer and potentially stronger models.

**Generation-based zero-shot learning.** Another line of ZSL studies uses generative models like GANs to generate representations of samples from unseen classes. Such generative approaches have been applied in two settings. Some studies assume they have access to test-class descriptions (attributes or text) during model training. Hence, they can train a classifier over test-class images, generated by leveraging the test-class descriptions [32, 54, 21]. Other studies assume access to test-class descriptions only at test time. Hence, they map the test-class descriptions to the shared space of training classes and apply a nearest-neighbor inference mechanism. In this work, we assume that any information about test classes is only available at test time. As a result, ZSL approaches that assume training-time access





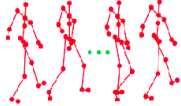
Dataset name and type	Sample data	Description type	Example description
ModelNet-40 [64] 3D Point Clouds CAD models		Class name	(1) <i>Airplane</i> (2) <i>Chair</i>
AwA [26] Animal images		Class name	(1) <i>Moose</i> (2) <i>Elephant</i>
		Long	(1) "An animal of the deer family with humped shoulders, long legs, and a large head with antlers.", (2) "A plant-eating mammal with a long trunk, large ears, and thick, grey skin."
		Negative	(1) "An animal without stripes and not gray", (2) "An animal without fur and without horns"
		Attribute	(1) "Animals with fur" (2) "Animals with long trunk"
SUN [41] Images of scenes and places		Short	(1) "Desert vegetation", (2) "Lecture room"
CUB [61] Images of bird species		Long	(1) "This bird is red with an orange beak and black eyes and eyebrow", (2) "a small yellow bird with a black chest and tail."
BABEL 120 [44] Sequences of 3D skeletal data		Short	(1) "Take off bag", (2) "Type on a keyboard"

Table 1: Overview of evaluation datasets and tasks.

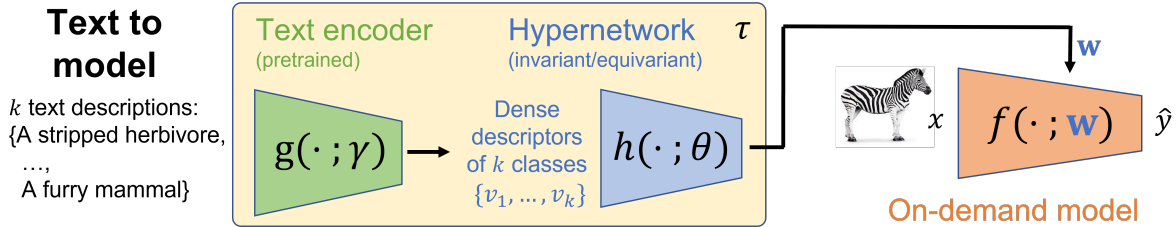


Figure 2: The text-to-model learning problem and our architecture. Our model (yellow box) receives a set of class descriptions as input and outputs weights  $w$  for a downstream on-demand model (orange). The model has two main blocks: A pretrained text encoder and a hypernetwork that obeys certain invariance and equivariance symmetries. The hypernetwork receives a set of dense descriptors to produce weights for the on-demand model.

to information about the test classes are beyond our scope.<sup>1</sup>

<sup>1</sup>While these algorithms could in principle be re-trained when new classes are presented at test-time (e.g. in a continual learning [51] setup), this would result in costly and inefficient inference mechanism, and possibly also in catastrophic forgetting [36] of previous class information. We hence do not include them in our experiments.

However, works that assume only test-time access to test-class information form some of our baselines [14, 24].

**Hypernetworks** [20] (HNs) were applied to many computer vision problems, including few-shot learning [68], federated learning [2], continual learning [60], weight pruning [33]. Here we use HNs for text-based ZSL. The work

by [29] also predicts model weights from textual descriptions, but differs in two key ways. (1) They learn a constant representation of each class; our method uses the context of all the classes in a task to predict data representation. (2) They predict weights of a linear architectures; our T2M-HN applies to deeper ones.

**ZSL with large vision-language models** Large-scale vision-language models, and notably CLIP [47] show remarkable zero-shot capabilities for vision-and-language tasks. A key difference between the CLIP approach and this paper is that CLIP was trained on *massive multimodal data*. In contrast, our approach leverages the semantic compositionality of *language models*, without requiring paired image-text data. As a result, we successfully applied T2M-HN directly to domains where no massive multimodal data exist, like 3D point cloud or skeleton sequence for action recognition. The downside is that the T2M-HN representation may be sensitive to language and semantic distinctions irrelevant to the visual modality.

### 3. Problem formulation

We describe the problem of text-to-model in the context of multiclass classification. It can be naturally extended to regression and ranking problems. Here, our objective is to learn a mapping  $\tau$  from a set of  $k$  natural language descriptions into the space of  $k$ -class classifiers. Here, we address the case where the architecture of the downstream classifier is fixed and given in advance, but this assumption can be relaxed as in [31].

Formally, let  $S^k = \{s_1, \dots, s_k\}$  be a set of  $k$  class descriptions drawn from a distribution  $\mathcal{P}_k$ , where  $s_j$  is a text description of the  $j^{\text{th}}$  class. The distribution  $\mathcal{P}_k$  can be characterized by a two-stage process: First, a set of  $k$  classes is drawn from a large set of classes. Then, a text description is drawn for each class.

Let  $\tau$  be a T2M model parameterized by a set of parameters  $\phi$ . It takes the text descriptors and produces a set of parameters  $W$  of a  $k$ -class classification model  $f(\cdot; W)$ . Therefore, we have  $\tau_\phi : \{s_1, \dots, s_k\} \rightarrow \mathbb{R}^d$ , where  $d$  is the dimension of  $W$ , that is, the number of parameters of the classification model  $f(\cdot; W)$ , and we denote  $W = \tau_\phi(S^k)$ .

Let  $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$  be a loss function, and let  $\{\mathbf{x}_i, y_i\}_{i=1}^n$  be a labeled dataset from a distribution  $\mathcal{P}$  over  $\mathcal{X} \times \mathcal{Y}$ . For  $k$ -class classification,  $\mathcal{Y} = \{1, \dots, k\}$ . We can explicitly write the loss in terms of  $\phi$  as follows.  $l(y_i, \hat{y}_i) = l(y_i, f(x_i; W)) = l(y_i, f(x_i; \tau_\phi(S^k)))$ . See also Figure 2 and note that  $\tau = h \circ g$ . The goal of T2M is to minimize

$$\phi^* = \arg \min_{\phi} \mathbb{E}_{S^k \sim \mathcal{P}^k} \mathbb{E}_{(x, y) \sim \mathcal{P}} [l(y, f(x; \tau_\phi(S^k)))]. \quad (1)$$

The training objective becomes  $\phi^* = \arg \min_{\phi} \sum_j \sum_i l(y_i, f(x_i; \tau_\phi(S^{k_j})))$ , where the sum

over  $j$  means summing over all descriptions from all sets in the training set.

## 4. Our approach

We first describe our approach, based on hypernetworks. We then discuss the symmetries of the problem, and an architecture that can leverage these symmetries.

We propose to address the text-to-model problem, using a hypernetwork architecture. A hypernetwork is a model that outputs the weights of another model [20]. In our case, it receives a set of textual descriptions of classes to be recognized, and outputs the weights of a classifier that can discriminate them. Figure 2 illustrates our architecture. It has two components. First, a text encoder  $g$  takes natural language descriptions and transforms them into dense descriptors; and second, a hypernetwork  $h$  takes these dense descriptors and emits weights for a downstream classifier. In this paper, we do not impose any special properties on the text encoder  $g$ . It can be any model trained using language data (no need for multi-modal data).

### 4.1. Symmetries of the T2M problem

Interestingly, the T2M setup imposes certain invariance and equivariance properties. Design an architecture that takes them into account can improve generalization. We now discuss these properties and then derive an architecture that captures them.

**Equivariance properties of the classifier layer.** As an illustrative example, consider a downstream multi-class classifier  $f_1$ , that is designed to distinguish *cats* from *dogs*, and another classifier  $f_2$ , designed to distinguish *dogs* from *cats*. Intuitively, at the optimum, the two classifiers should be identical except for a switch of two weight vectors at the last layer ( $w_1$  in  $f_1$  equal to  $w_2$  in  $f_2$ ). This has an important implication for the hypernetwork. Any permutation applied to its input class descriptions should be reflected in a parallel ordering of the weight vectors that it produces. Suppl. section B.1 provides a formal definition of this property.

**Invariance properties of intermediate layers.** Considering now the layers of the downstream classifier before the last (classifier) layer. In supplemental section B.1, we prove that using an equivariant transformation for the last layer and an invariant transformation for earlier layers is sufficient to ensure that the downstream classifier is equivariant to permutation over the descriptions.

**Invariant and equivariant Architectures.** Given the equivariance property discussed above, we wish to design a deep architecture that adheres to those symmetries, because that improves generalization. To ensure that certain

elements remain invariant permutation, they should be processed with a shared set of parameters [62, 48]. In our case, we need to share the parameters that process input descriptions, so the model is equivariant to permutations of those inputs.

Figure 3 gives the high-level structure of the equivariant architecture of T2M-HN. Schematics of equivariant layers and invariant layers are detailed in supplemental section B.1. Our experiments below show that using an equivariant architecture consistently improve generalization (Figure 7).

## 5. Experiments

The T2M setup is about producing a model that can be applied to data from new classes. Accordingly, the model trains on data from a set of training classes, alongside their text descriptions. Then, it is tested on data from new classes, given the text descriptions of these classes.

We evaluate T2M-HN in zero-shot classification, using three image datasets, one 3D point cloud dataset, and one action recognition dataset (see blow). We consider various forms of text description, including single-word class labels, few-word class names, and longer descriptions that could also include negative properties (i.e. properties that the images in the class do not have). Finally, we study one-class classification based on text attributes. Table 1 summarizes our tasks, datasets and descriptions.

**Baselines:** We compare our T2M-HN with four text-based zero-shot approaches for image recognition: (1) DE-VICE [18] projects images to a pre-trained language model space by adding a projection head to a pre-trained visual

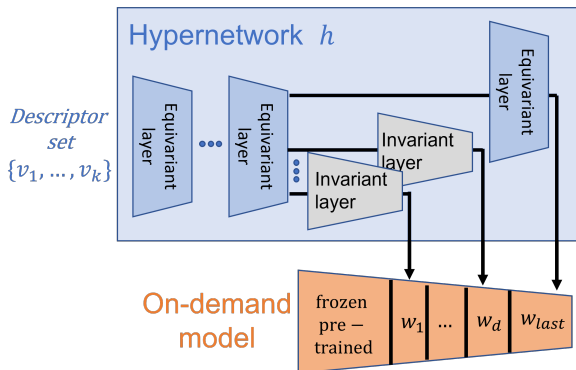


Figure 3: (a) The T2M-HN architecture for equivariant-invariant hypernetwork. The input is processed by equivariant layers, followed by a prediction head for each layer of the target on-demand classifier  $f$ . The prediction head for  $w_{last}$  is equivariant. Heads for earlier layers of  $f$ ,  $w_1, \dots, w_k$  are invariant. Also see schematics of the invariant and equivariant layers in supplemental Figure 8.

classification model; (2) Deep Embedding Model (DEM) [71] uses the visual space as the shared embedding space; (3) CIZSL [14] trains conditional GANs with a loss designed to generate samples from unseen classes without synthesizing unrealistic images. At inference time, the GAN is conditioned on test descriptions, generates synthetic image representations and test images are classified using kNN w.r.t. to the synthetic images.; and (4) GRaWD [24] trains a conditional GAN with a loss that helps to reach regions in space that are hard to classify as seen classes.

When relevant, we also computed the performance obtained when using the CLIP encoder followed by  $k$ -NN classification in the CLIP space [47]. Since CLIP was trained on 400 million (image, caption) pairs, it is reasonable to assume it has seen all classes studied here. It is hence not a zero-shot classifier and the results can be viewed as a “skyline” value that zero-shot approaches should aim at.

**Datasets:** We experiment with three image datasets: (1) **Animals with attributes (AWA)** [26]; (2) **SUN** [41]; and (3) **CUB** [61]; a 3D point-clouds dataset: (4) **ModelNet40** [64]; and an action recognition dataset: (5) **BABEL 120** [44], containing sequences of body skeletons.

**Implementation and architecture:** We encode single-word class names from the AWA dataset using GloVe [43] and longer descriptions, as well as class names, from ModelNet-40 using SBERT [50]. For images, the visual target model had a backbone based on a frozen ResNet-18 [22], pretrained on ImageNet with one or two fully connected layers, predicted by the HN. For 3D point-cloud data, the backbone was PointNet [45], again with one or two predicted fully-connected layers. For action recognition data, we follow [44] and use 2 stream-AGCN [56], with one or two predicted fully-connected layers as well.

**Experimental protocol:** We split the data in two dimensions: Classes and samples. For standardized comparisons the splitting classes into *seen classes* used for training and *unseen classes* used in evaluation, follows the split used by [65] for AWA, the split of [11] for Modelnet40 and the standard split of [61] for CUB. Since there is no official split for SUN and BABEL, we share our random split in the Supplemental material. As in other ZSL protocols, for each seen class we split out a set of evaluation images that are not presented during training, and used to evaluate the model on the seen classes. For AWA, CUB, SUN and BABEL 120 we randomly selected 10% of images for “seen” evaluation. For ModelNet40 we use the test split in [64]. We stress that “Seen” in our tables means *novel images* from *seen classes*.

**Workflow:** When training the whole architecture, we split the train seen classes. 80% of the classes were used for training the backbone. Then, we froze the weights of the backbone and use the remaining 20% to train the HN. This way, the HN learns to generalize to new classes. Finally, we evaluate the entire architecture on the evaluation

split of the seen classes, and on the unseen classes.

At test time, the model receives  $k$  class descriptions and predicts a standalone model to classify images drawn from the corresponding  $k$  classes. Unless otherwise specified, we experiment with the value of  $k = 2$ .

### 5.1. Zero-shot Classification using class names: Images and 3D point clouds

In the following experiment, we evaluate T2M-HN under two tasks: Zero-shot image classification and zero-shot 3D point cloud classification. We use single-word class names for both tasks as the textual class descriptions.

**Results:** Table 2 shows the average classification accuracy of all participating models. Our model reaches the highest accuracy in both experimental setups and datasets, indicating its effectiveness. We shed further light on the performance of our model on ModelNet-40 in Supplementary Section D.

### 5.2. Zero-shot classification using text descriptions: Images and sequences of 3D skeletons

Next, we evaluate T2M-HN when using richer text descriptions: **(1) For SUN**, we use short class descriptions provided by the original dataset. Specifically, SUN includes many multi-word class names like “parking garage indoor” or “control tower outdoor”. **(2) For BABEL 120** we use the action names provided by the original dataset. Many of the actions have multi-word, descriptive names such as “take of bag”. **(3) For AwA**, we use synthetic class descriptions generated by a GPT model. Specifically, we used GPT3 [8] to generate five different descriptions for each class of AwA. During training and evaluation, we randomly choose one description for each class in the batch, from its corresponding 5 class descriptions. See detailed examples in the Supplementary Section E. We will publish the full set of descriptions for reproducibility. **(4) For CUB**, we use the descriptions of each image in a given class as a possible description of the class.

In the CUB dataset, bird species from the same taxonomic family are harder to distinguish from each other than random pairs of species [59]. To investigate this further, we used the Datazone dataset of bird species [7] and annotated each species with its corresponding taxonomic family. Based on this information, we defined pairs of bird species from two different families as *easy* and pairs from the same family as *hard*.

**Results:** Table 3 presents the classification accuracy obtained using class descriptions, for the AWA, SUN, and BABEL datasets. T2M-HN outperforms all baselines.

Figure 4 shows the results for the CUB dataset with easy and hard tasks. To better understand the results, consider an important distinction between our approach and previous *shared-representation* approaches. These approaches

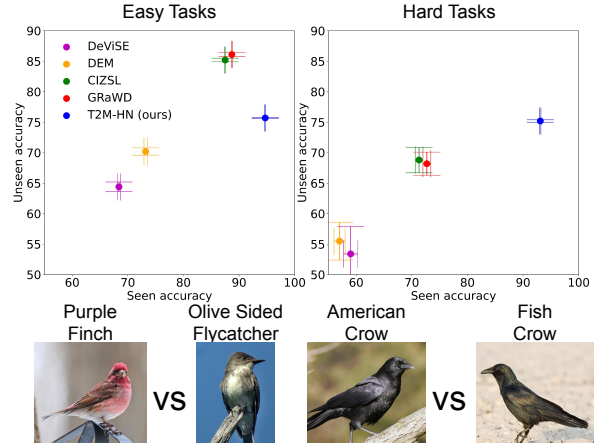


Figure 4: **Classifying easy and hard pairs of bird species from the CUB dataset.** Easy tasks are binary classification tasks, where pairs of birds are from different taxonomy families. In hard classification tasks, bird classes in each pair are from the same taxonomy family. Values are the mean classification accuracy on images from seen (x-axis) and unseen (y-axis) classes, averaged over all class pairs.

aim to learn class representations that would generalize to new classification tasks. In contrast, our approach aims to build task-specific representations and classifiers. For easy tasks, task-dependent representation may not be important because the input contains a sufficient signal for accurate classification. In contrast, in hard tasks, a model would benefit from task-dependent representation to focus on the few existing discriminative features of the input examples. Indeed, as demonstrated in Figure 4, in the easy tasks, although our model is superior on the seen classes, it is outperformed by the GAN-based baselines on unseen classes. In contrast, for the hard tasks, where task-specific class representation is more valuable, our model is superior on both seen and unseen classes.

### 5.3. Descriptions with negative terms

To this point, we have assumed that the descriptions correspond to properties of the class. However, descriptions could also state which properties the class does **not** have. For example, one may want to classify animals that “do not live in the water”, or animals that “do not fly”. To create such negative descriptions for the AwA data, we used the list of attributes provided for each class in AwA. For each class, we randomly sampled 4 attributes that do not apply to that class. For example, an elephant may be described as an “Animal that does/is not: fly, small, furry or white”.

**Results:** Table 4 presents our results. We tested two scenarios: Considering only negative descriptions (left side of the table) and including equal portions of positive and nega-

	AWA by class name			ModelNet40 by class name		
	Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
DeViSE [18]	78.1 $\pm$ 1.0	58.9 $\pm$ 1.4	67.2 $\pm$ 1.9	83.6 $\pm$ 2.7	58.6 $\pm$ 3.4	68.9 $\pm$ 3
DEM [71]	83.1 $\pm$ 1.6	75.1 $\pm$ 1.2	78.9 $\pm$ 2.0	86.7 $\pm$ 2.4	57.3 $\pm$ 3.3	69.0 $\pm$ 2.8
CIZSL [14]	97.0 $\pm$ 0.1	74.7 $\pm$ 3.2	84.20 $\pm$ 2.0	97.6 $\pm$ 0.6	50.1 $\pm$ 3.6	66.3 $\pm$ 3.3
GRaWD [24]	96.9 $\pm$ 0.1	81.6 $\pm$ 1.9	88.6 $\pm$ 1.1	97.8 $\pm$ 0.5	52.8 $\pm$ 3.3	68.3 $\pm$ 2.8
T2M-HN (ours)	<b>98.9 <math>\pm</math> 0.1</b>	<b>87.3 <math>\pm</math> 0.2</b>	<b>92.7 <math>\pm</math> 0.1</b>	<b>98.4 <math>\pm</math> 0.1</b>	<b>59.2 <math>\pm</math> 0.3</b>	<b>73.9 <math>\pm</math> 0.1</b>
CLIP [47]	98.9 $\pm$ 0.2	NA	NA	NA	NA	NA

Table 2: **Classification by single-word class names.** Mean classification accuracy on seen and unseen classes for AWA and ModelNet-40. Values are averages and SEM over all class pairs.

	SUN by short description			BABEL by short descriptions			AWA by GPT descriptions		
	Seen	Unseen	Harmonic	Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
CLIP	99.1 $\pm$ 0.4	NA	NA	NA	NA	NA	93.7 $\pm$ 0.2	NA	NA
DeViSE	52.0 $\pm$ 1.4	58.9 $\pm$ 1.1	55.2 $\pm$ 0.9	65.9 $\pm$ 4.4	51.1 $\pm$ 2.0	57.6 $\pm$ 2.8	91.8 $\pm$ 1.6	70.0 $\pm$ 3.7	79.4 $\pm$ 2.2
DEM	83.2 $\pm$ 1.1	83.2 $\pm$ 1.4	83.2 $\pm$ 0.9	56.6 $\pm$ 2.4	50.2 $\pm$ 1.1	53.2 $\pm$ 1.5	93.9 $\pm$ 1.2	73.0 $\pm$ 3.3	82.1 $\pm$ 1.8
CIZSL	94.0 $\pm$ 0.1	80.3 $\pm$ 0.6	86.6 $\pm$ 0.3	82.7 $\pm$ 2.1	62.5 $\pm$ 1.3	71.2 $\pm$ 1.2	96.6 $\pm$ 0.1	80.7 $\pm$ 2.2	87.9 $\pm$ 1.3
GRaWD	95.5 $\pm$ 0.1	84.7 $\pm$ 0.5	89.8 $\pm$ 0.3	83.7 $\pm$ 1.8	62.2 $\pm$ 1.1	71.3 $\pm$ 1.0	96.8 $\pm$ 0.1	81.1 $\pm$ 0.2	88.3 $\pm$ 1.2
T2M-HN (ours)	<b>95.8 <math>\pm</math> 0.1</b>	<b>88.4 <math>\pm</math> 0.1</b>	<b>92.0 <math>\pm</math> 0.1</b>	<b>95.3 <math>\pm</math> 0.1</b>	<b>77.6 <math>\pm</math> 0.1</b>	<b>85.5 <math>\pm</math> 0.1</b>	<b>98.7 <math>\pm</math> 0.1</b>	<b>83.3 <math>\pm</math> 0.1</b>	<b>90.3 <math>\pm</math> 0.1</b>

Table 3: **Classification using short and rich class descriptions.** Values are the mean ( $\pm$  s.e.m) classification accuracy averaged over 100 random class pairs (for SUN and BABEL 120) and all class pairs (for Awa).

AWA data	Negative descriptions			Negative and positive descriptions		
	Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
CLIP	19.9 $\pm$ 2.2	NA	NA	56.8 $\pm$ 2.9	NA	NA
DeViSE [18]	57.3 $\pm$ 4.9	54.5 $\pm$ 5.2	55.9 $\pm$ 5.0	79.5 $\pm$ 3.6	61.5 $\pm$ 4.5	69.4 $\pm$ 4.0
DEM [71]	81.7 $\pm$ 1.2	73.7 $\pm$ 1.6	77.5 $\pm$ 1.0	78.2 $\pm$ 1.7	69.1 $\pm$ 1.6	73.4 $\pm$ 1.2
CIZSL [14]	58.3 $\pm$ 0.8	56.6 $\pm$ 3.4	57.5 $\pm$ 1.8	93.9 $\pm$ 0.2	71.6 $\pm$ 2.3	81.2 $\pm$ 1.5
GRaWD [24]	54.9 $\pm$ 0.8	56.0 $\pm$ 3.2	55.3 $\pm$ 1.6	95.0 $\pm$ 0.2	73.9 $\pm$ 2.0	83.2 $\pm$ 1.5
T2M-HN (ours)	<b>90.0<math>\pm</math>0.2</b>	<b>77.1<math>\pm</math>0.3</b>	<b>83.0<math>\pm</math>0.2</b>	<b>96.6<math>\pm</math>0.2</b>	<b>82.9<math>\pm</math>0.2</b>	<b>89.2<math>\pm</math>0.1</b>

Table 4: **Classification with negative descriptions.** Mean accuracy on images from seen and unseen Awa classes. Values are averages over all class pairs. CLIP has seen all classes and therefore is marked as NA for unseen classes.

tive descriptions (right side of the table). For both scenarios, we keep the same ratio of positive and negative descriptions for training and testing.

T2M-HN outperforms all baselines by significant gaps. Presumably, the best baseline, GRaWD, which generates image features from the textual descriptions, fail to generate proper images given negative attributes. CLIP performance degrades dramatically in these setups, probably because the CLIP training set consisted of image captions, and these rarely contain negative descriptions.

#### 5.4. Identifying complex classes membership

Typically, zero-shot classification involves distinguishing “natural categories” [52] like “cats” and “dogs”. However, We may want to generate classifiers that follow more complex class boundaries, aggregating over multiple natural classes. For instance, “animals with horns” combine several classes from a rhino to a deer.

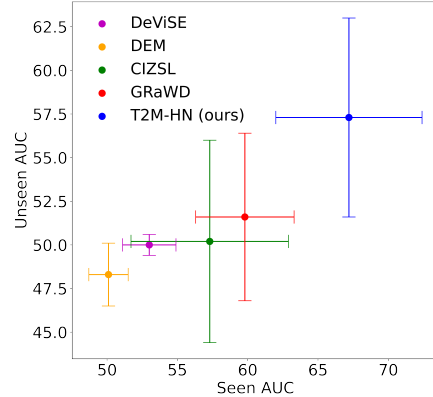


Figure 5: AUC of seen and unseen classes, in a one class task that crosses species boundaries: “Animals that have horns”. Shown are averages over 53 attributes.

To test 2M-HN in this scenario, we created a set of one-class classification tasks designed to recognize images based on properties that cut through class boundaries. To make the evaluation systematic, we used attributes from Awa, and eliminate non-visual attributes (like habitats and diet). The remaining 53 attributes, were split to 30/10/10. Details of the protocol are given in the supplemental material. We report the average Area Under the Recall-Precision Curve over seen classes and unseen classes.

**Results:** Figure 5 shows that T2M-HN captures the complex semantic distinctions of our task better than the

baselines. We attribute this to its ability to draw new classifiers for each new textual description.

### 5.5. T2M-HN classifiers depend on task context

Current leading text-based ZSL methods map class descriptions or images to a shared representation, but that mapping is constant for all classification tasks. Our T2M-HN is designed to use information about the classes of each specific classification task.

To demonstrate this effect, we use GradCam [55] and examine what image areas are used in different classification tasks. Figure 6 explores two such examples. The upper three panels show the image regions that are used for classifying the image as a *Dolphin*. When classifying dolphin vs. deer, the model gives most of its weight to the background (ocean water and waves), which is reasonable since an image of a deer probably will not contain those elements in the background. However, when classifying dolphin vs. killer whale, the model gives most of its weight to the dolphin itself, since the background of a dolphin image may be similar to the background of a whale image.

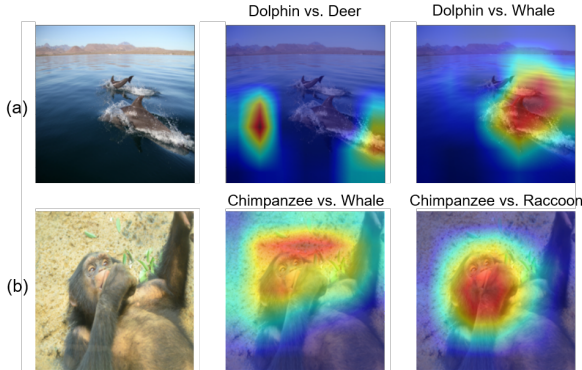


Figure 6: Class context affects the predicted classifier. **Top left:** An image of a dolphin. **Top middle:** gradcam heat map when classifying the dolphin image using a model trained for *dolphin vs deer*: The model is strongly affected by the background ocean water, presumably because the negative class lives on land. **Top right:** Recognition using a model for *dolphin vs. killer whale*: the model attends to the dolphin, since background would be similar for both classes. **Bottom:** A similar effect for a chimpanzee.

### 5.6. The Impact of Equivariance Design on HNs

To evaluate the effect of the equivariance property on our HN-based model performance, we compared variants with and without the equivariance design. We repeat the experiment for an on-demand model with one or two fully connected layers. Figure 7 shows the mean accuracy of the following variants: (1) **T2M-HN 1-layer** An equivariant HN

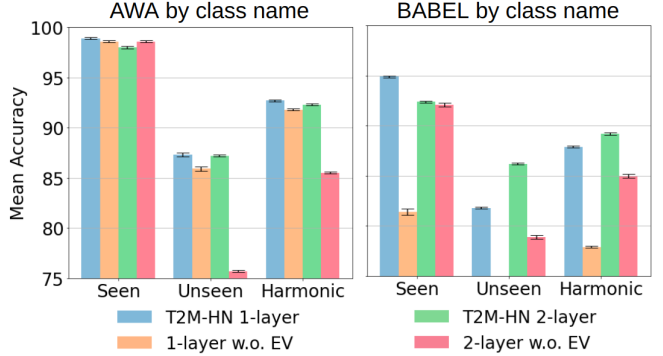


Figure 7: Ablation study. Mean classification accuracy (averaged across class pairs) on seen and unseen classes and their harmonic mean for the AWA and BABEL datasets.

that predicts one equivariant FC layer; (2) **1-layer w.o. EV** A FC HN that predicts one fully connected layer; (3) **T2M-HN 2-layers** An equivariant HN that predicts two FC layers for the on-demand model: The first is invariant and the second is equivariant; and (4) **2-layer w.o. EV** A FC HN that predicts two FC layers.

In all cases, the equivariant HN performs better than the simple fully connected. For Awa, T2M-HN 1-layer performs better than T2M-HN 2-layers. We believe this is because ResNet backbone separates the images so they are linearly separable. For BABEL, we used 2s-AGCN as a features extractor and in that case, T2M-HN 2-layer generalizes better to unseen classes. Since we use the accuracy over the seen classes to choose the model architecture, throughout the paper, we report the score achieved by T2M-HN 1-layer.

## 6. Conclusion

We presented T2M, a learning algorithm that generates a discriminative model "on demand", given test-time class descriptions only, such that class representations are task-dependant rather than fixed. We analyzed the group symmetries that a T2M model should obey, and characterized the proper invariance and equivariance properties that ensure these symmetries. We then proposed T2M-HN, a deep architecture T2M model based on HNs, which obeys the required symmetries. Next, we evaluated our T2M-HN approach in a series of recognition tasks, considering images, 3D point clouds, and action recognition setups. We experiment with descriptions at varying complexity: From single-word class names, through few-word class names and long text descriptions, all the way to "negative" and attribute descriptions. Our results clearly demonstrate the potential of the T2M modeling approach.

## References

- [1] Zeynep Akata, Scott Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. Evaluation of output embeddings for fine-grained image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2927–2936, 2015. 1, 2
- [2] Ohad Amosy, Gal Eyal, and Gal Chechik. Inference-time personalized federated learning. *arXiv preprint arXiv:2111.08356*, 2021. 3
- [3] Yashas Annadani and Soma Biswas. Preserving semantic relations for zero-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7603–7612, 2018. 2
- [4] Yuval Atzmon and Gal Chechik. Probabilistic and-or attribute grouping for zero-shot learning. *arXiv preprint arXiv:1806.02664*, 2018. 2
- [5] Yuval Atzmon and Gal Chechik. Adaptive confidence smoothing for generalized zero-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11671–11680, 2019. 2
- [6] Yuval Atzmon, Felix Kreuk, Uri Shalit, and Gal Chechik. A causal view of compositional zero-shot recognition. *Advances in Neural Information Processing Systems*, 33:1462–1473, 2020. 2
- [7] BirdLife, 2022. data retrieved from <http://datazone.birdlife.org/species/taxonomy>. 6
- [8] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *NeurIPS 2020*, 2020. 6, 13
- [9] Sebastian Bujwid and Josephine Sullivan. Large-scale zero-shot image classification from rich and diverse textual descriptions. *arXiv preprint arXiv:2103.09669*, 2021. 2
- [10] Soravit Changpinyo, Wei-Lun Chao, and Fei Sha. Predicting visual exemplars of unseen classes for zero-shot learning. In *Proceedings of the IEEE international conference on computer vision*, pages 3476–3485, 2017. 2
- [11] Ali Cheraghian, Shafin Rahman, Townim F Chowdhury, Dylan Campbell, and Lars Petersson. Zero-shot learning on 3d point cloud objects and beyond. *International Journal of Computer Vision*, pages 1–21, 2022. 2, 5, 15
- [12] Ali Cheraghian, Shafin Rahman, and Lars Petersson. Zero-shot learning of 3d point cloud objects. In *16th International Conference on Machine Vision Applications, MVA 2019, Tokyo, Japan, May 27-31, 2019*, pages 1–6. IEEE, 2019. 15
- [13] Taco S Cohen, Mario Geiger, and Maurice Weiler. A general theory of equivariant cnns on homogeneous spaces. *Advances in neural information processing systems*, 32, 2019. 2
- [14] Mohamed Elhoseiny and Mohamed Elfeki. Creativity inspired zero-shot learning. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 5783–5792. IEEE, 2019. 2, 3, 5, 7, 13
- [15] Mohamed Elhoseiny, Yizhe Zhu, Han Zhang, and Ahmed Elgammal. Link the head to the” beak”: Zero shot learning from noisy text description at part precision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5640–5649, 2017. 2
- [16] Marc Finzi, Samuel Stanton, Pavel Izmailov, and Andrew Gordon Wilson. Generalizing convolutional neural networks for equivariance to lie groups on arbitrary continuous data. In *International Conference on Machine Learning*, pages 3165–3176. PMLR, 2020. 2
- [17] Marc Finzi, Max Welling, and Andrew Gordon Wilson. A practical method for constructing equivariant multilayer perceptrons for arbitrary matrix groups. In *International Conference on Machine Learning*, pages 3318–3328. PMLR, 2021. 2
- [18] Andrea Frome, Greg S Corrado, Jon Shlens, Samy Bengio, Jeff Dean, Marc’Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. *Advances in neural information processing systems*, 26, 2013. 2, 5, 7, 13
- [19] Amir Globerson, Gal Chechik, Fernando Pereira, and Naftali Tishby. Euclidean embedding of co-occurrence data. *Advances in neural information processing systems*, 17, 2004. 1, 2
- [20] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016. 2, 3, 4
- [21] Zongyan Han, Zhenyong Fu, Shuo Chen, and Jian Yang. Contrastive embedding for generalized zero-shot learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 2371–2381. Computer Vision Foundation / IEEE, 2021. 2
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [23] Roei Herzig, Moshiko Raboh, Gal Chechik, Jonathan Berant, and Amir Globerson. Mapping images to scene graphs with permutation-invariant structured prediction. *Advances in Neural Information Processing Systems*, 31, 2018. 2
- [24] Divyansh Jha, Kai Yi, Ivan Skorokhodov, and Mohamed Elhoseiny. Imaginative walks: Generative random walk deviation loss for improved unseen learning representation. *arXiv preprint arXiv:2104.09757*, abs/2104.09757, 2021. 2, 3, 5, 7, 13
- [25] Risi Kondor and Shubhendu Trivedi. On the generalization of equivariance and convolution in neural networks to the action of compact groups. In *International Conference on Machine Learning*, pages 2747–2755. PMLR, 2018. 2
- [26] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-

- class attribute transfer. In *2009 IEEE conference on computer vision and pattern recognition*, pages 951–958. IEEE, 2009. [3](#), [5](#)
- [27] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute-based classification for zero-shot visual object categorization. *IEEE transactions on pattern analysis and machine intelligence*, 36(3):453–465, 2013. [2](#)
- [28] Leon Lang and Maurice Weiler. A wigner-eckart theorem for group equivariant convolution kernels. *arXiv preprint arXiv:2010.10952*, 2020. [2](#)
- [29] Jimmy Lei Ba, Kevin Swersky, Sanja Fidler, et al. Predicting deep zero-shot convolutional neural networks using textual descriptions. In *Proceedings of the IEEE international conference on computer vision*, pages 4247–4255, 2015. [2](#), [4](#)
- [30] Kai Li, Martin Renqiang Min, and Yun Fu. Rethinking zero-shot learning: A conditional visual classification perspective. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3583–3592, 2019. [2](#)
- [31] Or Litany, Haggai Maron, David Acuna, Jan Kautz, Gal Chechik, and Sanja Fidler. Federated learning with heterogeneous architectures using graph hypernetworks. *arXiv preprint arXiv:2201.08459*, 2022. [4](#)
- [32] Shichen Liu, Mingsheng Long, Jianmin Wang, and Michael I Jordan. Generalized zero-shot learning with deep calibration network. *Advances in neural information processing systems*, 31, 2018. [2](#)
- [33] Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Kwang-Ting Cheng, and Jian Sun. Metapruning: Meta learning for automatic neural network channel pruning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3296–3305, 2019. [3](#)
- [34] Ellen M Markman. Constraints children place on word meanings. *Cognitive science*, 14(1):57–77, 1990. [1](#)
- [35] Haggai Maron, Or Litany, Gal Chechik, and Ethan Fetaya. On learning sets of symmetric elements. In *International Conference on Machine Learning*, pages 6734–6744. PMLR, 2020. [2](#)
- [36] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989. [3](#)
- [37] Björn Michele, Alexandre Boulch, Gilles Puy, Maxime Bucher, and Renaud Marlet. Generative zero-shot learning for semantic segmentation of 3d point clouds. In *International Conference on 3D Vision, 3DV 2021, London, United Kingdom, December 1-3, 2021*, pages 992–1002. IEEE, 2021. [15](#)
- [38] Pedro Morgado and Nuno Vasconcelos. Semantically consistent regularization for zero-shot recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6060–6069, 2017. [2](#)
- [39] Aviv Navon, Aviv Shamsian, Idan Achituve, Ethan Fetaya, Gal Chechik, and Haggai Maron. Equivariant architectures for learning in deep weight spaces. *CoRR*, abs/2301.12780, 2023. [2](#)
- [40] Frederik Pahde, Mihai Puscas, Tassilo Klein, and Moin Nabi. Multimodal prototypical networks for few-shot learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2644–2653, January 2021. [1](#), [2](#)
- [41] Genevieve Patterson and James Hays. Sun attribute database: Discovering, annotating, and recognizing scene attributes. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2751–2758. IEEE, 2012. [2](#), [3](#), [5](#)
- [42] Tzuf Paz-Argaman, Yuval Atzmon, Gal Chechik, and Reut Tsarfaty. Zest: Zero-shot learning from text descriptions using textual similarity and visual summarization. *arXiv preprint arXiv:2010.03276*, 2020. [2](#)
- [43] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014. [5](#)
- [44] Abhinanda R. Punnakal, Arjun Chandrasekaran, Nikos Athanasiou, Alejandra Quiros-Ramirez, and Michael J. Black. BABEL: bodies, action and behavior with english labels. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 722–731. Computer Vision Foundation / IEEE, 2021. [3](#), [5](#)
- [45] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. [5](#)
- [46] Pengda Qin, Xin Wang, Wenhui Chen, Chunyun Zhang, Weiran Xu, and William Yang Wang. Generative adversarial zero-shot relational learning for knowledge graphs. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8673–8680. AAAI Press, 2020. [2](#)
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. [1](#), [2](#), [4](#), [5](#), [7](#)
- [48] Siamak Ravanbakhsh, Jeff Schneider, and Barnabas Poczos. Equivariance through parameter-sharing. In *International conference on machine learning*, pages 2892–2901. PMLR, 2017. [5](#)
- [49] Scott Reed, Zeynep Akata, Honglak Lee, and Bernt Schiele. Learning deep representations of fine-grained visual descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 49–58, 2016. [2](#)
- [50] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019. [5](#)
- [51] Mark B. Ring. *Continual learning in reinforcement environments*. PhD thesis, University of Texas at Austin, TX, USA, 1995. [3](#)
- [52] Eleanor H Rosch. Natural categories. *Cognitive psychology*, 4(3):328–350, 1973. [7](#)

- [53] Dvir Samuel, Yuval Atzmon, and Gal Chechik. From generalized zero-shot learning to long-tail with class descriptors. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 286–295, 2021. [2](#)
- [54] Edgar Schonfeld, Sayna Ebrahimi, Samarth Sinha, Trevor Darrell, and Zeynep Akata. Generalized zero-and few-shot learning via aligned variational autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8247–8255, 2019. [2](#)
- [55] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. [8](#)
- [56] Lei Shi, Yifan Zhang, Jian Cheng, and Hanqing Lu. Two-stream adaptive graph convolutional networks for skeleton-based action recognition. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 12026–12035. Computer Vision Foundation / IEEE, 2019. [5](#)
- [57] Jie Song, Chengchao Shen, Yezhou Yang, Yang Liu, and Mingli Song. Transductive unbiased embedding for zero-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1024–1033, 2018. [2](#)
- [58] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1199–1208, 2018. [1](#), [2](#)
- [59] Ramakrishna Vedantam, Samy Bengio, Kevin Murphy, Devi Parikh, and Gal Chechik. Context-aware captions from context-agnostic supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 251–260, 2017. [6](#)
- [60] Johannes Von Oswald, Christian Henning, João Sacramento, and Benjamin F Grewe. Continual learning with hypernetworks. *arXiv preprint arXiv:1906.00695*, 2019. [3](#)
- [61] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. [2](#), [3](#), [5](#)
- [62] Jeffrey Wood and John Shawe-Taylor. Representation theory and invariant neural networks. *Discrete applied mathematics*, 69(1-2):33–60, 1996. [5](#)
- [63] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32(1):4–24, 2020. [2](#)
- [64] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015. [3](#), [5](#)
- [65] Yongqin Xian, Tobias Lorenz, Bernt Schiele, and Zeynep Akata. Feature generating networks for zero-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5542–5551, 2018. [2](#), [5](#)
- [66] Cheng Xie, Hongxin Xiang, Ting Zeng, Yun Yang, Beibei Yu, and Qing Liu. Cross knowledge-based generative zero-shot learning approach with taxonomy regularization. *Neural Networks*, 139:168–178, 2021. [1](#), [2](#)
- [67] Guo-Sen Xie, Xu-Yao Zhang, Yazhou Yao, Zheng Zhang, Fang Zhao, and Ling Shao. Vman: A virtual mainstay alignment network for transductive zero-shot learning. *IEEE Transactions on Image Processing*, 30:4316–4329, 2021. [1](#), [2](#)
- [68] Li Yin, Juan M Perez-Rua, and Kevin J Liang. Sylph: A hypernetwork framework for incremental few-shot object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9035–9045, 2022. [3](#)
- [69] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30, 2017. [2](#)
- [70] Li Zhang, Tao Xiang, and Shaogang Gong. Learning a deep embedding model for zero-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2021–2030, 2017. [1](#), [2](#)
- [71] Li Zhang, Tao Xiang, and Shaogang Gong. Learning a deep embedding model for zero-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2021–2030, 2017. [5](#), [7](#), [13](#)
- [72] Ziming Zhang and Venkatesh Saligrama. Zero-shot learning via semantic similarity embedding. In *Proceedings of the IEEE international conference on computer vision*, pages 4166–4174, 2015. [1](#), [2](#)
- [73] Yizhe Zhu, Mohamed Elhoseiny, Bingchen Liu, Xi Peng, and Ahmed Elgammal. A generative adversarial approach for zero-shot learning from noisy texts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1004–1013, 2018. [2](#)

# Supplemental information

## A. Hyperparameter optimization

We tune hyperparameters using a held-out set described below.

For the HN optimizer, we tuned the learning rate  $\in \{0.001, 0.005, 0.01, 0.05, 0.1\}$ , momentum  $\in \{0.1, 0.3, 0.9\}$ , weight decay  $\in \{0.00001, 0.0001, 0.001, 0.1\}$ , and number of HN training epochs  $\in \{50, 70, 100\}$ .

For the on-demand target model, we fixed the optimizer to have a learning rate of 0.01, momentum of 0.9 and weight decay of 0.01. We tuned the batch size  $\in \{16, 32, 64, 128\}$  and the number of training epochs  $\{1, 2, 3, 5, 10\}$ .

We tried several sizes for the HN architecture with one hidden layer,  $\{30, 50, 120, 300\}$ . We also describe results with two layers in the ablation section 5.6.

Recall that we split the data across two dimensions: classes and samples. When training the backbone model, we held out 20% of training (seen) classes for training the HN on classes the backbone does not see. From those classes, we held out images to serve as a validation set. We used those images of seen classes to evaluate the architecture performance and chose the hyperparameters based on that estimation.

## B. Equivariant and invariant layers

As an illustrative example, consider a downstream multi-class classifier  $f_1$ , that is designed to distinguish *cats* from *dogs*, and another classifier  $f_2$ , designed to distinguish *dogs* from *cats*. Intuitively, at the optimum, the two classifiers should be identical except for a switch of two weight vectors at the last layer ( $w_1$  in  $f_1$  equal to  $w_2$  in  $f_2$ ). This has an important implication for the hypernetwork. Any permutation applied to its input class descriptions should be reflected in a parallel ordering of the weight vectors that it produces. We now show how to design a hypernetwork that obeys this property.

### B.1. Equivariance properties of the classifier layer.

Consider a downstream multiclass deep classifier whose last (classification) layer has a weight vector  $w_i \in \mathbb{R}^m$  for the output class  $i$ . The weight matrix of the last layer is  $W_{last} = \{w_1, \dots, w_k\}$  (See Figure 8a).

Let  $S^k = \{s_1, \dots, s_k\}$  be a set of  $k$  class descriptions drawn from a distribution  $\mathcal{P}_k$ , where  $s_j$  is a text description of the  $j^{th}$  class. The distribution  $\mathcal{P}_k$  can be characterized by a two-stage process: First, a set of  $k$  classes is drawn from a large set of classes. Then, a text description is drawn for each class.

Let  $\tau$  be a T2M model parameterized by a set of parameters  $\phi$ . It takes the text descriptors and produces a set

of parameters  $W$  of a  $k$ -class classification model  $f(\cdot; W)$ . Therefore, we have  $\tau_\phi : \{s_1, \dots, s_k\} \rightarrow \mathbb{R}^d$ , where  $d$  is the dimension of  $W$ , that is, the number of parameters of the classification model  $f(\cdot; W)$ , and we denote  $W = \tau_\phi(S^k)$ .

The HN receives  $k$  class descriptors and outputs their corresponding weights

$$W_{last} = \{w_1, \dots, w_k\} = R_{last}(\tau_\phi(\{s_1, \dots, s_k\})), \quad (2)$$

where  $R_{last}$  is a function that takes the output of  $\tau$  and resizes the last  $k \cdot m$  elements to the matrix  $W_{last}$ . If the input descriptions are permuted by a permutation  $\mathcal{P}$  the columns of the last layer weight should be permuted accordingly:

$$\mathcal{P}(f(x; \tau_\phi(S^k))) = f(x; \tau_\phi(\mathcal{P}(S^k))). \quad (3)$$

This is the equivariant property, and the HN must obey it.

### B.2. Invariance properties of intermediate layers

Considering now the layer of the downstream classifier before the last layer ( $w_d$  in Figure 8a). A similar argument holds for earlier (lower) intermediate layers. We now show that using an equivariant transformation for the last layer and an invariant transformation for the penultimate layer is sufficient to ensure that the downstream classifier is equivariant to permutation over the descriptions.

**Theorem B.1.** *Let  $f$  be a two-layer neural network  $f(x) = W^{last} \sigma(W^{pen} x)$ , whose weights are predicted from descriptors  $S^k = \{s_1, \dots, s_k\}$  such that  $[W^{last}, W^{pen}] = \tau(S^k)$ . If  $\tau(S^k)$  is equivariant to a permutation  $\mathcal{P}$  with respect to  $W^{last}$ , and invariant to  $\mathcal{P}$  with respect to  $W^{pen}$ , then  $f(x)$  is equivariant to  $\mathcal{P}$  with respect to the input of  $\tau(S^k)$ .*

*Proof.* From the equivariance of  $f(x)$  to a permutation  $\mathcal{P}$  over the input  $S^k$ , we have  $\mathcal{P}(f(x_i; \tau_\phi(S^k))) = f(x_i; \tau_\phi(\mathcal{P}(S^k)))$ . Denote by  $m$  the number of rows of  $W^{last}$  and  $z^{pen} = \sigma(W^{pen} x)$ . We have

$$\begin{aligned} \mathcal{P}(f(x; \tau_\phi(S^k))) &= \mathcal{P}(W^{last} \sigma(W^{pen} x)) = \mathcal{P}(W^{last} z^{pen}) \\ &= \mathcal{P}\left(\begin{bmatrix} W_1^{last} z^{pen} \\ \vdots \\ W_m^{last} z^{pen} \end{bmatrix}\right) \\ &= \begin{bmatrix} W_{\mathcal{P}(1)}^{last} z^{pen} \\ \vdots \\ W_{\mathcal{P}(m)}^{last} z^{pen} \end{bmatrix} = \mathcal{P}(W^{last}) z^{pen}. \end{aligned} \quad (4)$$

If  $\tau(S^k)$  is equivariant to  $\mathcal{P}$  with respect to  $W^{last}$ , and invariant to  $\mathcal{P}$  with respect to  $W^{pen}$ , then  $\tau(\mathcal{P}(S^k)) = [\mathcal{P}(W^{last}), W^{pen}]$ , so

$$\mathcal{P}(f(x; \tau_\phi(S^k))) = \mathcal{P}(W^{last}) z^{pen} = f(x; \tau_\phi(\mathcal{P}(S^k))). \quad (5)$$

	AwA triplets by class name		
	Seen	Unseen	Harmonic
DeViSE [18]	95.1 $\pm$ 0.7	55.6 $\pm$ 3.6	70.2 $\pm$ 1.2
DEM [71]	94.6 $\pm$ 0.7	64.3 $\pm$ 3.0	76.6 $\pm$ 1.1
CIZSL [14]	97.0 $\pm$ 0.4	62.0 $\pm$ 2.9	75.6 $\pm$ 2.1
GRaWD [24]	96.4 $\pm$ 0.5	68.5 $\pm$ 3.0	80.0 $\pm$ 2.0
T2M-HN (ours)	<b>98.1 <math>\pm</math> 0.1</b>	<b>75.3 <math>\pm</math> 0.1</b>	<b>85.2 <math>\pm</math> 0.1</b>

Table 5: Classification by class descriptions. Mean classification accuracy and SEM on images from seen and unseen classes. Averages are over 100 random class triplets

□

### B.3. Invariant and equivariant Architectures

Figure 8(b) shows the architecture of our equivariant layers. All inputs are fed into the same fully connected layer (vertical stripes). To take into account the context of each input, we sum all the inputs to obtain a context vector. We fed the context vector to a different fully connected layer (diagonal stripes) and add it to each one of the processed inputs. The invariant layer has a similar architecture (Figure 8(c)), but with additional summation over all equivariant outputs and another different fully connected layer (horizontal stripes).

Our HN uses several equivariant layers to process the input descriptions. We then use one prediction head for each layer of the output model. The last layer should be equivariant, so we use an equivariant prediction head. For the hidden layers, we use invariant layers (See Figure 8(a)).

### C. Multi-class classification

To demonstrate the flexibility of our approach to deal with multiple classes, we evaluated T2M-HN in 3-way classification tasks. In each task, the on-demand model classifies the image into one out of three classes. For example, such a task could be to classify whether an image is a dog, a cat, or an elephant. We use the same workflow as described in Section 5, with  $k = 3$ . Results are in Table 5. T2M-HN outperforms all baselines by a large margin.

### D. 3D point cloud multiclass classification

While T2M-HN is designed to excel in binary classification, it can be easily applied to multiclass problems. For comparison with previous models we evaluate its performance in multi-class settings, where T2M-HN predicts a model that classifies all seen and unseen classes, instead of two specific classes. Table 6 shows the results of this experiment. We report the result when classifying new samples from the seen classes (30-classes classification) and

	ModelNet40 by class name		
	Seen	Unseen	Harmonic
DeViSE [18]	47.2	14.5	22.2
DEM [71]	46.8	7.0	12.3
CIZSL [14]	75.6	6.0	11.0
GRaWD [24]	75.2	10.9	19.0
T2M-HN (ours)	<b>76.3</b>	<b>18.9</b>	<b>30.3</b>

Table 6: **3D point-cloud object recognition using single-word class names.** Multiclass accuracy on seen and unseen classes for ModelNet-40. The seen accuracy is between 30 classes, and the unseen accuracy is between 10 classes. Further description of the protocol can be found in Appendix D.

	AwA Super Sets		
	Seen	Unseen	Harmonic
DeViSE [18]	53.0 $\pm$ 1.9	50 $\pm$ 0.6	51.5 $\pm$ 0.9
DEM [71]	50.1 $\pm$ 1.4	48.3 $\pm$ 1.8	49.2 $\pm$ 1.6
CIZSL [14]	57.3 $\pm$ 5.6	50.2 $\pm$ 5.8	55.0 $\pm$ 4.0
GRaWD [24]	59.8 $\pm$ 3.5	51.6 $\pm$ 4.8	55.3 $\pm$ 3.1
T2M-HN (ours)	<b>67.2 <math>\pm</math> 5.2</b>	<b>57.3 <math>\pm</math> 5.7</b>	<b>61.9 <math>\pm</math> 5.4</b>

Table 7: **Classification using attributes.** Values denote the Area under the Recall-Precision curve averaged over the 13 test attributes  $\pm$  s.e.m. over these attributes. The seen results are new images from the seen classes, while the unseen results are images from unseen classes. Both are evaluated when classifying only the test attributes. The full protocol is in G.

from the unseen classes (10-classes classification). T2M-HN achieves SOTA results in this setup as well. It leverages the text generalization of the HN model to distinguish between unseen classes.

We further computed the top- $k$  accuracy achieved by running T2M-HN for the unseen classes. Figure 9 plots the accuracy as a function of  $k$ . T2M-HN provides superior accuracy for all tested values of  $k$ . To calculate the top- $k$  performance of the GAN-based models, after generating the images, we checked if any of  $K$  closest neighbors of an image is of the correct class.

### E. AwA GPT-3 descriptions

We use GPT3 [8] to generate 5 synthetic descriptions for each class of AwA. We use the API provided by OpenAI to ask "text-davinci-002" engine with a temperature of 0, max tokens of 512, and the prompt: "Suggest 5 definitions for an animal. Animal: {animal\_name}. Definitions:"

**Animal: moose**

Definitions:

1. A large, dark-colored deer with enormous antlers, native to North America and Europe.

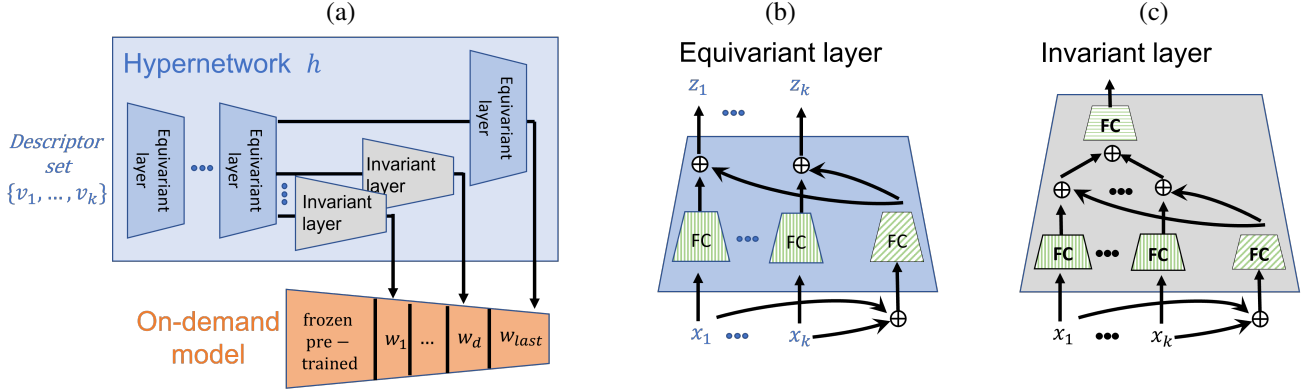


Figure 8: (a) The T2M-HN architecture for equivariant-invariant hypernetwork. The input is processed by equivariant layers, followed by a prediction head for each layer of the target on-demand classifier  $f$ . The prediction head for  $w_{last}$  is equivariant. Heads for earlier layers of  $f$ ,  $w_1, \dots, w_k$  are invariant. (b) An architecture for the equivariant layer. Every input is processed by a fully connected (FC) layer in a Siamese manner (shared weights). Inputs are also summed and processed by a second FC layer, whose output is added back to each output. (c) An architecture for an invariant layer, following a similar structure to b.

	AWA by class name			AWA by GPT descriptions			BABEL by class names		
	Seen	Unseen	Harmonic	Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
T2M-HN 1-layer	$98.9 \pm .1$	$87.3 \pm .2$	$92.7 \pm .1$	$98.7 \pm .1$	$83.3 \pm .1$	$90.3 \pm .1$	$94.9 \pm 0.1$	$81.8 \pm 0.1$	$87.9 \pm 0.1$
1-layer w.o. EV	$98.6 \pm .1$	$85.9 \pm .2$	$91.8 \pm .1$	$97.0 \pm .1$	$78.2 \pm .1$	$86.6 \pm .1$	$81.4 \pm 0.3$	$74.7 \pm 0.1$	$77.9 \pm 0.1$
T2M-HN 2-layers	$98.0 \pm .1$	$87.2 \pm .1$	$92.3 \pm .1$	$96.4 \pm .1$	$76.7 \pm .1$	$85.4 \pm .1$	$92.4 \pm 0.1$	$86.2 \pm 0.1$	$89.2 \pm 0.1$
2-layers w.o. EV	$98.6 \pm .1$	$75.7 \pm .1$	$85.5 \pm .1$	$74.0 \pm .3$	$57.4 \pm .1$	$64.6 \pm .2$	$92.1 \pm 0.2$	$78.9 \pm 0.2$	$85.0 \pm 0.2$

Table 8: **Ablation study.** Mean classification accuracy on seen and unseen classes and their harmonic mean for AWA and BABEL datasets. Values are averages over all class pairs. T2M-HN is the proposed method designed with equivariant and invariance properties. We evaluate two variants of T2M-HN, one that produces a single FC layer to the on-demand model, and a second variant that produces two FC layers. To demonstrate the importance of the equivariant design, we evaluate an HN that produces 1 and 2 layers without an equivariant design.

	Easy tasks			Hard tasks		
	Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
DeViSE	$68.4 \pm 0.3$	$64.4 \pm 0.8$	$66.3 \pm 0.5$	$59.0 \pm 1.3$	$53.4 \pm 4.5$	$55.9 \pm 2.6$
DEM	$73.2 \pm 0.3$	$70.2 \pm 0.6$	$71.7 \pm 0.3$	$57.0 \pm 1.0$	$55.5 \pm 3.1$	$56.3 \pm 1.7$
CIZSL	$87.5 \pm 0.1$	$85.2 \pm 0.3$	$86.3 \pm 0.2$	$71.3 \pm 0.7$	$68.8 \pm 2.1$	$70.0 \pm 1.2$
GRaWD	$88.7 \pm 0.1$	$86.1 \pm 0.3$	$87.4 \pm 0.2$	$72.7 \pm 0.7$	$68.2 \pm 1.9$	$70.4 \pm 1.1$
T2M-HN (ours)	$94.7 \pm 0.1$	$75.7 \pm 0.1$	$84.1 \pm 0.1$	$93.1 \pm 0.1$	$75.2 \pm 0.2$	$83.2 \pm 0.1$

Table 9: **Classification by class descriptions for the CUB dataset.** Easy tasks are binary classification tasks, where pairs of birds are from different taxonomy families. In hard classification tasks, bird classes in each pair are from the same taxonomy family. Values are the mean classification accuracy on images from seen and unseen classes, averaged over all class pairs.

2. An animal of the deer family with humped shoulders, long legs, and a large head with antlers.
3. A large, awkward-looking mammal with a long face and humped shoulders.
4. A very large deer with antlers that can spread six feet or more from tip to tip.
5. The largest member of the deer family, with males weighing up to 1,800 pounds and having antlers that

can spread up to six feet from tip to tip.

#### Animal: spider monkey

Definitions:

1. A type of monkey that has long legs and arms and a long tail.
2. A monkey that is found in the rainforests of Central and South America.

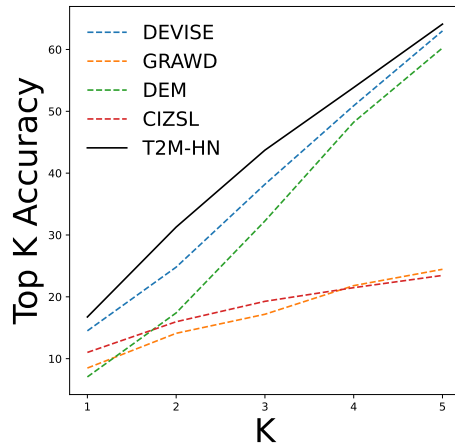


Figure 9: Accuracy at k for experiments with 3D point cloud from ModelNet-40. The solid line is our T2M model, dashed lines are for the baseline models.

3. A monkey that is known for its acrobatic abilities.
4. A monkey that is considered to be one of the most intelligent primates.
5. A monkey that is endangered in many parts of its range.

#### Animal: rhinoceros

Definitions:

1. A large, thick-skinned mammal with one or two horns on its snout, native to Africa and southern Asia.
2. An animal that is hunted for its horn, which is used in traditional Chinese medicine.
3. A large, herbivorous mammal with a single horn on its nose, found in Africa and southern Asia.
4. A mammal of the family Rhinocerotidae, having thick, grey or brown skin and one or two horns on the snout.
5. A very large, plant-eating mammal with one or two horns on its nose, found in Africa and southern Asia.

#### Elephant:

1. The largest land animal in the world, with males weighing up to six tons.
2. A plant-eating mammal with a long trunk, large ears, and thick, grey skin.
3. A mammal of the family Elephantidae, having a long trunk, large ears, and thick, grey skin.
4. An intelligent animal that is known for its memory and its ability to use its trunk for a variety of tasks.
5. An endangered species that is hunted for its ivory tusks.

## F. Data splits

**SUN unseen classes:** 'volcano', 'poolroom establishment', 'veterinarians office', 'reception', 'field wild', 'diner indoor', 'garbage dump', 'server room', 'vineyard', 'jewelry shop', 'drugstore', 'herb garden', 'lock chamber', 'temple east asia', 'marsh', 'cottage garden', 'cathedral outdoor', 'dentists office', 'pharmacy', 'hangar indoor', 'volleyball court indoor', 'lift bridge', 'synagogue outdoor', 'boathouse', 'ice shelf', 'boxing ring', 'rope bridge', 'electrical substation', 'auditorium', 'chalet', 'booth indoor', 'wine cellar barrel storage', 'greenhouse outdoor', 'badminton court indoor', 'thriftshop', 'cemetery', 'rainforest', 'courtyard', 'underwater coral reef', 'formal garden', 'ice skating rink outdoor', 'palace', 'movie theater indoor', 'dinetto home', 'sandbar', 'ball pit', 'amphitheater'

**SUN seen classes:** All remaining classes.

**ModelNet40:** We follow [11, 12, 37] and use the 10 classes included in ModelNet-10 as unseen classes, and the other 30 as seen.

**BABEL unseen classes:** 'a pose', 'action with ball', 'adjust', 'catch', 'clean something', 'communicate (vocalise)', 'crawl', 'get injured', 'hand movements', 'hop', 'limp', 'mix', 'play sport', 'press something', 'rolling movement', 'shuffle', 'side to side movement', 'sneak', 'spread', 'support', 'swing body part', 'trip', 'upper body movements', 'wait'

**BABEL seen classes:** All remaining classes.

**CUB unseen classes:** 'Acadian Flycatcher', 'American Crow', 'American Three Toed Woodpecker', 'Baltimore Oriole', 'Bank Swallow', 'Belted Kingfisher', 'Black Billed Cuckoo', 'Black Footed Albatross', 'Black Throated Sparrow', 'Boat Tailed Grackle', 'Bohemian Waxwing', 'Brandt Cormorant', 'Brewer Blackbird', 'Cape May Warbler', 'Cedar Waxwing', 'Chestnut Sided Warbler', 'Field Sparrow', 'Golden Winged Warbler', 'Grasshopper Sparrow', 'Gray Crowned Rosy Finch', 'Great Crested Flycatcher', 'Great Grey Shrike', 'Groove Billed Ani', 'Hooded Oriole', 'Horned Grebe', 'Indigo Bunting', 'Least Auklet', 'Least Tern', 'Marsh Wren', 'Mockingbird', 'Northern Flicker', 'Northern Waterthrush', 'Pacific Loon', 'Pied Billed Grebe', 'Pomarine Jaeger', 'Purple Finch', 'Red Legged Kittiwake', 'Rhinoceros Auklet', 'Sayornis', 'Scott Oriole', 'Tree Sparrow', 'Tree Swallow', 'Western Grebe', 'Western Gull', 'Western Wood Pewee', 'White Breasted Kingfisher', 'White Eyed Vireo', 'White Pelican', 'Wilson Warbler', 'Yellow Bellied Flycatcher', 'Yellow Billed Cuckoo'

**CUB seen classes:** All remaining classes.

## G. Attributes used for one-class classification

As mentioned in section 5.4, we use some of the attributes from the AWA dataset to define one-class classification.

cation tasks. First, we removed non-visual attributes. Then, we randomly split the remaining 53 attributes into 30 train, 10 validation, and 13 test attributes. We split both the images and the attributes, constructing 4 groups of images and attributes: (1) *Training images* from training attributes and training classes, used to train the hypernetwork; (2) *Validation images* from the training classes, with the validation attributes used to tune hyperparameters; (3) *Test images from seen classes*, new images of test attributes, whose class was seen during training (but not the specific images); and (4) *Test images from unseen classes*, new images of test attributes, whose class was not seen during training. We report the average Area under the Recall-Precision curve over seen (group (3)) and unseen classes (group (4)). The results are shown in Figure 5 and in Table 7. The attributes split is as follows:

**AwA train attributes:** 'orange', 'red', 'longneck', 'horns', 'tusks', 'flies', 'desert', 'cave', 'jungle', 'water', 'bush', 'lean', 'forest', 'gray', 'straianteeth', 'stripes', 'mountains', 'arctic', 'paws', 'hooves', 'pads', 'small', 'furry', 'ground', 'patches', 'white', 'fields', 'bipedal', 'toughskin', 'plains'.

**AwA validation attributes:** 'buckteeth', 'chewteeth', 'yellow', 'hairless', 'bulbous', 'big', 'flippers', 'tree', 'walks', 'coastal'.

**AwA test attributes:** 'quadrapedal', 'black', 'blue', 'ocean', 'longleg', 'spots', 'hands', 'claws', 'muscle', 'meatteeth', 'tail', 'brown', 'swims'.

## H. Results in Tables

Numeric results for Figure 7 and Figure 4 are in tables 8 and 9, respectively.