# Generative Model with Coordinate Metric Learning for Object Recognition Based on 3D Models

Yida Wang* and Weihong Deng*, *Member, IEEE*

*School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing, BJ 100876 China

*Abstract*—Collecting data for deep learning is so tedious which makes it hard to establish a perfect database. In this paper, we propose a generative model trained with synthetic images rendered from 3D models which can reduce the burden on collecting real training data and make the background conditions more sundry. Our architecture is composed of two sub-networks: semantic foreground object reconstruction network based on Bayesian inference and classification network based on multi-triplet cost training for avoiding over-fitting on monotone synthetic object surface and utilizing accurate informations of synthetic images like object poses and lightning conditions which are helpful for recognizing regular photos. Firstly, our generative model with metric learning utilizes additional foreground object channels generated from semantic foreground object reconstruction sub-network for recognizing the original input images. Multi-triplet cost function based on poses is used for metric learning which makes it possible training an effective categorical classifier purely based on synthetic data. Secondly, we design a coordinate training strategy with the help of adaptive noises applied on inputs of both of the concatenated sub-networks to make them benefit from each other and avoid inharmonious parameter tuning due to different convergence speed of two sub-networks. Our architecture achieves the state of the art accuracy of 50.5% on ShapeNet database with data migration obstacle from synthetic images to real photos. This pipeline makes it applicable to do recognition on real images only based on 3D models.

*Index Terms*—Bayesian rendering, triplet cost, synthetic image, semantic reconstruction, coordinate training, metric learning.

(a) Flowchart of generative model with coordinate metric learning using 3D models. **Top:** training stage. **Bottom:** testing stage.



(b) Reconstruction network    (c) Classification network

Fig. 1: Generative model with coordinate metric learning.

## I. INTRODUCTION

**D**EEP architectures based on convolutional operations such as AlexNet [1], Inception concept [2] and ResidualNet [3] work well on image classification trained on large scale real datasets such as ImageNet [4] and PASCAL [5]. Meanwhile, synthetic data could also be served as training data for specific tasks such as pose estimation [6], [7]. Triplet cost function [6] used in metric learning utilizes these accurate labels well by making arrangements for triplet set according to corresponding labels of synthetic data, so manifold learning is more applicable based on synthetic data than realistic data. Pose labels are also useful for exploiting geometric information, recent works on object rotation and deformation combining CNN with auto-encoder [8], [9] solve tasks such as face rotation and intrinsic transformations for objects based on additional pose information. Other supervision information such as remote codes [8] gives a help on image transformation
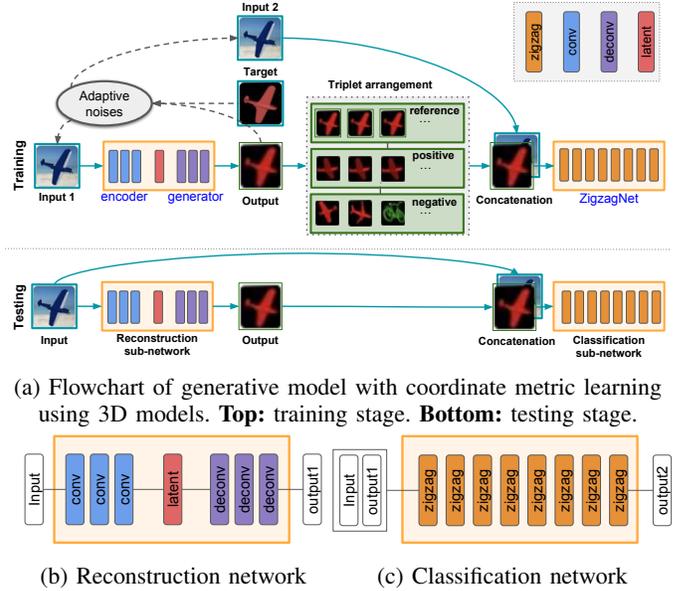
Corresponding author: W Deng (email: whdeng@bupt.edu.cn).

in regard to different purposes which modify auto-encoder to flexible reconstructor.

Though typical methods are effective in multi-task learning such as pose estimation [6] and categorical classification [10] with additional information recorded by special equipments, performance of testing on real photos with severe data migration problem compared to synthetic training data is still not satisfied. Works on foreground object segmentation [11] try to solve data migration problem from synthetic training data and real testing photos concatenating segmentation channel to original photo as input for classification architecture, but those concatenated architecture still have over fitting problem on original RGB channels without fully utilization on robust object segmentation result.

In this paper, we propose a deep architecture targeted on training with synthetic data and recognizing on real photos especially for classification. Such graph concatenated with two sub-networks which are reconstruction network based on Bayesian inference and classification network based on ZigzagNet. Our reconstruction architecture trained from two types of synthetic data rendered from models of ShapeNet [12] database reconstructs foreground objects from background with surface colors representing categories. Inspired by pos-

itive impact of depth images for pose estimation [6] and face rotation [8], three additional channels provided by object reconstruction sub-network are concatenated to original RGB channels altogether as a six-channel data feeding for classification sub-network. Pixel to pixel prediction is possible based on encoding and decoding architecture such as Fully Convolutional Network (FCN) [13], [14], Deconvolutional Network (DeconvNet) [15] and Variational Auto-encoder (VAE) [16], [17]. Inspired by good performance introduced by statistical latent variables in generative model [17], we modify conditional VAE [18], [19] to make the reconstruction sub-network able to generate semantic foreground objects robustly without feeding in additional supervision information for testing. Textured models from ShapeNet [12] database are used for avoiding over-fitting on monotone surface of texture-less models such as models in PASCAL 3D+ [20]. Our classification sub-network is optimized with a joint cost function composed of categorical softmax loss and pose assisted triplet loss. Such optimization tries to avoid over fitting on the unrealistic relationship between synthetic objects and background by means of learning discriminant geometric features according to different camera poses. Gradients from our joint cost function for object classification is back-propagated to object reconstruction sub-network which is the head of concatenated architecture as a supervision information like remote codes [8] to help reconstruct object in different categories accordingly. Our conjugate network aiming at reconstructing foreground objects and doing classification on categories in the same time is trained efficiently with adaptive noises for both inputs of two sub-networks. With the help of coupled noises calculated from variances of reconstructed output and ground truth mask applied on both inputs of two sub-networks, problem of gradient vanishing from classification network back to the output channels of reconstruction network is avoided. This means that the reconstructed object masks are more useful for training classification network as additional information.

Experiments on categorical prediction task for objects in real images based on photos provided in PASCAL 3D+ [20] show that our concatenated architecture trained merely on synthetic data solves severe data migration problem from synthetic training to realistic testing data well which achieves the state of the art accuracy of 50.5% on ShapeNet database without using Nearest Neighbor classifier in LDO [6]. Utilization on pose information with our joint cost function makes the classification accuracy much closer to the ones trained on real photos than some recent works like LDO [6] and SR [11].

In the following chapters, we describe our method for training deep architecture merely based on 3D models aiming for recognizing objects in real photos. Our theory could be roughly divided into those parts:

- General rendering strategy from 3D models to 2D images is described in Section III. Additional information prepared for metric learning such as poses are also recorded while generating images.
- Section IV-A describes our generative model targeted on pixel-wise reconstruction which provides additional information including segmentation masks with semantic colors and depth information.

- In Section V, we describe the second part of the conjugate architecture. Classification sub-network is modified into compact structure based on ZigzagNet in Section V-A to avoid over fitting.
- Furthermore, parametric model of classification sub-network is also optimized with metric learning in Section V-B supervised by category labels and pose information which also make it possible back-propagating effective gradients to reconstruction sub-network.
- Finally, paired adaptive noises contribute to a coordinate training process for both sub-networks by means of making adaptive corruptions on both of the RGB inputs channels which is described in Section VI.

## II. RELATED WORKS

Some methods related to training on synthetic data generated with rendering techniques [21], [22] are useful for tasks in specific conditions such as pose estimation [6], [7] in indoor conditions because camera positions are used for image rendering which could be recorded accurately. Although synthetic data could hardly be as realistic as real photos, some methods still utilize more information to overcome data migration problem for classification. For example, SR [11] utilize two types of synthetic data for training a concatenated architecture for final task of classification. One core task in this paper is utilizing 3D models with rich information for training deep parametric model which could be applied on real world visual recognition tasks. We train deep neural network directly from images rendered from 3D models and test on real photos, especially for classification.

Metric learning is helpful to learn and exploit implicit relationships among data samples. Descriptor-based metric learning methods such as Triplet training [6] and Siamese training [23] utilize additional informations for person re-identification, 3D object pose estimation and kinship verification by learning specific manifold in the embedded descriptor space. Metric learning on embeddings such as [24] uses pair-wise distances which is helpful for retrieval tasks. Joint Siamese and Triplet training is also applied on local descriptor optimization for patch matching [25]. Here we utilize a multi-triplet loss for optimizing a concatenated architecture to make it applicable for training on synthetic data and testing on real data.

## III. TRAINING DATA GENERATION

The first mission is representing 3D information in 2D space for training data generation based on 3D models. We construct a new database which consists of two types of synthetic images rendered from 3D models of ShapeNet [12] database which are the one including textured objects with background images crawled from Flickr and another one including textureless objects without any background. The second one will be used as semantic object masks in training stage, so we render the objects with monotone color in regard to specific category. Both types of synthetic images are rendered with the same camera parameters accordingly.

## A. Semantic Rendering

Inspired by the fact that depth images can provide complementary information for recognition on real images [6] by distinguishing foreground target from background, we embed an semantic foreground object reconstruction sub-network in the whole generative model based on Bayesian coding. The output of the reconstruction network is concatenated to the original images as a six channel data forwarded to the classification network. These concatenated channels reconstructed from realistic input images with background are expected to be similar to segmented object image with unique colors according to category.

As shown in Fig. 1a, our goal is training a robust semantic rendering sub-network which provides additional information for classification sub-network with the aim similar to the one of the self-restraint foreground object reconstruction network [11]. The reconstructed semantic masks with depth information is concatenated to the original RGB channels together as a 6 channel input to the next classification network. Such concatenation also makes the statistical gradient of the classifier able to be propagated back to the reconstructed network. Reconstructed channels for foreground objects are expected to contain segmentation information on contours, depth information and semantic color information which means that they do not only reflect informative details about the foreground objects, but also help to train the classification network with visual categorical information.

2D synthetic training images are rendered from 3D models in database with texture files like ShapeNet [12]. Foreground objects are shot from vertexes on a sphere net combined with triangles recursively generated from larger triangles which are initiated from a regular icosahedron [26] shown in Fig. 2b. we cut off the upper and lower parts of the whole sphere along with two meridians vertical to $z$ axis cause range of view positions are always limited in real world. For the coordinate $(V_x, V_y, V_z)$ of camera positions shown in Fig. 2a, $V_z \in [-0.1\min(V_x), 0.6\max(V_x)]$. Focal points are set randomly within a range around the center of each object to simulate photos taken as close shots. As we want to utilize realistic data to train a semantic foreground objects reconstruction network afterwards, background images are crawled from *Flickr* where target objects are not included.

For fine-tuning with more data, we also render another collection of close shot images with shifted focal points which are moved from the center of object to the head of objects. Shifting distance from the original focal points in the center could be represent as $(F'_x, F'_y, F'_z) = (F_x, F_y, F_z) + 0.2((C_x, C_y, C_z) - P_{axis})$ where $F$ represents the focal points used for the camera shots, $C$ represents the camera position on the sphere and $P_{axis}$ represents the intersection point of the axis of front view and the sphere net of camera positions.

Mean and normalized standard deviation images from ShapeNet and ImageNet in Fig. 3 show that gap between synthetic ones and real ones is still too large to train a recognition model directly from synthetic data. A challenging task is overcoming data migration problem between synthetic data and real photos because we want to get rid of the over-



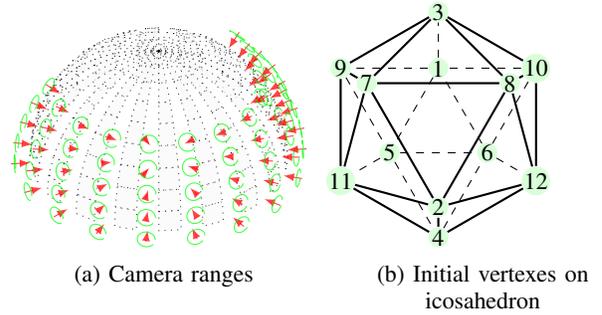(a) Camera ranges     (b) Initial vertexes on icosahedron

Fig. 2: Synthetic camera positions on the vertexes of a sphere.

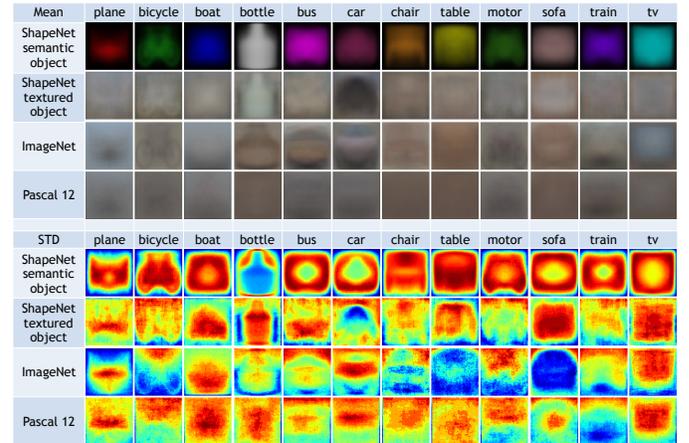fitting problem introduced by our realistic data.



Fig. 3: Mean value and normalized heat map for standard deviation of images including 2 types of synthetic images (upper 2 rows) and 2 types of real photos (lower 2 rows) in large scale database.

## IV. SEMANTIC RECONSTRUCTION

Our semantic reconstruction removes background environments from foreground objects and renders them with specific color according to specific category for further classification, so images with textured foreground objects will have the same monotone surface color if they are coming from the same category. Reconstruction sub-network is trained from paired synthetic images which are realistic images with background pictures and semantic segmentation masks of foreground objects. We apply two methods to make the reconstruction network capable on semantic rendering for foreground objects. Firstly, we design a generative model using Bayesian inference to make the generator more capable on perceiving global information on objects compared to traditional pixel-wise coding methods like FCN [13] and DeconvNet [15]. Other supervision assistants for reconstruction like remote code described in face rotating [8] for multi-task learning according to the identities of objects are discarded while our gradients of joint loss acting as supervision in classification sub-network is back-propagated into the reconstruction network. Then adaptive noises are applied automatically on input of the reconstruction sub-network to fully utilize the back-propagated gradients from

classification sub-network sequentially concatenated with the generated output of reconstruction sub-network, so our whole architecture is called generative model with coordinate metric learning. Both of the conjugate sub-networks are coupled with each others by utilizing the adaptive noises.

### A. Pixel-wise Generative Model

Let us suppose $X$ to represent the input samples. Our probabilistic variational framework is composed of two parts: an encoder $q(X)$, that compresses each input $X$ into latent variables $z$, and a generator or decoder $p(z)$ that maps $z$ into the output $Y$. Generally speaking, our model could be regarded as a variational coder (VC) specifically designed for image related tasks.

We represent the prediction task $p(q(X))$ as a stochastic process modeled by the two probability distributions $P(Y|z)$ and $P(z)$[1].

Inspired by VAE, a Gaussian distribution $P(z) \sim N(0, I)$ with the identity matrix $I$ is adopted for modeling the latent variables $z$, this meaning that the likelihood $P(Y|z)$ also follows a Gaussian distribution: $Y \sim N(p(z), \sigma^2 * I)$ with covariance defined by the scalar $\sigma$ having the same dimension as the number of latent variables.

One of the main ideas behind our approach is to optimize both the encoder $p(X)$ and the generator $p(z)$. The negative log probability of the output $Y$ is proportional to the squared Euclidean distance between the output $p(z)$ and the expected target $Y$. Therefore, the main task is how to optimize $p(X)$. Since in theory the latent variables $z$ should be able to reproduce any output $Y$, the target of the optimization procedure for the encoder function $p(X)$ is set to be the matching posterior $P(z|Y)$, which means that the latent variables are more likely to produce expected outputs.

Firstly, we need to apply Bayesian rule

$$P(z|Y) = \frac{P(Y|z)P(z)}{P(Y)} = \frac{P(Y|z)P(z)}{\int P(Y|z)P(z)dz} \quad (1)$$

to calculate the true posterior $P(z|Y)$ from training samples. The problem is that the range of $z$ is not limited at all, so the likelihood $P(Y|z)$ is almost zero for most $z$. Thus, we can not calculate this likelihood directly from the samples. To solve this problem, we set some constraints on the latent variables $z$ to make them more likely to generate $Y$. Meanwhile, since the optimization of the encoder $p(X)$ suffers a similar limitation induced by the limited output range, we redefine the encoder function as $q(X)$, an alternative for $p(X)$, such that the associated probability density function $Q(z|X)$ is expected to match the posterior $P(z|Y)$.

Under these conditions, the expectation of the likelihood $E_{z \sim Q} P(Y|z)$ conditioned on the latent variables should be close to the true probability $P(Y)$. Kullback-Leibler (KL) divergence between $P(z|Y)$ and $Q(z|X)$ can be used to evaluate the capability of the encoding network to generate

latent variables which are likely to produce the expected target $Y$

$$D[Q(z|X)||P(z|Y)]$$
$$= E_{z \sim Q}[\log Q(z|X) - \log P(Y|z) - \log P(z)] + \log P(Y). \quad (2)$$

When we expand the KL divergence in (2) to optimize $Q(z|X)$ based on the Bayes rule in (1), $P(Y)$ still exists for calculating $P(z|Y)$ which is not tractable due to the same reason mentioned above.

So we rewrite (2) by combining the KL divergence with $\log P(Y)$ and represent them as the new optimization target $B(P, Q)$, which is the evidence lower bound of $P(Y)$, to get rid of the process for calculating $P(Y)$ directly for optimizing $Q(z|X)$ in divergence

$$B(P, Q) = \log P(Y) - D[Q(z|X)||P(z|Y)] . \quad (3)$$

So our target becomes maximizing the evidence lower bound $B$ with respect to $\log P(Y)$ to make $D[Q(z|X)||P(z|Y)]$ as small as possible. To implement our target via deep learning, we define the cost function $\mathcal{L}_{\text{enc-gen}} = -B$ of the generative model.

By expanding the divergence $D[Q(z|X)||P(z|Y)]$ in (3) and combining $\log Q(z|X)$ and $-\log P(z)$ together as another divergence, we can obtain the final form of the cost function as our encoder and decoder

$$\mathcal{L}_{\text{enc-gen}} = E_{z \sim Q}[\log P(Y|z) + \log P(z) - \log Q(z|X)]$$
$$= \underbrace{D[Q(z|X)||P(z)]}_{\mathcal{L}_{\text{enc}}} \underbrace{-E_{z \sim q}[\log P(Y|z)]}_{\mathcal{L}_{\text{gen}}} . \quad (4)$$

The final loss function in (4) can be interpreted as follows: we replace the information component $\log P(X|z)$ used in VAEs as $\log P(Y|z)$ to make the decoding process suitable for different targets $Y$.

Traditionally, foreground object reconstruction is treated as pixel-wise segmentation using one-hot coding like DCN [15], so the number of output channel is no fewer than the number of categories of training data. In our work, semantic reconstruction results are represented in RGB channels which could be easily applied on different databases because the number of channels are fixed as 3.

### V. Categorical Classification with Metric Learning

Object recognition is dependent on categorical information, so we concatenate a classification network after the generative sub-network. Assuming that a descriptor of the reconstructed channels from one sample is represented as $r = [pixels_r, pixels_g, pixels_b]$ where $pixels$ are feature for a single channel, the effectiveness of back-propagated Euclidean loss directly depends on the assigned absolute RGB values of masks which means that greater disparity between two classes matters more than the minor one. This is a disaster when there are so many categories that we can't attribute effective RGB colors anymore for a discriminative rendering.

We make features of our training data and real photos similar to each other by exploiting common factors to minimize

---

[1]We use capital letters for representing probability density function and lower case letters for normal functions which is deterministic from input to output.

over-fitting problem triggered by absence of textures. Based on additional information such as object poses, ZigzagNet is used as architecture for pre-training which is optimised by a multi-triplet cost function together with an optional pairwise term for speeding up convergence. Then additional synthetic data with different focal points is used for fine-tuning with softmax loss. Our training method makes the descriptors in classification sub-network suitable for Nearest Neighbour (NN) classifier without softmax and support vector machine classifier due to data migration problem between synthetic data and real images. We evaluate the performance of our work on output of the penultimate layer of the classification network before the last inner product matrix. Our discriminant descriptor satisfies NN classifier better than models trained from real images and is even similar to the performance of AlexNet [1] descriptors which is also trained from realistic data using larger parametric model.

### A. Compact Classification Architecture

With reconstruction sub-network 14 MB on disk, we apply ZigzagNet for the conjugate classification which has $30\times$ fewer parameters than AlexNet and is just 6.2 MB on disk to make the whole architecture as compact as possible. So the overall model size of our conjugating model is just about 20.2 MB. Over-fitting problem from synthetic data to real photos is also overcame by applying Zigzag module instead of simply using convolutional layer.
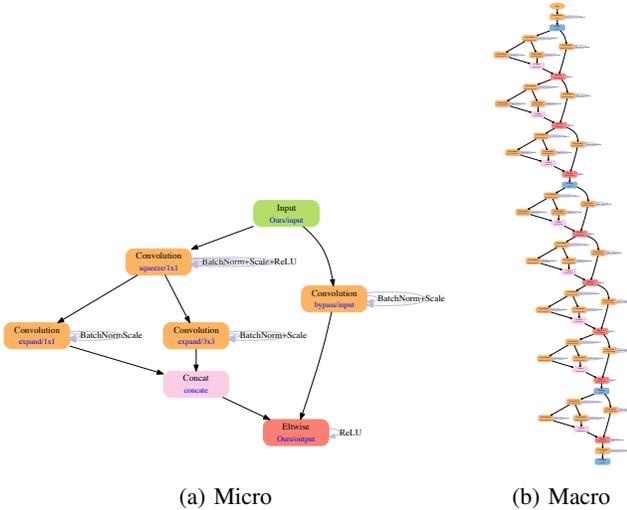


(a) Micro    (b) Macro

Fig. 4: Micro Zigzag module and macro ZigzagNet architecture.

We try to solve data migration problem between synthetic training data and real testing data by making the whole architecture more compact. Due to the fact that we want to use the concatenated RGB channels and the generated object reconstruction channels which differ a lot directly as 6-channel input for classification sub-network, so we adopt channel-wise compression between micro architectures for deep architecture to avoid over fitting problem. We modify convolutional Fire module in SqueezeNet [27] into Zigzag module to form ZigzagNet which performs better in classification compared

to SqueezeNet. ZigzagNet combines SqueezeNet with residual concept [3] without using any fully connected layers. Macro architecture of ZigzagNet is concatenated of repetitive micro Zigzag modules composed of three $1 \times 1$ and one $3 \times 3$ convolutional layers in a zigzag style of 3 steps: channel wise squeezing for principal representation, reception fields expanding by parallel connection of multi-scale convolution kernels and consolidation by adding $1 \times 1$ convolutional layer to bypass the original information before squeezing to keep the learning process stable while remaining the size of parametric model. Input to a micro Zigzag module is compressed by a $1 \times 1$ convolutional layer for channel-wise linear projection to make the following parameter model compact and representative. This is pretty important for our data generated from reconstruction sub-network with concatenated synthetic reconstruction channels and realistic RGB channels because such $1 \times 1$ convolutional operation behaves like a trainable pre-processing procedure before spatial convolutions afterwards. Output in a single module within macro architecture is joined with bypassed input information by a convolutional layer as the input of next module to keep the learning process stable. Our micro module also differs from Fire module of bypassed SqueezeNet shown in Fig. 4a in nonlinear operation, ReLU operation on the expanding layer before the element wise summation layer is moved to the output of Fire module to eliminate the scale difference between two inputs of the element wise operation layer and make information from both branches compressed at the same time. For convenience of comparison experiments, depth of macro architecture of ZigzagNet is set as the same as AlexNet and SqueezeNet in account of micro architectures with channel-wise representation. The final macro architecture of ZigzagNet is shown in Fig. 4b which is a concatenated structure of several micro modules of Zigzag module shown in Fig. 4a.

### B. Metric Learning with Multi-triplet Set

Realistic images generated by two camera modes separately are both well utilised using multi-triplet cost and softmax cost. As background images are attached behind the rendered objects for realistic images, training with category labels based on texture-less models leads to over-fitting on distinguishing edges of objects against background. If we just train classification model directly based on those synthetic data, objects without clear boundary against background in real images won't be recognized well. We introduce a triplet cost to utilize information contained in poses, lighting conditions of synthetic data which are also common attributes for real photos.

Our multi-triplet cost function which is modified from work of [6] together with an optional pairwise term is used in the initial training process to improve the performance for recognizing physical characteristics of foreground objects no matter whether there is background or not. Fig. 5 shows that triplet training based on poses makes descriptors of synthetic images form in a sphere-like distribution and Nearest Neighbor Prediction with 4 candidates in Fig. 5e indicates that trained parametric model has ability on both identification on object categories and distribution regression on object poses.

(a) Triplet set with 5 samples



(b) Softmax cost distrbution     (c) Triplet cost distribution



(d) Nearest Neighbor      (e) Nearest Neighbor
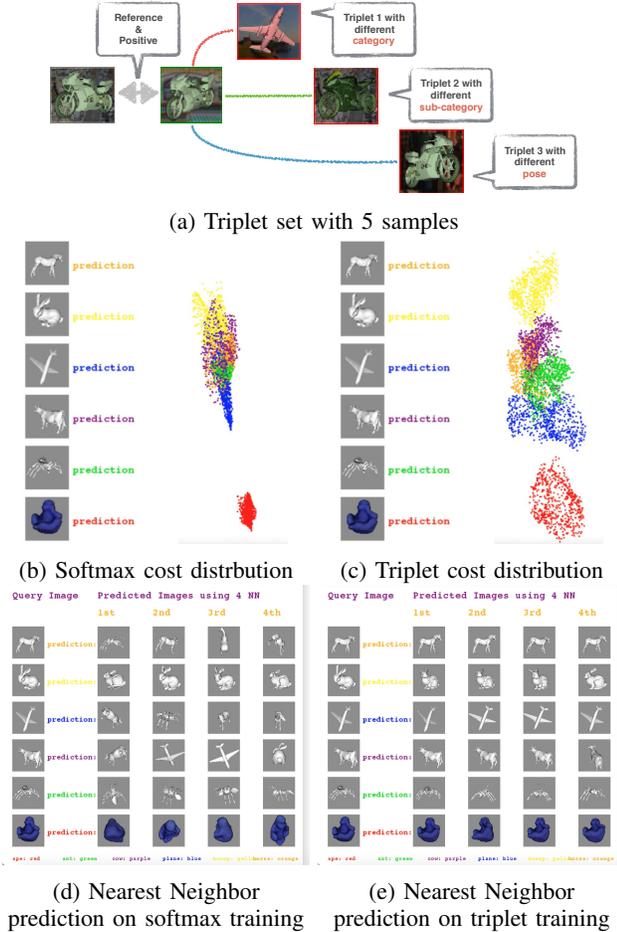prediction on softmax training    prediction on triplet training

Fig. 5: Triplet set arrangement and feature visualization with PCA matrix calculated from descriptors.

Fig. 5b shows that training directly with softmax loss based on categorical information over fits the blue surface of "Ape" which is bad for distinguishing the other 5 categories with white surface. Descriptors trained by our multi-triplet loss shown in Fig. 5c have a more reasonable distribution related to poses. Triplet loss makes it possible to optimise the parametric model with common characteristics including object poses, lighting conditions and camera modes which are accurately recorded during rendering procedure of synthetic images. Fully exploiting the information in poses can help to avoid learning bad relationships between objects and background which make those synthetic images not so realistic. One multi-triplet set is composed of a reference sample, a positive sample and several negative samples:

$$\mathcal{L}_{triplet} = \mathcal{L}_{pair}(x_i, x_j) + \sum_{(x_i, x_j, x_k) \in \mathcal{T}} \mathcal{L}_{tri}(x_i, x_j, x_k) , \quad (5)$$

where $\mathcal{L}_{pair}(x_i, x_j) = ||f(x_i) - f(x_j)||_2^2$. As shown in Fig. 5a, our multi-triplet set is composed of 5 samples, the positive sample is the one with closest pose from reference sample or only has a different lighting condition which is the same pose and 3 negative samples are selected as one differs more in pose from the same class and two from other classes or sub-classes. Our triplet loss makes the whole training process

going smoothly and effectively by solving gradient vanishing and exploding. Firstly, the problem of gradient vanishing in traditional triplet [28] training where positive sample is similar to negative sample is solved by adopting basic form of triplet loss as the one in [6]. Feature distance is modified from Euclidean distance to its squared form to make distribution of learned descriptors has a manifold of sphere as shown in Fig. 5c which is related to geometric camera positions on a sphere. This also makes the learning process more stable for 3 samples which differs little in a triplet set, so tiny difference between three samples without background with similar poses will not lead to a not-a-number gradient which might appear in [6]. Secondly, large variance in background images from Flickr makes the gradient explodes easily which is not so helpful for understanding hidden information of the foreground objects, so we further apply natural logarithms on loss of a triplet:

$$\mathcal{L}_{tri} = \ln(\max(1, 2 - \frac{||f(x_i) - f(x_k)||_2^2}{||f(x_i) - f(x_j)||_2^2 + m_{\text{tri}}})) , \quad (6)$$

where $f(x)$ is the input of the loss for sample $x$ and $m_{\text{tri}}$ is the margin for triplet. Denote that $D_{ij} = ||f(x_i) - f(x_j)||_2^2$ and $D_{ik} = ||f(x_i) - f(x_k)||_2^2$, so the partial differential equations for the input of triplet loss layer are:

$$\frac{\partial \mathcal{L}_{tri}}{\partial f(x_i)} = \frac{D_{ik}(f(x_i) - f(x_j)) - (D_{ij} + m_{\text{tri}})(f(x_i) - f(x_k))}{\mathcal{L}_{tri}(D_{ij} + m_{\text{tri}})^2}$$
$$\frac{\partial \mathcal{L}_{tri}}{\partial f(x_j)} = \frac{D_{ik}(f(x_j) - f(x_i))}{\mathcal{L}_{tri}(D_{ij} + m_{\text{tri}})^2}$$
$$\frac{\partial \mathcal{L}_{tri}}{\partial f(x_k)} = \frac{f(x_i) - f(x_k)}{\mathcal{L}_{tri}(D_{ij} + m_{\text{tri}})} \quad (7)$$

For convenience of an easier training data preparation procedure, we set the positive and the reference samples fixed in a multi-triplet set and use them directly as the pairwise term if it is needed. As training based on multi-triplet sets costs more time for convergence, samples without background are used at first stage to make parametric model converge faster. The pairwise term is only applied to make the descriptor robust to small variance of objects rather than the more complicated backgrounds images which are added in the following training stage.

### C. Fine-tuning with Additional Data

Triplet training based on pose information uses synthetic samples where focal points are all in the center of every objects. We add softmax loss for fine-tuning based on close shot images with shifting focal points. Distribution of descriptors projected with 3 dimensional PCA matrix shows that special samples from real test images could be better clustered with other normal samples in the same class after fine-tuning and all samples are set apart better according to categories. Additional rendered data gives a help on reducing intra-class variability while keeping the geometry relationship for descriptors of regular samples based on pre-trained model.

## VI. COORDINATE TRAINING FOR CONJUGATE STRUCTURE

As shown in Fig. 1a, the concatenated reconstruction and classification sub-networks both have inputs of realistic images which means that both sub-networks are relatively coordinate in training stage. Here coordinate training means that output of reconstruction sub-network always help to optimize classification sub-network and the gradients back-propagated from the classification sub-network are always effective. Our coordinate training uses coupled noising ratio $R_{rec}$ and $R_{cls}$ to corrupt realistic inputs for the aim of feeding effective information back and forth between two sub-networks. The following part introduces problems in concatenated network and introduce how our coupled noise ratios on each input of sub-networks help to efficiently train both sub-networks.

### A. Adaptive Noise Assistant for Information Feeding

We makes two sub-networks trained more effectively by applying adaptive noises for realistic input channels. Two sub-networks should not be simply concatenated together because different convergence speed of sub-networks will make it hard training the second network according to the output of the first network and utilizing the back-propagated gradients for the first network. As shown in Fig. 9a, reconstruction sub-network is only able to reconstruct average images of all semantic masks in the beginning, so these three additional channels are not so discriminant that the classification network can not utilize them compared to the realistic channels. Assume that synthetic training data is evenly distributed in categories and two types of images are concatenated together as inputs for classification sub-network, the classification results will be similar no matter whether using the reconstruction network as supporting architecture or not. The variance between original images and reconstructed images are so large that the statistically back-propagated gradients could not be effectively passed back to the reconstructed channels according to loss functions designed for classification network. We present synthetic images as object channels $O$, semantic depth masks as mask channels $M$, reconstruction network with input $I_{rec} = [O]$ produces reconstruction channels $M'$ similar to $M$ and the data fed in classification network is represented as a channel-wise concatenated tensor as $I_{cls} = [O, M']$. Direct training for classification network based on $[O, M']$ will lead to a severe problem of over-fitting on reconstructed channels because the variance of $M'$ is far smaller than the other three object channels which means that the classification network utilizes few information from reconstructed channels, so the classification result differs little even without $M'$ which means that $cls(I_{cls}) \approx cls(I_{rec})$. For perspective of statistical gradient decent for the reconstruction sub-network, gradients feeding back through $I_{cls}$ are not effective due to the absence of correlation between the classification loss and information in $I_{cls}$. This means that the two parts of the whole network are not so relevant to each other and also can not take benefits from each other's forward outputs and back-propagated gradients.

We propose an efficient strategy using adaptive noise on the realistic data separately on the input of reconstruction network and classification network to solve this problem. We utilize coupled noising function $n_{rec}$ and $n_{cls}$ for inputs of reconstruction network $I_{rec}$ and classification network $I_{cls}$ to restrict information feeding in the original realistic images separately which are defined as

$$
\begin{aligned}
n_{rec}(I_{rec}) &= R_{rec} \; Ran + (1 - R_{rec}) \; I_{rec} \\
n_{cls}(I_{cls}) &= R_{cls} \; Ran + (1 - R_{cls}) \; I_{cls} \;,
\end{aligned} \tag{8}
$$

where $Ran$ is Gaussian noise, $R_{rec}$ is noising ratio for synthetic input data fed in reconstruction sub-network and $R_{cls}$ is noising ratio for classification sub-network. Here we ensure that information from the realistic synthetic images are partially utilized according to the variation of the reconstructed samples and ground truth semantic masks. Noising ratio $R_{rec}$ is set as result of hyperbolic tangent function applied on ratio of variance between reconstructed images and semantic masks in current batch while noising ratio $R_{cls}$ is obtained by hyperbolic tangent function applied on the result of $1 - \gamma(\alpha, \beta)R_{rec}$, so we can see that $R_{rec}$ and $R_{cls}$ are strongly correlated with each other through $\gamma(\alpha, \beta)$. Assuming that hyper parameters $\alpha$ and $\beta$ are set to be the same, $\gamma$ will be 1 which means that we can easily represent $R_{cls}$ using $R_{rec}$. Both ratios could be represented in (9).

$$
\begin{aligned}
R_{rec} &= \tanh(\alpha \frac{Var(M')}{Var(M) + m_{noise}}) \\
R_{cls} &= \tanh(1 - \tanh(\beta \frac{Var(M')}{Var(M) + m_{noise}}))
\end{aligned} \tag{9}
$$

Here we usually set hyper parameters $(\alpha, \beta)$ as 0.25 and 2 according to results of experimental tests. Generally speaking, the noising ratio for classification network $R_{cls}$ and the one for reconstruction network $R_{rec}$ are coupled as (10).

$$
R_{cls} = \tanh(1 - \tanh(\frac{\beta}{\alpha} arctanh(R_{rec}))) \tag{10}
$$

Such strategy on training procedure ensures that information in realistic synthetic images is well utilized together with reconstructed images because classification sub-network tries best to learn from reconstructed channels without too much discriminant information in image channels in the beginning, meanwhile the orderly concatenated sub-networks and spatially concatenated six-channel data which is the junction between two sub-networks are trained in balance with effective gradients. Our coordinate training strategy makes the whole network converges easier than generative adversarial networks (GAN) [29], [30] where the generator and discriminator are sometimes trained separately according to the loss or discriminative accuracy, our method could also utilize metric learning to do multi-task learning compared to GAN with a flexible classification network. Similarly, parametric model in our network benefits from each others when foreground object reconstruction and classifier are always supervising each others in the training stage. If we regard them as generator and discriminator like concept in GAN, the parametric model are updated with different paces which means that it's always updating parametric model even with zero gradients defined by the variance of generated channels.

## VII. OVERALL MODEL AND IMPLEMENTATION

The overall loss function for our model can be thus defined as:

$$\mathcal{L}_{total} = \mathcal{L}_{\text{enc-gen}} + \mathcal{L}_{\text{pair}} + \mathcal{L}_{\text{tri}} + \mathcal{L}_{\text{s}} \ , \tag{11}$$

The variational prediction loss ($\mathcal{L}_{\text{enc-gen}}$) is summed with the pair-wise loss ($\mathcal{L}_{\text{pair}}$) and triplet loss ($\mathcal{L}_{\text{tri}}$). If we want to do classification for real photos, then the softmax loss ($\mathcal{L}_{\text{s}}$) is also added.

Similar to optimization process of VAEs [16], since the probability density function of the encoder, $Q(z|X)$, is still restricted to being a normal distribution $N(z|\mu,\Sigma)$, where $\mu(X;\theta)$ and $\Sigma(X;\theta)$ are arbitrary deterministic functions. The formulation of the KL divergence that we employ is:

$$\begin{aligned}
&D[N(\mu(p), \Sigma(p))||N(0, I)] \\
=&\frac{1}{2}(\text{tr}(\Sigma(p)) + (\mu(p))^T \mu(p) - k - \log \det (\Sigma(p)))) \ ,
\end{aligned} \tag{12}$$

where tr() is the trace and det() is the determinant. Here $\Sigma(p)$ is represented as $\log \Sigma(p)$.

One sample of $z$ for $P(Y|z)$ could be an approximation of $E_{z \sim Q} [\log P(Y|z)]$. So, the full equation to be optimized is

$$E_X[E_{z \sim Q}[\log P(Y|z)] - D[Q(z|X)||P(z)]] \ . \tag{13}$$

The gradient symbol of this equation can be moved into the expectations. Therefore, we can sample values of $z^{(l)}$ from the distribution $q(z|X)$, and compute the gradient of

$$\mathcal{L}_r = \frac{1}{L} \sum_{l=1}^{L} \log P(Y|z^{(l)}) - D[Q(z|X)||P(z)] \ . \tag{14}$$

We can then average the gradient of this function over samples of $X$ and $z$, so the result converges to the gradient of (13). The problem in (13) is that $E_{z \sim Q}[\log P(Y|z)]$ depends on both parameters $p$ and $q$. We solve this problem in (14) by means of the re-parameterization trick used for VAE and explained in [16].

In conclusion, our coordinate training makes training process more efficient with one loss functions of (11) for concatenated architectures. In functional representation of our deep architecture, the whole parametric model is tuned in the same time. Suppose that we have paired synthetic training data of realistic images $O$ and semantic depth masks $M$, reconstruction network $rec()$ with mapping network $map()$ and rendering network $ren()$, classification network $cls()$ lead by multi-triplet cost function, paired adaptive noising $noi()$ and $noi'()$, the concatenated network could be represented as $cls(ren(map(noi(O))), \ noi'(O))$ in training stage and $cls(ren(map(O)), \ O)$ in testing stage. The generative model is embedded in the mapper $map()$ and render $ren()$ while metric learning is applied in the classifier $cls()$. In the end, we make all those functions work as a whole architecture in perspective of effective parametric model training with paired noising functions $noi()$ and $noi'()$.

## VIII. EXPERIMENTS

### A. Experiments on ShapeNet Database

*1) Visualization of Reconstruction:* Both training and testing data are synthetic realistic images with little data migration problem rendered from 3D models in 12 categories provided by ShapeNet for basic experiments to prove that our generative model has a strong fitting capability and semantic embeddings. Reconstruction results of synthetic training images rendered from ShapeNet database show that our architecture is able to reconstruct ideal pixel-wise semantic foreground objects with expected colors. Objects with 12 unique RGB colors representing for different classes are used as semantic reconstruction targets for training shown in Fig. 6b which means that the reconstructed channels are represented by RGB values rather than channel-wise one-hot coding mask [31] used in AAE [32] with unfixed number of channels. There are two main advantages contributed by reconstruction using RGB channels. Firstly, reconstructed RGB channels could be widely used for representing segmentation masks in real database because itself is a kind of visualisation. Secondly, semantic segmentation masks represented using one-hot coding have large amount of channels which make the pixel-wise representation so sparse that can hardly be stored efficiently on disks. Our reconstruction cost function defined using Euclidean distance suits for optimising the decoder of the variational generative model. It does not have dependency on uniform distribution on three RGB channels which means that semantic masks in two different categories with similar colors are not less discriminative compared to ones in another pair of colors with bigger Euclidean distance. Our variational deep architecture is more discriminant and robuster in perspective of the reconstruction result. Reconstruction target in form of semantic depth information shown in Fig. 6b could be retrieved well in Fig. 6d while suppressing the background at the same time. Fig. 6c is the outputs of latent variables sampled in specific continuous range which shows that our generative reconstruction sub-network contains rich information about monotone surface colors of foreground object without background. Similar architecture without variational inference is simply represented as FCN [13] in Fig. 9e, it shows that contours are not reconstructed as clear as our generative method and the colors in single object also differ more which means that there are not enough meaningful semantic information according to the categories of foreground objects.



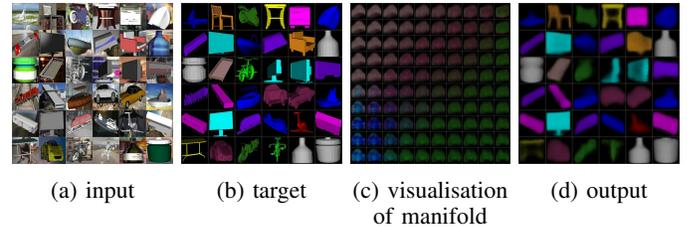(a) input    (b) target    (c) visualisation of manifold    (d) output

Fig. 6: Experiments on synthetic training data from ShapeNet models.

As for the training process, reconstruction procedure shown in Fig. 9a satisfies for our expectation on stability. In initial

training stage, the reconstruction sub-network is only able to reveal average color of all masks together with convex hull of objects by suppressing surrounding background. Then the outputs of the reconstruction network according to input samples in categories with less intraclass variance such as bottles are reconstructed clearly. Finally, objects are extracted from background together with relative depth information for samples in categories with larger intraclass variance like planes and bikes.

For analyzing how cost functions optimizing the reconstruction sub-network, visualizations of reconstructed channels during the whole training process are shown sequentially in Fig. 9a. The average color and contour of samples are reconstructed quickly in this initial stage. In this stage, loss of generator $\mathcal{L}_{gen}$ based on Euclidean distance shown in Fig. 8b drops faster than Loss of encoder $\mathcal{L}_{enc}$ represented as KL divergence shown in Fig.8a which means that the reconstruction sub-network behaves similarly to reconstruction based on principle component analysis [33], [34] to ensure reconstruction error as small as possible. This means that samples in categories with little intraclass variance is reconstructed better in this stage. Then encoder loss of KL divergence $\mathcal{L}_{enc}$ takes effect for rendering details to make sure that input samples from different categories with larger intraclass and smaller interclass variance can be reconstructed better contributed to latent variables. Two components of the reconstruction loss play their own roles at different stages of training which means that $\mathcal{L}_{gen}$ make the generator quickly fit the mean value of the training data and simple samples, then KL divergence component makes a finer fitting for interclass variance of samples while keeping previous samples reconstructed stably by utilizing latent space well with encoder.

*2) Semantic Embeddings Trained by Triplet Sets:* For semantic rendering task shown in Fig. 6, there are specific colors for each kind of objects in 12 categories. As categorical information is included for setting up triplet sets, triplet training could also be applied on latent variables. In such situation, clusters of latent variables from every categories could be set apart from each others. Fig. 7 shows the visualization of output from latent samples $z$ of our reconstruction sub-network around cluster centers of every categories. In each manifold, the poses of visualized outputs are continuous according to latent variables in specific range. It is obvious that all of the rightmost samples at the bottom line have the same pose which means that our multi-triplet cost function helps to establish hidden semantic embeddings.

*3) Analysis for Optimizers in Training:* Here we explain how our generative model works in regarding to both corruption function (9) and loss function (11). The conjugate network works as a whole in the training stage rather than gets trained separately. Both KL divergence loss in Fig. 8a and the Euclidean loss in the reconstruction network in Fig. 8b are optimized in similar trends and converge in similar training stage. There is a rapid drop of the KL divergence loss at the very beginning of training, then the divergence keeps in a particular range for a long time while the Euclidean loss keeping going smaller. This phenomenon indicates that the encoder in our generative model helps to construct a latent



(a) plane  (b) bottle  (c) chair
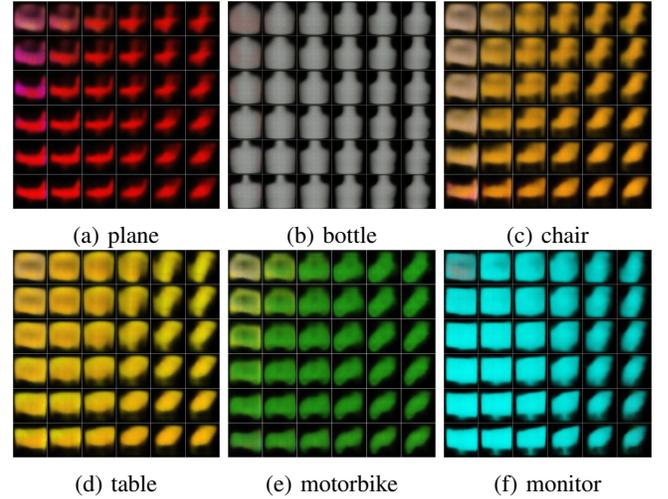


(d) table  (e) motorbike  (f) monitor

Fig. 7: Visualized outputs by sampling cluster centers of latent variables from synthetic images.

space where the output of the encoder can provide discriminant information for the generator stably which makes it easier for training a reconstructor without over-fitting.

Noising ratio plays an important role for training concatenated network, analysis on relationship of losses and noising ratio explains how adaptive noises work for tuning the whole parametric model. Variation ratio between reconstructed images and semantic masks in Fig. 8f gradually increases which means that discriminant information in reconstructed images increases during training because the information in expected semantic masks is almost the same. Noising ratios used for processing realistic inputs of two sub-networks are shown in Fig. 8g for reconstruction network and Fig. 8h for classification network. They are in negative and positive correlation with the variation ratio in Fig. 8f accordingly which means that difficulty for training reconstruction network increases while training for classification network becoming easier during training. Coupled noising ratios make the reconstruction sub-network utilize the back-propagated gradient from classification sub-network more efficiently. At the same time, the classification sub-network utilizes all six input channels evenly without over-fitting on realistic RGB channels more easily than the reconstructed channels in the beginning. By deceasing information of input images for classification sub-network, information in all six channels are always balanced during training. Gradient from all six channels in classification network could be fully utilized in the beginning rather than only from three of them in synthetic images with background. The intersection of three lines in Fig. 8d means that two noising ratios applied on reconstruction and classification networks equal to each other at this time. Here it shows that loss of classification in Fig. 8c is already decreased substantially and classification parametric model converges in a much smoother way in training afterwards without unstable loss turbulence meanwhile. The variance ratio in Fig. 8f keeps increasing monotonically with reconstruction ratio in Fig. 8g which means that realistic RBG images with background are not well-fitted on purpose before there are enough discriminant

information given by the reconstructed channels. Once the noising ratio of realistic images for classification network is smaller than noising ratio for reconstruction network, the reconstruction network tends to feed more information to classification network. Although corruption for the input of reconstruction network in Fig. 8g keeps increasing, reconstruction loss in Fig. 8b is always decreasing which means that the generator keeps learning discriminant information even if the realistic input is more and more confusing.
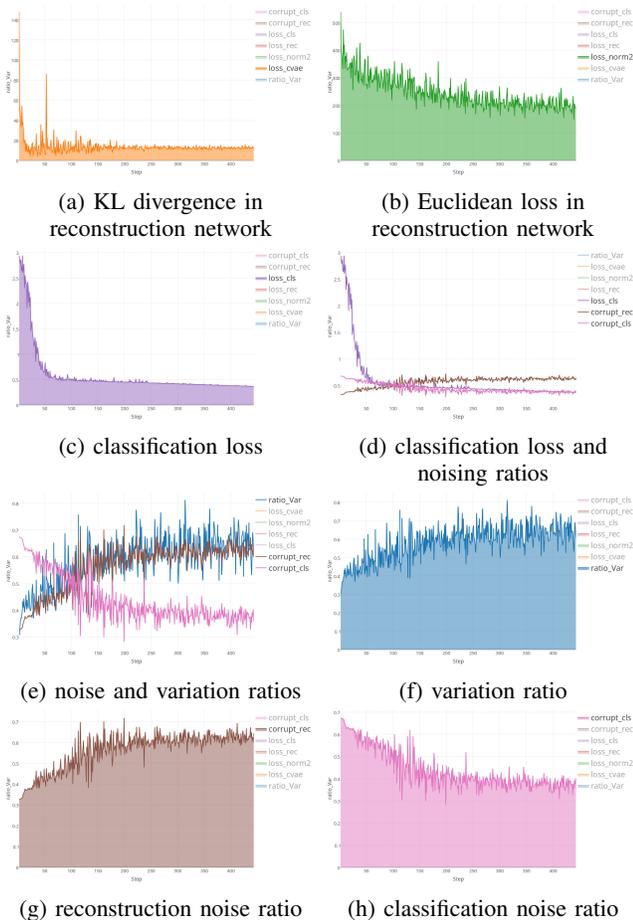


(a) KL divergence in reconstruction network

(b) Euclidean loss in reconstruction network

(c) classification loss

(d) classification loss and noising ratios

(e) noise and variation ratios

(f) variation ratio

(g) reconstruction noise ratio

(h) classification noise ratio

Fig. 8: Line charts of training parameters for experiments based on ShapeNet.

*4) Advantages to Other Pixel-wise Reconstruction Methods:* As reconstruction sub-network of our conjugate architecture depends on variational inference, the final reconstruction result is similar to the one learned from CVAE [18], [19]. The advantage of our architecture compared to simple variational coder (VC) is shown in the learning process in Fig. 9a and Fig. 9c where reconstruction sub-network learns faster than VC due to the classification sub-network. Our conjugate generative model also performs better than other popular pixel-wise regression models such as FCN [13], [14] and DeconvNet [15] if the final one-hot coding layer is replaced by pixel-wise regression layer with 3 channels. As shown in Fig. 9, our generative model with metric learning (GM-ML) has advantages both on the speed of learning and the precision of reconstruction accuracy. Results of methods like DeconvNet

shows that directly doing a pixel-wise coding from realistic images to semantic masks does not work well because FCN and DeconvNet treat semantic segmentation tasks without probabilistic variables in latent space, direct learning from input makes it less robust to noises for FCN. Firstly, FCN reconstruction result shown in the first two images of Fig. 9e shows over fitting problem for synthetic data where the local contours are learned at first without forming global shapes. Secondly, final result of FCN for the last image in Fig. 9e can not render monotone RGB colors according to unique categories where a single object has more than one colors on the surface. Our GM-ML makes it possible to train a pixel-wise coding strategy based on paired input and expected output with hidden meaning of category by forming clusters in the latent space. Our method shown in Fig. 9a indicates that it fits for the learning target shown in Fig. 6b well. If noises are fixed rather than adaptively changed according to output of reconstruction channels like result in the first two images of Fig. 9b, the learning progress is much slower.
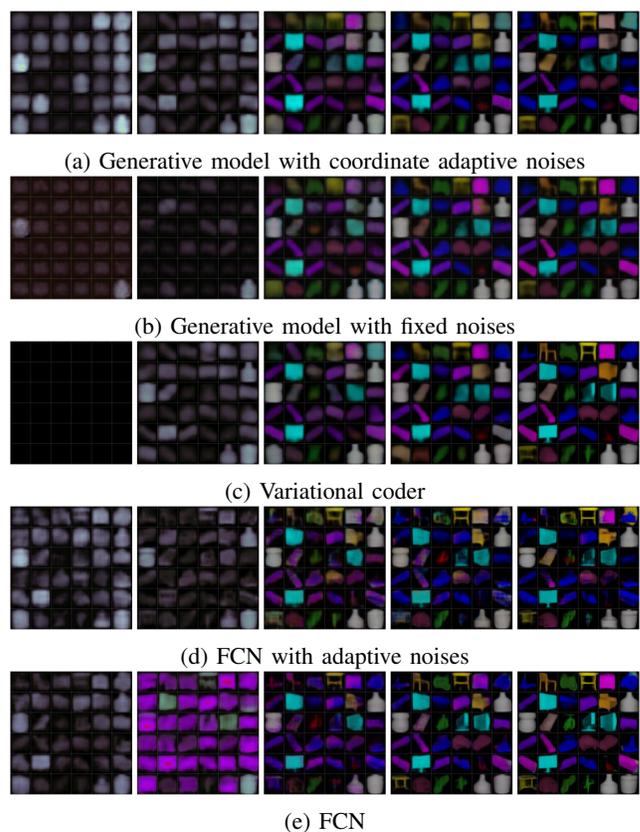


(a) Generative model with coordinate adaptive noises

(b) Generative model with fixed noises

(c) Variational coder

(d) FCN with adaptive noises

(e) FCN

Fig. 9: Visualization of reconstructed samples according to inputs and randomly selected latent variables.

*5) Feature Distributions, Reconstruction and Classification:* We visualize the distribution of descriptors by projecting the output of classification sub-network to two dimensions using principal component analysis (PCA) [35] and independent component analysis (ICA). Distribution shown in Fig. 11d indicates that our generative model trained with adaptive noises has a desired distribution of descriptors for classification where descriptors in different class are separated evenly.

Projected features of VC in Fig. 11a can not identify samples according to their categories as good as ours shown in Fig. 11c where samples from the same category are better clustered. This comparison experiment shows that directly concatenating a classification network to generative model without back-propagating gradients from classification sub-network makes it hard to fully utilize the reconstructed channels due to the absence of categories information. Here the reconstruction layer could be regarded as both a reconstruction target and a bridge providing additional supervision information. The variational variables also matter much in the whole pipeline, reconstructed foreground objects using FCN [13] concatenated with a classification network shown in Fig. 10c have more blended colors than ours in Fig. 10d where foreground objects are rendered with monotone colors. This means that the pixel-wise reconstructor can not render foreground objects stably according to categories without using generative model even when it can also reconstruct contours well.



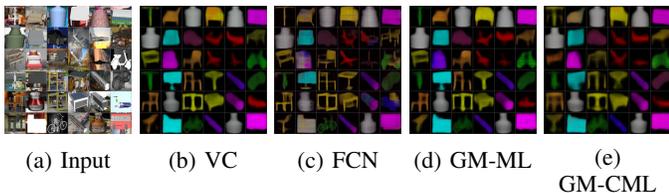(a) Input    (b) VC    (c) FCN    (d) GM-ML    (e) GM-CML

Fig. 10: Reconstruction samples of synthetic images rendered from 3D models of ShapeNet. VC is variational coder with flexible output. GM-CML (Generative Model with Coordinate Metric Learning) indicates that the conjugate generative model is trained with adaptive noises and GM-CML represent the same architecture without adaptive noises.



(a) VC      (b) FCN
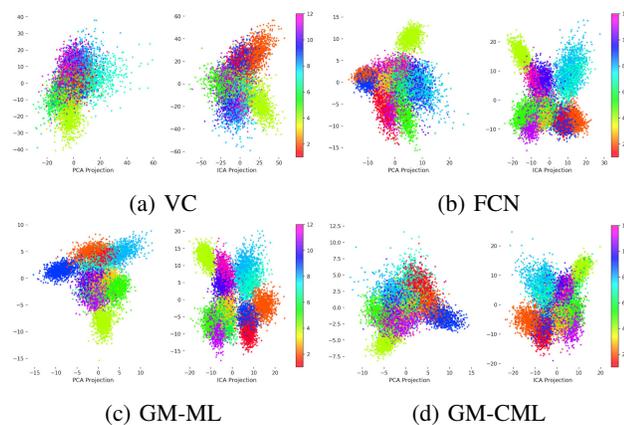
(c) GM-ML      (d) GM-CML

Fig. 11: Two-dimensional projection using PCA and ICA.

Table. I shows classification results for testing on the same type of synthetic realistic images rendered from ShapeNet [12] with different backgrounds and object poses compared to training data. Our conjugate deep architecture is also more capable for classification tasks than other popular deep architectures due to the help of the foreground reconstruction sub-network. Table. I shows that our generative model (GM-ML_1) of which output of the reconstruction sub-network is concatenated as a part of the input of ZigzagNet for classification achieves

the accuracy of 97.9% which is 7.8% higher than the one trained from ZigzagNet. Such improvement shows that the additional 3 channels of semantic foreground object extraction helps a lot for improving the classification network when data migration problem is not severe. Comparison experiments on different reconstruction sub-networks show that generative model retrieves the hidden category information with 5.1% improvement on classification accuracy compared to FCN.

TABLE I: Classification on rendered samples from ShapeNet database. Both training and testing data are synthetic images rendered from models of ShapeNet database in 12 categories. $GM\text{-}CML_1$ means that our conjugate generative model is trained with adaptive noises and the output of reconstruction network is concatenated with the input of ZigzagNet. $GM\text{-}ML_2$ means that the classification network is AlexNet.

| Method \ Result | Accuracy | Model size |
|---|---|---|
| $GM\text{-}CML_1$ | 96.3% | 20.2 MB |
| $GM\text{-}ML_1$ | **97.9**% | 20.2 MB |
| $GM\text{-}ML_2$ | 97.1% | 20.2 MB |
| FCN [13] | 91.2% | 20.2 MB |
| DeconvNet [15] | 92.7% | 20.2 MB |
| AlexNet [1] | 86.2% | 240 MB |
| SqueezeNet [36] | 84.7% | 4.8 MB |
| ZigzagNet [37] | 90.1% | 6.2 MB |

*B. Experiments on ImageNet Database*



(a) Input    (b) VC    (c) FCN    (d) GM-ML    (e) GM-CML

Fig. 12: Reconstruction samples of real images from ImageNet.



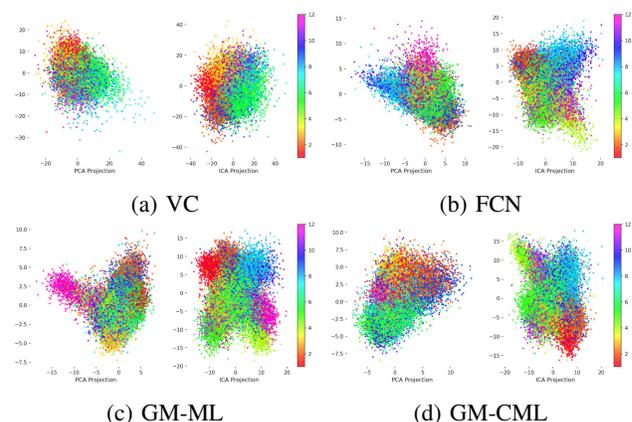(a) VC      (b) FCN

(c) GM-ML      (d) GM-CML

Fig. 13: Two-dimensional projection using PCA and ICA.

Here real photos of PASCAL and ImageNet are part of the original database which are provided by PASCAL 3D+ [20]

database with 12 categories altogether. Those experiments are designed for proving that our method works well for testing on real photos with the parametric model trained by realistic data rendered by 3D models. We do classification using outputs from softmax classifier which is the last layer in classification network to evaluate the performance. We render training data from models of PASCAL 3D+ with rendering method of SR [11] for basic experiments and models of ShapeNet database with rendering method in ZigzagNet. Some experiment results are collected from original papers if our testing results are similar to previously declared ones. Experiment in Table. II on trained network for regular photos from ImageNet attached in PASCAL 3D+ [4] shows that conjugating generative model with the help of adaptive noises could extract semantic foreground object information from background better compared to other popular and recent networks, achieving the state of the art accuracy of 50.5% for real photos from ImageNet. Comparison experiments on other reconstruction methods such as FCN [13] and DCN [15] show that latent variables in our generative model are helpful for extracting categorical informations in the foreground semantic masks which contribute to about 10% accuracy rising. We use the same macro reconstruction architecture for our conjugate generative model, FCN and DeconvNet, the main difference in other reconstruction sub-network for comparison is that they do not have the latent variables. Visualization of the reconstructed channels in Fig. 12e shows that our method has advantage in revealing details of objects compared to object segmentation results without texture information and can also suppressing close-up background compared to depth images. Here object reconstruction is applicable in extracting objects in real images by distinguishing object from background while suppressing highlights regions in background. Although foreground objects reconstructed in form of 1-channel depth images in SR [11] could give a help on doing classification on real photos by means of concatenating reconstructed channel to RGB channels, semantic colors are missing.

Our compact generative model based on ZigzagNet could be stored on disk occupying about 20.2 MB including both reconstruction network and classification network and achieves highest accuracy compared to both SqueezeNet [36] and AlexNet.

### C. Experiments on Pascal 12 Database

We further test our generative model on a difficult test set provided by PASCAL 3D+ database. Test images selected from PASCAL 2012 [5] are much more complicated than ones from ImageNet where foreground objects are not centered and shot as a whole part. This means that parts of object without complete contour are captured in sample photos. Basic result is carried out on texture-less models provided in PASCAL 3D+ database which shows that a severe data migration problem makes it too hard to get a good performance no matter which network we choose where all accuracies are not higher than 20% in Table. III. It's so necessary to use synthetic images rendered from textured model provided in ShapeNet database instead of ones from texture-less models in PASCAL 3D+

TABLE II: Classification on ImageNet samples attached in PASCAL 3D+ database. Different types of data used for training are real photos (1st column) and synthetic images which are rendered from PASCAL 3D+ (2nd column) and ShapeNet (3rd column) models. GM-CML$_1$ means that our conjugate generative model is trained with adaptive noises and the output of reconstruction network is concatenated with the input of ZigzagNet. GM-ML$_2$ means that the classification network is AlexNet.

| Data / Method | Real photo | PASCAL 3D | ShapeNet |
|---|---|---|---|
| GM-CML$_1$ | **63.1%** | **27.8%** | **50.5%** |
| GM-ML$_1$ | 58.6% | 26.1% | 49.2% |
| GM-ML$_2$ | 60.9% | 25.1% | 48.2% |
| FCN [13] | — | 23.9% | 38.4% |
| DeconvNet [15] | — | 26.2% | 40.2% |
| SR [11] | — | 25.1% | 40.1% |
| AlexNet [1] | 55.8% | 24% | 39.6% |
| SqueezeNet [36] | 45.5% | 21.4% | 35% |
| ZigzagNet [37] | 55.9% | 25.1% | 42.7% |
| Inception V1 [2] | 59% | 21.9% | 29.9% |
| Inception V4 [38] | — | 23.7% | 37.2% |
| ResidualNet [3] | — | 19.1% | 28.5% |
| LDO [6] | — | 14.8% | 25.5% |

in such condition with incomplete foreground objects waiting to be recognized. Parametric network trained from ShapeNet models has 10% improvement on accuracy with the help of reconstructed foreground objects. As objects are not centered in test photos, three reconstructed channels help to reconstruct foreground objects and improve recognition accuracy to 29.7% which is already much higher than that from deep models like Inception V4 [38] which is only 20.1%. This phenomenon means that deep residual architectures like ResidualNet [3] and Inception V4 [38] with strong fitting capability still have their weakness when training on synthetic data and testing on real photos with data migration problem without the help of pose information. The reason why direct classification networks like AlexNet without reconstruction sub-network fail to achieve acceptable accuracy as our model in complex condition is that our coordinate generative is assisted with metric learning targeting on utilizing common knowledge rather than over-fitting all information. This means that our generative model has its own limitation on fitting with input data, but the whole parametric model is tuned on fundamental knowledge like relative pose information which could be easily collected from synthetic data. As show in Table. III, the most amazing point is that the accuracy of our pipeline is nearly twice as high as ones trained without generative model which are 15.2% for FCN [13] concatenated with a classification network and 15.1% for DeconvNet [15]. Obvious advantage in classification against traditional pixel-wise reconstruction methods could be explained by reconstruction results shown in Fig. 12 where our method avoids problem of tearing foreground objects up into several different objects in form of rendering it with different colors in Fig. 12c. Close accuracy between our generative model and ZigzagNet indicates that pose information helps

to enhance the ability of perceiving ground truth relative information which avoid over-fitting for uncommon or random category information in test samples. Here we can draw conclusion that our generative model with coordinate metric learning could try best to reconstruct foreground objects in a global view rather than a pixel-wise view like FCN where the surface colors are always pure, so over-fitting problem is avoided by not reconstructing clear foreground object object channels.

TABLE III: Classification on Pascal 2012 samples attached in PASCAL 3D+ database. Different types of data used for training are real photos (1st column) and synthetic images which are rendered from PASCAL 3D+ (2nd column) and ShapeNet (3rd column) models. GM-CML$_1$ means that our conjugate generative model is trained with adaptive noises and the output of reconstruction network is concatenated with the input of ZigzagNet. GM-ML$_2$ means that the classification network is AlexNet.

| Data / Method | Real photo | PASCAL 3D | ShapeNet |
|---|---|---|---|
| GM-CML$_1$ | *46.7* | **18.1**% | **29.7**% |
| GM-ML$_1$ | *42.1* | 18% | 29.4% |
| GM-ML$_2$ | *44.5* | 17.2% | 28% |
| FCN [13] | — | 13.4% | 15.2% |
| DeconvNet [15] | — | 12.1% | 15.1% |
| SR [11] | — | 13.3% | 16.9% |
| AlexNet [1] | 39.2% | 15.2% | 20.7% |
| SqueezeNet [36] | 32.9% | 14.8% | 19.6% |
| ZigzagNet [37] | **42.9**% | 16.7% | 27.9% |
| Inception V1 [2] | 42.3% | 14.7% | 17.3% |
| Inception V4 [38] | — | 16.7% | 20.1% |
| ResidualNet [3] | — | 12.3% | 16.2% |
| LDO [6] | — | 11.2% | 18.7% |

As we have achieve the state of the art performance based for training with synthetic data alone and testing on real photos. We further carry out experiments for our methods in Table II and Table III to prove that our pre-trained model can also improve the performance of fine-tuning with real data. Accuracies shown in *italic* font in Table II and Table III represent that our synthetic data are used for pre-training and real data are used for fine-tuning. Most of those results are higher than results trained directly by real photos which shows that our parametric model pre-trained by synthetic data could improve the final performance if real images are used as training data for fine-tuning afterwards.

## IX. CONCLUSION

We designed a unified deep generative model with two concatenated sub-networks to do foreground object reconstruction and categorical classification with metric learning at the same time. We try to solve the ultimate target of training deep parametric model barely on synthetic images rendered from 3D models and testing on real images. A coordinate training strategy based on variance ratio makes our conjugate network work effectively where two concatenated sub-networks converge quicker and benefit more from each others. Experiments

on real photos from ImageNet database attached in ShapeNet database shows that our generative model with coordinate metric learning achieves the state of the art classification accuracy of 50.5% trained on rendered data. The generative model also shows more flexible reconstruction ability compared to VAE while removing additional supervision dependency of CVAE. From perspective of generative model, it gains faster converging speed compared to variational coder.
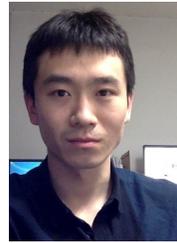
## REFERENCES

[1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc., 2012, pp. 1097–1105.

[2] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," *CoRR*, vol. abs/1409.4842, 2014. [Online]. Available: http://arxiv.org/abs/1409.4842

[3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[4] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and F. Li, "Imagenet large scale visual recognition challenge," *CoRR*, vol. abs/1409.0575, 2014. [Online]. Available: http://arxiv.org/abs/1409.0575

[5] M. Everingham, S. M. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes challenge: A retrospective," *International Journal of Computer Vision*, vol. 111, no. 1, pp. 98–136, 2015.

[6] P. Wohlhart and V. Lepetit, "Learning descriptors for object recognition and 3d pose estimation," in *Proc. IEEE CVPR*, 2015.

[7] H. Su, C. R. Qi, Y. Li, and L. J. Guibas, "Render for CNN: viewpoint estimation in images using cnns trained with rendered 3d model views," *CoRR*, vol. abs/1505.05641, 2015. [Online]. Available: http://arxiv.org/abs/1505.05641

[8] J. Yim, H. Jung, B. I. Yoo, and C. Choi, "Rotating your face using multi-task deep neural network," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 676–684.

[9] J. Yang, S. E. Reed, M.-H. Yang, and H. Lee, "Weakly-supervised disentangling with recurrent transformations for 3d view synthesis," in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 1099–1107.

[10] B. Pepik, R. Benenson, T. Ritschel, and B. Schiele, "What is holding back convnets for detection?" *CoRR*, vol. abs/1508.02844, 2015. [Online]. Available: http://arxiv.org/abs/1508.02844

[11] Y. Wang and W. Deng, "Self-restraint object recognition by model based cnn learning," in *2016 IEEE International Conference on Image Processing (ICIP)*, Sept 2016, pp. 654–658.

[12] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu, "ShapeNet: An Information-Rich 3D Model Repository," Stanford University — Princeton University — Toyota Technological Institute at Chicago, Tech. Rep. arXiv:1512.03012 [cs.GR], 2015.

[13] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.

[14] L. Wang, W. Ouyang, X. Wang, and H. Lu, "Visual tracking with fully convolutional networks," in *IEEE International Conference on Computer Vision*, 2015, pp. 3119–3127.

[15] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1520–1528.

[16] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[17] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, "Variational inference: A review for statisticians," *Journal of the American Statistical Association*, vol. 112, no. 518, pp. 859–877, 2017.

[18] K. Sohn, H. Lee, and X. Yan, "Learning structured output representation using deep conditional generative models," in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 3483–3491.

[19] D. P. Kingma, D. J. Rezende, S. Mohamed, and M. Welling, "Semi-supervised learning with deep generative models," *Advances in Neural Information Processing Systems*, vol. 4, pp. 3581–3589, 2014.

[20] Y. Xiang, R. Mottaghi, and S. Savarese, "Beyond pascal: A benchmark for 3d object detection in the wild," in *IEEE WACV*, 2014.

[21] M. Aubry, D. Maturana, A. A. Efros, B. C. Russell, and J. Sivic, "Seeing 3d chairs: exemplar part-based 2d-3d alignment using a large dataset of cad models," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 3762–3769.

[22] X. Peng, B. Sun, K. Ali, and K. Saenko, "Exploring invariances in deep convolutional neural networks using synthetic images," *CoRR*, vol. abs/1412.7122, 2014. [Online]. Available: http://arxiv.org/abs/1412.7122

[23] J. Lu, X. Zhou, Y.-P. Tan, Y. Shang, and J. Zhou, "Neighborhood repulsed metric learning for kinship verification," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 2, pp. 331–345, 2014.

[24] H. Oh Song, Y. Xiang, S. Jegelka, and S. Savarese, "Deep metric learning via lifted structured feature embedding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4004–4012.

[25] B. Kumar, G. Carneiro, I. Reid *et al.*, "Learning local image descriptors with deep siamese and triplet convolutional networks by minimising global loss functions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5385–5394.

[26] S. Hinterstoisser, S. Benhimane, V. Lepetit, P. Fua, and N. Navab, "Simultaneous recognition and homography extraction of local patches with a simple linear classifier," in *Proceedings of the BMVC*. BMVA Press, 2008, pp. 10.1–10.10.

[27] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size," *CoRR*, vol. abs/1602.07360, 2016. [Online]. Available: http://arxiv.org/abs/1602.07360

[28] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning fine-grained image similarity with deep ranking," *CoRR*, vol. abs/1404.4661, 2014. [Online]. Available: http://arxiv.org/abs/1404.4661

[29] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 2672–2680.

[30] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *CoRR*, vol. abs/1511.06434, 2015. [Online]. Available: http://arxiv.org/abs/1511.06434

[31] D. K. Taleshmekaeil, A. Safari, and Y. Kong, "Using one hot residue number system(ohrns) for digital image processing," in *The 16th CSI International Symposium on Artificial Intelligence and Signal Processing (AISP 2012)*, May 2012, pp. 064–067.

[32] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey, "Adversarial autoencoders," *arXiv preprint arXiv:1511.05644*, 2015.

[33] I. T. Jolliffe, "Principal component analysis," *Springer Berlin*, vol. 87, no. 100, pp. 41–64, 1986.

[34] H. Abdi and L. J. Williams, "Principal component analysis," *Wiley Interdisciplinary Reviews Computational Statistics*, vol. 2, no. 4, pp. 433–459, 2010.

[35] I. T. Jolliffe, "Principal component analysis," *Technometrics*, 2014.

[36] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size," *CoRR*, vol. abs/1602.07360, 2016. [Online]. Available: http://arxiv.org/abs/1602.07360

[37] Y. Wang, C. Cui, X. Zhou, and W. Deng, "Zigzagnet: Efficient deep learning for real object recognition based on 3d models," in *Asian Conference on Computer Vision*. Springer, 2016, pp. 456–471.

[38] C. Szegedy, S. Ioffe, and V. Vanhoucke, "Inception-v4, inception-resnet and the impact of residual connections on learning," *CoRR*, vol. abs/1602.07261, 2016. [Online]. Available: http://arxiv.org/abs/1602.07261

**Yida Wang** is Ph.D candidate in Technical University of Munich, Munich, Germany. He received B.E. and M.E. degrees from Beijing University of Posts and Telecommunications, Beijing, China in 2014 and 2017, respectively. His research interests include pattern recognition and computer vision. He was invited by Microsoft Research Seattle for Microsoft Faculty Summit in 2016 for project of "CNTK on Mac: 2D Object Restoration and Recognition Based on 3D Model" which was awarded as the second prize for Microsoft Open Source Challenge 2016. He is currently sponsored by Bleenco Research Fellowship in 2018 and was also sponsored by Google Summer of Code project twice for deep learning projects together with OpenCV organization. He was named of Excellent Graduate Student of Beijing City twice in 2014 and 2017 and been awarded for National Scholarship for Graduate Students in 2016.



**Weihong Deng** received the B.E. degree in information engineering and the Ph.D. degree in signal and information processing from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 2004 and 2009, respectively. From Oct. 2007 to Dec. 2008, he was a postgraduate exchange student in the School of Information Technologies, University of Sydney, Australia. He is currently an professor in School of Information and Telecommunications Engineering, BUPT. His research interests include statistical pattern recognition and computer vision, with a particular emphasis in face recognition. He has published over 100 technical papers in international journals and conferences, such as IEEE TPAMI and CVPR. He serves as associate editor for IEEE Access, and guest editor for Image and Vision Computing Journal and the reviewer for dozens of international journals, such as IEEE TPAMI / TIP / TIFS / TNNLS / TMM / TSMC, IJCV, PR / PRL. His Dissertation titled Highly accurate face recognition algorithms was awarded the Outstanding Doctoral Dissertation Award by Beijing Municipal Commission of Education in 2011. He has been supported by the program for New Century Excellent Talents by the Ministry of Education of China in 2013 and Beijing Nova Program in 2016.