

# ChatGPT-guided Semantics for Zero-shot Learning

Fahimul Hoque Shubho<sup>\*</sup>, Townim Faisal Chowdhury<sup>†</sup>, Ali Cheraghian<sup>‡</sup>,  
Morteza Saberi<sup>§</sup>, Nabeel Mohammed<sup>\*</sup>, Shafin Rahman<sup>\*</sup>

<sup>\*</sup>Dept. of Electrical and Computer Engineering, North South University, Bangladesh,

<sup>†</sup>Australian Institute for Machine Learning, University of Adelaide, Australia,

<sup>‡</sup>Data61, Commonwealth Scientific and Industrial Research Organisation, Australia

<sup>§</sup>School of Computer Science and DSI, University of Technology Sydney, Australia

{hoque.shubho, nabeel.mohammed, shafin.rahman}@northsouth.edu,

townim.chowdhury@adelaide.edu.au, ali.cheraghian@data61.csiro.au, morteza.saberi@uts.edu.au

**Abstract**—Zero-shot learning (ZSL) aims to classify objects that are not observed or seen during training. It relies on class semantic description to transfer knowledge from the seen classes to the unseen classes. Existing methods of obtaining class semantics include manual attributes or automatic word vectors from language models (like word2vec). We know attribute annotation is costly, whereas automatic word vectors are relatively noisy. To address this problem, we explore how ChatGPT, a large language model, can enhance class semantics for ZSL tasks. ChatGPT can be a helpful source to obtain text descriptions for each class containing related attributes and semantics. We use the word2vec model to get a word vector using the texts from ChatGPT. Then, we enrich word vectors by combining the word embeddings from class names and descriptions generated by ChatGPT. More specifically, we leverage ChatGPT to provide extra supervision for the class description, eventually benefiting ZSL models. We evaluate our approach on various 2D image (CUB and AwA) and 3D point cloud (ModelNet10, ModelNet40, and ScanObjectNN) datasets and show that it improves ZSL performance. Our work contributes to the ZSL literature by applying ChatGPT for class semantics enhancement and proposing a novel word vector fusion method.

**Index Terms**—Zero-shot learning, 3D point cloud, ChatGPT, Word vectors, Language models

## I. INTRODUCTION

The goal of Zero-shot Learning (ZSL) is to classify unseen objects not observed in training. A more generic version termed Generalized ZSL (GZSL) attempts to predict a class from seen and unseen classes together. Researchers have started exploring ZSL and GZSL with 2D image datasets. Later, considering the availability of depth-sensing cameras, exploring ZSL on 3D point cloud data got considerable attention [1]–[3]. For both 2D and 3D cases, semantic descriptions of classes play a pivotal role in transferring knowledge from seen to unseen classes. Class semantics are designed to describe all objects with a common set of features or components, working as a bridge between seen and unseen worlds. The literature shows that addressing ZSL tasks in 3D has more challenges than its 2D counterpart. Therefore, improving class semantics may help to address some challenges. This paper attempts to improve class semantic descriptions for 2D and 3D ZSL tasks.

Class semantics can be obtained manually (attribute vectors) or automatically (word vectors). Attributes are identifiable

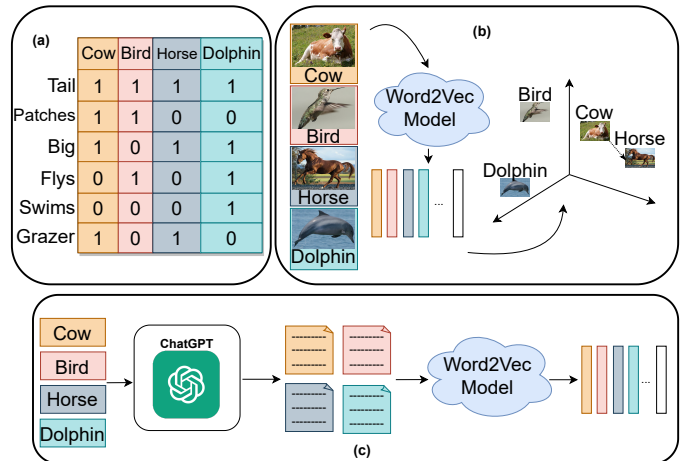


Fig. 1: Class semantics can be obtained from (a) manually annotated attributes or (b) automatically calculated from language models like word2vec. Attributes need hard manual labor, whereas word vectors are relatively noisy. In this paper (c), we attempt to improve the quality of automatic word vectors by using extra supervision from ChatGPT.

features to describe a class (see Fig. 1(a)) that require laborious human annotation to obtain and are not readily available for many large-scale or 3D point cloud datasets. In contrast, automatic word vectors are the output (as vectors) of language models (like word2vec [4], GloVe [5], Fasttext [6], Bert [7], etc.), given class names as input (see Fig. 1(b)). These models are usually trained using billions of text corpus from Wikipedia, news articles, etc. Compared to attributes, the automatic extraction of word vectors makes them more realistic for real-world applications. However, embeddings from word vectors are noisier than manual attribute vectors resulting in poorer ZSL performance than attributes. This issue becomes more challenging, especially for ZSL on 3D point cloud objects, because of pre-trained models, poor quality features, dataset size, etc. [2]. In this paper, we investigate the use of ChatGPT [8] to improve class semantics.

Usually, word vectors are calculated using a single class name as input. However, many related words and definitions

associated with that single class name are also necessary to improve the representativeness of the word vector. Considering related semantics can improve the semantic description of the given class [9]. Again, such related semantics can be obtained by hard manual annotations [10], a costly process, or noisy web crawling annotation [11]. To address this problem, recent ChatGPT can be a valuable source for describing a class with related semantics and attributes (see Fig. 1(c)). It is both automatic and less noisy. Therefore, given a class name, we ask ChatGPT to describe that class with a paragraph of text containing related semantics and attributes. Then, we extract a word embedding of that paragraph using a language model(word2vec). Finally, we linearly combine word embeddings from class names and ChatGPT to calculate an improved word vector. Different from [3], this approach does not need any prompt engineering, so it can be used in any existing ZSL models to increase accuracy. We test our approach on seven methods DEM [10], LATEM [12], SYNC [13], GDAN [14], f-CLSWGAN [15], TF-VAEGAN [16] and CADA-VAE [17] covering both 2D image (CUB [18] and AwA [19]) and 3D point cloud datasets (ModelNet10, ModelNet40 [20] and ScanObjectNN [21]). Note that our method is also applicable for both synthetic ModelNet40 and real-world scanned ScanObjectNN cases. We consistently improve performance in experiments across different ZSL and GZSL problem setups.

Our contributions are as follows:

- Employing ChatGPT guided semantics for solving ZSL problems,
- Proposing an approach that uses texts provided by ChatGPT as extra supervision to enhance word vectors while representing class semantics that readily be plugged into any ZSL method to improve accuracy,
- Extensive experiments on 2D images and 3D point cloud object (synthetic and real-world scanned) datasets.

## II. RELATED WORKS

**ZSL with 2D images:** Initial work of ZSL [22] used image-related information attributes to classify images in the unseen test set, and the model was trained by utilizing the seen attributes. In another work, [23] proposed using labeled image data and semantic information gleaned from unlabeled text using a deep visual-semantic embedding model. The model learns the semantic relationship between labels from the textual data and directly maps images into a deep semantic embedding space. In a later work, Akata et al. [24] used image features with supervised attributes and unsupervised output embedding derived from unlabeled text corpora to project image features in attribute space and measure similarity with class description to classify images. In [25], Kodirov et al. proposed learning in the reverse direction of the previous work by using an unsupervised domain adaption approach that employs the target domain class labels’ projections in the semantic space to regularise the learned target domain projection. In contrast to using a global linear embedding like previous works, [12] learns various linear models and allows each image-class combination to be chosen from them. Recent

works use a generative model-based approach to mitigate bias and domain shift problems observed in the ZSL. In [26], Mishra et al. proposed procreating samples using the attributes with conditional variational autoencoder and used the created samples to classify the unseen classes. In [27], the authors leveraged the semantic relationship between the seen and unseen classes to generate visual features for the unseen courses while maintaining the same semantic relationship between both classes. Hayat et al. [28] proposed synthesizing the features for unseen classes using a generative model with their class semantics for ZSL in object detection. This paper aims to ZSL and GZSL tasks based on improved class semantics.

**ZSL with 3D point cloud object:** The works around 3D ZSL are relatively new compared to 2D images. In a pioneer work, Cheraghian et al. [29] proposed using 3D point clouds objects for zero-shot tasks following an inductive process with transductive ZSL to classify 3D cloud points and show that it can achieve competitive performance for 3D point cloud classification. Later, [30] proposed a loss function comprised of a regression and skewness term to reduce the hubness problems described in [29]. The authors used transductive ZSL and Generalized Zero-Shot Learning (GZSL) approaches for classifying 3D point data. They also introduced a triplet loss function for these approaches. [2] proposed a unified approach to address the hubness and domain shift problems using a novel loss function. [31] introduced a generative approach to classify and segment 3D data for ZSL and GZSL. In [32], the authors used NLP based prompt to augment new data that improves the performance of ZSL on 3D datasets. This paper investigates the performance of both embedding and generative methods using improved class semantics.

**LLM in computer vision tasks:** Large Language Models (LLM) like GPT have been used widely in natural language processing tasks. Recent works have used LLM for computer vision tasks also. [33] leveraged LLM to generate narration from video input. [34] proposed that LLM models such as GPT-3.5 can annotate unlabeled data. They also showed that the annotation from LLM works better for user input and keyword relevance assessment compared to using crowd-sourced annotation data. Later, [35] used 2 LLM, GPT-3.5 and GPT-4, for automatic annotation of texts with specific categories of linguistic information and compared it with human annotator to show that it is feasible to use LLM with the annotation for local grammar analysis. In [36], the authors used several LLMs, such as Dolly-v2, StableVicuna, ChatGPT, and GPT-4, for data augmentation in cross-lingual commonsense reasoning datasets to overcome the imitation of training data and showed that it has a positive impact on the overall performance of the trained models. [37] applied a language model to procreate class-conditioned texts attuned by prompts and utilized them as the training data for fine-tuning a pre-trained bidirectional language model used for zero-shot learning of Natural Language Understanding tasks. This paper uses ChatGPT to ensure extra supervision for class semantics during ZSL and GZSL tasks.

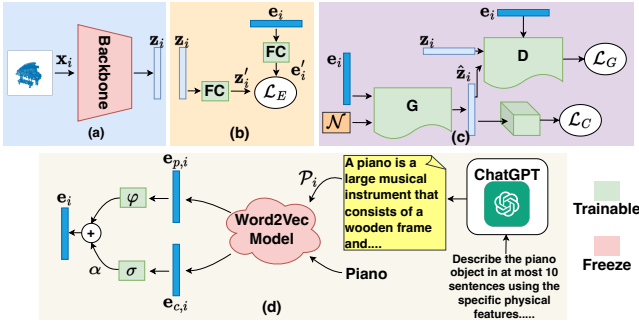


Fig. 2: We demonstrate the profound impact of integrating ChatGPT-guided semantic information into traditional ZSL techniques. (a) The input data  $\mathbf{x}_i \in \mathbb{R}^K$  (e.g., images or point clouds) is initially processed through a Backbone, producing a feature embedding  $\mathbf{z}_i \in \mathbb{R}^m$ . (b) ZSL methods often utilize a shared common space, where the distance between the mapped feature embeddings,  $\mathbf{z}_i \in \mathbb{R}^m$  and  $\mathbf{e}_i \in \mathbb{R}^d$ , is minimized through a loss function. (c) In generative ZSL methods, a generator model receives a semantic embedding  $\mathbf{e}_i \in \mathbb{R}^d$ , along with Gaussian noise, as input to produce synthetic feature embeddings  $\hat{\mathbf{z}}_i \in \mathbb{R}^m$ . (d) In our proposed method, we leverage ChatGPT to generate additional descriptions for classes belonging to the seen category. Both the ChatGPT output and the class names are fed into a Word2vec model to extract their respective semantic feature embeddings. These extracted feature embeddings are then passed through fully connected layers and merged in the subsequent stage to generate the final comprehensive semantic feature embeddings,  $\mathbf{e}_i \in \mathbb{R}^d$ .

### III. METHOD

We demonstrate the potential of leveraging the generation capability of LLMs to enhance ZSL methods for recognizing unseen classes. LLMs have gained significant attention due to numerous advantages across various applications. To harness this power, we specifically utilize the ChatGPT model to generate supplementary descriptions for the class names during the training phase of ZSL. This description includes attributes and semantics related to a class that can enhance its discriminative ability from other classes. Encoding this description with a word vector by forwarding it through a language model (word2vec) can augment additional supervision to ZSL models.

#### A. Problem formulation

Let  $\mathbf{x} \in \mathbb{R}^K$  represent the input data, corresponding to either an image or point cloud data. We have two sets of class labels,  $\mathcal{Y}^s = \{y_1^s, \dots, y_S^s\}$  for seen classes and  $\mathcal{Y}^u = \{y_1^u, \dots, y_U^u\}$  for unseen classes. Here, the seen and unseen labels are disjoint, i.e.,  $\mathcal{Y}^s \cap \mathcal{Y}^u = \emptyset$ . Additionally,  $\mathcal{E}^s = \{\phi(y_1^s), \dots, \phi(y_S^s)\}$  and  $\mathcal{E}^u = \{\phi(y_1^u), \dots, \phi(y_U^u)\}$  represent the sets of semantic feature embeddings obtained using the embedding function  $\phi(\cdot)$ , where  $\phi(y) \in \mathbb{R}^d$ . To proceed, we define the set of  $n_s$  seen samples as  $\mathcal{D}^s = \{(\mathbf{x}_i^s, y_i^s, \mathbf{e}_i^s)\}_{i=1}^{n_s}$ , where  $\mathbf{x}_i^s$  represents

the  $i^{\text{th}}$  instance from the seen set, with ground truth label  $y_i^s \in \mathcal{Y}^s$  and corresponding semantic vector  $\mathbf{e}_i^s = \phi(y_i^s) \in \mathcal{E}^s$ . Similarly, the set of  $n_u$  unseen samples is defined as  $\mathcal{D}^u = \{(\mathbf{x}_i^u, y_i^u, \mathbf{e}_i^u)\}_{i=1}^{n_u}$ , where  $\mathbf{x}_i^u$  represents the  $i^{\text{th}}$  sample from the unseen set, with ground truth label  $y_i^u \in \mathcal{Y}^u$  and corresponding semantic vector  $\mathbf{e}_i^u = \phi(y_i^u) \in \mathcal{E}^u$ . This paper addresses two main tasks: Zero-Shot Learning (ZSL) and Generalized Zero-Shot Learning (GZSL). In ZSL, we assign unseen class labels  $\hat{\mathbf{y}} \in \mathcal{Y}^u$  to unseen instances  $\mathbf{x} \in \mathcal{D}^u$ . In GZSL, we assign class labels  $\hat{\mathbf{y}} \in \mathcal{Y}^s \cup \mathcal{Y}^u$  to instances  $\mathbf{x} \in \mathcal{D}^s \cup \mathcal{D}^u$ , encompassing both seen and unseen classes. The semantic vectors  $\mathbf{e}_i^s$  and  $\mathbf{e}_i^u$  are usually implemented with attribute sets or word vectors. This paper proposes a method to produce an improved version of word2vec-based word vectors.

#### B. Preliminaries

Two main types ZSL methods are embedding [20]–[24] and generative [26], [28], [38]. Embedding methods learn to map both the visual features of the input data and the semantic attributes to a common space. They can recognize unseen classes in this space by comparing their attribute representations with the input’s features. In contrast, generative methods synthesize samples that look like unseen classes by using a mix of labeled data from seen classes and auxiliary information, such as class attributes or textual descriptions.

Overall, in any ZSL methodology, an input backbone transfers the input data  $\mathbf{x}_i \in \mathbb{R}^K$  into a meaningful feature embedding, denoted as  $\mathbf{z}_i \in \mathbb{R}^m$ , see Fig 2 (a). Next, we will delve deeper into the embedding and synthetic Zero-Shot Learning (ZSL) methods.

**Embedding ZSL:** Embedding-based ZSL aims to learn functions that map input data (images or point clouds) and semantic information (attributes or class labels) to a shared space. In this space, the model can link input features with semantic information, allowing it to identify and categorize classes it has not seen before. In general, There are two branches in the embedding approach Fig 2 (b). In the first branch, the input feature embedding  $\mathbf{z}_i \in \mathbb{R}^m$  is forwarded into a fully connected layer, resulting in  $\mathbf{z}'_i \in \mathbb{R}^q$ . Simultaneously, the corresponding semantic representation  $\mathbf{e}_i \in \mathbb{R}^d$  is fed into a fully connected layer, which generates  $\mathbf{e}'_i \in \mathbb{R}^q$  to map the semantic embedding to a common space, where the Euclidean distance between the semantic and feature embeddings is minimized. This is achieved by optimizing the following objective function:

$$\mathcal{L}_E = \frac{1}{n_s} \sum_{i=1}^{n_s} \|\mathbf{z}'_i - \mathbf{e}'_i\|_2^2 + \lambda \mathcal{R}(\theta) \quad (1)$$

where the  $\theta$  is the trainable parameters, and  $\mathcal{R}$  refers to the regularization loss function. The hyperparameter  $\lambda$  is crucial in controlling the trade-off between the regularization and embedding losses.

**Generative ZSL:** In generative Zero-Shot Learning (ZSL) methods [14], [15], the objective is to generate synthetic samples for the unseen classes based on their semantic attributes.

Input Prompt	Describe the <b>piano</b> object in at most ten sentences using the specific physical features and do not need to mention the features that are not available in the object. Also, do not use any numeric in descriptions; instead, use words.
ChatGPT Generate Description	A piano is a large musical instrument that consists of a wooden frame and various physical components. The main body of the piano contains the strings and soundboard. The strings are stretched horizontally across the length of the piano, with different lengths and thicknesses producing various pitches when struck. The soundboard is a large wooden panel located underneath the strings, responsible for resonating and amplifying the vibrations produced by the strings. The piano has a keyboard consisting of black and white keys.

TABLE I: A sample input prompt to ChatGPT and generated text description. The word vector generated from the text description of a class name (**piano**) can encode auxiliary semantics inside the word vector of **piano**.

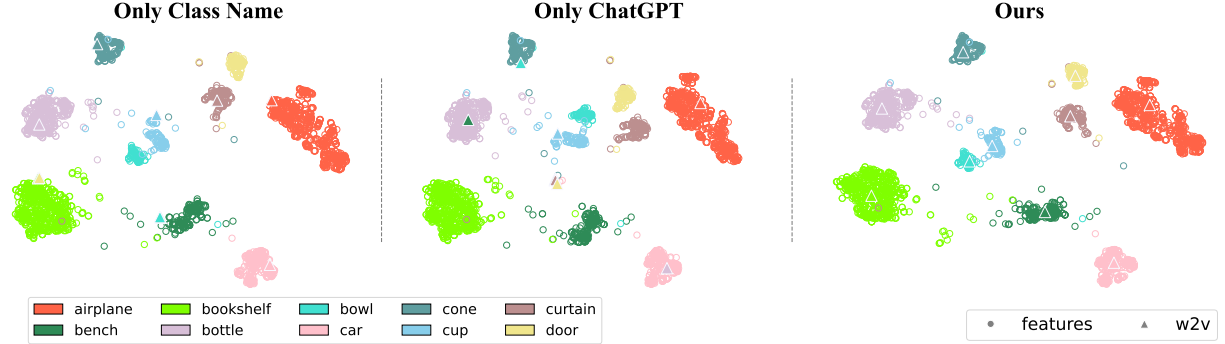


Fig. 3: tSNE visualization of features and semantics of 10 classes from ModelNet40 dataset. Our method has achieved better feature-semantic alignment than only Class name and ChatGPT-based embedding vectors.

These methods enable the model to generate realistic and representative samples of unseen classes by leveraging the semantic information associated with those classes.

To be more specific, our goal is to learn a conditional Generative Adversarial Network (GAN) model, denoted as  $G$ , which takes random Gaussian noise  $\mathbf{h} \sim \mathcal{N}(0, 1)$  and the semantic class embedding  $\mathbf{e}_i$  as inputs to generate the feature representation of class  $i$ , denoted as  $\hat{\mathbf{z}}_i \in \mathbb{R}^m$ . Concurrently, we train a discriminator model,  $D$ , to classify real features  $\mathbf{z}_i \in \mathbb{R}^m$  against synthetic features  $\hat{\mathbf{z}}_i \in \mathbb{R}^m$ .

The objective function for generating synthetic feature samples is defined as follows:

$$\mathcal{L}_G = \mathbb{E}_{\mathbf{z}, \mathbf{e}} [D(\mathbf{z}, \mathbf{e})] - \mathbb{E}_{\hat{\mathbf{z}}, \mathbf{e}} [D(\hat{\mathbf{z}}, \mathbf{e})] - \eta \mathbb{E}_{\tilde{\mathbf{z}}, \mathbf{e}} \left[ (\|\nabla_{\tilde{\mathbf{z}}} D(\tilde{\mathbf{z}}, \mathbf{e})\|_2 - 1)^2 \right], \quad (2)$$

where  $\tilde{\mathbf{z}} = \beta \mathbf{z} + (1 - \beta) \hat{\mathbf{z}}$ , with  $\beta \sim \mathcal{U}(0, 1)$ , and  $\eta$  is the penalty coefficient. This objective function guides the training of our conditional GAN model, enabling it to generate realistic and diverse feature representations based on the input class embeddings and Gaussian noise.

Also, a classification loss is required to ensure that the synthetic samples are suitable for the classifier. This loss is defined as follows:

$$\mathcal{L}_C = -E_{\hat{\mathbf{z}} \sim p_{\hat{\mathbf{z}}}} [\log P(\mathbf{y} | \hat{\mathbf{z}}; \Theta)], \quad (3)$$

The provided loss function is computed using a linear softmax classifier parameterized by  $\Theta$ , which undergoes pre-training on the real features  $\mathbf{z} \in \mathcal{D}^s$  from seen classes. To elaborate, this loss function serves as a regularizer, motivating the generator to construct discerning features in its generated samples.

### C. Improved class semantics

We utilize the ChatGPT model by providing it with the class name of a seen class,  $\mathbf{y}^s \in \mathcal{Y}^s$ , to generate descriptive sentences. Specifically, we input the class names into the GPT-3.5 [8], using prompts such as ‘‘Describe the [CLASS] in at most ten sentences, focusing on specific physical features and excluding any unavailable features.’’ The generated sentences offer detailed descriptions of the class names, aiding in improving class semantics for our ZSL tasks. We denote the prompts created for each sample, formulated as follows:

$$\mathcal{P} = \text{ChatGPT}(\text{Commands}). \quad (4)$$

An illustrative example of a class name description obtained from ChatGPT is presented in the following Table I.

In our innovative approach, visualized in Fig. 2 (d), we have devised a multi-step process to enrich our understanding of class characteristics. The class name is initially processed by ChatGPT, generating a comprehensive description denoted as  $\mathcal{P}_i$ . This description includes multiple sentences closely related to the class name, as discussed earlier.

Subsequently, both the class name’s semantic representation, denoted as  $\mathbf{e}_{c,i} \in \mathbb{R}^d$ , and the ChatGPT output’s semantic representation, denoted as  $\mathbf{e}_{p,i} \in \mathbb{R}^d$ , are extracted using a Word2vec model. These semantic representations offer valuable insights into the meaning and context of the class name and its associated description. Afterwards, these semantic descriptions are passed through separate fully connected layers, represented by  $\sigma$  and  $\alpha$  respectively, before being merged using the following equation:

$$\mathbf{e}_i = \sigma(\mathbf{e}_{c,i}) + \alpha \cdot \varphi(\mathbf{e}_{p,i}) \quad (5)$$

The merging process through addition allows the combination of the class-specific information with the contextual

Method	Variations	AwA					CUB				
		ZSL	GZSL				ZSL	GZSL			
		Acc	Acc <sub>s</sub>	Acc <sub>u</sub>	HM	BC	Acc	Acc <sub>s</sub>	Acc <sub>u</sub>	HM	BC
DEM [10]	Only Class Name	45.79	83.15	10.95	19.35	0	10.68	18.47	2.62	4.59	0
	Only ChatGPT	<b>54.74</b>	56.09	8.71	15.09	1	11.69	<b>47.78</b>	1.62	3.14	1
	Ours	52.26	<b>85.45</b>	<b>15.33</b>	<b>25.99</b>	<b>3</b>	<b>15.17</b>	30.91	<b>2.63</b>	<b>4.85</b>	<b>3</b>
LATEM [12]	Only Class Name	46.21	-	-	-	0	11.35	-	-	-	0
	Only ChatGPT	50.47	-	-	-	0	10.19	-	-	-	0
	Ours	<b>52.31</b>	-	-	-	<b>1</b>	<b>14.50</b>	-	-	-	<b>1</b>
SYNC [13]	Only Class Name	48.83	-	-	-	0	12.71	-	-	-	0
	Only ChatGPT	<b>55.99</b>	-	-	-	<b>1</b>	11.90	-	-	-	0
	Ours	55.10	-	-	-	0	<b>15.28</b>	-	-	-	<b>1</b>
GDAN [14]	Only Class Name	-	<b>60.07</b>	7.88	13.93	1	-	21.02	<b>6.54</b>	<b>9.97</b>	<b>2</b>
	Only ChatGPT	-	56.92	20.29	29.91	0	-	<b>27.09</b>	1.12	2.15	1
	Ours	-	56.35	<b>23.48</b>	<b>33.15</b>	<b>2</b>	-	20.72	6.44	9.83	0
f-CLSWGAN [15]	Only Class Name	42.48	<b>91.4</b>	0.62	1.24	1	12.56	<b>59.01</b>	2.00	3.87	1
	Only ChatGPT	45.05	90.81	0.71	1.41	0	11.84	58.34	1.64	3.19	0
	Ours	<b>45.74</b>	91.00	<b>0.76</b>	<b>1.50</b>	<b>3</b>	<b>15.37</b>	57.41	<b>3.20</b>	<b>6.05</b>	<b>3</b>
TF-VAEGAN [16]	Only Class Name	56.3	72.44	<b>45.34</b>	55.37	1	16.55	30.77	14.01	19.25	0
	Only ChatGPT	<b>64.31</b>	<b>78.72</b>	40.53	53.51	<b>2</b>	15.51	26.51	<b>14.68</b>	18.89	1
	Ours	57.14	73.67	44.54	<b>55.47</b>	1	<b>20.46</b>	<b>48.39</b>	14.64	<b>22.48</b>	<b>3</b>
CADA-VAE [17]	Only Class Name	36.79	<b>90.02</b>	19.16	31.6	1	16.17	64.71	4.65	8.68	0
	Only ChatGPT	<b>45.01</b>	82.93	20.51	32.89	1	11.71	63.44	4.02	7.56	0
	Ours	40.17	89.25	<b>20.69</b>	<b>33.59</b>	<b>2</b>	<b>16.50</b>	<b>64.97</b>	<b>5.30</b>	<b>9.79</b>	<b>4</b>

TABLE II: ZSL and GZSL results on 2D image datasets.

details from ChatGPT, resulting in a more comprehensive and enriched representation, denoted as  $\mathbf{e}_i \in \mathbb{R}^d$ . In Fig. 3, we visualize  $\mathbf{e}_c$ ,  $\mathbf{e}_p$  and  $\mathbf{e}$  of each unseen class and their corresponding visual features. We notice that with our method class semantic word vectors ( $\mathbf{e}$ ) and features form clusters, indicating adequate feature-semantic alignment, which is essential for ZSL.

#### IV. EXPERIMENTS

**2D Image dataset setup:** For ZSL on 2D images, two well-known datasets are used in this work: (1) The Caltech-UCSD Birds (CUB-200-2011) [18] contains 11,788 images of 200 bird species. We have used the seen-unseen split setting from [19] with 150 classes as seen data for training and 50 as unseen data for testing. (2) Animals with Attributes (AwA) [19] dataset consists of 30,745 images of 50 classes with the seen-unseen split of 40 seen classes for training and ten unseen classes for testing.

**3D Point cloud dataset setup:** We have used three well-known datasets, ScanObjectNN [39], ModelNet10 and ModelNet40 [20]. ScanObjectNN dataset consists of 2,902 3D real-world point cloud data of objects of 15 categories. The ModelNet40 dataset contains 12,311 CAD-generated 3D meshes of 40 categories, and ModelNet10 is a part of it. For the Point Cloud datasets, we have used the splits from [2].

**Train-Test-Split:** The dataset split used for ZSL divides the data into a seen set and an unseen set. The seen set consists of examples from the classes the model is trained on, while the unseen set contains samples from the classes the model was unknown to. This strategy assesses the model’s ability to generalize to unseen classes. In this study, the split configuration from [19] is used, and the authors emphasized that the choice of dataset split could significantly impact the results and should be selected carefully. The accuracy increases whenever

the classes from the test split overlap with the train classes. Therefore, the authors proposed a new split for the train and test splits alongside the standard for all datasets used in this study for ZSL. For 3D datasets, the ModelNet10 and ModelNet40 [20] datasets are combined, and the 30 classes from the ModelNet40 which are not present in the ModelNet10 are considered as seen classes, and unseen classes are considered from the test set of ModelNet10 following the settings from [2]. For the ScanObjectNN [39] dataset, we consider 26 classes from the ScanObjectNN dataset that are not overlapped with the ModelNet40 dataset as seen classes, and the overlapped classes are considered seen.

**Language models:** Several LLM models have been used in recent works. Considering the open access and the performance on the NLP domain, we have used the GPT-3.5 turbo<sup>1</sup> model, which is based on the autoregressive language model GPT-3 from [8]. The GPT-3.5 turbo model generates additional human-made like text descriptions for each class in datasets. This gives the advantage of having additional annotations of each class without captioning them manually. This description is used along with the available attributes and class names. We also use Word2vec [4] to generate the word vector representation of the text descriptions obtained through GPT-3.5. Word2vec learns word embeddings using a neural network from datasets containing billions of words and has millions of words in its vocabulary. And given a text input, it produces a set of vectors representing the words, and it can be used as input alongside the visual features. This gives the advantage of producing high-quality word vectors from the text descriptions using a pre-trained word2vec model to train ZSL models instead of using the attributes vectors produced by manual work alongside the visual features. This work uses the pre-trained word2vec model through Gensim

<sup>1</sup><https://platform.openai.com/docs/api-reference/chat>

[40] to produce the word vectors. This pre-trained model is trained on a partial Google News Dataset, which contains around 100 billion words. And the trained model contains 300-dimensional vectors for 3 million words and phrases. Passing the obtained class descriptions from GPT-3.5 produces a 300-dim. vector representation for each class.

**Evaluation:** We have used the top-1 accuracy metric to measure the performance of our models for ZSL, where only the unseen class’s information is used to predict an unseen class. For GZSL, the class name is predicted based on both seen and unseen class information. In addition to that, the Harmonic Mean (HM) [19] of the accuracy of seen and unseen classes in GZSL is also calculated.  $HM = \frac{2 \times acc_s \times acc_u}{acc_s + acc_u}$  where  $acc_s$  and  $acc_u$  are the top-1 accuracy of the seen and unseen classes, respectively. To understand the better-performing variations in all methods, we have used Borda Count (BC) [41]. Each of the variations is only assigned 1 point for ranking highest for each evaluation metric, and the total points are counted to find the best variation.

**Validation strategy:** To find the best hyperparameters, we have varied the learning rate, number of epochs, loss weighting factor, and value of  $\alpha$ . For the 2D datasets with embedding methods, we varied the number of epochs from 1000 to 3000 with the combination of learning rate ranging from 0.001 to 0.00001 and loss weighting factor from 0.01 to 0.001. For the 3D datasets and generative methods, we varied the number of epochs from 100 to 500 with a combination of learning rates ranging from 0.01 to 0.00001. We have varied the value of  $\alpha$  from 0.1 to 1.0 with the word vectors. For all the datasets,  $\alpha$  values are selected from the set  $\{0.1, 0.3, 0.5, 0.7, 1.0\}$  for our proposed method.

**Compared methods:** We compare our work with embedding (DEM [10], LATEM [12], SYNC [13], GDAN [14]) and generative (f-CLSWGAN [15], TF-VAEGAN [16] and CADA-VAE [17]) methods. These methods mostly reported results on manual attribute vectors for each class. Instead of attributes, we have used word vectors in this paper to eliminate manual annotation. We compare three different variations of experiments for each method: **(1)** Only Classname: Here, the word vector is generated as class semantics using the word2vec model using the class name only. This is the traditional way and serves as a baseline for the ZSL method. **(2)** Only ChatGPT: Instead of merely a class name, the text descriptions for each class generated by ChatGPT are used to calculate word vectors using the word2vec model. **(3)** Ours: This is our final proposal, where word vectors produced by **(1)** and **(2)** are linearly combined to fuse them and produce an improved word vector. Auxiliary supervision from ChatGPT helps to improve (noisy) only name-based descriptions.

**Implementation details<sup>2</sup>:** We have used ChatGPT to generate text descriptions through OpenAI’s API. The text descriptions, as well as the class names, are passed through a pre-trained word2vec model using Gensim [40] to generate word vectors from class names and ChatGPT text descriptions. The word

vectors are then combined to create our variation of the word vector. These word vector variations are used with the visual features extracted from the datasets’ 2D images and 3D point clouds. For the 2D images, ResNet101 is used to extract the visual features similar to [19], which produces visual features of 2048 dimensions. For the 3D point clouds, we used PointNet [42], PointConv [43], and EdgeConv [44] to extract visual features from the 3D point clouds where the first two produce visual features of 1024 dimensions and the later one of 2048 dimensions.

#### A. Performance on 2D Image datasets

The performances are measured for ZSL and GZSL settings and shown in Table II. The following observations can be made from the results. **(1)** For all the methods, the performance consistently improves while using the image features with the combined word vector of the class names and descriptions. **(2)** The results are better for the AWA datasets than those on CUB and 3D datasets. One reason is that the animal classes in the AWA dataset are more distinguishable in terms of visual and descriptive features. In contrast, the CUB dataset contains only bird classes with similar physical and descriptive features. So describing the birds in natural language will be more similar in words, and it will also be reflected in the word vectors generated from the descriptions. **(3)** Using only ChatGPT-based word vectors from the description produces better results than only class name-based ones. For AWA datasets, using word vectors for descriptions gives better ZSL results than the other two variations. But overall, our proposed combined word vector of class name and description holds top BC score on both datasets over all methods. It tells us that our proposed method can improve the quality of class semantics for ZSL and GZSL.

#### B. Performance on 3D Point cloud datasets

We extract 3D point cloud features from the PointConv backbone and apply and align them with the word vectors generated using the Word2vec model. Table III we report 3D ZSL results. Here are our observations: **(1)** Our proposed combined word vector of class name and description improves performance across all the methods. For both ZSL and GZSL settings, accuracy is higher than using word vectors of only labels or only ChatGPT descriptions. The BC score of our combined vector of class name and description also affirms that it improves performance for all the methods. **(2)** For embedding-based methods like DEM and GDAN, the performance of only ChatGPT description word vector is better than the word vector of only the class name. The combined word vector of the class name and description beats the previous ones. For the generative models like TF-VAEGAN and CADA-VAE, the performance with the word vector of the class name is better than with the word vector of the ChatGPT description. And again, the performance of our proposed combined word vector is better than the previous ones. **(3)** The methods used for 2D ZSL can be used for 3D point clouds by using the extracted features of the 3D point cloud as input in place of

<sup>2</sup>Codes and models are available in <https://github.com/FHShubho/CGS-ZSL>

Method	Variations	ModelNet40					ScanObjectNN				
		ZSL	GZSL				ZSL	GZSL			
		Acc	Acc <sub>s</sub>	Acc <sub>u</sub>	HM	BC	Acc	Acc <sub>s</sub>	Acc <sub>u</sub>	HM	BC
DEM [10]	Only Class Name	17.33	<b>87.69</b>	6.35	11.84	1	17.99	88.88	12.32	21.64	0
	Only ChatGPT	<b>27.55</b>	77.20	8.24	14.88	1	13.93	88.20	4.61	8.76	0
	Ours	22.57	87.43	<b>10.30</b>	<b>18.43</b>	<b>2</b>	<b>18.48</b>	<b>89.48</b>	<b>14.20</b>	<b>24.51</b>	<b>4</b>
LATEM [12]	Only Class Name	26.29	-	-	-	0	11.88	-	-	-	0
	Only ChatGPT	27.50	-	-	-	0	10.61	-	-	-	0
	Ours	<b>28.11</b>	-	-	-	<b>1</b>	<b>13.42</b>	-	-	-	<b>1</b>
SYNC [13]	Only Class Name	21.17	-	-	-	0	17.43	-	-	-	0
	Only ChatGPT	20.22	-	-	-	0	15.90	-	-	-	0
	Ours	<b>25.79</b>	-	-	-	<b>1</b>	<b>18.35</b>	-	-	-	<b>1</b>
GDAN [14]	Only Class Name	-	86.80	<b>2.70</b>	<b>5.23</b>	<b>2</b>	-	<b>88.93</b>	16.60	27.97	1
	Only ChatGPT	-	86.87	2.25	4.38	0	-	88.18	17.33	28.97	0
	Ours	-	<b>87.07</b>	2.36	4.60	1	-	<b>88.50</b>	<b>20.31</b>	<b>33.04</b>	<b>2</b>
f-CLSWGAN [15]	Only Class Name	18.15	88.83	1.42	2.79	0	22.51	89.17	11.83	20.91	0
	Only ChatGPT	19.11	<b>89.33</b>	1.10	2.17	1	17.12	89.26	11.58	20.50	0
	Ours	<b>26.77</b>	88.73	<b>5.90</b>	<b>11.06</b>	<b>3</b>	<b>22.84</b>	<b>89.85</b>	<b>13.01</b>	<b>22.71</b>	<b>4</b>
TF-VAEGAN [16]	Only Class Name	28.13	67.17	20.11	30.95	0	28.17	75.85	25.70	38.40	0
	Only ChatGPT	24.34	<b>76.57</b>	18.71	30.07	1	25.86	<b>86.10</b>	18.95	31.06	1
	Ours	<b>30.37</b>	73.47	<b>25.45</b>	<b>37.81</b>	<b>3</b>	<b>28.76</b>	78.60	<b>26.54</b>	<b>39.68</b>	<b>3</b>
CADA-VAE [17]	Only Class Name	21.50	88.40	5.71	10.73	0	-	89.54	18.18	30.59	0
	Only ChatGPT	15.61	87.03	8.21	15.01	0	-	89.17	15.43	26.31	0
	Ours	<b>21.74</b>	<b>88.93</b>	<b>13.16</b>	<b>22.89</b>	<b>4</b>	-	<b>89.92</b>	<b>18.68</b>	<b>30.93</b>	<b>3</b>

TABLE III: ZSL and GZSL result on 3D image datasets using PointConv backbone.

Backbone	Variations	ModelNet40					ScanObjectNN				
		ZSL	GZSL				ZSL	GZSL			
		Acc	Acc <sub>s</sub>	Acc <sub>u</sub>	HM	BC	Acc	Acc <sub>s</sub>	Acc <sub>u</sub>	HM	BC
PointNet	Only Class Name	18.51	<b>56.8</b>	15.13	23.90	1	20.48	18.06	<b>16.80</b>	17.41	1
	Only ChatGPT	23.50	46.63	<b>20.80</b>	<b>28.75</b>	<b>2</b>	19.78	17.92	15.80	16.80	0
	Ours	<b>25.64</b>	41.53	14.44	21.43	1	<b>21.57</b>	<b>18.77</b>	16.25	<b>17.42</b>	<b>3</b>
EdgeConv	Only Class Name	29.00	44.47	29.53	35.49	0	23.20	<b>76.03</b>	11.40	19.87	1
	Only ChatGPT	23.17	59.93	20.05	30.05	0	<b>25.41</b>	75.95	12.06	20.82	1
	Ours	<b>37.86</b>	<b>71.33</b>	<b>29.62</b>	<b>41.86</b>	<b>4</b>	24.71	75.35	<b>12.79</b>	<b>21.86</b>	<b>2</b>
PointConv	Only Class Name	28.13	67.17	20.11	30.95	0	28.17	75.85	25.70	38.40	0
	Only ChatGPT	24.34	<b>76.57</b>	18.71	30.95	1	25.86	<b>86.10</b>	18.95	31.06	1
	Ours	<b>30.37</b>	73.47	<b>25.45</b>	<b>37.81</b>	<b>3</b>	<b>28.76</b>	78.60	<b>26.54</b>	<b>39.68</b>	<b>3</b>

TABLE IV: ZSL and GZSL performance of TF-VAEGAN [16] method using our word vector and different backbones.

the features extracted from the images of the 2D dataset. And the results show that the methods from 2D ZSL also perform well for the 3D point cloud datasets. (4) Using our proposed combined word vector representation consistently improved the ZSL performance and the unseen class in the GZSL setting. (5) For the ScanObjectNN dataset case, the training is done on synthetic ModelNet40 objects, but testing is done on real-world-scanned (not synthetic) ScanObjectNN objects. This is a complex setup because the domain of train and testing is different. Working consistently in this challenging setting proves the robustness of our approach.

### C. Ablation studies

**Different 3D backbones:** We have experimented with three different backbones for 3D datasets, PointNet [42], PointConv [43], and EdgeConv [44]. These three backbones extract point cloud features from the 3D object, later used as the input with all the variants of the word vector representations. As shown in Table IV, the features extracted through the PointConv backbone yield better performance than PointNet and PointConv backbones. We experimented with all the backbones for all seven methods and observed that the features from PointConv showed better performance. We notice that using the word

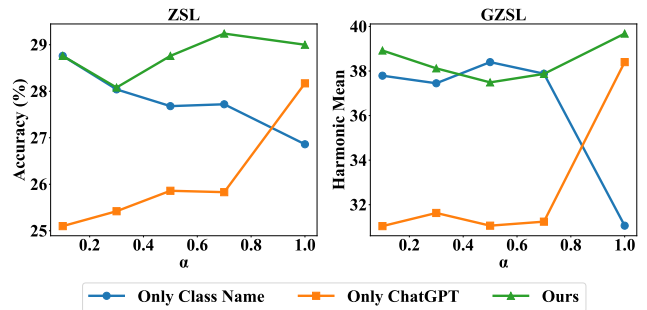


Fig. 4: ZSL and GZSL for different values of  $\alpha$  for TF-VAEGAN [16] on ScanObjectNN dataset

vectors from ChatGPT and our representation of word vectors improves performance consistently across all the backbones.

**Varying  $\alpha$ :** While combining the word vectors generated using the class name and text description from ChatGPT, we varied the value of  $\alpha$ . We experimented with the  $\alpha$  values of 0.1, 0.3, 0.5, 0.7, and 1.0. Fig. 4 shows the accuracy of ZSL and Harmonic Mean in the GZSL setting over varying values of  $\alpha$ . We can observe that increasing the value of  $\alpha$  yield performance improvements for ZSL and GZSL with our vector

representation. We can also observe that the word vector of the text representation from the ChatGPT also increases with the increasing values of  $\alpha$ . Still, the performance with the word vector of only class name decreases as the value of  $\alpha$  increases.

## V. CONCLUSION

This paper investigates the possible use of ChatGPT to improve class semantics for zero-shot learning tasks. ChatGPT-based word vectors fused with traditional class name-based vectors can achieve better ZSL and GZSL performance on 2D image and 3D point cloud recognition tasks. In experiments, we apply our proposal to multiple embedding-based and generative ZSL methods and notice that our technique consistently improves those methods' existing performance. We perform extensive experiments on 2D image and 3D point cloud datasets. It tells that ChatGPT could be a suitable annotation tool that can provide automatic and less noisy information without manual labor. In the future, we will investigate the use of ChatGPT-guided embedding vectors for object detection and segmentation tasks.

## REFERENCES

- [1] A. Cheraghian, S. Rahman, and L. Petersson, "Zero-shot learning of 3d point cloud objects," in *MVA*, 2019.
- [2] A. Cheraghian, S. Rahman, T. F. Chowdhury, D. Campbell, and L. Petersson, "Zero-shot learning on 3d point cloud objects and beyond," *IJCV*, 2022.
- [3] R. Zhang, Z. Guo, W. Zhang, K. Li, X. Miao, B. Cui, Y. Qiao, P. Gao, and H. Li, "Pointclip: Point cloud understanding by clip," *arXiv preprint arXiv:2112.02413*, 2021.
- [4] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," 2013.
- [5] J. Pennington, R. Socher, and C. Manning, "GloVe: Global vectors for word representation," in *EMNLP*, Oct. 2014.
- [6] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, "Enriching word vectors with subword information," *arXiv preprint arXiv:1607.04606*, 2016.
- [7] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," 2019.
- [8] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *NeurIPS*, vol. 33, 2020.
- [9] S. Rahman, S. Khan, and N. Barnes, "Polarity loss: Improving visual-semantic alignment for zero-shot detection," *IEEE TNNLS*, 2022.
- [10] L. Zhang, T. Xiang, and S. Gong, "Learning a deep embedding model for zero-shot learning," in *CVPR*, 2017.
- [11] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.
- [12] Y. Xian, Z. Akata, G. Sharma, Q. Nguyen, M. Hein, and B. Schiele, "Latent embeddings for zero-shot classification," in *CVPR*, 2016.
- [13] S. Changpinyo, W.-L. Chao, B. Gong, and F. Sha, "Synthesized classifiers for zero-shot learning," in *CVPR*, 2016.
- [14] H. Huang, C. Wang, P. S. Yu, and C.-D. Wang, "Generative dual adversarial network for generalized zero-shot learning," in *CVPR*, 2019.
- [15] Y. Xian, T. Lorenz, B. Schiele, and Z. Akata, "Feature generating networks for zero-shot learning," in *CVPR*, 2018.
- [16] S. Narayan, A. Gupta, F. S. Khan, C. G. M. Snoek, and L. Shao, "Latent embedding feedback and discriminative features for zero-shot classification," in *ECCV*, August 2020.
- [17] E. Schonfeld, S. Ebrahimi, S. Sinha, T. Darrell, and Z. Akata, "Generalized zero and few-shot learning via aligned variational autoencoders," in *CVPR*, 2019.
- [18] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "Caltech-ucsd birds-200-2011 (cub-200-2011)," California Institute of Technology, Tech. Rep. CNS-TR-2011-001, 2011.
- [19] Y. Xian, C. H. Lampert, B. Schiele, and Z. Akata, "Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly," *IEEE TPAMI*, vol. 41, 2019.
- [20] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *CVPR*, 2015.
- [21] M. A. Uy, Q.-H. Pham, B.-S. Hua, T. Nguyen, and S.-K. Yeung, "Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data," in *ICCV*, 2019.
- [22] C. H. Lampert, H. Nickisch, and S. Harmeling, "Attribute-based classification for zero-shot visual object categorization," *IEEE TPAMI*, vol. 36, no. 3, 2014.
- [23] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, "Devise: A deep visual-semantic embedding model," *NIPS*, vol. 26, 2013.
- [24] Z. Akata, S. Reed, D. Walter, H. Lee, and B. Schiele, "Evaluation of output embeddings for fine-grained image classification," in *CVPR*, 2015.
- [25] E. Kodirov, T. Xiang, Z. Fu, and S. Gong, "Unsupervised domain adaptation for zero-shot learning," in *ICCV*, December 2015.
- [26] A. Mishra, S. Krishna Reddy, A. Mittal, and H. A. Murthy, "A generative model for zero shot learning using conditional variational autoencoders," in *CVPR*, 2018.
- [27] M. R. Vyas, H. Venkateswara, and S. Panchanathan, "Leveraging seen and unseen semantic relationships for generative zero-shot learning," in *ECCV*, 2020.
- [28] N. Hayat, M. Hayat, S. Rahman, S. Khan, S. W. Zamir, and F. S. Khan, "Synthesizing the unseen for zero-shot object detection," in *ACCV*, 2020.
- [29] A. Cheraghian, S. Rahman, and L. Petersson, "Zero-shot learning of 3d point cloud objects," in *MVA*. IEEE, 2019.
- [30] A. Cheraghian, S. Rahman, D. Campbell, and L. Petersson, "Mitigating the hubness problem for zero-shot learning of 3d objects," in *BMVC*, 2019.
- [31] B. Michele, A. Boulch, G. Puy, M. Bucher, and R. Marlet, "Generative zero-shot learning for semantic segmentation of 3d point clouds," in *3DV*, 2021.
- [32] M. Nasiri, A. Cheraghian, T. F. Chowdhury, S. Ahmadi, M. Saberi, and S. Rahman, "Prompt-guided scene generation for 3d zero-shot learning," in *DICTA*, 2022.
- [33] Y. Zhao, I. Misra, P. Krähenbühl, and R. Girdhar, "Learning video representations from large language models," in *CVPR*, June 2023.
- [34] X. He, Z. Lin, Y. Gong, A.-L. Jin, H. Zhang, C. Lin, J. Jiao, S. M. Yiu, N. Duan, and W. Chen, "Anollm: Making large language models to be better crowdsourced annotators," 2023.
- [35] D. Yu, L. Li, H. Su, and M. Fuoli, "Using llm-assisted annotation for corpus linguistics: A case study of local grammar analysis," 2023.
- [36] C. Whitehouse, M. Choudhury, and A. F. Aji, "Llm-powered data augmentation for enhanced crosslingual performance," 2023.
- [37] Y. Meng, J. Huang, Y. Zhang, and J. Han, "Generating training data with language models: Towards zero-shot language understanding," in *NeurIPS*, vol. 35, 2022.
- [38] E. Kodirov, T. Xiang, Z. Fu, and S. Gong, "Unsupervised domain adaptation for zero-shot learning," in *ICCV*, Dec 2015.
- [39] M. A. Uy, Q.-H. Pham, B.-S. Hua, D. T. Nguyen, and S.-K. Yeung, "Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data," in *ICCV*, 2019.
- [40] R. Řehůřek and P. Sojka, "Software Framework for Topic Modelling with Large Corpora," in *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, May 2010.
- [41] D. G. Saari, *Basic geometry of voting*. Springer Science & Business Media, 1995, vol. 12.
- [42] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," 2017.
- [43] W. Wu, Z. Qi, and L. Fuxin, "Pointconv: Deep convolutional networks on 3d point clouds," in *CVPR*, 2019.
- [44] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *arXiv preprint arXiv:1801.07829*, 2018.