# ZSCRGAN: A GAN-based Expectation Maximization Model for Zero-Shot Retrieval of Images from Textual Descriptions*

Anurag Roy
Department of CSE, IIT Kharagpur, India

Vinay Kumar Verma
Department of ECE, Duke University, USA

Kripabandhu Ghosh
Department of CSA, IISER Kolkata, India

Saptarshi Ghosh
Department of CSE, IIT Kharagpur, India

## ABSTRACT

Most existing algorithms for cross-modal Information Retrieval are based on a supervised train-test setup, where a model learns to align the mode of the query (e.g., text) to the mode of the documents (e.g., images) from a given training set. Such a setup assumes that the training set contains an exhaustive representation of all possible classes of queries. In reality, a retrieval model may need to be deployed on previously unseen classes, which implies a *zero-shot* IR setup. In this paper, we propose a novel GAN-based model for zero-shot text to image retrieval. When given a textual description as the query, our model can retrieve relevant images in a zero-shot setup. The proposed model is trained using an Expectation-Maximization framework. Experiments on multiple benchmark datasets show that our proposed model comfortably outperforms several state-of-the-art zero-shot text to image retrieval models, as well as zero-shot classification and hashing models suitably used for retrieval.

## KEYWORDS

Zero-shot; Text to image retrieval; GAN; E-M

## 1 INTRODUCTION

Today, information is generated in several modes, e.g., text, image, audio, video, etc. Thus, for a query in one mode (e.g., text), the relevant information may be present in a different mode (e.g., image). Cross-modal Information Retrieval (IR) algorithms are being developed to cater to such search requirements.

**Need for Zero-shot Information Retrieval (ZSIR):** A train-test setup of an IR task comprises parameter learning for various classes / categories of queries. Standard cross-modal retrieval methods require training data of all classes of queries to train the retrieval models. But such methods can fail to retrieve data for queries of *new* or *unseen* classes. For instance, suppose the retrieval model has been trained on images and textual descriptions of various classes of vehicles, such as 'car', 'motorbike', 'aeroplane', and so on. Now, given a query 'bus', the model is expected to retrieve images and textual descriptions of buses (for which the model has not been trained). Such a situation conforms to the "zero-shot" setup [14, 17, 20] which focuses on recognizing new/unseen classes with limited training classes.

Such situations are relevant in any modern-day search system, where new events, hashtags, etc. emerge every day. So, contrary to the conventional IR evaluation setup, the zero-shot paradigm needs to be incorporated in an IR setting. Specifically, zero-shot

cross-media retrieval intends to achieve retrieval across multiple modes (e.g., images to be retrieved in response to a textual query) where there is no overlap between the query-classes in training and test data. Zero-Shot IR (ZSIR) is especially challenging since models need to handle not only different semantics across seen and unseen query-classes, but also the heterogeneous features of data across different modes.

**Present work and differences with prior works:** Though lot of research has been reported on general multimodal and cross-modal retrieval [26], to our knowledge, only a few prior works have attempted cross-modal IR in a zero-shot setting [2, 3, 11, 16, 22, 33] (see Section 2 and Section 4.3 for details of these methods). Some of these prior works assume additional information about the class labels (e.g., a measure of semantic similarity between class labels) whcih may not always be available.

In this paper, we propose a novel model for cross-modal IR in zero-shot setting, based on Conditional Generative Adversarial Networks (GANs) [18], that can retrieve *images* relevant to a given *textual* query. Our model – which we name **ZSCRGAN** (Zero-Shot Cross-modal Retrieval GAN) – relies only on the textual data to perform the retrieval, and *does not need additional information about class labels*. Though prior ZSIR models [2, 3, 33] also use GANs, the main novel contributions of the proposed model can be summarized as follows: (1) We propose **use of wrong classes** to enable the generator to generate features that are unique to a specific class, by distinguishing it from other classes. (2) We develop a **Common Space Embedding Mapper (CSEM)** to map both the image embeddings and the text embeddings to a common space where retrieval can be performed. *This is the key step that enables our model to perform retrieval without relying on additional semantic information of class labels.* (3) We develop an **Expectation-Maximization (E-M) based method** for efficiently training the retrieval model, where the GAN and the CSEM are trained alternately. We show that this E-M setup enables better retrieval than jointly training the GAN and the CSEM.

We experiment on several benchmark datasets for zero-shot retrieval – (1) the Caltech-UCSD Birds dataset, (2) the Oxford Flowers-102 dataset, (3) the North America Birds (NAB) dataset, (4) the Wikipedia dataset, and (5) the Animals With Attribute (AWA) dataset. Our proposed ZSCRGAN comfortably out-performs several strong and varied baselines on all the datasets, including ZSIR models [3, 22, 33], state-of-the-art Zero-Shot classification models [25, 28] suitably adapted for the Text-to-Image retrieval setting, as well as state-of-the-art Zero-Shot Hashing models [11, 16, 30]. Also note that our proposed model can be used not only with

textual queries, but also with other forms of queries such as attribute vectors, as demonstrated by its application on the AWA dataset. We make the implementation of ZSCRGAN publicly available at https://github.com/ranarag/ZSCRGAN.

## 2 RELATED WORK

There are many works on multimodal retrieval (see [26] for a survey). However, most of these works are *not* in zero-shot setup on which we focus in this paper.

**Zero-Shot Learning (ZSL):** The initial models for ZSL mostly focused on learning a similarity metric in the joint attribute space or feature space [23, 32]. With the recent advancements of the generative model [9, 12], models based on Variational Autoencoders (VAE) [25] and GAN-based [28] approaches have attained the state-of-the-art results for ZSL.

**Multimodal IR in Zero-Shot setup (ZSIR):** There have been several recent works on multimodal ZSIR. For instance, some works have attempted zero-shot *sketch-based* image retrieval [5, 7, 13], most of which use GANs to perform the retrieval. Note that sketch-based image retrieval is different from the text-to-image retrieval that we consider in this paper – while the input query is a sketch in the former, we assume the query to be a textual description (or its vector representation).

There have also been works on zero-shot text-to-image retrieval. Reed *et al.* [22] minimized empirical risk function to create a joint embedding for both text and images. Chi *et al.* proposed two very similar models [3] and [2]. The DADN model developed in [3] was shown to out-perform the one in [2]. These models adopt a dual GAN approach where a forward GAN is used to project the embeddings to a semantic space, and a reverse GAN is used to reconstruct the semantic space embeddings to the original embeddings. Zhu *et al.* [33] developed a GAN-based model named ZSL-GAN for retrieving images from textual queries, where the GAN is trained with classification loss as a regularizer. Additionally, several Zero-Shot Hashing models have been developed [11, 15, 16, 30] that can also be used for ZS text-to-image retrieval.

All the above-mentioned prior models for text-to-image ZSIR are considered as baselines in this work. Further details of the prior models are given in Section 4, where the primary differences of these models with our proposed model are also explained.

## 3 PROPOSED APPROACH

We start by formally defining the zero-shot text to image retrieval problem, and then describe our proposed model ZSCRGAN (implementation available at https://github.com/ranarag/ZSCRGAN). Table 1 gives the notations used in this paper.

### 3.1 Problem Definition

In the zero-shot Text to Image retrieval setup, we consider training data $\mathcal{D} = \{I_r^k, \varphi_{t_r}^k, y^k\}_{k=1}^K$ with $K$ samples. $I_r^k$ is the real image embedding of the $k^{th}$ image. $\varphi_{t_r}^k$ is the real text embedding of the text accompanying the $k^{th}$ image. $y^k \in \mathcal{Y}$ is a class label, and $\mathcal{Y} = \{1, 2, \ldots, S\}$ is the set of *seen* classes, where $S$ is the number of seen classes (only seen classes are available at training time).

| | Description | | | Description |
|---|---|---|---|---|
| $G()$ | Generator | | $D()$ | Discriminator |
| $L_D$ | Discriminator Loss Function | | $\mathcal{L}_G$ | Generator Loss Function |
| $\varphi_t$ | Text embedding of unseen class | | $I_i$ | Image embedding of unseen class |
| $\theta_t$ | Common Space embedding generated from $\varphi_t$ | | $\theta_i$ | Common space embedding generated from $I_i$ |
| $I_r$ | Real Image Embedding | | $I_w$ | Wrong Image Embedding |
| $i_r$ | Real representative Embedding | | $i_w$ | Wrong representative Embedding |
| $\varphi_{t_r}$ | Real Text Embedding | | $\varphi_{t_w}$ | Wrong Text Embedding |
| $\hat{c_{t_r}}$ | Real Latent Embedding | | $\hat{c_{t_w}}$ | Wrong Latent Embedding |
| $\psi_C$ | trainable parameters of CSEM | | $\psi_G$ | trainable parameters of Generator |
| $\theta_{t_r}$ | Common Space embedding generated from $\varphi_{t_r}$ | | $\theta_{t_w}$ | Common space embedding generated from $\varphi_{t_w}$ |
| $\psi_C$ | Trainable parameters of CSEM | | $\psi_G$ | Trainable parameters of Generator |
| $\hat{c_t}$ | Latent Embedding Generated from $\varphi_t$ | | z | noise vector sampled from a normal distribution $p_z$ |

**Table 1: Notations used in this paper.**

Let $\mathcal{U} = \{1, 2, \ldots, U\}$ be the set of unseen classes, where $U$ is the number of unseen classes (not available at training time). For each unseen class query $u \in \mathcal{U}$ a relevant set of images are present that we have to retrieve. At test-time, for each unseen class $u \in \mathcal{U}$, we use a textual embedding $\varphi_t$ from $u$ as query. Textual embedding $\varphi_t$ and unseen class image are projected into joint space to perform the retrieval. In the zero-shot setup $\mathcal{U} \cap \mathcal{Y} = \Phi$, i.e. training and test classes are disjoint.
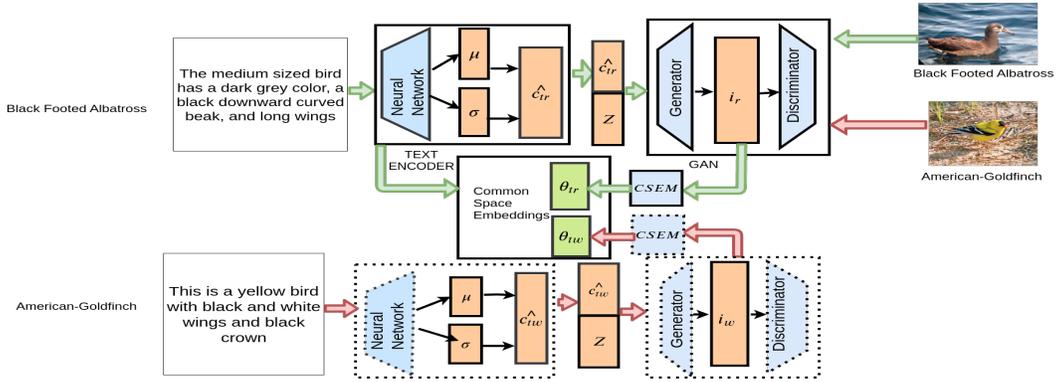
### 3.2 Overview of our approach

One of the main challenges of zero-shot cross-modal retrieval is the *representation problem* of *novel (unseen) class* data. The text and image are of different modalities and the challenge is to represent them in a common embedding space. We train a model to perform this mapping. However training such a model becomes difficult when there is a huge overlap among embeddings of different classes, e.g., there may be cases where the images of one class have high overlap with images from other classes, making it difficult to map it to a common embedding space. As part of addressing this challenge, we use a Generative Adversarial Network (GAN) [9] to generate a per class 'representative embedding' for all image embeddings of a particular class. This generated embedding is desired to have *high* similarity with image embeddings from the same class and *low* similarity with image embeddings from other classes.

We chose a generative model rather than a discriminative one for this purpose, because generative models are most robust to visual imagination, and this helps the model to map the images from unseen classes more robustly. Two popular generative models are Variational Autoencoders (VAE) [12] and Generative Adversarial Networks (GAN) [9]. Samples from VAE models tend to be blurry (i.e., with less information content) as compared to GANs, because of the probability mass diffusion over the data space [24]. Therefore, we found GANs to be the most suitable for this problem.

While there have been prior works using GANs for zero-shot retrieval tasks [2, 3, 22], they rely on class labels for training. Some of the prior models [2, 3] make use of word-embeddings of class labels as class level information. However, in many cases, class label information may not be available. In contrast, our proposed model does *not need* class labels to perform its task. This is achieved by learning to map image embeddings and text embeddings to a common embedding space.

Our proposed model ZSCRGAN (shown in Figure 1) works as follows. We take two text embeddings, one belonging to class $y \in \mathcal{Y}$ and the other belonging to some other class $\hat{y} \in \mathcal{Y}$. We pass the two text embeddings through a Text Encoder (TE) which generates

**Figure 1: [color online] The proposed `ZSCRGAN` architecture used for generating a common space embedding for text and image. The example here demonstrates training of the architecture for learning a common space embedding for the class 'Black Footed Albatross', using textual descriptions and images of that class. To this end, textual descriptions and images of some *other* class (here 'American-Goldfinch') are used. We refer to this other class as the 'wrong' class. The blocks in dashed lines (through which textual embeddings of the wrong class are passed) are identical copies of the corresponding blocks in solid line (through which textual embeddings of the correct class are passed).**

(i) a latent embedding $\hat{c_{tr}}$ (a.k.a real text embedding) for class $y$ and (ii) $\hat{c_{tw}}$ (a.k.a wrong text embedding) for class $\hat{y}$. For each class, we generate a per-class representative embedding $i$ for all images of that class. We train the Text Encoder jointly with the GAN. We use the representative embeddings to train a *Common Space Embedding Mapper* (CSEM) which learns to map all the image and text embeddings to a common space where these common space embeddings will have a high similarity if they belong to the same class.

We formulate an E-M paradigm in which we train the CSEM and the GAN alternatively. In a latter section we also justify our choice of such a paradigm by comparing the model performance with this E-M formulation and without it (i.e., on jointly training the CSEM and GAN).

## 3.3 Details of our proposed model `ZSCRGAN`

Let $Q(I, \phi)$ be the joint probability distribution denoting the probability of text embeddings $\phi$ and relevant image embeddings $I$. Maximizing this probability is expected to ensure high similarity between $\phi$ and relevant image embeddings $I$ therefore leading to better performance of the retrieval model. We plan to do this maximization using a machine learning model having $\psi_C$ as its parameters. Hence, our aim will be to maximize the probability distribution $Q(I, \phi | \psi_C)$. For simplicity, we will maximize the log of this distribution $\log Q(I, \phi | \psi_C)$. .Let $\psi_C$ be the random variable representing the values that can be taken by the trainable parameters of the neural network model. Thus, our log probability function becomes $\log Q(I, \phi | \psi_C)$. However, when we try to maximize $\log Q(I, \phi | \psi_C)$ directly using a machine learning model , we see very poor retrieval performance. The reason being, images from one class are very similar to images from another class . For example, *'Crested Auklet'* and *'Rhinoceros Auklet'* in the CUB dataset are two very similar looking birds and are almost indistinguishable to the human eye. Due to these overlaps among images from classes, training the neural network model becomes difficult, as it encounters very similar positive and negative examples during training. Thus , we introduce a latent variable $I'$ – an unobserved

embedding which would be a representative for the images of a class. The embedding will have high similarity with the images of a particular class and very low similarity with images from all other classes. Using these representative embeddings instead of the actual image embeddings will solve the training problem of our model . We adopt an Expectation-Maximization (E-M) framework to perform the maximization in presence of the latent variable.

*3.3.1* **E-M formulation:** As stated above, our objective is to maximize the following expression:

$$\log Q(I, \phi | \psi_C) = \log \sum_{i \in I'} Q(i, I, \phi | \psi_C) \tag{1}$$

Let $P(I' = i | I, \phi, \psi_G)$ denote the probability of generating $i$ given $I, \phi$ and $\psi_G$. Here is $\psi_G$ are the trainable parameters of the model used to generate $i$. Thus we have:

$$\log \sum_{i \in I'} Q(i, I, \phi | \psi_C) = \log \sum_{i \in I'} \frac{P(i | I, \phi, \psi_G) \cdot Q(i, I, \phi | \psi_C)}{P(i | I, \phi, \psi_G)}$$

$$\geq \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log \frac{Q(i, I, \phi | \psi_C)}{P(i | I, \phi, \psi_G)} \quad \text{by Jensen's Inequality}$$

$$= \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log Q(i, I, \phi | \psi_C) - \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log P(i | I, \phi, \psi_G)$$

$$= \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot [\log Q(i, \phi | I, \psi_C) + \log Q(I | \psi_C)]$$

$$- \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log P(i | I, \phi, \psi_G)$$

since $I$ is conditionally independent of $i$ and $\phi$ given $\psi_C$

$$= \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log Q(i, \phi | \psi_C) + \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log Q(I | \psi_C)$$

$$- \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log P(i | I, \phi, \psi_G)$$

$$= \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log Q(i, \phi | \psi_C) + 1 \cdot \log Q(I)$$

$$- \sum_{i \in I'} P(i | I, \phi, \psi_G) \cdot \log P(i | I, \phi, \psi_G)$$

$$\tag{2}$$

where the last step holds since the distribution of $I$ is independent of $\psi_C$. Now, after applying Jensen's Inequality Eqn. 2 implies the following:

$$\log \sum_{i \in I'} Q(i, I, \phi | \psi_C) - \log Q(I)$$
$$\geq \sum_{i \in I'} P(i|I, \phi, \psi_G) \cdot \log Q(i, \phi|\psi_C) - \sum_{i \in I'} P(i|I, \phi, \psi_G) \cdot \log P(i|I, \phi, \psi_G)$$
$$= F(P, \psi_C) \tag{3}$$

$Q(i, \phi|\psi_C)$ denotes the joint probability of $i$ and $\phi$ being similar. Thus, $F(P, \psi_C)$ is a lower bound for $\log Q(I, \phi|\psi_C)$, and maximizing the lower bound will ensure a high minimum value of $\log Q(I, \phi|\psi_C)$.[1] We train a neural network architecture using the E-M algorithm to maximize $F(P, \psi_C)$ where the E and M steps at iteration $it$ are:

$$\text{E-step}: \quad P^{it}(I' = i|I, \phi) = \underset{P}{\operatorname{argmax}} F(P, \psi_C^{it-1}) \tag{4}$$

$$\text{M-step}: \quad \psi_C^{it} = \underset{\psi_C}{\operatorname{argmax}} F(P^{it}, \psi_C) \tag{5}$$

So, the challenge now is to maximize $F(P, \psi_C)$. To this end, we propose to use neural networks as follows.

3.3.2 **Using neural networks to approximate** $F(P, \psi_C)$**:** Eqn. 2 can be re-written as:

$$F(P, \psi_C) = \mathbb{E}_{i \sim P(I'=i|I, \phi, \psi_G)}[\log Q(i, \phi|\psi_C)]$$
$$- \mathbb{E}_{i \sim P(I'=i|I, \phi, \psi_G)}[\log P(i|I, \phi, \psi_G)] \tag{6}$$

We take the help of two neural networks to approximate the two parts of the function $F$ as shown in Eqn. 6 – (1) Common Space Embedding Mapper (CSEM) to represent the first term, and (2) GAN to represent the second term in Eqn. 6.

**Common Space Embedding Mapper (CSEM):** This module is a feed-forward neural network trained to maximize the probability $Q(i, \phi|\psi_C)$ which denotes the joint probability of $i$ and $\phi$ being similar. We define such a $Q$ as:

$$Q(i, \phi|\psi_C) = \frac{e^{v_p}}{e^{v_p} + e^{v_n}} \tag{7}$$

where $v_p$ and $v_n$ are scores calculated as:

$$v_p = CosineSim(CSEM(i), \hat{c_{tr}})$$
$$v_n = CosineSim(CSEM(G(z, \hat{c_{tw}})), \hat{c_{tr}}) \tag{8}$$

The CSEM is trained using the cost function Triplet Loss [6] $\mathcal{L}_T$, which can be written as:

$$\mathcal{L}_T = -\mathbb{E}_{(i \sim P(I'=i|I, \phi, \psi_G), \hat{c_{tr}} \sim TE(\phi)}[\log Q(i, \phi|\psi_C)]$$
$$= \mathbb{E}_{(i \sim P(I'=i|I, \phi, \psi_G), (\hat{c_{tr}}, \hat{c_{tw}}) \sim TE(\phi)}[-\log \frac{e^{v_p}}{e^{v_p} + e^{v_n}}] \tag{9}$$
$$= \mathbb{E}_{(i \sim P(I'=i|I, \phi, \psi_G), (\hat{c_{tr}}, \hat{c_{tw}}) \sim TE(\phi)}[\log(1 + e^{v_n - v_p})]$$

We call this module the *Common Space Embedding Mapper* because it learns to map the image embeddings to a space where the resulting embeddings will have high cosine similarity among themselves if they are from the same class and low cosine similarity with images from different classes. We train it by using the triplet loss $\mathcal{L}_T$ considering $\hat{c_{tr}}$ as the pivot, $i$ as the positive example and $G(z, \hat{c_{tw}})$ as the negative example. Here, $\hat{c_{tr}}$ and $\hat{c_{tw}}$ are generated by $TE(\phi_{tr})$

---

[1]The function $F(P, \psi_C)$ is related to the Free Energy Principle that is used to bound the 'surprise' on sampling some data, given a generative model (see https://en.wikipedia.org/wiki/Free_energy_principle).

and $TE(\phi_{tw})$ respectively, $TE()$ being the Text Encoder (described in Section 3.3.2).

**GAN based Learning:** $\mathbb{E}_{i \sim P(I'=i|I, \phi)}[\log P(i|I, \phi)]$ is calculated using a Generative Adversarial Network. Generative adversarial networks (GAN) [9] are composed of two neural network models – the *Generator* and the *Discriminator*, which are trained to compete with each other. The generator (G) is trained to generate samples from a data distribution, $p_g$ that are difficult to distinguish from real samples for the discriminator. Meanwhile, the discriminator (D) is trained to differentiate between real samples sampled from the true data distribution $p_{data}$ and images generated by G. The GAN training is often very unstable; To stabilize the GAN training, Arjovsky et al. [1] proposed Wasserstein-GAN (WGAN). WGAN optimizes the Wasserstein (a.k.a Earth Mover's) distance between $p_g$ and $p_{data}$. The loss function to optimize in WGAN is:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}, z \sim p_g}[D(x) - D(G(z)))]$$
$$s.t. \|D_w\| \leq k \tag{10}$$

where $D_w$ indicates all the parameters of the discriminator and $k$ is a constant.

The WGAN architecture does not have control over the data generation of a particular class, which is a requirement in our application. Hence, to enable generation of data of a particular class, we use the *Conditional-WGAN* [18] (CWGAN), where the WGAN is conditioned on some latent embedding. The proposed model uses the CWGAN as the base architecture, conditioned on the latent text embedding (which helps to achieve robust sample generation).

The discriminator and generator losses are as follows:

$$\mathcal{L}_D = \frac{1}{2}( \mathbb{E}_{(I_w, \varphi_{tr}) \sim p_{data}}[D(I_w, \varphi_{tr})]$$
$$+ \mathbb{E}_{z \sim p_z, \varphi_{tr} \sim p_{data}, \hat{c_{ir}} \sim TE(\phi)}[D(G(z, \hat{c_{ir}}), \varphi_{tr})] ) \tag{11}$$
$$- \mathbb{E}_{(I_r, \varphi_{tr}) \sim p_{data}}[D(I_r, \varphi_{tr})]$$

$$\mathcal{L}_G = -\mathbb{E}_{z \sim p_z, \varphi_{tr} \sim p_{data}, \hat{c_{ir}} \sim TE(\phi)}[D(G(z, \hat{c_{ir}}), \varphi_{tr})]$$
$$+ \alpha * (D_{JS}(\mathcal{N}(\mu(\varphi_{tr}), \Sigma(\varphi_{tr})) || \mathcal{N}(0, I)$$
$$+ D_{JS}(\mathcal{N}(\mu(\varphi_{tw}), \Sigma(\varphi_{tw})) || \mathcal{N}(0, I)) \tag{12}$$
$$+ \beta * -\mathbb{E}_{i \sim P(I'=i|I, \phi, \psi_G), (I_r, I_w) \sim p_{data}}[R(I' = i, I)]$$

where $I_r$ and $\varphi_{tr}$ are the image and text embeddings that belong to the same class. Any image embedding *not* belonging to this corresponding class will be a candidate wrong image embedding $I_w$ and any text embedding not belonging to this class will be a candidate wrong text embedding $\varphi_{tw}$. The latent embeddings $\hat{c_{tr}}$ and $\hat{c_{tw}}$ are sampled from the Normal Distribution with $\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t))$, that is the latent space representation of the text embedding. This allow a small perturbations in $\varphi_{tr}$ and $\varphi_{tw}$ which is required to increase the generalizability of the model. The function $R(I' = i, I)$ is a regularizer in order to ensure that the embeddings generated by the generator can be used to retrieve all the relevant images. The equations are as follows:

$$R(I' = i, I) = -\log \frac{p(I' = i, I)}{1 - p(I' = i, I)} \tag{13}$$

where $p$ denotes the joint probability of $i$ and $I$ which we define as follows:

$$p(i, I) = \frac{e^{-d_p}}{e^{-d_p} + e^{-d_n}} \tag{14}$$

where $d_p = |I_r - i|$ and $d_n = |I_w - i| - \lambda$ are the Manhattan distances between the generated embedding $i$ and $I_r$ and $I_w$ respectively. We have also provided a margin $\lambda$ in order to ensure that $i$ is atleast $\lambda$ separated from $I_w$ in terms of Manhattan distance. Our formulation of probability $p(i, I)$ also ensures that $p \in (0, 1)$, since for $p$ to attain the value 0 or 1, $d_p$ or $d_n$ will have to be tending to $\infty$. However, in our formulation, $d_p$ and $d_n$ always attain finite values. This ensures that the term $R(I' = i, I)$ is not undefined anywhere during our optimization. The log of odds ratio function can be further simplified as:

$$\log \frac{p(i, I)}{1 - p(i, I)} = \log \frac{\frac{e^{-d_p}}{e^{-d_p} + e^{-d_n}}}{\frac{e^{-d_n}}{e^{-d_p} + e^{-d_n}}} = \log \frac{e^{-d_p}}{e^{-d_n}} = d_n - d_p \quad (15)$$

Here the discriminator $(D)$ and generator $(G)$ are trained by alternatively maximizing $L_D$ and minimizing $L_G$. Also, $\alpha$ and $\beta$, are hyper-parameters to be tuned to achieve optimal results.

**Discriminator:** In this adversarial game between the $D()$ and the $G()$, the $D()$ is trained to separate a real image from a fake one. In our task of generating representative embedding $i \in I'$ from text embeddings of a particular class, an image embedding from a different class $I_w$ should also be identified by the discriminator as fake given the text embedding $\hat{c_{tr}}$. For example as in Figure 1, if the text embedding of *Black-Footed-Albatross* is given and $G()$ generates an embedding of any other class, say *American Goldfinch*, then $D()$ should label it as fake and force $G()$ to generate $i$ for Black-Footed-Albatross from text-embedding of Black-Footed-Albatross. Hence we added another condition to the discriminator loss of classifying the image $I_w$ as fake and gave it equal weightage as compared to discriminator loss of classifying the generated embedding as fake. Hence, as shown in Eqn. 16, we add another discriminator loss to the CWGAN discriminator loss:

$$\begin{aligned}
\mathcal{L}_D &= (\mathbb{E}_{z \sim p_z, \varphi_{tr} \sim p_{data}, \hat{c_{tr}} \sim TE(\phi)}[D(G(z, \hat{c_{tr}}), \varphi_{tr})] \\
&\quad - \mathbb{E}_{(I_r, \varphi_{tr}) \sim p_{data}}[D(I_r, \varphi_{tr})]) * 0.5 \\
&\quad + (\mathbb{E}_{(I_w, \varphi_{tr}) \sim p_{data}}[D(I_w, \varphi_{tr})] \\
&\quad - \mathbb{E}_{(I_r, \varphi_{tr}) \sim p_{data}}[D(I_r, \varphi_{tr})]) * 0.5
\end{aligned} \quad (16)$$

**Generator:** As shown in Figure 1, we design a $G$ (with the loss function stated in Eqn. 12) which takes the latent embedding vector as input and tries to generate image embedding of the same class. In order to achieve this goal, we add two regularizers to $G$ – (1) Negative Log of odds regularizer (Eqn. 13),(2) Jensen-Shannon Divergence (Eqn. 17). The intention for jointly learning the text encoder model instead of training the text encoder separately is to ensure that the text encoder model learns to keep the important features required for the generator to generate image embedding.

**Learning text encodings:** As shown in Figure 1, the text embedding $\varphi_t$ is first fed into an encoder which encodes $\varphi_t$ into a Gaussian distribution $\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t))$, where mean $\mu$ and standard deviation $\Sigma$ are functions of $\varphi_t$. This helps to increase the samples of the text encoding and hence reduces the chance of overfitting. This kind of encoding provides small perturbations to $\varphi_t$, thus yielding more train pairs given a small number of image-text pairs. We train this model by optimizing the following condition as a regularization term while training the generator.

$$D_{JS}(\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t)) \,||\, \mathcal{N}(0, I)) \quad (17)$$

---

**Algorithm 1** Training the proposed model

1: **for** $it$ in $1 \ldots$n **do**
2:     /* **Update P (E-step)** */
3:     **for** $j$ in $1 \ldots it$ **do**                         /* train GAN */
4:         **for** $l$ in $1 \ldots 5$ **do**
5:             minimize $\mathcal{L}_D$
6:         **end for**
7:         $\underset{\psi_G}{\text{minimize}} \mathcal{L}_G$
8:     **end for**
9:     /* **Update $\psi_C$ (M-step)** */
10:    **for** $j$ in $1 \ldots it$ **do**                     /* train CSEM */
11:       Take text embedding $\varphi_{tr}$ and $\varphi_{tw} \sim \varphi$
12:       Obtain the $\hat{c_{tr}}$ and $\hat{c_{tw}}$ using the TE
13:       Obtain $i_r$ and $i_w$ using Generator
14:       $\underset{\psi_C}{\text{minimize}} \mathcal{L}_T$
15:    **end for**
16: **end for**

---

where $D_{JS}$ is the Jensen-Shanon divergence (JS divergence) between the conditioning distribution and the standard Gaussian distribution. Unlike the previous conditional GANs [18] our model does not append $z$ with the conditioned $c$ directly. Instead it learns the distribution over each embedding and thereby learns a new embedding space. The samples of the learned space are passed on to generator with $z$. This approach is more discriminative in the latent space, and helps to separate two classes in the original generated space.

*3.3.3* **Training setup and implementation details.** The values of the hyper-parameters in the model are set to $\alpha = 0.5$, $\beta = 2$, and $\lambda = 2$ using grid-search. The generator and the discriminator are trained in an iterative manner with the given objective functions (Eqns. 11 and 12). Both the generator and the discriminator are optimized with the root mean squared propagation(RmsProp) optimizer. The generator is a fully connected neural network with two hidden layers. The first hidden layer has $2,048$ units and the second hidden layer has $4,096$ units. *Leaky ReLU* activation function has been used in the hidden layer, and the output of the generator is passed through a *ReLU* function to ensure that there are no negative values. The discriminator is also a fully connected neural network with $1,024$ hidden layer units and 1 output unit. *Leaky ReLU* activation function is used in the hidden layer and no activation function is used for the output layer. The Text Encoder is a fully connected neural network with 2048 output units – 1024 dimension for the mean ($\mu$) and 1024 dimension for the standard deviation ($\Sigma$). The CSEM is implemented as a single layer feed forward network with *ReLU* activation. The weights and biases of both the generator and the discriminator are initialized with a random normal initializer having mean 0.0 and standard deviation 0.02.

## 3.4 Applying model for retrieval

Once the model ZSCRGAN is trained, the retrieval of images for a novel/unseen class proceeds as shown in Algorithm 2. The query for retrieval is the text embedding $\varphi_t$ of a novel/unseen class. Given $\varphi_t$, $\hat{c_t}$ is generated by the *text encoder*. Then $G()$ produces the image embedding $i_t$ which is passed through CSEM() to produce $\theta_t$. Now, for each image $I_{us}$ from an unseen class in the test set, we obtain

**Algorithm 2** Retrieval Algorithm

---
1: **procedure** GETRELEVANTIMAGES($\varphi_t$ , $k$)
2:                                  ▷ $\varphi_t$: the text embedding of a query-class
3:                        ▷ $k$: number of images to be retrieved
4:     $\hat{c}_t \leftarrow$ TEXTENCODER($\varphi_t$)
5:     $z \leftarrow$ GETRANDOMNORMALNOISE
6:     $i_t \leftarrow$ G($z, \hat{c}_t$)
7:     $\theta_t \leftarrow$ CSEM($i_t$)
8:     $simList \leftarrow []$
9:     **for** $I_{us}$ in test_image_set **do**
10:        $I_i \leftarrow$ FETCHIMGEMBEDDING($I_{us}$)
11:        $\theta_i \leftarrow$ CSEM($I_i$)
12:        $sim_{it} \leftarrow$ COSINESIM($\theta_t, \theta_i$)
13:        APPEND($simList, < sim_{it}, I_{us} >$)
14:     **end for**
15:     sort $simList$ in descending order of $sim_{it}$
16:     $imageList \leftarrow$ images in first k indices of simList
17:     **return** $imageList$
18: **end procedure**

---

the corresponding image embedding $I_i$ which is also passed through CSEM() to get $\theta_i$. Let $sim_{it} = cosineSim(CSEM(I_i), CSEM(G(z, \hat{c}_t)))$ be the cosine similarity between $\theta_t$ and $\theta_i$, where $z$ is a noise vector sampled from a random normal distribution. Thereafter, a 2-tuple $< sim_{it}, I_{us} >$ is formed and appended to a list called $simList$. The list is then sorted in descending order of the $sim_{it}$ values. The top $k$ images are extracted from the sorted $simList$ and stored in $imageList$ which is returned as the ranked list of retrieved images.

## 4 EXPERIMENTS AND ANALYSIS

This section details our experiments through which we compare the performance of the proposed model with that of several state-of-the-art baselines, over several standard datasets.

### 4.1 Datasets

We use the following five datasets for the experiments. Statistics of the datasets are summarized in Table 2. For each dataset, the models work on image embeddings and text/attribute embeddings (e.g., of the image captions). Table 2 also states the sources for the various embeddings for each dataset. For fair evaluation, every model, including the proposed and the baseline models, use the same text and image embeddings.

**(1) Oxford Flowers-102 (Flowers) dataset** contains 8,189 images of flowers from 102 classes [19]. Each class consists of between 40 and 258 images. Each image is accompanied by 10 sentences, each describing the image [22]. The data is split into 82 training and 20 test classes, with *no overlap* among the training and test classes [19].

As text embeddings, we use charCNN-RNN embeddings of the image captions provided by [22]. We use the image embeddings provided by [29], which are generated by passing each image through ResNet 101 architecture [10] pre-trained on the Imagenet [4], and taking the pre-trained ResNet's final layer embeddings.

**(2) Caltech-UCSD Birds (CUB) dataset**, which contains 11,788 images of birds, divided into 200 species (classes), with each class containing approx. 60 images [27]. Each image has 10 associated sentences describing the image [22]. Following [29], the images

in CUB are split into 150 training classes and 50 test classes, such that there is *no overlap* among the classes in training and test sets. Similar to the Flowers dataset, we use the text embeddings provided by [22] and image embeddings provided by [29].

**(3) North American Birds (NAB) dataset** is a larger version of the CUB dataset, with 48,562 images of birds categorized into 1,011 classes. The dataset was extended by Elhoseiny et. al. [8] by adding a Wikipedia article for each class. The authors also merged some of the classes to finally have 404 classes. Two types of split have been proposed in [8] in terms of how the seen (S) and unseen (U) classes are related – (1) Super-Category-Shared (SCS), and (2) Super-Category-Exclusives (SCE). In SCS, unseen classes are chosen such that for each unseen class, there exists at least one seen class that have the *same parent category*. However, in SCE split, there is no overlap in the parent category among the seen and unseen classes. Hence, retrieval is easier for the SCS split than for the SCE split.

We use the text and image embeddings provided by [33]. The 13, 217-dimensional text embeddings are actually TF-IDF vectors obtained from the Wikipedia articles (suitably preprocessed) corresponding to the classes. The image embeddings are obtained by feeding the images into VPDE-net [31], and extracting the activations of the part-based FC layer of the VPDE-net. The NAB dataset has six semantic parts ('head', 'back', 'belly', 'breast', 'wing', and 'tail'). A 512-dimension feature vector is extracted for each part and concatenated in order. The resulting 3072-dimensional vectors are considered as image embeddings.

**(4) Animals with Attributes (AWA) dataset** consists of 30, 475 images of animals from 50 different classes. 85 attributes are used to characterize each of the 50 categories, giving an class-level attribute embedding of 85 dimensions. The 2048-dimension image embeddings are taken from the final layer ResNet-101 model trained on ImageNet [4]. There exists two kinds of splits of the data according to [29] – *Standard-split* and *proposed split*. The standard-split does not take into consideration the ImageNet classes, while the proposed-split consists test classes which have no overlap with the ImageNet classes. We used the proposed-split [29] of the dataset for the ZSIR models (including the proposed model, DADN, ZSL-GAN).

**(5) Wikipedia (Wiki) dataset** consists of 2, 866 image-text pairs taken from Wikipedia documents [21]. The image-text pairs have been categorized into 10 semantic labels (classes). Images are represented by 128-dimensional SIFT feature vectors while the textual data are represented as probability distributions over the 10 topics, which are derived from a Latent Dirichlet Allocation (LDA) model. Following [16], the classes are split randomly in an 80:20 train:test ratio, and average of results over 10 such random splits is reported.

**Ensuring zero-shot setup for some of the datasets:** The Flowers, CUB, and AWA datasets use image embeddings from ResNet-101 that is pretrained over ImageNet. Hence, for these datasets, it is important to ensure that the test/unseen classes should *not* have any overlap with the ImageNet classes, otherwise such overlap will violate the ZSL setting [29]. Therefore, for all these three datasets, we use the same train-test class split as proposed by [29], which ensures that there is *no overlap among the test classes and the ImageNet classes on which ResNet-101 is pretrained*.

| Dataset | Dimensions of T/A/I | # T/A/I | # Seen / Unseen Classes | Text Embedding Source | Image Embedding Source | Interpretation of Query |
|---|---|---|---|---|---|---|
| **FLO** | T: 1024, I: 2048 | T: 102, I: 8,189 | 82/20 | https://github.com/reedscot/cvpr2016 [22] | http://datasets.d2.mpi-inf.mpg.de/xian/ImageNet2011_res101_feature.zip [29] | Textual description of a particular category of flowers |
| **CUB** | T: 1024, I: 2048 | T: 200, I: 11,788 | 150/50 | https://github.com/reedscot/cvpr2016 [22] | http://datasets.d2.mpi-inf.mpg.de/xian/ImageNet2011_res101_feature.zip [29] | Textual description of a particular category of birds |
| **NAB** | T: 13217, I: 3072 | T: 404, I: 48,562 | 323/81 | https://github.com/EthanZhu90/ZSL_GAN [33] | https://github.com/EthanZhu90/ZSL_GAN [33] | Textual description of a particular category of birds |
| **AWA** | A: 85, I: 2048 | A: 50, I: 30,475 | 40/10 | http://datasets.d2.mpi-inf.mpg.de/xian/ImageNet2011_res101_feature.zip [29] | http://datasets.d2.mpi-inf.mpg.de/xian/ImageNet2011_res101_feature.zip [29] | Human annotated attributes of a particular category of animals |
| **Wiki** | T: 10, I: 128 | T: 2,866, I: 2,866 | 8/2 (following split in [16]) | ttp://www.svcl.ucsd.edu/projects/crossmodal/ [21] | http://www.svcl.ucsd.edu/projects/crossmodal/ [21] | Textual part of Wikipedia articles |

**Table 2: Statistical description of the datasets (T: text, A: attribute, I: image). All models, including the proposed and baseline models, use the same embeddings for fair comparison (details in Section 4.1.**

## 4.2 Evaluation Metrics

For each class in the test set, we consider as the *query* its per-class text embedding $\varphi_t$ (or the per-class attribute embedding in case of the AWA dataset). The physical interpretation of the query for each dataset is explained in Table 2 (last column). For each class, we retrieve 50 top-ranked images (according to each model). Let the number of queries (test classes) be $Q$. We report the following evaluation measures:

**(1) Precision@50, averaged over all queries:** For a certain query $q$, Precision@50$(q) = \frac{k}{50}$, where $k$ is the number of relevant images among the top-ranked 50 images retrieved for $q$. We report Precision@50 averaged over all $Q$ queries.

**(2) mean Average Precision (mAP@50):** mAP@50 is the mean of *Average Precision at rank* 50, where the mean is taken over all $Q$ queries. $mAP@50 = \frac{1}{Q} \sum_{q=1}^{Q} \text{AveP}_{50}(q)$ where $\text{AveP}_{50}(q)$ for query $q$ is $\frac{\sum_{r \in R} \text{Precision@r}(q)}{|R|}$ where $R$ is the set of ranks (in $[1, 50]$) at which a relevant image has been found for $q$.

**(2) Top-1 Accuracy (Top-1 Acc):** This metric measures the fraction of queries (unseen classes) for which the *top-ranked* retrieval result is relevant [22, 33]. In our experiments we report the average Top-1 Accuracy of the models over the set of all unseen classes.

## 4.3 Baselines

We compare the proposed ZSCRGAN model with different kinds of baseline models, as described below.

**Zero-Shot Information Retrieval (ZSIR) models:** We consider three state-of-the-art models for ZSIR:

(1) Reed *et al.* developed the **DS-SJE [22]** model that jointly learns the image and text embeddings using a joint embedding loss function. We use the codes and pre-trained models provided by the authors (at https://github.com/reedscot/cvpr2016).

(2) Chi *et al.* proposed two very similar models for zero-shot IR [2, 3]. The **DADN model [3]** out-performs the one in [2]. Hence, we consider DADN as our baseline. DADN uses dual GANs to exploit the category label embeddings to project the image embeddings and text embeddings to have a common representation in a semantic space [3]. We use the codes provided by the authors (at https://github.com/PKU-ICST-MIPL/DADN_TCSVT2019).

(3) *Image generation models* are optimized to generate high-fidelity images from given textual descriptions. These models can be used for text-to-image retrieval as follows – given the textual query, we

can use such a model to generate an image, and then retrieve images that are 'similar' to the generated image as answers to the query. Zhu *et al.* [33] developed such a GAN-based approach for retrieval of images from textual data. The GAN (called **ZSL-GAN**) is trained with classification loss as a regularizer. We use the codes provided by the authors (at https://github.com/EthanZhu90/ZSL_GAN).

**Zero-Shot Classification (ZSC) models:** We consider two state-of-the-art ZSC models [25, 28] (both of which use generative models), and adopt them for the retrieval task as described below.

(1) **f-CLSWGAN [28]** used the GAN to synthesize the unseen class samples (image embeddings) using the class attribute. Then, using the synthesized samples of the unseen class, they trained the softmax classifier. We used the codes provided by the authors (at http://datasets.d2.mpi-inf.mpg.de/xian/cvpr18xian.zip). To use this model for retrieval, we used the image embeddings generated by the generator for the unseen classes, and ranked the image embeddings of the unseen classes using their cosine similarity with the generated image embedding.

(2) **SE-ZSL [25]** used the Conditional-VAE (CVAE) with feedback connection for Zero-Shot Learning (implementation obtained on request from the authors). We adopt the same architecture for the text to image retrieval as follows. Using the CVAE, we first trained the model over the training class samples. At test time, we generated the image embedding of the unseen classes conditioned on the unseen text query. In the image embedding space, we perform the nearest neighbour search between the generated image embedding and the original image embedding.

**Zero-Shot Hashing models:** Hashing based retrieval models have been widely studied because of their low storage cost and fast query speed. We compare the proposed model with some Zero-Shot Hashing models, namely DCMH [30], SePH [15], and the more recent AgNet [11] and CZHash [16]. The implementations of most of these hashing models are not available publicly; hence we adopt the following approach. We apply our method to the same datasets (e.g., AwA, Wiki) for which results have been reported in the corresponding papers [11, 16], and use the same experimental setting as reported in those papers.

Note that different prior works have reported different performance metrics. For those baselines whose implementations are available to us, we have changed the codes of the baselines minimally, to report Precision@50, mAP@50, and Top-1 Accuracy for all models.

| Retrieval Model | Prec@50(%) | mAP@50(%) | Top-1 Acc(%) |
|---|---|---|---|
| Zero-shot classification models adopted for retrieval | | | |
| SE-ZSL [25] | 29.3% | 45.6% | 59.6% |
| fCLSWGAN [28] | 36.1% | 52.3% | 64% |
| Zero-shot retrieval models | | | |
| DS-SJE [22] | 45.6% | 58.8% | 54% |
| ZSL-GAN [33] | 42.2% | 59.2% | 60% |
| DADN [3] | 48.9% | 62.7% | 68% |
| **ZSCRGAN (proposed)** | **52%**$^{SFJGD}$ | **65.4%**$^{SFJGD}$ | **74%**$^{SFJGD}$ |

**Table 3: Zero-Shot Retrieval on CUB dataset. The proposed model outperforms all baselines (bold-font indicates the best results in all tables). The super-scripts S, F, J, G, and D indicate that the proposed method is statistically significantly better at 95% confidence level (p < 0.05) than SE-ZSL, fCLSWGAN, DS-SJE, ZSL-GAN and DADN respectively.**

| Retrieval Model | Prec@50(%) | mAP@50(%) | Top-1 Acc(%) |
|---|---|---|---|
| Zero-shot classification models adopted for retrieval | | | |
| SE-ZSL [25] | 41.7% | 63.1% | 66.4% |
| fCLSWGAN [28] | 44.1% | 67.2% | 71.2% |
| Zero-shot retrieval models | | | |
| DS-SJE [22] | 55.1% | 65.7% | 63.7% |
| ZSL-GAN [33] | 38.7% | 46.6% | 45% |
| DADN [3] | 20.8% | 28.6% | 25% |
| **ZSCRGAN (proposed)** | **59.5%**$^{SFJGD}$ | **69.4%**$^{SFJGD}$ | **80%**$^{SFJGD}$ |

**Table 4: Zero-Shot Retrieval on Flower dataset. The proposed model outperforms all the baselines. Superscripts S, F, J, G, D show significant improvements (see Table 3).**

For comparing with the hashing models (whose implementations are not available), we report only mAP which is the only metric reported in [11, 16].

## 4.4 Comparing performances of models

Table 3 and Table 4 respectively compare the performance of the various models on the CUB and Flower datasets. Similarly, Table 5 and Table 6 compare performances over the NAB dataset (SCE and SCS splits respectively). Table 7 compares performances over the AwA and Wiki datasets. The main purpose of Table 7 is to compare performance with the Zero-Shot Hashing models whose implementations are not available to us; hence we report only mAP which is the only metric reported in [11, 16].

The proposed ZSCRGAN considerably outperforms almost all the baselines across all datasets, the only exception being that DADN outperforms ZSCRGAN on the Wiki dataset (Table 7). We performed Wilcoxon signed-rank statistical significance test at a confidence level of 95%. The superscripts S, F, J, G, and D in the tables indicate that the proposed method is statistically significantly better at 95% confidence level ($p < 0.05$) than SE-ZSL, fCLSWGAN, DS-SJE, ZSL-GAN and DADN respectively. We find that the results of the proposed model are significantly better than most of the baselines. Note that we could not perform significance tests for the Hashing methods owing to the unavailability of their implementations.

We now perform a detailed analysis of why our model performs better than the baselines.

**ZSCRGAN vs. DS-SJE:** DS-SJE [22] uses discriminative models to create text and image embeddings.[2] This discriminative approach

---

[2]We could not run DS-SJE over some of the datasets, as we were unable to modify the Lua implementation to suit these datasets. Also, the Prec@50 of DS-SJE on the Flower dataset (in Table 4) is what we obtained by running the pre-trained model provided by the authors, and is slightly different from what is reported in the original paper.

| Retrieval Model | Prec@50(%) | mAP@50(%) | Top-1 Acc(%) |
|---|---|---|---|
| Zero-shot classification models adopted for retrieval | | | |
| SE-ZSL [25] | 7.5% | 3.6% | 7.2% |
| Zero-shot retrieval models | | | |
| ZSL-GAN [33] | 6% | 9.3% | 6.2% |
| DADN [3] | 4.7% | 7.3% | 2.5% |
| **ZSCRGAN (proposed)** | **8.4%**$^{GD}$ | **11.8%**$^{SGD}$ | **7.4%**$^{GD}$ |

**Table 5: Zero-Shot Retrieval on NAB dataset (SCE split). The proposed model outperforms all the baselines. Superscripts S, F, J, G, D show significant improvements (see Table 3).**

| Retrieval Model | Prec@50(%) | mAP@50(%) | Top-1 Acc(%) |
|---|---|---|---|
| Zero-shot classification models adopted for retrieval | | | |
| SE-ZSL [25] | 25.3% | 34.7% | 11.4% |
| Zero-shot retrieval models | | | |
| ZSL-GAN [33] | 32.6% | 39.4% | 34.6% |
| DADN [3] | 26.5% | 28.6% | 17.3% |
| **ZSCRGAN (proposed)** | **36%**$^{SGD}$ | **43%**$^{SGD}$ | **49.4%**$^{SGD}$ |

**Table 6: Zero-Shot Retrieval on NAB dataset on (SCS split). The proposed model outperforms all the baselines. Superscripts S, G, D show significant improvements (see Table 3).**

| Retrieval Model | mAP(%) on AwA | mAP(%) on Wiki |
|---|---|---|
| Zero-shot Hashing models | | |
| DCMH [30] | 10.3% (with 64-bit hash) | 24.83% (with 128-bit hash) |
| AgNet [11] | 58.8% (with 64-bit hash) | 25.11% (with 128-bit hash) |
| SePH [15] | – | 50.44% (with 128-bit hash) |
| CZHash [16] | – | 25.87% (with 128-bit hash) |
| Zero-shot Retrieval models | | |
| ZSL-GAN [33] | 12.5% | - |
| DADN [3] | 27.9% | **58.94%** |
| **ZSCRGAN (proposed)** | **62.2%**$^{GD}$ | 56.9% |

**Table 7: Zero-shot retrieval on (i) AwA dataset, and (ii) Wiki dataset. Results of hashing models reproduced from [11] and [16]. Other metrics could not be reported due to unavailability of the implementations of the hashing models.**

has limited visual imaginative capability, whereas the generative approach of the proposed ZSCRGAN does not suffer from this problem. Specifically, the low performance of DS-SJE on the CUB datast (Table 3) is due to the lack of good visual imaginative capability which is required to capture the different postures of birds in the CUB dataset (which is not so much necessary for flowers).

**ZSCRGAN vs. DADN:** DADN [3] uses semantic information contained in the class labels to train the GANs to project the textual embeddings and image embeddings in a common semantic embedding space. Specifically, they use 300-dimensional Word2vec embeddings pretrained on the Google News corpus[3], to get embeddings of the class labels. A limitation of DADN is that unavailability of proper class label embeddings can make the model perform very poorly in retrieval tasks. For instance, DADN performs extremely poorly on the Flowers dataset, since out of the 102 class labels in the Flowers dataset, the pretrained Word2vec embeddings are not available for as many as 23 labels. Similarly, out of 404 class labels in the NAB dataset, the pretrained Word2vec embeddings are not available for 8 labels. ZSCRGAN does *not* rely on class labels, and performs better than DADN on both these datasets. On the other hand, DADN performs better than ZSCRGAN on the Wiki dataset (Table 7 since pretrained embeddings are available for all class labels.

---

[3]https://code.google.com/archive/p/word2vec/

**ZSCRGAN vs. ZSL-GAN:** ZSL-GAN uses a GAN to generate image embeddings from textual descriptions.[4] In their model the discriminator branches to two fully connected layers – one for distinguishing the real image embeddings from the fake ones, and the other for classifying real and fake embeddings to the correct class. We believe that the retrieval efficiency is lower due to this classification layer of the discriminator. This layer learns to classify the real and generated embeddings to the same class, which is contradictory to its primary task of distinguishing between the real and generated embeddings (where ideally it should learn to assign different labels to the real and generated embeddings).

**ZSCRGAN vs. Hashing models (e.g., AgNet, CZHash):** Hashing-based retrieval methods are popular due to their low storage cost and fast retrieval. However, achieving these comes at a cost of retrieval accuracy – to generate the hash, the models sacrifice some information content of the embeddings, and this results in the loss in accuracy. As can been seen from Table 7, all the hashing methods have substantially lower mAP scores than our proposed method.

## 4.5 Analysing two design choices in ZSCRGAN

In this section, we analyse two design choices in our proposed model – (1) why we adopted an Expectation-Maximization approach, instead of jointly training the GAN and the CSEM, and (2) why we chose to select wrong class embeddings *randomly*.

*4.5.1* **E-M vs Joint Optimization:** In the proposed model, the CSEM and the GAN are trained alternately using an E-M setup (see Sec. 3.3). However, it is also possible to train both these networks jointly; i.e., when the Generator is trained, the CSEM loss is also optimized. We performed a comparison between these two approaches, and observed that the performance of the model drops by a significant amount when jointly trained. Figure 2 compares the performance of the two approaches in terms of Precision@50 over the CUB dataset, and shows that the EM setup results in better retrieval. Similar observations were made for other performance metrics, across all datasets (details omitted for brevity).

The reason why the jointly training approach, where the CSEM loss is optimized along with the generator loss, does not work well is as follows. Using the triplet Loss $\mathcal{L}_T$ (defined in Eqn. 9), the CSEM learns maximize similarity with the relevant embeddings and minimize similarity with the irrelevant embeddings. Thus, during backpropagation, the relevant representative embeddings have a different gradient than the irrelevant representative embeddings $G(z, \hat{c_{t_w}})$. When the CSEM is jointly trained with the generator, the weights of the CSEM get updated after every iteration, causing different gradients (or weights) for each of the wrong class embeddings of the generator, thus causing a distorted space of wrong class embeddings, and thus causing hindrance to the learning of the Generator. This problem, however, has been removed in the EM setup where the parameters of the generator is frozen while training the CSEM.

*4.5.2* **Choice of wrong class embeddings:** In the proposed model, for given $I_r$ and $\phi_{tr}$ for a certain target class $c_{target}$, we learn the
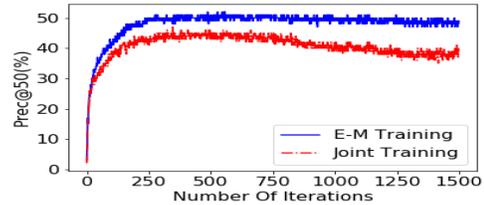
**Figure 2: [color online] Comparing performance of proposed E-M setup (solid blue curve) with joint training of the GAN and the CSEM (dashed red curve). Shown is how Prec@50 varies with the number of iterations for which the model is trained, over the CUB dataset. The E-M setup achieves consistently better performance.**
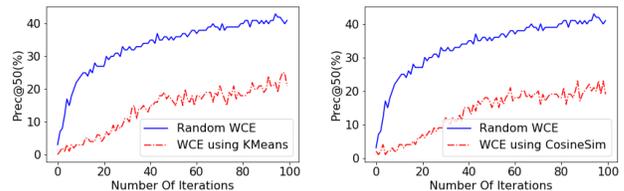


**Figure 3: [color online] Comparing random selection of Wrong Class Embeddings (WCE) with other ways of selecting WCE: (a) using KMeans clustering of images, and (b) using cosine similarity (details in text). Both figures show how Prec@50 varies with number of iterations for which the model is trained, over the CUB dataset. Random selection of WCE performs the best.**

representative embedding for the images of that class. To this end, we use wrong class embeddings (WCE) selected randomly from among all other classes. One might think that, instead of randomly selecting wrong classes, we should employ some intelligent strategy of selecting wrong classes. For instance, one can think of selecting WCE such that they are most similar to $\phi_{tr}$, which may have the apparent benefit that the model will learn to distinguish between confusing (similar) classes well. However, we argue otherwise.

Restricting the choice of the wrong classes distorts the space of the wrong embeddings, and hence runs the risk of the model identifying embeddings from classes outside the space of distorted wrong embeddings as relevant. In other words, though the model can learn to distinguish the selected wrong classes well, it fails to identify other classes (that are *not* selected as the wrong classes) as non-relevant.

To support our claim, we perform two experiments – (1) The wrong class is selected as the class whose text embedding has the highest cosine similarity to $\phi_{tr}$, and (2) The images are clustered using K-Means clustering algorithm, and the wrong class is selected as that class whose images co-occur with the maximum frequency in the same clusters as the images from $c_{target}$. Figure 3 compares the performance of the proposed model (where WCE are selected randomly) with these two modified models. Specifically Precision@50 is compared over the CUB dataset. In both cases, the accuracy drops drastically when WCE are chosen in some way other than randomly. Observations are similar for other performance metrics and other datasets (omitted for brevity).
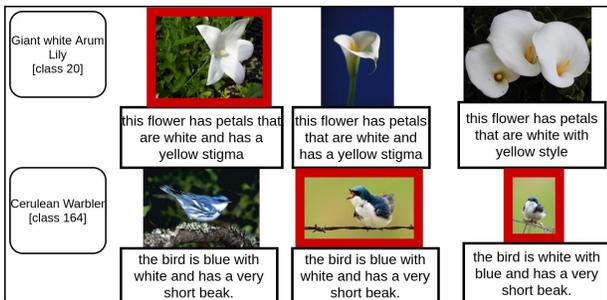
**Figure 4: [color online] Images from the top three classes retrieved by ZSCRGAN, for the query-classes shown on the left. Top panel for Flower dataset, bottom panel for CUB dataset. The images with thick red boundaries are from some class that is *not* the query-class (hence not considered relevant), but they are very similar to some images in the query-class.**

## 4.6 Error Analysis of ZSCRGAN

We analyse the failure cases where ZSCRGAN retrieves an image from a different class, compared to the query-class (whose text embedding has been issued as the query). Figure 4 shows some such examples, where the images enclosed in thick red boxes are not from the query-class. In general, we find that the wrongly retrieved images are in fact very similar to some (correctly retrieved) images in the query-class. For instance, in the CUB dataset, for the query-class *Cerulean Warbler*, the textual description of an image from this class (*this bird is blue with white and has a very short beak*) matches exactly with that of an image from a different class (which was retrieved by the model). Other cases can be observed where an image of some different class that has been retrieved, matches almost exactly with the description of the query-class. For instance, in the Flower dataset, for the query class *Giant White Arum Lily*, the wrongly retrieved class also has flowers with white petals and a yellow stigma, which matches exactly with many of the flowers in the *Giant White Arum Lily* class.

## 4.7 Ablation Analysis of ZSCRGAN

Table 8 reports an ablation analysis, meant to analyze the importance of different components of our proposed architecture. For brevity, we report only Prec@50 for the two datasets CUB and Flowers (observations on other datasets are qualitatively similar).

The largest drop in performance occurs when the wrong class embedding is not used. As stated earlier, this use of wrong class embeddings is one of our major contributions, and an important factor in the model's performance. Another crucial factor is the generation of representative embeddding $i \in I'$ for each class using a GAN. Removing this step also causes significant drop in performance. The Triplet loss and the Log of odds ratio regularizer (R) are also crucial – removal of either leads to significant degradation in performance. Especially, if R is removed, the performance drop is very high for the Flower dataset. R is more important for the Flower dataset, since it is common to find different flowers having similar shape but different colors, and R helps to distinguish flowers based on their colors.

| Retrieval Model | Prec@50 CUB | Prec@50 Flowers |
|---|---|---|
| Complete proposed model | 52% | 59.5% |
| w/o use of wrong class embedding | 24.7% | 27.2% |
| w/o R (regularizer) and Triplet Loss (CSEM) | 23.8% | 33.7% |
| w/o Triplet Loss | 36.2% | 41.4% |
| w/o R (regularizer) | 48.4% | 35.2% |
| w/o GAN (i.e. the representative embedding for a class) | 25.9% | 32% |

**Table 8: Results of ablation analysis on the proposed model. Precision@50 reported on CUB and Flower datasets.**

## 5 CONCLUSION

We propose a novel model for zero-shot text to image (T → I) retrieval, which outperforms many state-of-the-art models for ZSIR as well as several ZS classification and hashing models on several standard datasets. The main points of novelty of the proposed model ZSCRGAN are (i) use of an E-M setup in training, and (ii) use of wrong class embeddings to learn the representation of classes. The implementation of ZSCRGAN is publicly available at https://github.com/ranarag/ZSCRGAN.

In future, we look to apply the proposed model to other types of cross-modal retrieval (I → T), as well as to the zero -shot multi-view setup (TI → I, TI → T, I → TI, etc.) where multiple modes can be queried or retrieved together.

## REFERENCES

[1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein Generative Adversarial Networks. In *Proc. ICML*. 214–223.
[2] Jingze Chi and Yuxin Peng. 2018. Dual Adversarial Networks for Zero-shot Cross-media Retrieval. In *Proc. IJCAI*. 663–669.
[3] J. Chi and Y. Peng. 2019. Zero-shot Cross-media Embedding Learning with Dual Adversarial Distribution Network. *IEEE Transactions on Circuits and Systems for Video Technology* 30, 4 (2019), 1173–1187.
[4] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *Proc. IEEE CVPR*.
[5] Sounak Dey, Pau Riba, Anjan Dutta, Josep Llados, and Yi-Zhe Song. 2019. Doodle to Search: Practical Zero-Shot Sketch-Based Image Retrieval. In *Proc. IEEE CVPR*.
[6] Xingping Dong and Jianbing Shen. 2018. Triplet Loss in Siamese Network for Object Tracking. In *The European Conference on Computer Vision (ECCV)*.
[7] Anjan Dutta and Zeynep Akata. 2019. Semantically Tied Paired Cycle Consistency for Zero-Shot Sketch-based Image Retrieval. In *Proc. IEEE CVPR*.
[8] M. Elhoseiny, Y. Zhu, H. Zhang, and A. Elgammal. 2017. Link the Head to the "Beak": Zero Shot Learning from Noisy Text Description at Part Precision. In *Proc. IEEE CVPR*. 6288–6297.
[9] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio. 2014. Generative adversarial nets. In *Proc. NIPS*.
[10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proc. IEEE CVPR*.
[11] Zhong Ji, Yunxin Sun, Yunlong Yu, Yanwei Pang, and Jungong Han. 2020. Attribute-Guided Network for Cross-Modal Zero-Shot Hashing. *IEEE Transactions on Neural Networks and Learning Systems* 31, 321–330.
[12] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. In *Proc. ICLR*.
[13] Sasi Kiran Yelamarthi, Shiva Krishna Reddy, Ashish Mishra, and Anurag Mittal. 2018. A Zero-Shot Framework for Sketch based Image Retrieval. In *Proc. ECCV*.
[14] Hugo Larochelle, Dumitru Erhan, and Yoshua Bengio. 2008. Zero-data Learning of New Tasks. In *Proc. AAAI - Volume 2*.
[15] Z. Lin, G. Ding, J. Han, and J. Wang. 2017. Cross-View Retrieval via Probability-Based Semantics-Preserving Hashing. *IEEE Transactions on Cybernetics* 47, 12 (2017), 4342–4355.
[16] X. Liu, Z. Li, J. Wang, G. Yu, C. Domenicon, and X. Zhang. 2019. Cross-Modal Zero-Shot Hashing. In *Proc. IEEE ICDM*.
[17] Tomas Mikolov, Andrea Frome, Samy Bengio, Jonathon Shlens, Yoram Singer, Greg S Corrado, Jeffrey Dean, and Mohammad Norouzi. 2013. Zero-Shot Learning by Convex Combination of Semantic Embeddings. In *Proc. ICLR*.

[18] M. Mirza and S. Osindero. 2014. Conditional generative adversar-ial nets. In *arXiv:1411.1784*.

[19] M-E. Nilsback and A. Zisserman. 2008. Automated Flower Classification over a Large Number of Classes. In *Proc. ICCVGIP*.

[20] Mark Palatucci, Dean Pomerleau, Geoffrey E Hinton, and Tom M Mitchell. 2009. Zero-shot Learning with Semantic Output Codes. In *Proc. NIPS*. 1410–1418.

[21] N. Rasiwasia, J. Costa Pereira, E. Coviello, G. Doyle, G.R.G. Lanckriet, R. Levy, and N. Vasconcelos. 2010. A New Approach to Cross-Modal Multimedia Retrieval. In *ACM International Conference on Multimedia*. 251–260.

[22] A. Reed, Z. Akata, B. Schiele, and H. Lee. 2016. Learning deep representations of fine-grained visual descriptions. In *Proc. IEEE CVPR*.

[23] Richard Socher, Milind Ganjoo, Christopher D Manning, and Andrew Ng. 2013. Zero-shot learning through cross-modal transfer. In *Proc. NIPS*.

[24] L. Theis, A. van den Oord, and M. Bethge. 2016. A note on the evaluation of generative models. In *Proc. ICLR*.

[25] Vinay Kumar Verma, Gundeep Arora, Ashish Mishra, and Piyush Rai. 2018. Generalized zero-shot learning via synthesized examples. In *Proc. IEEE CVPR*.

[26] Kaiye Wang, Qiyue Yin, Wei Wang, Shu Wu, and Liang Wang. 2016. A Compre-hensive Survey on Cross-modal Retrieval. *CoRR* abs/1607.06215 (2016).

[27] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. 2010. *Caltech-UCSD Birds 200*. Technical Report CNS-TR-2010-001. California Institute of Technology.

[28] Yongqin Xian, Tobias Lorenz, Bernt Schiele, and Zeynep Akata. 2018. Feature generating networks for zero-shot learning. In *Proc. IEEE CVPR*.

[29] Yongqin Xian, Bernt Schiele, and Zeynep Akata. 2017. Zero-Shot Learning – The Good, the Bad and the Ugly. In *Proc. IEEE CVPR*.

[30] Erkun Yang, Cheng Deng, Wei Liu, Xianglong Liu, Dacheng Tao, and Xinbo Gao. 2017. Pairwise Relationship Guided Deep Hashing for Cross-Modal Retrieval. In *Proc. AAAI*. 1618âĂŞ1625.

[31] Han Zhang, Tao Xu, Mohamed Elhoseiny, Xiaolei Huang, Shaoting Zhang, Ahmed Elgammal, and Dimitris Metaxas. 2016. SPDA-CNN: Unifying semantic part detection and abstraction for fine-grained recognition. In *Proc. IEEE CVPR*.

[32] Li Zhang, Tao Xiang, and Shaogang Gong. 2017. Learning a deep embedding model for zero-shot learning. In *Proc. IEEE CVPR*.

[33] Yizhe Zhu, Mohamed Elhoseiny, Bingchen Liu, Xi Peng, and Ahmed Elgammal. 2018. A Generative Adversarial Approach for Zero-Shot Learning from Noisy Texts. In *Proc. IEEE CVPR*.