

Adversarial Training of Variational Auto-encoders for Continual Zero-shot Learning(A-CZSL)

Subhankar Ghosh
Aerospace Engineering Department
Indian Institute of Science
Bengaluru, India
subhankarg@alum.iisc.ac.in

Abstract—Most of the existing artificial neural networks(ANNs) fail to learn continually due to catastrophic forgetting, while humans can do the same by maintaining previous tasks’ performances. Although storing all the previous data can alleviate the problem, it takes a large memory, infeasible in real-world utilization. We propose a continual zero-shot learning model(A-CZSL) that is more suitable in real-case scenarios to address the issue that can learn sequentially and distinguish classes the model has not seen during training. Further, to enhance the reliability, we develop A-CZSL for a single head continual learning setting where task identity is revealed during the training process but not during the testing. We present a hybrid network that consists of a shared VAE module to hold information of all tasks and task-specific private VAE modules for each task. The model’s size grows with each task to prevent catastrophic forgetting of task-specific skills, and it includes a replay approach to preserve shared skills. We demonstrate our hybrid model outperforms the baselines and is effective on several datasets, i.e., CUB, AWA1, AWA2, and aPY. We show our method is superior in class sequentially learning with ZSL(Zero-Shot Learning) and GZSL(Generalized Zero-Shot Learning). Our code is available at https://github.com/CZSLwithCVAE/CZSL_CVAE

I. INTRODUCTION

Conventional supervised machine learning (ML) and, more recently, deep learning algorithms have shown remarkable performance on image classification, Computer Vision and Natural Language Processing. Despite the recent success of supervised machine learning algorithms, they have two significant limitations: 1) Conventional machine learning models have restricted themselves to training limited classes. When any example from the novel/unseen class occurs at the test time, such examples can not be correctly classified. 2) Conventional machine learning models cannot continually learn over time by accommodating new knowledge from streaming data while retaining learned information. We still need to improve our existing algorithms to achieve human-level performance the way humans learn [1] sequentially without forgetting previous tasks throughout their life. Can we build an ANN model that can learn sequentially and simultaneously works for zero-shot learning(distinguish the classes it has not seen during training)? The name of such methods is continual zero-shot learning.

ZSL(zero-shot learning), where a trained model sees data only from unseen classes during testing, and

GZSL(generalized zero-shot learning), where data may come from both seen and unseen classes for prediction. Both mentioned tasks are challenging for a model to perform continually. Researchers have addressed both ZSL [5], [20] and continual learning [9], [13], [18] approaches separately. Therefore, a more preferable and feasible approach needs to tackle continual learning and unseen object problems simultaneously. This model aims to leverage the advantages of both continual learning and zero-shot learning in a single framework. Consequently, we propose a model that can predict the unseen class object and adapt to a new task without completely forgetting the knowledge about the previous task. References [7], [22] combine ZSL with continual learning before, though our approach is entirely different from that, but is more realistic and beats their baselines results on the same datasets.

A traditional zero-shot learning(ZSL) method identifies unseen classes using seen class samples and class embeddings. Despite showing a promising performance of traditional ZSL, the method is unable to learn from streaming data. Generative approaches have received quite an attention over embedding based approaches due to the synthesized ability of unseen class features and a significant improvement in the model’s performance. In this paper, we propose a novel adversarial training of variational autoencoders(VAEs) for continual zero-shot learning. Here, the model composes a task-specific private module that is learned for each task, and a task-invariant shared module that is learned for all tasks. Task-specific in a sense, the 1st private module learns from the 1st task’s real data only, whereas the $t^{th}(t = 2, 3, \dots, T)$ private module gets trained by real data from t^{th} task and replay synthesized data from all previous (t-1) tasks generated from the shared module. We tackle both ZSL and continual learning together by using CVAE(conditional variational autoencoders) [12] that transfer knowledge from seen to unseen classes through class embeddings [24] to counter ZSL problems. As visual data is not available during training time, knowledge transfer from seen to unseen classes is formed through side information that makes semantic relationships between classes and class-embeddings. Our approach is motivated by the fact that processing and synthesizing images

are time taking for continual learning when the number of classes is high. Therefore, instead of images, we train and test our model with the features of the same images generated using a pre-trained model. Another thing that motivates us is that the human brain structure is complex and contains billions of neurons [10], so we may need to eventually make complicated networks in the coming future containing a huge number of neurons to learn sequentially. The main contributions of this work are summarized as follows:

- We develop an generative replay-based and structure-based continual zero-shot learning method using CVAE.
- The proposed method is developed for a single head setting that is more convenient to solve real case scenarios.
- We present results for four ZSL benchmark datasets for continual zero-shot learning.
- We propose two types of modules: one, task-invariant holds information for all tasks, another is task-specific. If there are T tasks, our proposed architecture consists of one task-invariant VAE and T task-specific VAEs.

II. RELETED WORK

A. Continual learning

The main problem in continual learning is catastrophic forgetting. McCloskey and Cohen first introduced the term catastrophic forgetting or catastrophic interference in the 1980s [17]. They claimed that catastrophic interference is a fundamental limitation of neural networks and a downside of its high generalization ability. While the cause of catastrophic forgetting has not been studied analytically, it is known that the neural networks parameterize the internal features of inputs, and training the networks on new samples causes alteration in already established representations. There are three types of continual learning to tackle forgetting: regularization-based, memory-based, and structure-based methods.

Regularization methods

Here, the learning capacity is fixed [13], [27], and continual learning is performed by penalizing a network’s parameters. Researchers use a new regularizer for a novel method. In [27], the essential parameters are computed online. They keep track of how the loss function changes due to a specific parameter change and accumulate this information during training. There should be a weight importance process for parameters selection to prioritize parameters usage—the way elastic weight consolidation(EWC) [13] gives importance to parameters based on the Fisher information matrix. The usages of these methods are limited because they cannot perform well for a large number of tasks.

Memory-based methods

Methods in this category try to prevent forgetting by either storing [14], [15] or synthesizing data from previous classes. The first one needs memory for rehearsal, whereas the latter is a generative model like GAN [8] or VAE [12], or both synthesize data of previous tasks to perform pseudo-rehearsal. The number of examplers stored decreases with the increase

in classes if the memory budget is limited. Researchers have recently proposed using a tiny memory [14] to store a few examples per class for old tasks.

Structure-based methods

The third approach to mitigate forgetting is structure-based methods [19]. The size of a network grows with each task to prevent catastrophic forgetting. Previous tasks’ performance is maintained by freezing the learned module while accommodating new tasks by augmenting the network with new modules. The computational cost for this method is inevitable if the number of tasks is high.

B. Zero-shot learning and Generalized zero-shot learning

Recently, zero-shot learning(ZSL) [5], [20] has attracted a lot of attention because the model can distinguish unseen classes during testing. ZSL models are able to do so by transferring knowledge from seen to unseen levels through a semantic relationship between classes and their attributes. We can transform a ZSL problem into a supervised machine learning problem using generative models like GAN or VAE, or both. Once a generative model gets trained, it can synthesize data for unseen classes, and the data is useful for training a classifier like a conventional supervised problem. Another modification of ZSL is GZSL, a more practical approach, where data come from both seen and unseen classes.

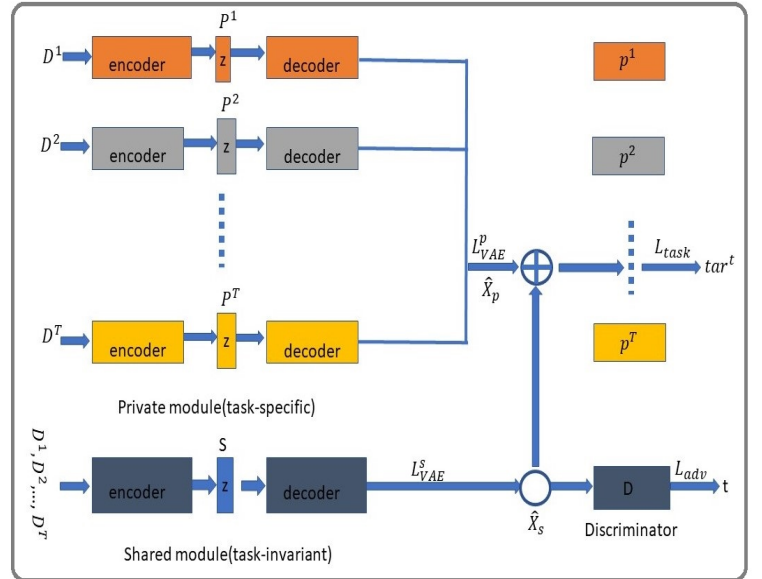


Fig. 1: shows our model at training time where the shared module plays a minimax game with the discriminator to generate task-invariant \hat{X}_s^s features. In contrast, the discriminator attempts to assign task labels to the synthesized features (\hat{X}_s^s). The discriminator tries to maximize the probability of task label. Architecture grows because of the task-specific modules denoted as P^t and p^t , task-specific perceptron networks assign class labels during each task’s training. To prevent forgetting shared module is trained with replay examples of previously learned classes generated from the decoder of the shared module using $z \sim \mathcal{N}(\mu = 0, \Sigma = 1)$.

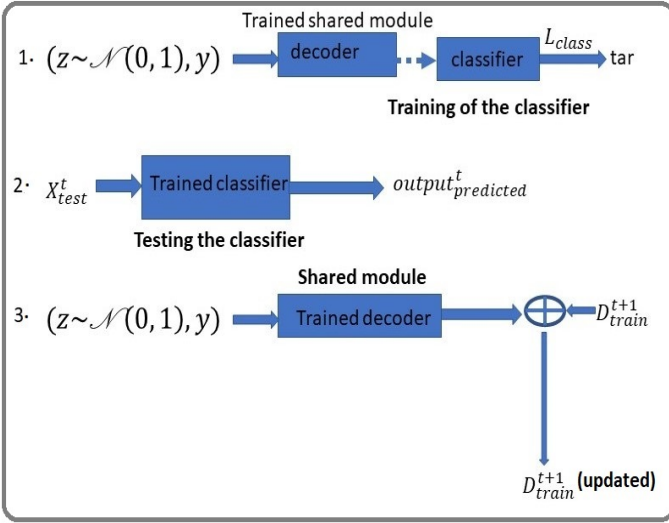


Fig. 2: has three parts:

- 1) **Training of the classifier**(separate) with generated data(of both seen and unseen classes) using the shared module's trained decoder.
- 2) **Testing the classifier** for both seen and unseen classes(ZSL).
- 3) The **shared module's** decoder is used to generate replay samples of previous tasks that get concatenated to the next task to update the model with new and old tasks.

C. Continual Zero-shot Learning

Reference [7] proposed a continual zero-shot learning model called Generalized continual zero-shot learning(GCZSL), a single head CZSL where the task identity is revealed during training but not during testing. To mitigate forgetting, they used knowledge distillation and stored few actual data per class using a small episodic memory. The GCZSL and reference [22] used the single-head setting, where the task identity is disclosed during training but not at test time. There are a few works [3], [26] that used the multi-head setting. Reference [3] proposed an average gradient episodic memory(A-GEM) based CZSL, and [26] offered a generative model-based CZSL.

D. Adversarial learning

Adversarial learning has usages in many domains such as generative models [8], object composition [2], representation learning [16], domain adaptation [23], active learning [21] etc. In adversarial training, a model learns the parameters through the minimax game, where a module wants to maximize the cost function, and another wants to minimize the same. This paper shows shared play the minimax game with discriminator, where shared tries to minimize the loss function, and the discriminator wants to maximize.

III. ADVERSARIAL TRAINING OF CONDITIONAL VARIATIONAL AUTOENCODERS FOR CONTINUAL ZERO-SHOT LEARNING(A-CZSL)

We study the problem of learning a sequence of T data distributions denoted as $D_{tr} = \{D_{tr}^1, D_{tr}^2, \dots, D_{tr}^T\}$, where $D_{tr}^t = \{(X_i^t, Y_i^t, tar_i^t, T_i^t)_{i=1}^{n_t}\}$ is the data distribution for

the task t with n_t sample tuples of input($X^t \in \mathcal{X}$), target label($tar^t \in tar$), attributes of classes ($Y^t \in \mathcal{Y}$), and task label($T^t \in \mathcal{T}$). D_{tr}^t contains seen class information. Apart from this, class embeddings of unseen classes($\mathcal{U}_c = \{(a_i)_{i=1}^{n_{uc}}\}$) are also available, where n_{uc} is the number of unseen labels. The goal is to learn a sequential function, $f : (z \sim \mathcal{N}(0,1), \mathcal{Y}) \rightarrow \hat{\mathcal{X}}_s$, for each task, where $\hat{\mathcal{X}}_s$ is synthesized data generated from the shared module. The synthetic data can be used to train a supervised classifier. We aim to learn another function, $f_{\theta_c} : \hat{\mathcal{X}} \rightarrow tar$, after training each task, that can map input(from seen or unseen or both classes) into it's target output without affecting the previous model's performance on prior works. We try to achieve our goal by training two separate modules: shared and private, to enhance a better knowledge transfer from seen to unseen classes and better forget avoidance of prior knowledge. The model prevents catastrophic forgetting in shared and private spaces separately and begins learning f_{θ}^t where $\theta \in (\theta_S, \theta_P)$ as mapping function from $(\mathcal{X}_{tr}^t, \mathcal{Y}_{tr}^t)$ to tar^t . We use some n samples per class to be synthesized prior to t^{th} task and accumulate the generated data to the current task(t^{th}) to train the model.

$$D_{tr}^t \leftarrow D_{tr}^t \cup D_{gen}^{1:(t-1)}$$

The cross-entropy loss function for the f_{θ}^t mapping corresponds to:

$$L_{task=t}(f_{\theta}^t, D_{tr}^t) = - \mathbb{E}_{(\mathcal{X}^t, \mathcal{Y}^t, tar^t) \sim D^t} \left[\sum_{c=1}^C \mathbb{1}_{(c=tar^t)} \log(\sigma(f_{\theta}^t(\mathcal{X}^t, \mathcal{Y}^t))) \right] \quad (1)$$

Where σ is the softmax function, in learning a sequence of tasks, an ideal f^t maps the input features X^t to two independent feature spaces: X_s^t a shared features space among all tasks and X_p^t remains private for each task. Both X_s^t and X_p^t get concatenated and fed to a task-specific multi-layer perceptron network to get desired output labels.

We introduce a mapping named shared ($S_{\theta_s} : X \rightarrow \hat{\mathcal{X}}_s$) and train it to generate features by feeding noise into the shared module's decoder to fool a discriminator D. In contrast, the $D(D_{\theta_d} : \hat{\mathcal{X}}_s \rightarrow T)$ try to assign the synthesized features to their corresponding task labels($T^t \in \{0,1,2,\dots,T\}$). The decoder and the discriminator can do so when the D gets trained to maximize the probability of assigning correct task labels to the features generated from the shared module. Simultaneously, the shared tries to minimize the same probability.

The corresponding cross-entropy adversarial loss for the minimax game:

$$L_{adv}(D, S, D_{tr}^t) = \min_S \max_D \sum_{t=0}^T \mathbb{1}_{t=t} [\log(D(S(\mathcal{X}^t, \mathcal{Y}^t)))] \quad (2)$$

The extra-label zero is there for fake data generated from the Gaussian distribution with mean = 0 and std = 1. In most

cases, we use adversarial training in a generative adversarial network that tries to learn the input data distribution in order to synthesize more data from the same distribution. Here we do the same by utilizing generative models task-invariant shared(VAE), and task-specific private(VAE); both try to learn input data distribution.

To facilitate adversarial training for S, we use the Gradient Reversal Layer [6] that directly tries to maximize the discriminator’s loss. The layer acts like an identity function during forward-propagation but multiplies the loss with a negative one during backpropagation in order to maximize the cost function for the discriminator. The adversarial training between the discriminator and the shared is complete when the discriminator can no longer predict the correct task label for features generated from the shared module. The private module, however, merely learns any task-invariant features.

Variational autoencoders

Autoencoders can effectively learn feature space and representation [2], [27]. A variational Autoencoder(VAE) is a generative model that follows an encoder-latent vector-decode architecture of classical autoencoder, which places a prior distribution on the feature space and uses an expected lower bound to optimize the learned posterior. Conditional VAE is an extension of the VAE, where data are fed to network with class properties such as labels, attributes, etc. The VAE is a fundamental building block of our approach. Variational distribution aims to find a true conditional probability distribution over the latent variables z through minimizing their distance using a variational lower bound limit. The loss function for a VAE is:

$$L_{VAE} = \mathbb{E}_{q_\phi(z|x)} [\log(p_\theta(x|z))] - D_{KL}(q_\phi(z|x) || p_\theta(z)) \quad (3)$$

Where the first term is the reconstruction loss, and the second one is the KL divergence between $q(z|x)$ and $p(z)$. The encoder predicts μ and Σ such that $q_\phi(z|x) = \mathcal{N}(\mu, \Sigma)$, from which a latent vector is synthesized via reparametrization process.

The final objective function of the model for the t^{th} task is:

$$L^{(t)} = \lambda_1 L_{adv} + \lambda_2 L_{task} + \lambda_3 L_{VAE}^s + \lambda_4 L_{VAE}^p \quad (4)$$

Where, $\lambda_1, \lambda_2, \lambda_3$, and λ_4 are regularizers to control the effect of each component. The full algorithm of the model is given in Algorithm [1].

A. Avoid forgetting

Catastrophic forgetting occurs because of the imbalance between old and new classes that results in a bias of the network towards the newest ones. One insight of our approach is to decouple the single representation learned for all tasks continually into two parts: private and shared. Though knowledge is transferred for ZSL and GZSL mostly from the shared module from seen to unseen classes. The critical approach is experience replay that gets concatenated to the current task’s data during training of the model with the same task to avoid forgetting sequentially.

B. Datasets

We evaluate our model on four benchmark datasets used for ZSL: Attribute Pascal and Yahoo(aPY) [20], Animals With Attributes(AWA1, AWA2) [20], and Caltech-UCSD-Birds 200-2011(CUB) [25]. Statistics of the datasets are presented in Table [1].

C. Continual Zero-shot learning(CZSL) setting

The dataset we use follows the setting used in [22]. It explains whether a class is seen or unseen is decided based on the number of tasks a model has been trained so far. If a model goes trained continually up to the t^{th} task, the classes are assumed to be seen till the t^{th} task, and the rest of the whole dataset’s classes are accepted unseen for the model while training.

D. evaluation matrices

We evaluate the resulting model on all previous tasks similar to [4], [15] after training for each new task. We use ACC as the average test classification accuracy across all classes for GZSL, seen classes, and unseen classes for GSL to measure our model’s performance. To measure forgetting, we calculate backward transfer, BWT that says how much learning new tasks have influenced previous tasks’ performance. While $BWT > 0$ indicates catastrophic forgetting and $BWT < 0$, learning new tasks has helped improve performance on previous tasks. We calculate forgetting measure for seen classes only.

$$BWT = \frac{1}{T-1} \sum_{t=1}^{T-1} [R_{t,t}^{seen} - R_{T,t}^{seen}] \quad (5)$$

$$mSA = \frac{1}{T} \sum_{t=1}^T R_{t,t}^{seen} \quad (6)$$

mSA is the mean seen classification accuracy across all tasks.

$$mUA = \frac{1}{T-1} \sum_{t=1}^{T-1} R_{t,t}^{unseen} \quad (7)$$

mUA is the measure of zero-shot learning performance for the model.

$$mOA = \frac{1}{T} \sum_{t=1}^T R_{t,t}^{overall} \quad (8)$$

mOA is the measure of generalized zero-shot learning performance.

$$mH = \frac{1}{T-1} \sum_{t=1}^{T-1} \left[\frac{2 * R_{t,t}^{seen} * R_{t,t}^{unseen}}{R_{t,t}^{seen} + R_{t,t}^{unseen}} \right] \quad (9)$$

mH is the harmonic mean classification accuracy.

Dataset	Semantic Dim	#Images	#SC	#UC
CUB	312	11788	150	50
aPY	64	15339	20	12
AWA1	85	30475	40	10
AWA2	85	37322	40	10

TABLE I: Datasets and their statistics, Where SC and UC are seen and unseen classes respectively.

Algorithm 1 Continual Zero-shot Learning

Input: $(\mathcal{X}, \mathcal{Y}, tar) \sim D^{all}$

Parameters: $\theta_S, \theta_P, \theta_D, \theta_c$

Output: \hat{X}_S, \hat{X}_P

```

1:  $D^{gen} \leftarrow \{\}$ 
2: for  $t \leftarrow 1$  to  $T$  do
3:   for  $e \leftarrow 1$  to epochs do
4:     for  $k \leftarrow 1$  to  $S_{steps}$  do
5:       Compute  $L_{task}$  using  $(\mathcal{X}^t, \mathcal{Y}^t, tar^t) \in D^t$ 
6:       Compute  $L_{adv}$  using  $(\mathcal{X}^t, \mathcal{Y}^t, t) \in D^t$ 
7:       Compute  $L_{VAE}^S$  for shared module using
          $(\mathcal{X}^t, \mathcal{Y}^t) \in D^t$ 
8:       Compute  $L_{VAE}^P$  for private module using
          $(\mathcal{X}^t, \mathcal{Y}^t) \in D^t$ 
9:        $L^{(t)} = \lambda_1 L_{adv} + \lambda_2 L_{task} + \lambda_3 L_{VAE}^S + \lambda_4 L_{VAE}^P$ 
10:       $\theta'_S \leftarrow \theta_S - \alpha_S \nabla L^{(t)}$ 
11:       $\theta'_P \leftarrow \theta_P - \alpha_P \nabla L^{(t)}$ 
12:    end for
13:    for  $j \leftarrow 1$  to  $D_{steps}$  do
14:      Compute  $L_{adv}$  for  $D$  using  $(S(x)^t, tar^t)$  and  $(z' \sim \mathcal{N}(\mu = 0, \Sigma = 1), tar = 0)$ 
15:       $\theta'_D \leftarrow \theta_D - \alpha_D \nabla L^{(t)}$ 
16:    end for
17:  end for
18:  Generate data from the shared module for seen and
  unseen classes to train a separate classifier.
19:   $D \leftarrow D_{seen} \cup D_{unseen}$ 
20:  for  $C_e \leftarrow 1$  to  $C_{epochs}$  do
21:    Compute  $L_{class}$  using  $(\mathcal{X}, tar) \in D$ 
22:     $\theta'_c \leftarrow \theta_c - \alpha_c \nabla L_{class}$ 
23:  end for
24:  Test the classifier for seen data.
25:  Test the classifier for unseen data(ZSL).
26:  Test the classifier for all seen and unseen data(GZSL).
27:  for  $c \leftarrow 1$  to  $C$  do
28:     $C$  is the replay classes.
29:    for  $i \leftarrow 1$  to  $n$  do
30:       $n$  is the number of samples to be generated per
      class for the experience replay.
31:       $(\mathcal{X}_i, \mathcal{Y}_i, tar_i) \sim D^{gen}$ 
32:    end for
33:  end for
34:   $D^{t+1} \leftarrow D^{t+1} \cup D^{gen}$ 
35: end for

```

Where $R_{j,i}$ is the test classification accuracy on task i after sequentially finishing learning the j^{th} task.

IV. EXPERIMENTS

This section consists of baselines and implementation details we used in our experiment.

A. Baselines

The research on continual zero-shot learning(CZSL) has been less explored. References [7], [22] has investigated the work before on a single-head setting that we represent in this paper. Reference [7] used the following baselines, so we do the same in this paper.

- **AGEM + CZSL** [22]: It is an average gradient episodic memory-based continual zero-shot learning. The authors of [22] have mentioned the harmonic mean of the CUB dataset only.
- **SEQ + CVAE** [7]: The authors train CVAE sequentially without considering any continual learning strategy. After SEQ+CVAE is trained on the current task, synthetic are generated using noise and class embeddings for all classes to train a separate classifier.

B. Other Methods

We also compare our results with CZSL-CV+mof [7], CZSL-CV+rb [7], and CZSL-CV+res [7].

C. Implementation Details

We use Pytorch as our framework. We train our model with a hundred epochs and a classifier for thirty epochs for each task on all datasets except CUB. We use the same number of epochs for the CUB dataset till task fifteen and then reduce to fifty for the model and ten for the classifier till task eighteen and again decrease to twenty for the model and five for the classifier. The Adam [11] optimizer is used in all experiments, and the learning rate for the classifier and others is 0.0001 and 0.001, respectively. We use weight decay 0.0001 as a regularizer for the classifier. We use 500 hidden units for both shared and private modules. Latent dimension is 50, and batch size for both model and classifier is 61. We take $\lambda_1 = \lambda_2 = \lambda_3 = 1$, and $\lambda_4 = 0.5$.

V. RESULTS AND DISCUSSION

In the first set of experiments, we measure mSA, mUA, mH, and compare it against state-of-the-art methods on CUB, aPY, AWA1, and AWA2 datasets.

A. Performance on CUB Dataset

We divide the CUB dataset into 20 tasks, where each consists of ten classes. We compare our results with several methods in Table III. Generalized CZSL [7] used memory to store real data as the replay, and it achieved mH = 20.15. In contrast, our model achieves mH = 17.41 when each task gets trained for 100 epochs till the 15th task, then from task 16th to 18th for 50 epochs and rest of the tasks for 20 epochs. Similarly, the classifier gets trained for 30 epochs till the task 15th, then from task 16th to 18th for 10 epochs and rest of the tasks for 5 epochs. The performance of our model per task is given in Figures [3, 4].

CUB				
Methods	mSA	mUA(ZSL)	mH	mOA(GZSL)
AGEM+CZSL [22]	-	-	13.20	-
Seq-CVAE [7]	24.66	8.57	12.18	-
CZSL-CV+mof [7]	43.73	10.26	16.34	-
CZSL-CV+rb [7]	42.97	13.07	19.53	-
CZSL-CV+res [7]	44.89	13.45	20.15	-
ours				
A-CZSL**	34.25	12.42	17.41	22.40
A-CZSL*	34.47	12.00	17.15	21.72

TABLE II: Results for the CUB dataset, where mSA: Mean Seen Accuracy, mUA: Mean Unseen Accuracy, mH: Hermonic Mean Accuracy, mOA = Mean Overall Accuracy. The best results in the table are presented in boldface. (**) and (*) denote that the model has been trained without and with adversarially, respectively.

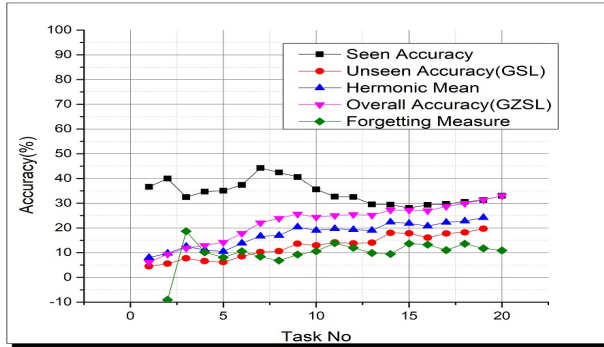


Fig. 3: Results without adversarial training for the CUB dataset.

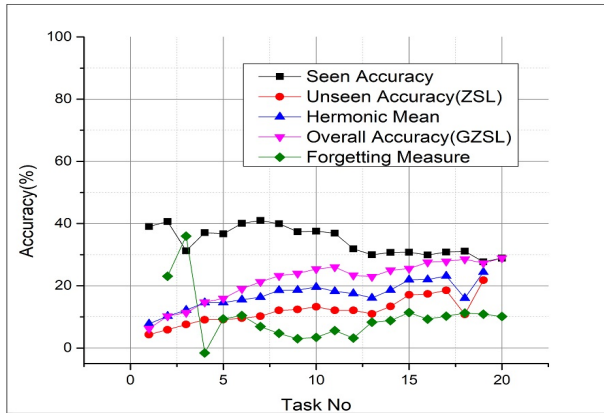


Fig. 4: Results with adversarial training for the CUB dataset.

B. Performance on aPY Dataset

We divide the aPY dataset into 4 tasks, where each consists of eight classes. We compare our results with several methods in Table [III]. Generalized CZSL [7] used memory to store real data as the replay, and it achieved mH = 23.90. In contrast, our model achieves mH = 23.05 when each task gets trained for 100 epochs. Similarly, the classifier gets trained for 30 epochs for all tasks. The performance of our model per task is given in Figures [5, 6].

aPY				
Methods	mSA	mUA(ZSL)	mH	mOA(GZSL)
AGEM+CZSL [22]	-	-	-	-
Seq-CVAE [7]	51.57	11.38	18.33	-
CZSL-CV+mof [7]	64.91	10.79	18.27	-
CZSL-CV+rb [7]	64.45	11.00	18.60	-
CZSL-CV+res [7]	64.88	15.24	23.90	-
ours				
A-CZSL**	58.14	15.91	23.05	38.20
A-CZSL*	55.46	11.2	18.63	35.97

TABLE III: Results for the aPY dataset, where mSA: Mean Seen Accuracy, mUA: Mean Unseen Accuracy, mH: Hermonic Mean Accuracy, mOA = Mean Overall Accuracy. The best results in the table are presented in boldface. (**) and (*) denote that the model has been trained without and with adversarially, respectively.

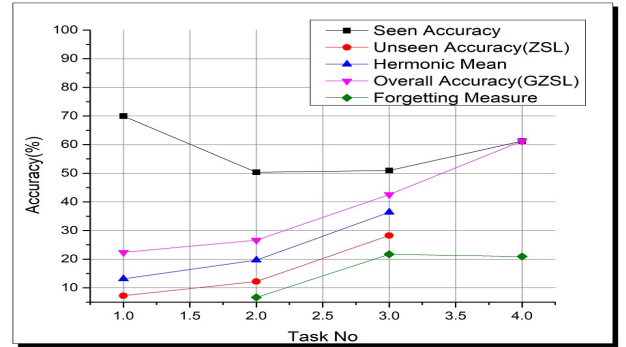


Fig. 5: Results without adversarial training for the aPY dataset.

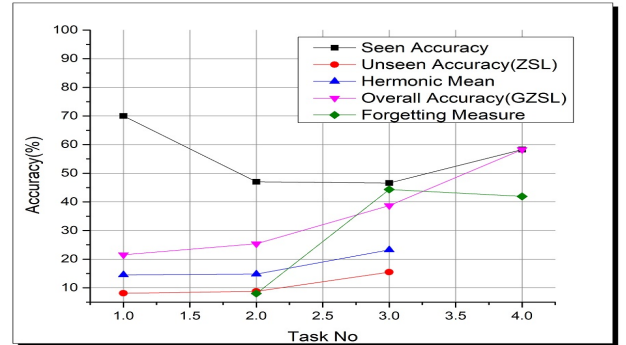


Fig. 6: Results with adversarial training for the aPY dataset.

C. Performance on AWA1 Dataset

We divide the AWA1 dataset into 5 tasks, where each consists of ten classes. We compare our results with several methods in Table [IV]. Generalized CZSL [7] used memory replay to alleviate the catastrophic forgetting and used memory to store real data, and it achieved mH = 35.51. In contrast, our model achieves mH = 35.75 when each task gets trained for 100 epochs. Similarly, the classifier gets trained for 30 epochs for all tasks. The performance of our model per task is given in Figures [7, 8].

D. Performance on AWA2 Dataset

We divide the aPY dataset into 5 tasks, where each consists of ten classes. We compare our results with several methods

AWA1				
Methods	mSA	mUA(ZSL)	mH	mOA(GZSL)
AGEM+CZSL [22]	-	-	-	-
Seq-CVAE [7]	59.27	18.24	27.14	-
CZSL-CV+mof [7]	76.77	19.26	30.46	-
CZSL-CV+rb [7]	77.85	21.90	33.64	-
CZSL-CV+res [7]	78.56	23.65	35.51	-
ours				
A-CZSL**	67.98	20.58	30.83	43.08
A-CZSL*	71.00	24.26	35.75	50.9

TABLE IV: Results for the AWA1 dataset, where mSA: Mean Seen Accuracy, mUA: Mean Unseen Accuracy, mH: Hermonic Mean Accuracy, mOA = Mean Overall Accuracy. The best results in the table are presented in boldface. (**) and (*) denote that the model has been trained without and with adversarially, respectively.

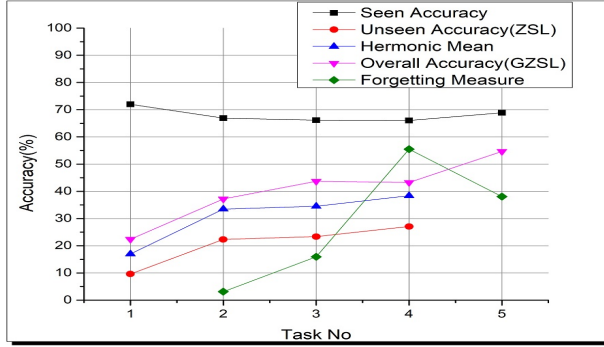


Fig. 7: Results without adversarial training for the AWA1 dataset.

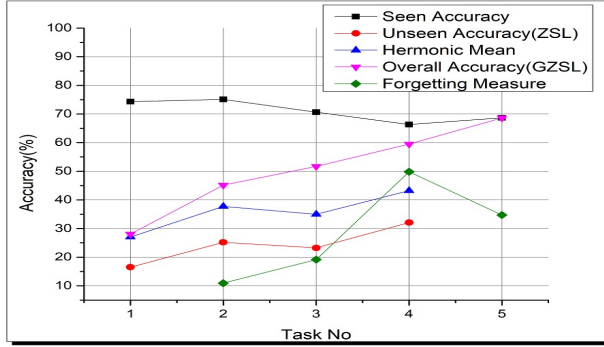


Fig. 8: Results with adversarial training for the AWA1 dataset.

AWA2				
Methods	mSA	mUA(ZSL)	mH	mOA(GZSL)
AGEM+CZSL [22]	-	-	-	-
Seq-CVAE [7]	61.42	19.34	28.67	-
CZSL-CV+mof [7]	79.11	24.41	36.60	-
CZSL-CV+rb [7]	80.92	24.82	37.32	-
CZSL-CV+res [7]	80.97	25.75	38.34	-
ours				
A-CZSL**	70.05	22.85	32.98	44.97
A-CZSL*	70.16	25.93	37.19	51.55

TABLE V: Results for the AWA2 dataset, where mSA: Mean Seen Accuracy, mUA: Mean Unseen Accuracy, mH: Hermonic Mean Accuracy, mOA = Mean Overall Accuracy. The best results in the table are presented in boldface. (**) and (*) denote that the model has been trained without and with adversarially, respectively.

in Table [V]. Generalized CZSL [7] used memory to store real data as the replay, and it achieved $mH = 38.34$. In contrast, our model achieves $mH = 37.19$ when each task gets trained for 100 epochs. Similarly, the classifier gets trained for 30 epochs for all tasks. The performance of our model per task is given in Figures [9, 10].

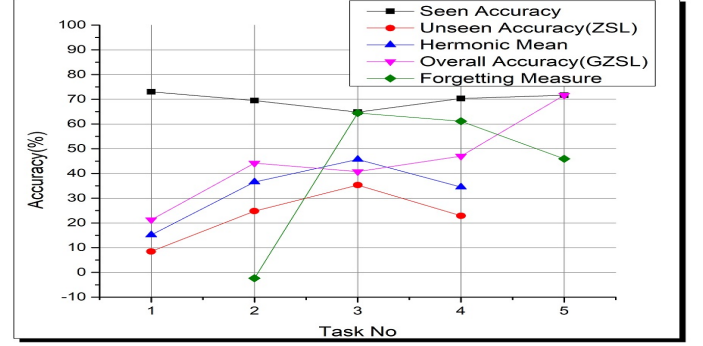


Fig. 9: Results without adversarial training for the AWA2 dataset.

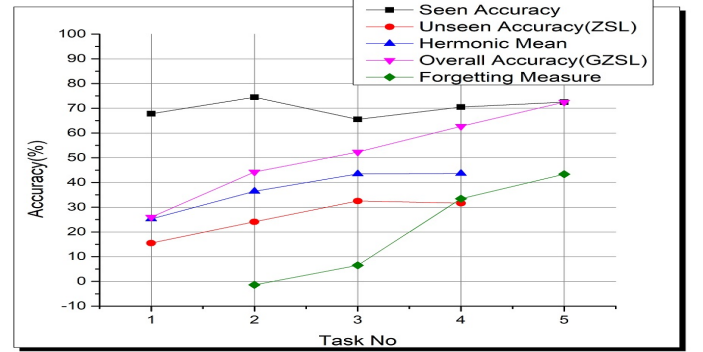


Fig. 10: Results with adversarial training for the AWA2 dataset.

VI. CONCLUSION

In this work, we proposed a novel hybrid algorithm. The novelty of our work is that we use adversarial learning. Here the model needs generative replay and it grows for task incremental learning. The private module barely shares knowledge from seen to unseen classes that can be future work to optimize the private module for ZSL. Another future work might be to develop task-free continual zero-shot learning. How can we build a continual zero-shot learning model for object detection? What should be the optimum latent dimension, hidden-layer size?

ACKNOWLEDGEMENT

We thank prof Suresh Sundaram and Dr. Chandan Gautam from the Artificial Intelligence(AI) lab in the Aerospace Engineering Department, Indian Institute of Science(IISc), Bangalore, for the valuable discussion in the initial phase of this work.

REFERENCES

- [1] Wickliffe C Abraham and Anthony Robins. Memory retention—the synaptic stability versus plasticity dilemma. *Trends in neurosciences*, 28(2):73–78, 2005.
- [2] Samaneh Azadi, Deepak Pathak, Sayna Ebrahimi, and Trevor Darrell. Compositional gan: Learning image-conditional binary composition. *International Journal of Computer Vision*, 128(10):2570–2585, 2020.
- [3] Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- [4] Sayna Ebrahimi. Continual learning with neural networks. 2020.
- [5] Rafael Felix, Ian Reid, Gustavo Carneiro, et al. Multi-modal cycle-consistent generalized zero-shot learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 21–37, 2018.
- [6] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016.
- [7] Chandan Gautam, Sethupathy Parameswaran, Ashish Mishra, and Suresh Sundaram. Generalized continual zero-shot learning. *arXiv preprint arXiv:2011.08508*, 2020.
- [8] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014.
- [9] Tyler L Hayes, Nathan D Cahill, and Christopher Kanan. Memory efficient experience replay for streaming learning. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 9769–9776. IEEE, 2019.
- [10] S Herculano-Houzel. The human brain in numbers: a linearly scaled-up primate brain. *front. hum. neurosci.* 3, 31 (2009), 2009.
- [11] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [12] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [13] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [14] Yuanpeng Li, Liang Zhao, Kenneth Church, and Mohamed Elhoseiny. Compositional language continual learning. In *International Conference on Learning Representations*, 2019.
- [15] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *arXiv preprint arXiv:1706.08840*, 2017.
- [16] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.
- [17] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989.
- [18] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.
- [19] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [20] Aashish Sheshadri, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes. *Computer Vision and*, pages 1778–1785, 2012.
- [21] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5972–5981, 2019.
- [22] Ivan Skorokhodov and Mohamed Elhoseiny. Normalization matters in zero-shot learning. *arXiv preprint arXiv:2006.11328*, 2020.
- [23] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017.
- [24] Vinay Kumar Verma and Piyush Rai. A simple exponential family framework for zero-shot learning. In *Joint European conference on machine learning and knowledge discovery in databases*, pages 792–808. Springer, 2017.
- [25] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [26] Kun Wei, Cheng Deng, and Xu Yang. Lifelong zero-shot learning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 551–557, 2020.
- [27] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pages 3987–3995. PMLR, 2017.