# Learning to Generate with Memory

**Chongxuan Li**                                                    LICX14@MAILS.TSINGHUA.EDU.CN
**Jun Zhu**                                                              DCSZJ@MAIL.TSINGHUA.EDU.CN
**Bo Zhang**                                                           DCSZB@MAIL.TSINGHUA.EDU.CN
Dept. of Comp. Sci. & Tech., State Key Lab of Intell. Tech. & Sys., TNList Lab,
Center for Bio-Inspired Computing Research, Tsinghua University, Beijing, 100084, China

## Abstract

Memory units have been widely used to enrich the capabilities of deep networks on capturing long-term dependencies in reasoning and prediction tasks, but little investigation exists on deep generative models (DGMs) which are good at inferring high-level invariant representations from unlabeled data. This paper presents a deep generative model with a possibly large external memory and an attention mechanism to capture the local detail information that is often lost in the bottom-up abstraction process in representation learning. By adopting a smooth attention model, the whole network is trained end-to-end by optimizing a variational bound of data likelihood via auto-encoding variational Bayesian methods, where an asymmetric recognition network is learnt jointly to infer high-level invariant representations. The asymmetric architecture can reduce the competition between bottom-up invariant feature extraction and top-down generation of instance details. Our experiments on several datasets demonstrate that memory can significantly boost the performance of DGMs on various tasks, including density estimation, image generation, and missing value imputation, and DGMs with memory can achieve state-of-the-art quantitative results.

## 1. Introduction

Deep learning models are able to extract abstract representations from low-level inputs by adopting a deep architecture with explicitly designed nonlinear transformations (Bengio et al., 2013a). Among many types of deep models, deep generative models (DGMs) learn abstract

representations from unlabeled data and can perform a wide range of tasks, including density estimation, data generation and missing value imputation. Depending on the building blocks, various types of DGMs exist, including undirected models (Salakhutdinov & Hinton, 2009), directed models (Neal, 1992; Hinton et al., 2006), autoregressive models (Larochelle & Murray, 2011; Gregor et al., 2014), and Markov chain based models (Bengio et al., 2014). Recently, DGMs have attracted much attention on developing efficient and (approximately) accurate learning algorithms, such as stochastic variational methods (Kingma & Welling, 2014; Rezende et al., 2014; Bornschein & Bengio, 2015; Burda et al., 2015) and Monte Carlo methods (Adams et al., 2010; Gan et al., 2015; Du et al., 2015).

Although current DGMs are able to extract high-level abstract representations, they may not be sufficient in generating high-quality input samples. This is because more abstract representations are generally *invariant* or less sensitive to most specific types of local changes of the input. This bottom-up abstraction progress is good for identifying predictive patterns, especially when a discriminative objective is optimized (Li et al., 2015); but it also loses the detail information that is necessary in the top-down generating process. It remains a challenge for DGMs to generate real data, especially for images that have complex structures. Simply increasing the model size is apparently not wise, as it may lead to serious over-fitting without proper regularization as well as heavy computation burden. Some recent progress has been made to improve the generation quality. For example, DRAW (Gregor et al., 2015) iteratively constructs complex images over time through a recurrent encoder and decoder together with an attention mechanism and LAPGAN (Denton et al., 2015) employs a cascade of generative adversarial networks (GANs) (Goodfellow et al., 2014) to generate high quality natural images through a Laplacian pyramid framework (Burt & Adelson, 1983). However, no efforts have been made on enriching the capabilities of probabilistic DGMs by designing novel building blocks in the generative model.

In this paper, we address the above challenges by presenting a new architecture for building probabilistic deep generative models with a possibly large external memory and an attention mechanism. Although memory has been explored in various deep models for capturing long-term dependencies in reasoning and prediction tasks (See Section 2 for a review), our work represents a first attempt to leverage external memory to enrich the capabilities of probabilistic DGMs for better density estimation, data generation and missing value imputation. The overall architecture of our model is an interleave between stochastic layers and deterministic layers, where each deterministic layer is associated with an external memory to capture local variant information. An attention mechanism is used to record information in the memory during learning and retrieve information from the memory during data generation. This attention mechanism can be trained because the invariant information and local variant information are correlated, e.g., both containing implicit label information. Both the memory and attention mechanisms are parameterized as differentiable components with some smooth nonlinear transformation functions. Such a design allows us to learn the whole network end-to-end by developing a stochastic variational method, which introduces a recognition network without memory to characterize the variational distribution. Different from (Kingma & Welling, 2014; Burda et al., 2015), our recognition network is asymmetric to the generative network. This asymmetric recognition network is sufficient for extracting invariant representations in bottom-up inference, and is compact in parameterization. Furthermore, this asymmetry can help reduce the competition between bottom-up invariant feature extraction (using the recognition network) and top-down input generation (using the deep generative model with memory).

We quantitatively and qualitatively evaluate our method on several datasets in various tasks, including density estimation, data generation and missing value imputation. Our results demonstrate that an external memory together with a proper attention mechanism can significantly improve DGMs to obtain state-of-the-art performance.

## 2. Related Work

Memory has recently been leveraged in deep models to capture long-term dependencies for various tasks, such as algorithm inference (Graves et al., 2014), question answering (Weston et al., 2015; Sukhbaatar et al., 2015) and neural language transduction (Grefenstette et al., 2015). The external memory in these models provides a way to record information stably and interact with the environment, and hence extends the capability of traditional learning models. Typically the interaction, e.g., reading from and writing on the memory, is done through an associated attention mechanism and the whole system is trained with supervision. The attention mechanism can be differentiable and trained in an end-to-end manner (Graves et al., 2014; Sukhbaatar et al., 2015), or really discrete and trained by a Reinforcement Learning algorithm (Zaremba & Sutskever, 2015).

In addition to the memory-based models mentioned above, attention mechanisms have been used in other deep models for various tasks, such as image classification (Larochelle & Hinton, 2010; Ba et al., 2015), object tracking (Mnih et al., 2014), conditional caption generation (Xu et al., 2015), machine translation (Bahdanau et al., 2015) and image generation (Graves, 2013; Gregor et al., 2015). Recently, DRAW (Gregor et al., 2015) introduces a novel 2-D attention mechanism to decide "where to read and write" on the image and does well in generating objects with clear track, such as handwritten digits and sequences of real digits.

Compared with previous memory-based networks (Graves et al., 2014; Weston et al., 2015), we propose to employ an external hierarchical memory to capture variant information at different abstraction levels trained in an unsupervised manner. Besides, our memory cannot be written directly like (Graves et al., 2014; Weston et al., 2015); instead it is updated through optimization. Compared with previous DGMs with visual attention (Tang et al., 2014; Gregor et al., 2015), we make different assumptions about the data, i.e., the main object (such as faces) has massive local features, which cannot be modeled by a limited number of latent factors. We employ an external memory to capture this and the associated attention mechanism is used to retrieve the memory, not to learn "what-where" combination on the images. Besides, the external memory used in our model and the memory units of LSTMs used in DRAW (Gregor et al., 2015) can complement each other (Graves et al., 2014). Further investigation on DRAW with external memory is our future work.

Considering the bottom-up inference procedure and top-down generation procedure together, additional memory mechanisms can help to reduce the competition between invariant feature extraction and local variant reconstruction, especially when label information is provided (e.g., in supervised or semi-supervised setting). Similar idea is highlighted in the Ladder Network (Valpola, 2014; Rasmus et al., 2015), which reconstructs the input hierarchically using an extension of denoising autoencoders (dAEs) (Vincent et al., 2010) with the help of lateral connections and achieves excellent performance on semi-supervised learning (Rasmus et al., 2015). Though it is possible to interpret the Ladder Network probabilistically as in (Bengio et al., 2013b), we model the data likelihood directly with the help of external memory instead of explicit lateral edges. Our method can also be extended to do supervised or semi-

supervised learning as in (Kingma et al., 2014), which is our future work.

# 3. Probabilistic DGMs with Memory

We present a probabilistic deep generative model (DGM) with a possibly large external memory as well as a soft attention mechanism.

## 3.1. Overall Architecture

Formally, given a set of training data $\mathcal{D}$, we assume each $\mathbf{x} \in \mathcal{D}$ is independently generated with a set of hierarchically organized latent factors $\mathbf{z}_L, \ldots, \mathbf{z}_1$ as follows:

- Draw the top-layer factors $\mathbf{z}_L \sim \mathcal{N}(0, I)$.
- For $l = L - 1, \ldots, 0$, calculate the mean parameters $\boldsymbol{\mu}_l = g_l(\mathbf{z}_{l+1}; M_l)$ and draw the factors $\mathbf{z}_l \sim P_l(\boldsymbol{\mu}_l)$,

where each $g_l$ is a nonlinear function, often assumed to be smooth for the ease of learning. To connect with observations, the bottom layer is clamped at $\mathbf{z}_0 = \mathbf{x}$. Each $\mathbf{z}_l$ is randomly sampled from a Gaussian distribution except $\mathbf{z}_0$ whose distribution depends on the properties of the data (e.g., Gaussian for continuous data or Multinomial for discrete case). All the distributions $P_l$ are assumed to be of an exponential family form, with mean parameters $\boldsymbol{\mu}_l$.

Here, we define $g_l$ as a feed-forward deep neural network with $I_l$ deterministic layers and a set of associated memories $\{M_l^{(i)}\}_{i=0}^{I_l-1}$, one per layer. We parameterize each memory as a trainable matrix with dimension $d_s \times n_s$, where $d_s$ is the number of slots in the memory and $n_s$ is the dimension of each slot. Then, the network is formally parameterized as follows:

- Initialize the top-layer factors $\mathbf{h}_l^{(I_l)} = \mathbf{z}_{l+1}$.
- For $i = I_l - 1, \ldots, 0$, do the transformation $\mathbf{h}_l^{(i)} = \phi(\mathbf{h}_l^{(i+1)}; M_l^{(i)})$,

where $\phi$ is a proper (e.g., smooth) function for linear or nonlinear transformation. The bottom layer is our output $\boldsymbol{\mu}_l = \mathbf{h}_l^0$, which is called a stochastic layer as it computes the mean parameters for a distribution to get samples from. All the other layers are called deterministic layers.

Compared with previous DGMs, one key feature of our model is that it incorporates an external memory at each deterministic layer, as detailed below. The overall architecture is a stack of multiple such layers interleaved with stochastic layers as above. In such a DGM architecture, memory $M_l^{(i)}$ can recover the information that is missing in higher-layers $\mathbf{h}_l^{(>i)}$. In other words, the higher layers do not need to represent all details, but focusing on representing abstract invariant features if they seem more relevant to the task at hand than the more detailed information.

## 3.2. General Memory Mechanism for a Single Layer

We now present a single layer with memory generally, which is our building block for the above DGM. For notation simplicity, we omit the sub-script $l$ in the following text. Formally, let $\mathbf{h}_{in}$ denote the input information, and $\mathbf{h}_{out}$ denote the output after some deterministic transformation with memory. In our model, $\mathbf{h}_{in}$ can be either the samples of latent factors or the output from a higher-level deterministic layer; and similarly $\mathbf{h}_{out}$ can be used as the input of either a stochastic layer or a lower-level deterministic layer.

A layer of standard DGMs without memory generates the low-level generative information $\mathbf{h}_g$ based on $\mathbf{h}_{in}$ through a proper transformation, which can be generally put as:

$$\mathbf{h}_g = \phi(\mathbf{h}_{in}; W_g, b_g),$$

where $W_g$ and $b_g$ are the weights and biases of the transformation respectively and uses it as the final output, i.e. $\mathbf{h}_{out} = \mathbf{h}_g$.

In our DGM with memory $M$, we first compute the low-level generative information $\mathbf{h}_g$ in the same way as a standard layer, and then retrieve the memory with some proper attention mechanism to get knowledge $\mathbf{h}_m$. Finally, we combine $\mathbf{h}_g$ and $\mathbf{h}_m$ to get the output $\mathbf{h}_{out}$. Formally, the memory retrieval process is parameterized as

$$\mathbf{h}_m = f_m(\mathbf{h}_a; M),$$

where $\mathbf{h}_a = f_a(\mathbf{h}_g; A, b_A)$ is the information used to access the memory and computed by an attention mechanism parameterized by a controlling matrix $A$ and a bias vector $b_A$. The attention mechanism takes the generative information $\mathbf{h}_g$, which is the final output of a vanilla layer described previously, as input. $f_a$ is the mapping function in the attention mechanism and $f_m$ is the mapping function in the memory mechanism, which are deterministic transformations to be specified. The final output $\mathbf{h}_{out}$ is the combination of $\mathbf{h}_g$ and $\mathbf{h}_m$ as follows:

$$\mathbf{h}_{out} = f_c(\mathbf{h}_g, \mathbf{h}_m; C),$$

where $C$ is a set of trainable parameters in the combination function $f_c$, which is another deterministic transformation to be specified. We visualize the computation flow of these two types of layers in Figure 1, where each component will be specified next.

## 3.3. Concrete Examples with Hierarchical Memory Mechanisms

With the above building blocks, we can stack multiple layers to build a DGM as in Section 3.1. For simplicity, here we consider a generative model with only one stochastic
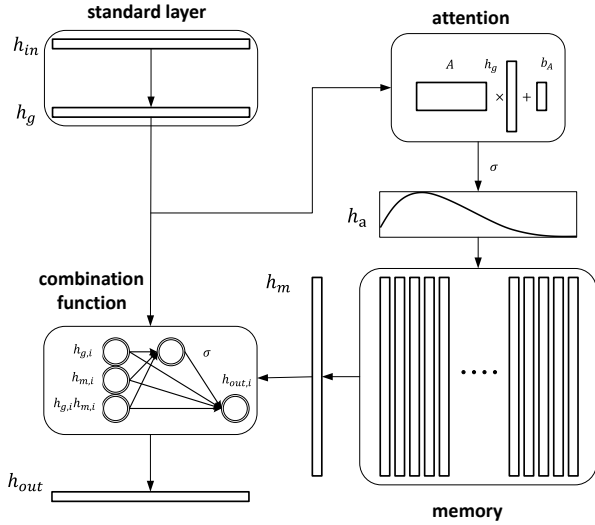
*Figure 1.* Architecture comparison between a standard layer (top-left part) and a layer with memory (the whole figure).

layer and $I$ deterministic layers to explain our memory mechanism, which can be straightforwardly extended to cases with multiple stochastic layers.

Let the top most information to be the random samples from the prior, i.e., $\mathbf{h}^{(I+1)} = \mathbf{z}$. Using permutation invariant architecture as an example, we compute the low-level generative information $\mathbf{h}_g^{(i)}$ based on the input $\mathbf{h}^{(i+1)}$ as:

$$\mathbf{h}_g^{(i)} = \phi(W_g^{(i)}\mathbf{h}^{(i+1)} + b_g^{(i)}).$$

We further retrieve the knowledge $\mathbf{h}_m^{(i)}$ from memory. Though various strategies exist, we consider the simple one that adopts a linear combination of slots in memory as

$$\mathbf{h}_m^{(i)} = f_m(\mathbf{h}_a^{(i)}) = M^{(i)}\mathbf{h}_a^{(i)},$$

where the coefficients $\mathbf{h}_a^{(i)}$ are computed as

$$\mathbf{h}_a^{(i)} = f_a(\mathbf{h}_g^{(i)}) = \sigma(A^{(i)}\mathbf{h}_g^{(i)} + b_A^{(i)}),$$

and $\sigma(x) = 1/(1 + \exp(-x))$ is the sigmoid function. Therefore, each element of $\mathbf{h}_a^{(i)}$ is a real value in the interval $(0, 1)$, which represents the preference of $\mathbf{x}$ to the corresponding memory slot. An alternative soft attention function used in our experiment is the softmax function, which normalizes the summation of the preference values on all slots to be one (Bahdanau et al., 2015). A hard attention mechanism trained with Reinforcement Learning (Xu et al., 2015) can be further investigated in the future work.

The most straightforward choice of the composition function is the element-wise summation:

$$\mathbf{h}^{(i)} = \mathbf{h}_g^{(i)} + \mathbf{h}_m^{(i)},$$

where the memory encodes the residual between the true target $\mathbf{h}^{(i)}$ and the generative information $\mathbf{h}_g^{(i)}$. However, in practise, we found that a more flexible composition function can lead to a better result. Inspired by the Ladder Network (Valpola, 2014; Rasmus et al., 2015), we specify the combination function of $\mathbf{h}_m^{(i)}$ and $\mathbf{h}_g^{(i)}$ as element wise multiple layer perceptron with optionally final nonlinearity $\phi$:

$$\mathbf{h}^{(i)} = f_c(\mathbf{h}_g^{(i)}, \mathbf{h}_m^{(i)}) = \phi(a^{(i)} + b_1^{(i)}c^{(i)}),$$

where the inside linear part $a^{(i)}$ is the summation of scaled inputs and cross terms as well as biases:

$$\begin{aligned} a^{(i)} &= a_1^{(i)} + a_2^{(i)} \odot \mathbf{h}_m^{(i)} + a_3^{(i)} \odot \mathbf{h}_g^{(i)} \\ &+ a_4^{(i)} \odot \mathbf{h}_g^{(i)} \odot \mathbf{h}_m^{(i)}, \end{aligned}$$

and the inside nonlinear part $c^{(i)}$ is computed similarly but goes through a sigmoid function:

$$\begin{aligned} c^{(i)} &= \sigma(c_1^{(i)} + c_2^{(i)} \odot \mathbf{h}_m^{(i)} + c_3^{(i)} \odot \mathbf{h}_g^{(i)} \\ &+ c_4^{(i)} \odot \mathbf{h}_g^{(i)} \odot \mathbf{h}_m^{(i)}), \end{aligned}$$

where $\odot$ is the element wise product. The output in our model only depends on the top-down signals $\mathbf{h}_g$ initially, instead of the auxiliary information as in the Ladder Network, which will be discussed in the experiment setting. $(W_g^{(i)}, b_g^{(i)}, M^{(i)}, A^{(i)}, b_A^{(i)}, a_{1,2,3,4}^{(i)}, b_1^{(i)}, c_{1,2,3,4}^{(i)})$ are trainable parameters in single layer. We illustrate each component in Figure. 1.

## 4. Inference and Learning

Learning a DGM is generally challenging due to the highly nonlinear transformations in multiple layers plus a stochastic formalism. To develop a variational approximation method, it is important to have a rich family of variational distributions that can well-characterize the nonlinear transforms. Significant progress has been made recently on stochastic variational inference methods with a sophisticated recognition model to parameterize the variational distributions (Kingma & Welling, 2014; Rezende et al., 2014). In this section, we develop such an algorithm for our DGM with memory.

Let $\boldsymbol{\theta}_g$ be the collection of parameters in the DGM. Then the joint distribution of each data $\mathbf{x}$ and the corresponding latent factor $\mathbf{z}$ can be generally put in a factorized form:

$$p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}_g) = p(\mathbf{z}; \boldsymbol{\theta}_g)p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}_g),$$

where the prior is often of a simple form, such as spherical Gaussian in our experiments, and the form of the conditional distribution $p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}_g)$ is chosen according to the data and its mean parameters depend on the external memories through a deep architecture as stated above.
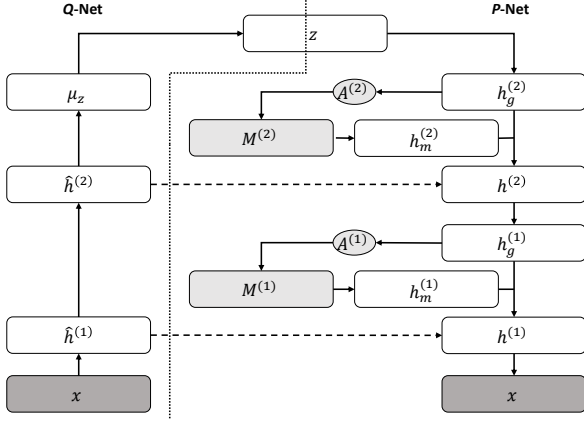
*Figure 2.* A model with one stochastic layer and two deterministic layers, where $\mathbf{z}$ is shared by the $P$-net and $Q$-net.

As in (Kingma & Welling, 2014), we adopt deep neural networks to parameterize a recognition model as the approximate posterior distribution $q(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}_r)$, where $\boldsymbol{\theta}_r$ is the collection of the parameters in the recognition model (denoted by $Q$-Net, as it characterizes distribution $q$). Since the $Q$-Net implements the bottom-up abstraction process to identify invariant features, it is unnecessary to have an external memory. Furthermore, the $Q$-Net without memory is compact in parameterization. The overall architecture is asymmetric, as illustrated in Figure 2, where the components at the left side of the dot line together with sampling $\mathbf{z}$ from $q(\mathbf{z}|\mathbf{x})$ is the $Q$-Net and the components at the right side of the dot line with $\mathbf{z}$ sampled from the prior is the generative model (denoted by $P$-Net, as it characterizes model distribution $p$). The solid arrow means the corresponding component is used as input of next component and the dash arrow means the corresponding component is used as the training target of next component, as explained below. The components representing external memory and associated attention mechanisms are filled with shallow gray. We omit the components corresponding to the combination functions for better visualization.

We define the $Q$-Net as follows. Following the example with one stochastic layer and $I$ deterministic layers in the previous section, we extract the high-level features $\hat{\mathbf{h}}^{(i+1)}$ as follows:

$$\hat{\mathbf{h}}^{(i+1)} = \phi(V^{(i)}\hat{\mathbf{h}}^{(i)} + b_r^{(i)}),$$

where $\phi$ is a proper nonlinear function and $(V^{(i)}, b_r^{(i)})$ are trainable parameters. The bottom layer is the input data, i.e. $\hat{\mathbf{h}}^{(0)} = \mathbf{x}$ and the top layer is still factorized Gaussian distribution. The mean of $\mathbf{z}$ is computed by linear transformation of $\hat{\mathbf{h}}^{(I)}$ and the variance of $\mathbf{z}$ is computed similarly but with a final exponential nonlinearity.

A variational lower bound of log-likelihood for per data $\mathbf{x}$ can be formulated as:

$$\mathcal{L}(\boldsymbol{\theta}_g, \boldsymbol{\theta}_r; \mathbf{x}) \triangleq \mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}_r)}[\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}_g) - \log q(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}_r)].$$

We add local reconstruction error terms as an optional regularizer, and jointly optimize the parameters in the generative model and the recognition model:

$$\min_{\boldsymbol{\theta}_g, \boldsymbol{\theta}_r} \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \left( \mathcal{L}(\boldsymbol{\theta}_g, \boldsymbol{\theta}_r; \mathbf{x}) + \sum_{i=1}^{I} \lambda^{(i)} ||\mathbf{h}^{(i)} - \hat{\mathbf{h}}^{(i)}||_2^2 \right),$$

where the relative weights $\lambda^{(i)}$ are prefixed hyperparameters. We optimize the objective with a stochastic gradient variational Bayes (SGVB) method (Kingma & Welling, 2014). Note that we cannot send the message of a intermediate layer in the recognition model to a layer in the generative model through a lateral connection as in Ladder Network (Valpola, 2014; Rasmus et al., 2015) because that indeed changes the distribution of $p(\mathbf{x}|\mathbf{z})$ according to the data $\mathbf{x}$. However, we do not use any information of $\mathbf{x}$ in the generative model explicitly and the correctness of the variational bound can be verified.

We employ batch normalization layers (Ioffe & Szegedy, 2015) in both the recognition model and generative model to accelerate the training procedure, and the intermediate features in local reconstruction error terms are replaced by a corresponding normalized version. To compare with stat-of-the-art results, we also train our method as in importance weighted autoencoders (IWAE) (Burda et al., 2015), which uses importance weighting estimate of log likelihood with multiple samples in the training procedure to achieve a strictly tighter variational lower bound.

## 5. Experiments

We now present both quantitative and qualitative evaluations of our method on the real-valued MNIST, OCR-letters and Frey faces datasets for various tasks. The MNIST dataset (Lecun et al., 1998) consists of 50,000 training, 10,000 validation and 10,000 testing images of handwritten digits and each image is of $28 \times 28$ pixels. The OCR-letters dataset (Bache & Lichman, 2013) consists of 32,152 training, 10,000 validation and 10,000 testing letter images of size $16 \times 8$ pixels. The Frey faces dataset consists of 1,965 real facial expression images of size $28 \times 20$ pixels. We model MNIST and OCR-letters datasets as Bernoulli distribution and model Frey faces dataset as Gaussian distribution at data level.

Our basic competitors are VAE (Kingma & Welling, 2014) and IWAE (Burda et al., 2015). We add the memory mechanisms to these methods and denote our models as MEM-VAE and MEM-IWAE, respectively. In all experiments except the visualization in Appendix D, MEM-VAE employs

the sigmoid function and element-wise MLP as the attention and composition functions respectively.

Our implementation is based on Theano (Bastien et al., 2012).[1] We use ADAM (Kingma & Ba, 2015) in all experiments with parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$ (decay rates of moving averages) and $\epsilon = 10^{-4}$ (a constant that prevents overflow). As a default, the global learning rate is fixed as $10^{-3}$ for 1,000 epochs and annealed by a factor 0.998 for 2,000 epochs with minibatch size 100. Initially, We set $a_3^i$ and $c_3^i$ as vectors filled with ones and $(a_{1,2,4}^{(i)}, b_1^{(i)}, c_{1,2,4}^{(i)})$ as vectors filled with zeros to avoid poor local optima. This means that we initialize the output as signals from top-down inference, which is different from the Ladder Network (Rasmus et al., 2015). We initialize the memory matrix as Gaussian random variables and other parameters following (Glorot & Bengio, 2010). We specify $\phi$ as rectified linear units (ReLu) (Nair & Hinton, 2010) in both the generative model and the recognition model.

We do not tune the hyper-parameters of our method heavily. We choose a model with one stochastic layer and two deterministic layers as the default setting. The values of $\lambda^{(1)}$ and $\lambda^{(2)}$ are fixed as 0.1 following Ladder Network (Rasmus et al., 2015). We do not include a local reconstruction error term at data level since the variational lower bound penalizes the reconstruction error of data already. The dimension of slots in memory $d_s$ is the same as that of the corresponding generative information $\mathbf{h}_g$ because we use element-wise combination function $f_c$. We employ the memory mechanism in both deterministic layers and make the total number of slots $n_s^{(1)} + n_s^{(2)}$ to be 100 to keep the number of additional parameters relatively small. We choose a 70-30 architecture according to the validation performance on the MNIST dataset and then make it default for all experiments if not mentioned.

## 5.1. Density Estimation

We follow (Burda et al., 2015) to split the MNIST dataset into 60,000 training data and 10,000 testing data after choosing the hyper-parameters. We train both the baselines and our models with 1, 5 and 50 importance samples respectively and evaluate the test likelihood with 5,000 importance samples as in (Burda et al., 2015). In each training epoch, we binarize the data stochastically as the input. The results of VAE, IWAE-5 (trained with 5 importance samples) and IWAE-50 (trained with 50 importance samples) with one stochastic layer in (Burda et al., 2015) are -86.76, -85.54 and -84.78 nats respectively. However, we use 500 hidden units in the deterministic layers and 100 latent variables in the stochastic layer to achieve a stronger baseline result with a different architecture and more parameters.

[1]Source code at https://github.com/zhenxuan00/MEM_DGM

*Table 1.* Log likelihood estimation on MNIST and OCR-letters datasets. Results are from [1] (Murray & Salakhutdinov, 2009), [2] (Burda et al., 2015), [3] (Bornschein & Bengio, 2015), [4] (Larochelle & Murray, 2011) and [5] (Gregor et al., 2014). Results with * are evaluated on binarised MNIST dataset.

| MODELS | MNIST | OCR-LETTERS |
|---|---|---|
| *VAE* | -85.67 | -30.09 |
| *MEM-VAE(ours)* | -84.41 | -29.09 |
| *IWAE-5* | -84.49 | -28.69 |
| *MEM-IWAE-5(ours)* | -83.26 | -27.65 |
| *IWAE-50* | -83.67 | -27.60 |
| *MEM-IWAE-50(ours)* | **-82.84** | **-26.90** |
| *DBN*[1] | -84.55 | - |
| *S2-IWAE-50*[2] | **-82.90** | - |
| *RWS-SBN/SBN*[3]* | -85.48 | -29.99 |
| *RWS-NADE/NADE*[3]* | -85.23 | **-26.43** |
| *NADE*[4]* | -88.86 | -27.22 |
| *DARN*[5]* | **-84.13** | -28.17 |

We present our likelihood results in Table 1. We can see that our methods improve the results of baselines (both VAE and IWAE) significantly and achieve state-of-the-art results on the real-valued MNIST dataset with permutation invariant architectures. DRAW (Gregor et al., 2015) achieves -80.97 nats by exploiting the spatial information. Our method MEM-IWAE-50 even outperforms *S2-IWAE-50*, which is the best model in (Burda et al., 2015) with two stochastic layers and four deterministic layers.

To compare with a broader family of benchmarks, we further quantitatively evaluate our model on the OCR-letters dataset. We use 200 hidden units in the deterministic layers and 50 latent variables in the stochastic layer as the dimension of the input is much smaller. The test log-likelihood is evaluated with 100,000 importance samples as in (Bornschein & Bengio, 2015) and shown in Table 1. Again, our methods outperform the baseline approaches significantly and are comparable with the best competitors, which often employ autoregressive connections (Larochelle & Murray, 2011; Gregor et al., 2014) that are effective on small images with simple structures. Note that these sophisticated structures are not exclusive to our memory mechanisms. A systematic investigation of using memory with such structures is our future work.

## 5.2. Analysis of Our Model

We now present a careful analysis of our model to investigate the possible reasons for the outstanding performance.

**Classification:** We investigate the effect of external memory on the training of the recognition model by classifica-

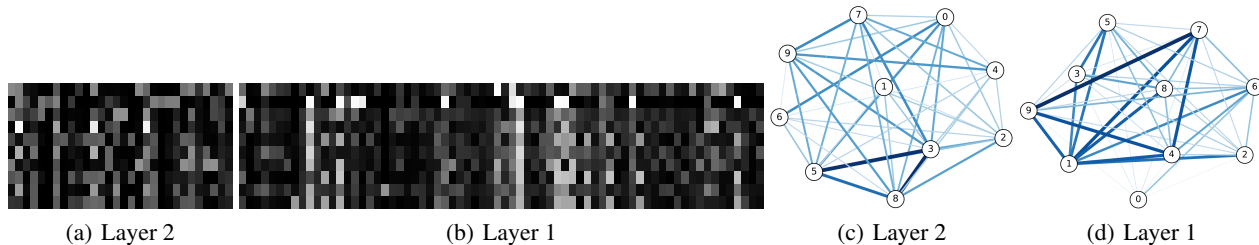(a) Layer 2      (b) Layer 1      (c) Layer 2      (d) Layer 1

*Figure 3.* (a-b): Averaged activations on each memory slot over different classes of testing data on MNIST dataset in layer 2 and layer 1 respectively. (c-d): 2-D visualization for correlation between classes for layer 2 and layer 1 respectively (best view in color).

tion and MEM-VAE outperforms VAE (See details in Appendix A).

**A larger baseline:** We test VAE with a 530-530-100 architecture, which has almost the same number of parameters as MEM-VAE. The log-likelihood trained with 1, 5 and 50 importance samples on MNIST are -85.69, -84.43 and -83.58 respectively. We can see that using our memory leads to much better results than simply increasing the model size. A comparison of number of parameters used in all of the models can be seen in Appendix B.

**Importance of memory:** We test the relative importance of the memory mechanism and local reconstruction error regularizer. MEM-VAE in the default settings but without local reconstruction error regularizer achieves a test log density estimation of -84.44 nats. VAE with additional local reconstruction regularizer achieves test log density estimation of -85.68 nats. These experiments demonstrate that the memory mechanism plays a central role in the recovery of detailed information. The local reconstruction error regularizer may help more provided supervision.

**Preference of memory slots over classes:** We investigate the preference of memory slots over different classes in MEM-VAE. We average $\mathbf{h}_a$ and normalize the activations for each class and visualize the matrices in Figure 3(a-b), where each column represents a slot and each row represents a class (0-9 in top-down order). The averaged and normalized activations are used as the intensities for the corresponding positions in the matrices. Furthermore, we compute the correlation coefficients between activations of different classes and visualize them in a 2-D graph in Figure 3(c-d), where each node represents a class and each edge represents the correlation between two endpoints. The larger the correlation is, the wider and darker the edge is. We observe that the trained attention model can access the memory based on the implicit label information in the input, which accords with our assumption. The activations are correlated for those digits that share similar structures such as "7" and "9". Furthermore, different layers of memory focus on different patterns. For example, layer 1 has a strong activation of a vertical line pattern which is shared among digits "1", "4", "7" and "9", while layer 2 activates
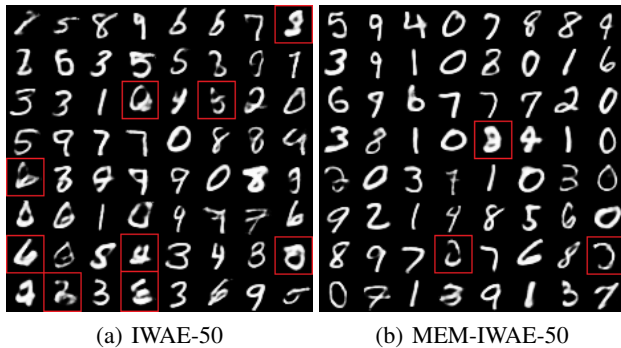


(a) IWAE-50      (b) MEM-IWAE-50

*Figure 4.* (a-b): Random generation from IWAE-50 and MEM-IWAE-50 on MNIST dataset respectively.

most to a semi-circle pattern which is shared among digits "3", "5" and "8". Besides, layer 1 has almost the same 2D-visualization result as the raw data.

**Visualization:** We visualize the generative information $\mathbf{h}_g$ and memory information $\mathbf{h}_m$ by mapping these vectors to images (See details in Appendix C and D respectively).

### 5.3. Random Generation

We further evaluate the random generations from the baseline and our model empirically on MNIST and Frey faces datasets, which is shown in Figure 4 and Figure 5 respectively. We label unclear or meaningless images with red rectangles. This is done by majority voting of several volunteers. We do not select any pictures for both datasets.

For the MNIST dataset, the setting is same as in Section 5.1. We observe that the memory mechanism helps a lot to get clear and meaningful samples as in Figure 4.

For Frey faces dataset, we randomly split into 1,865 training data and 100 testing data. We use a single deterministic layer with 200 hidden units and a stochastic layer with 10 latent factors and set $n_s^{(1)}$ to be 20 as the number of training samples is small. We use one sample of the recognition model in both of the training and testing procedure as in (Kingma & Welling, 2014). We find that the minibatch size effects the results a lot, and the quality of visualization
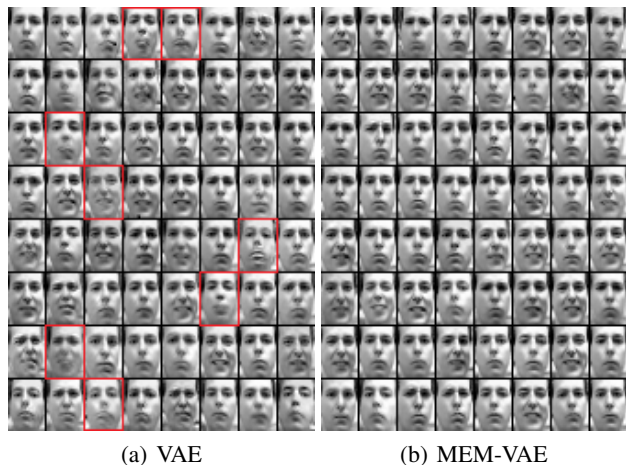
(a) VAE          (b) MEM-VAE

*Figure 5.* (a-b): Random generation from VAE and MEM-VAE on Frey faces dataset respectively.

and the averaged test log density are inconsistent (Theis et al., 2016). Specifically, setting the minibatch size to be 100, VAE achieves test log density of 1308 nats, which reproduces the result with same architectures in (Kingma & Welling, 2014), but the visualization is somehow unclear; while setting the minibatch size to be 10, VAE achieves test log density of 1055 nats, but the visualization is much better. All of the parameters are set referred to (Kingma & Welling, 2014) or based on the performance of test log density of VAE. We also find that MEM-VAE outperforms VAE in both cases in terms of the quantitative test likelihood and qualitative visualization — the corresponding log density of MEM-VAE are 1330 and 1240 nats respectively. The random samples given minibatch size 100 is shown in Figure 5, where we can see that all samples of MEM-VAE are clear but some of VAE cannot present all details in the facial expression successfully.

### 5.4. Missing Value Imputation

Finally, we evaluate our method on the task of missing value imputation with three different types of noise, including (1) *RECT*-12 means that a centered rectangle of size $12 \times 12$ is missing; (2) *RAND*-0.6 means that each pixel is missing with a prefixed probability 0.6; and (3) *HALF* means that the left half of the image is missing. For both VAE and MEM-VAE, the missing values are randomly initialized and then inferred by a Markov chain that samples latent factors based on the current guess of missing values and then refines the missing values based on the current latent factors. We compare the mean square error (MSE) results after 100 epochs of inference as in Table 2 on the MNIST dataset. The results demonstrate that DGM with external memory can capture the underlying structures of data better than vanilla methods under different types of

*Table 2.* MSE results on MNIST dataset with different types of noise.

| NOISE TYPE | VAE | MEM-VAE |
|---|---|---|
| *RECT*-12 | 0.1403 | **0.1362** |
| *RAND*-0.6 | 0.0194 | **0.0187** |
| *HALF* | 0.0550 | **0.0539** |

noise. Besides, MEM-VAE has better qualitative results (See Appendix E).

## 6. Conclusions and Future Work

In this paper, we introduce a novel building block for deep generative models (DGMs) with an external memory and an associated soft attention mechanism. In the top-down generative procedure, the additional memory helps to recover the local detail information, which is often lost in the bottom-up abstraction procedure for learning invariant representations. Various experiments on handwritten digits and letters as well as real faces datasets demonstrate that our method can substantially improve the vanilla DGM on density estimation, random generation and missing value imputation tasks, and we can achieve state-of-the-art results among a broad family of benchmarks.

There are three possible extensions of our method:

- The use of other types of memory and attention mechanisms in DGMs can be further investigated. Particularly, the combination of external memory and visual attention as well as recurrent networks (Gregor et al., 2015) may achieve better results in generative tasks.

- A class conditional DGM (Kingma et al., 2014) with memory can potentially achieve better performance on both classification and generation because the external memory helps to reduce the competition between the invariant feature extraction and detailed generation, and explicit label information can make the whole system be easier to train.

- Our method can be further applied to convolutional neural networks by sharing parameters across different channels and then employed in non-probabilistic DGMs such as LAPGAN (Denton et al., 2015) to refine generation on high-dimensional data.

# References

Adams, R., Wallach, H., and Ghahramani, Z. Learning the structure of deep sparse graphical models. In *AISTATS*, 2010.

Ba, J. L., Mnih, V., and Kavukcuoglu, K. Multiple object recognition with visual attention. In *ICLR*, 2015.

Bache, K. and Lichman, M. UCI machine learning repository. 2013. URL http://archive.ics.uci.edu/ml.

Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. In *ICLR*, 2015.

Bastien, F., Lamblin, P., Pascanu, R., Bergstra, J., Goodfellow, I., Bergeron, A., Bouchard, N., Warde-Farley, D., and Bengio, Y. Theano: new features and speed improvements. In *Deep Learning and Unsupervised Feature Learning NIPS Workshop*, 2012.

Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35:8:1798–1828, 2013a.

Bengio, Y., Yao, L., Alain, G., and Vincent, P. Generalized denoising autoencoders as generative models. In *NIPS*, 2013b.

Bengio, Y., Thlhodeau-Laufer, E., Alain, G., and Yosinski, J. Deep generative stochastic networks trainable by backprop. In *ICML*, 2014.

Bornschein, J. and Bengio, Y. Reweighted wake-sleep. In *ICLR*, 2015.

Burda, Y., Grosse, R., and Salakhutdinov, R. Importance weighted autoencoders. In *arXiv:1509.00519*, 2015.

Burt, P. J. and Adelson, E. H. The Laplacian pyramid as a compact image code. *IEEE Transactions on Communications*, 31:532ÍC540, 1983.

Denton, E., Chintala, S., Szlam, A., and Fergus, R. Deep generative image models using a laplacian pyramid of adversarial networks. In *NPIS*, 2015.

Du, C., Zhu, J., and Zhang, B. Learning deep generative models with doubly stochastic MCMC. In *arXiv:1506.04557*, 2015.

Gan, Z., Henao, R., Carlson, D. E., and Carin, L. Learning deep sigmoid belief networks with data augmentation. In *AISTATS*, 2015.

Glorot, X. and Bengio, Y. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010.

Goodfellow, I. J., Abadie, J. P., Mirza, M., Xu, B., Farley, D. W., S.ozair, Courville, A., and Bengio, Y. Generative adversarial nets. In *NIPS*, 2014.

Graves, A. Generating sequences with recurrent neural networks. In *arXiv:1308.0850*, 2013.

Graves, A., Wayne, G., and Danihelka, I. Neural Turing machines. In *arXiv:1410.5401*, 2014.

Grefenstette, E., Hermann, K. M., Suleyman, M., and Blunsom, P. Learning to transduce with unbounded memory. In *NIPS*, 2015.

Gregor, K., Danihelka, I., Mnih, A., Blundell, C., and Wierstra, D. Deep autoregressive networks. In *ICML*, 2014.

Gregor, K., Danihelka, I., Graves, A., Rezende, D. J., and Wierstra, D. DRAW: A recurrent neural network for image generation. In *ICML*, 2015.

Hinton, G. E., Osindero, S., and Teh, Y. A fast learning algorithm for deep belief nets. *Neural Computation*, 18, 2006.

Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.

Kingma, D. P. and Ba, J. L. Adam: A method for stochastic optimization. In *ICLR*, 2015.

Kingma, D. P. and Welling, M. Auto-encoding variational Bayes. In *ICLR*, 2014.

Kingma, D. P., Rezende, D. J., Mohamed, S., and Welling, M. Semi-supervised learning with deep generative models. In *NIPS*, 2014.

Larochelle, H. and Hinton, G. Learning to combine foveal glimpses with a third-order Boltzmann machine. In *NIPS*, 2010.

Larochelle, H. and Murray, I. The neural autoregressive distribution estimator. In *AISTATS*, 2011.

Lecun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *In Proceedings of the IEEE*, 1998.

Li, C., Zhu, J., Shi, T., and Zhang, B. Max-margin deep generative models. In *NIPS*, 2015.

Mnih, V., Heess, N., Graves, A., and Kavukcuoglu, K. Recurrent models of visual attention. In *NIPS*, 2014.

Murray, I. and Salakhutdinov, R. Evaluating probabilities under high-dimensional latent variable models. In *NIPS*, 2009.

Nair, V. and Hinton, G. E. Rectified linear units improve restricted Boltzmann machines. In *ICML*, 2010.

Neal, R. M. Connectionist learning of belief networks. In *Artificial intelligence*, 1992.

Rasmus, A., Berglund, M., Honkala, M., Valpola, H., and Raiko, T. Semi-supervised learning with ladder networks. In *NIPS*, 2015.

Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014.

Salakhutdinov, R. and Hinton, G. E. Deep Boltzmann machines. In *AISTATS*, 2009.

Sukhbaatar, S., Szlam, A., Weston, J., and Fergus, R. End-to-end memory networks. In *NIPS*, 2015.

Tang, Y., Srivastava, N., and Salakhutdinov, R. Learning generative models with visual attention. In *NIPS*, 2014.

Theis, L., Oord, A., and Bethge, M. A note on the evaluation of generative models. In *arXiv:1511.01844*, 2016.

Valpola, H. From neural PCA to deep unsupervised learning. In *arXiv:1411.7783*, 2014.

Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., and Manzagol, P. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *JMLR*, 11:3371ÍC340, 2010.

Weston, J., Chopra, S., and Bordes, A. Memory networks. In *ICLR*, 2015.

Xu, K., Ba, J. L., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., Zemel, R. S., and Bengio, Y. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.

Zaremba, W. and Sutskever, I. Reinforcement learning neural Turing machines. In *arXiv:1505.00521*, 2015.

*Table 3.* Number of parameters used in practice.

| MODEL | NUMBER OF PARAMETERS |
|---|---|
| *MNIST-500-500-100* | 1,441K |
| *MNIST-530-530-100* | 1,559K |
| *MNIST-500-500-100-MEM* | 1,550K |
| *OCR-letters-200-200-50* | 164K |
| *OCR-letters-200-200-50-MEM* | 208K |

*Table 4.* Average preference values of selected slots.

| "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8" | "9" |
|---|---|---|---|---|---|---|---|---|---|
| 0.27 | 0.82 | 0.33 | 0.11 | 0.34 | 0.15 | 0.49 | 0.27 | 0.09 | 0.28 |
| 0.24 | 0.09 | 0.06 | 0.11 | 0.30 | 0.13 | 0.12 | 0.27 | 0.09 | 0.21 |
| 0.18 | 0.05 | 0.06 | 0.11 | 0.07 | 0.07 | 0.05 | 0.11 | 0.09 | 0.18 |



*Figure 6.* Generation from MEM-VAE when memory is disabled.



*Figure 7.* Visualization of selected slots.

## A. Recognition Model

We feed the output of the last deterministic layer in the recognition models into a linear SVM to classify the MNIST digits to examine the invariance in features. We achieve slightly better classification accuracy (97.90% for VAE and 98.03% MEM-VAE), which means that additional memory mechanisms do not hurt or even improve the invariance of the features extracted by the recognition model.

## B. Number of Parameters Used

As we employ external memory and attention mechanisms, the number of parameters in our building block is larger than that in a standard layer. However, the total number of parameters in the whole model is controlled given a limited number of slots in the memory. See Table 3 for a comparison, and we do not observe that our method suffers from overfitting.

## C. Disabling Memory

We investigate the performance of MEM-VAE when the memory is disabled (setting $\mathbf{h}_m$ as a vector filled with ones) as in Figure 6. The top row shows original samples; the middle row shows samples with memory of the first layer disabled; and the bottom row shows samples with memory of both layers disabled. It can be seen that, without information from memory, the main pattern of the generation does not change much but the local details are lost in some sense, which supports our assumption.

## D. Visualizing the Memory Slots Directly

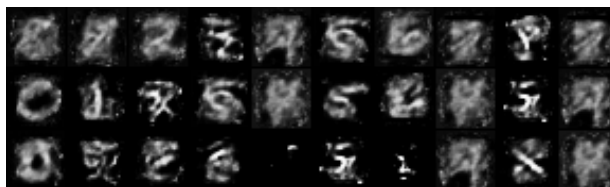As mentioned in Section 5, MEM-VAE employs the sigmoid function and element-wise MLP as the attention and composition functions respectively. MEM-VAE is complicated and flexible, which has better quantitative results on log-likelihood estimation but is hard to be visualized directly as it has much nonlinearity. We introduce VIS-MEM-VAE, which uses the softmax function and element-wise summation as the attention and composition functions respectively. VIS-MEM-VAE achieves a slightly worse density estimation result (-84.68 nats on MNIST, still better than -85.67 nats of VAE) but the memory slots can be mapped to data level for visualization due to the simple composition function and sparse attention mechanism.

We average the preference values $\mathbf{h}_a$ of the test data and select top-3 memory slots that has the highest activations for each class. Note that the activations are normalized, i.e., the summation of the preference values on all memory slots equals to one for each class. We set the generative information to be a vector filled with zeros and the memory information to be one of the selected slots and then generate a corresponding image. The activations and images are shown in Table 4 and Figure 7, where each column represents a class (0-9 in left-right order). For example, the image at the first column and the second row in Figure 7 corresponds to the memory slot that has the second-highest averaged activation of class "0" and the value is 0.24.

It can be seen that most of the selected slots respond to one class or similar classes (some slots are shared by similar classes such as "4", "7" and "9") and the corresponding image contains a blurry sketch of the digit (or mixture of some digits) with different local styles, which indicates that the external memories can encode local variants of objects and can be retrieved based on generative information $\mathbf{h}_g$. A few images are less meaningful but the corresponding activations are relatively small (smaller than 0.08 or so).

(a) Data

(b) Noisy data
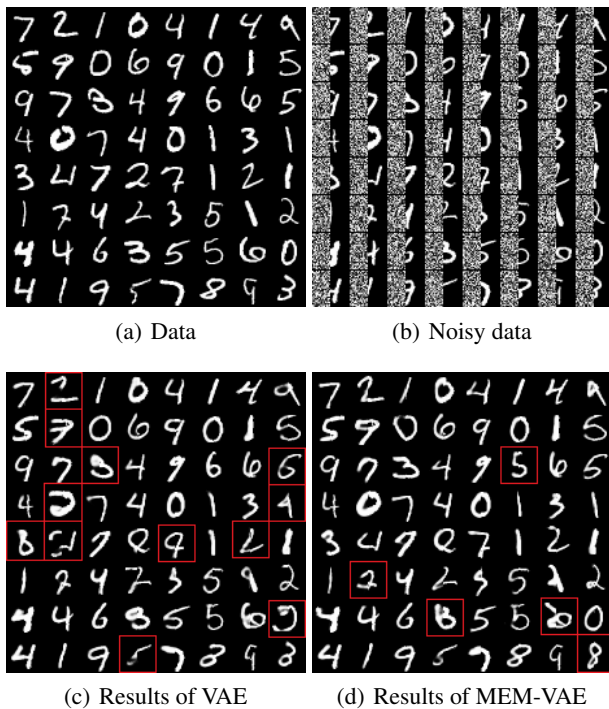
(c) Results of VAE

(d) Results of MEM-VAE

*Figure 8.* (a-b): Original test data on MNIST and perturbed data with left half missing respectively; (c-d): Imputation results of VAE and MEM-VAE respectively. Red rectangles mean that the corresponding model fails to infer digits correctly but the competitor succeeds.

## E. Missing Value Imputation

We visualize the images recovered by VAE and MEM-VAE in Figure 8 given incomplete test data. It can be seen that MEM-VAE has better visualization than VAE with fewer meaningless images, clearer digits and more accurate inference.