

Style Memory: Making a Classifier Network Generative

Rey Wiyatno
Mechatronics Engineering
University of Waterloo
Waterloo, ON, N2L 3G1, Canada
Email: rrwiyatn@edu.uwaterloo.ca

Jeff Orchard
Cheriton School of Computer Science
University of Waterloo
Waterloo, ON, N2L 3G1, Canada
Email: jorchard@uwaterloo.ca

Abstract—Deep networks have shown great performance in classification tasks. However, the parameters learned by the classifier networks usually discard stylistic information of the input, in favour of information strictly relevant to classification. We introduce a network that has the capacity to do both classification and reconstruction by adding a “style memory” to the output layer of the network. We also show how to train such a neural network as a deep multi-layer autoencoder, jointly minimizing both classification and reconstruction losses. The generative capacity of our network demonstrates that the combination of style-memory neurons with the classifier neurons yield good reconstructions of the inputs when the classification is correct. We further investigate the nature of the style memory, and how it relates to composing digits and letters. Finally, we propose that this architecture enables the bidirectional flow of information used in predictive coding, and that such bidirectional networks can help mitigate against being fooled by ambiguous or adversarial input.

I. INTRODUCTION

Deep neural networks now rival human performance in many complex classification tasks, such as image recognition. However, these classification networks are different from human brains in some basic ways. First of all, the mammalian cortex has many feed-back connections that project in the direction opposite the sensory stream [1]. Moreover, these feed-back connections are implicated in the processing of sensory input, and seem to enable improved object/background contrast [2], and imagination [3]. Feed-back connections are also hypothesized to be involved in generating predictions in the service of perceptual decision making [4], [5].

Humans (and presumably other mammals) are also less susceptible to being fooled by ambiguous or adversarial inputs. Deep neural networks have been shown to be vulnerable to adversarial examples [6], [7]. Slight modifications to an input can cause the neural network to misclassify it, sometimes with great confidence! Humans do not get fooled as easily, leading us to wonder if the feed-back, generative nature of real mammalian brains contributes to accurate classification.

In pursuit of that research, we wish to augment classification networks so that they are capable of both recognition (in the feed-forward direction) and reconstruction (in the feed-back direction). We want to build networks that are both classifiers and generative.

The nature of a classifier network is that it throws away most of the information, keeping only what is necessary to make accurate classifications. Simply adding feed-back connections to the network will not be enough to generate specific examples of the input – only a generic class archetype. But what if we combine the features of a classifier network and an autoencoder network by adding a “style memory” to the top layer of the network? The top layer would then consist of a classification component (eg. a softmax vector of neurons) as well as a collection of neurons that are not constrained by any target classes.

We hypothesized that adding a style memory to the top layer of a deep autoencoder would give us the best of both worlds, allowing the classification neurons to contribute the class of the input, while the style memory would record additional information about the encoded input – presumably information not encoded by the classification neurons. The objective of our network is to minimize both classification and reconstruction losses so that the network can perform both classification and reconstruction effectively. As a proof of concept, we report on a number of experiments with MNIST and Extended MNIST (EMNIST) [8] that investigate the properties of the network augmented with this style memory.

II. RELATED WORK

Others have developed neural architectures that encode both the class and style of digits to enable reconstruction. For example, a group of researchers introduced a method called bidirectional backpropagation [9]. Their network is generative because it has feed-back connections that project down from the top (softmax) layer. A digit class can be chosen at the top layer, and the feed-back connections render an instance of the desired digit in the bottom layer (as an image). However, the network always renders the same, generic sample of each class, and is incapable of generating different versions of each digit.

Networks that have the capacity to generate images have been shown to learn meaningful features. Previous work showed that in order to recognize images, the network needs to first learn to generate images [10]. Other work also showed that a network consisting of stacked Restricted Boltzmann Machines (RBMs) learns good generative models, effective

for pre-training a classifier network [11]. RBMs are stochastic in nature, so while they can generate different inputs, they are not used to generate a specific sample of input. Furthermore, others demonstrated that autoencoders pre-trained in a greedy manner also lead to better classifier networks [12]. Both [12] and [13] use tied weights, where the feed-back weight matrices are simply the transpose of the feed-forward weights; this solution is not biologically feasible. These findings have inspired other models such as stacked denoising autoencoders [14], which learn to reconstruct the original input image given a noise-corrupted input image.

Another method was proposed to map an input to a lower dimensional space that minimizes within-class distance of the input [15]. They first pre-trained a network as RBMs, and then “unrolled” the network to form a deep autoencoder. The network was then fine-tuned by performing nonlinear neighbourhood component analysis (NCA) between the low-dimensional representations of inputs that have the same class. They were able to separate the class-relevant and class-irrelevant parts by using only 60% of the lower-dimensional code units when performing nonlinear NCA, but all the codes were used to perform reconstruction. As a result, their network was able to minimize within-class distance in the lower-dimensional space while maintaining good reconstruction. Inference was then performed by using K-nearest neighbours in that lower-dimensional space. Our method is similar, but our top layer includes an explicit classification vector alongside the class-agnostic style memory.

Finally, recent work has shown that humans are also vulnerable to adversarial examples when in a time-constrained context [16]. Human subjects we given a limited time to perform an image-classification task, which was similar to the experimental settings in [17]. The hypothesis driving the experimental design is that the brain may only have time to perform a single forward pass through the sensory stream; this mimics the feed-forward behavior typical of deep neural networks. Surprisingly, human subjects were also fooled by adversarial examples generated by a basic iterative method [18], [19]. This further strengthens our motivation to create a network that is both a classifier and generative.

III. METHOD

A. Model Description

Our bidirectional network consists of an input layer, convolutional layers, fully connected layers, and an output layer. However, the output layer is augmented; in addition to classifier neurons (denoted by y in Fig. 1), it also includes style-memory neurons (denoted m in Fig. 1). A standard classifier network maps $x \in X \rightarrow y \in Y$, where the dimension of Y is usually much smaller than the dimension of X . The feed-forward connections of our augmented network map $x \in X \rightarrow (y, m) \in Y \times M$. The output y is the classification vector (softmax). The output m is the style memory, meant to encode information about the particular form of an input. For the example of MNIST, the classification vector might represent that the digit is a ‘2’, while the style memory records

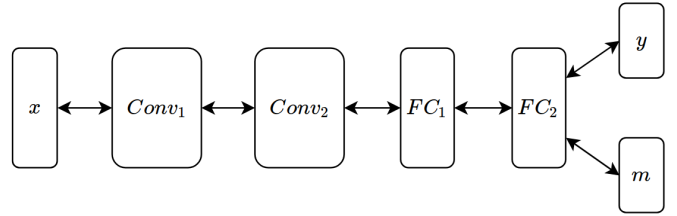


Fig. 1: Our bidirectional network with a style memory in the output layer. Here, x denotes the input ($x \in X$), while $Conv_i$ and FC_i denote convolutional layer and fully connected layer i , respectively. Lastly, y denotes output label ($y \in Y$), and m denotes the style memory ($m \in M$).

that the ‘2’ was written on a slant, and with a loop in the bottom-left corner.

A classifier network can be trained as a deep autoencoder network. However, the decoder will only be able to generate a single, generic element of a given class. Instead, by adding a style memory in the output layer, the network will be able to learn to generate a variety of different renderings of a particular class.

B. Training

We trained the network following a standard training procedure for deep autoencoders, depicted in Fig. 2. For the input layer, we follow the work from [14] by injecting small additive Gaussian noise to the input.

The objective for our network’s top layer is to jointly minimize two loss functions. The first loss function is the classifier loss L_y , which is a categorical cross-entropy loss function,

$$L_y(y_t, y) = - \sum_x y_t \log(y), \quad (1)$$

where y_t is the target label, and y is the predicted label. The second loss function is the reconstruction loss between the input and its reconstruction. This reconstruction loss, denoted L_r , is the squared Euclidean distance between the input to the network, and the reconstruction of that input,

$$L_r(\hat{x}, x) = \|\hat{x} - x\|^2, \quad (2)$$

where \hat{x} is the reconstruction of the input x , as shown in Fig. 2.

Our goal is to find connection weights, W^* , that minimize the combination of both loss functions in the last layer,

$$W^* = \arg \min_W \sum_{x \in X} L_y(y_t, y) + \alpha(L_r(\hat{x}, x)), \quad (3)$$

where W represents the parameters of the network, and α adjusts the weight of L_r .

IV. EXPERIMENTS

We performed all experiments in this paper using digits and letters from MNIST and EMNIST datasets, respectively. Both datasets have an input dimensionality of 28×28 pixels. The networks used for the experiments have two convolutional layers and two fully connected layers. The first and second

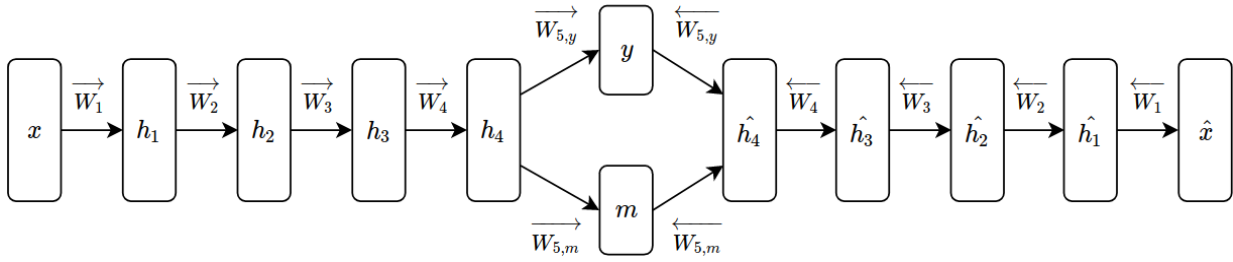


Fig. 2: The “unrolled” network. Learning consists of training the network as a deep autoencoder, where h_i denotes the hidden layer representation of layer i .



Fig. 3: Reconstruction of MNIST digits using the network’s predictions and style memories. The top row shows the original images from the MNIST test set, and the bottom row shows the corresponding reconstructions produced by the network.

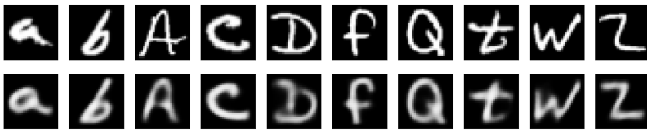


Fig. 4: Reconstruction of EMNIST letters using the network’s predictions and style memories.

convolutional layers are made of 32 and 64 filters, respectively. The receptive fields of both convolutional layers are 5×5 with a stride of 2, using ReLU activation functions. The fully connected layers FC_1 and FC_2 have 256 and 128 ReLU neurons, respectively.

The style memory consists of 16 logistic neurons, and the classifier vector contains either 10 or 26 softmax neurons depending on the datasets (MNIST or EMNIST). The reconstruction loss weight (α) was set to 0.05, and the optimization method used to train the network was Adam [20] with a learning rate η of 0.00001 for 250 epochs. The network achieved 98.48% and 91.27% classification accuracy on the MNIST and EMNIST test sets, respectively.

A. Reconstruction Using Style Memory

The reconstructions produced by our network show that the network has the capacity to reconstruct a specific sample, rather than just a generic example from a specific class. Figures 3 and 4 show examples of digit and letter reconstructions. Notice how the network has the ability to reconstruct different styles of a class, like the two different ‘4’s, two different ‘9’s, and two different ‘A’s. For each sample, the reconstruction mimics the style of the original character. Also note that the digits and letters in both figures were correctly classified by the network.

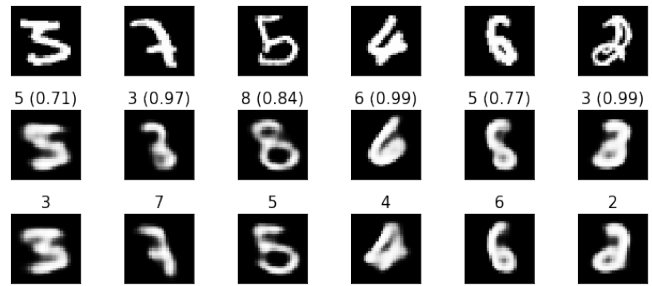


Fig. 5: Comparison of MNIST digit reconstruction using the prediction from the network versus ground-truth label. The top row shows the original images from the MNIST test set that the network misclassified. The middle row shows the reconstruction of the images using the predictions of the network, along with the predicted class and confidence score (i.e. the first image in the middle row was the reconstruction of the first image in the top row, where the network predicted the input to be a digit ‘5’ with 71% confidence). The bottom row shows the reconstructions using the corrected one-hot labels where the ground truth labels are printed on top of the images in the bottom row.

B. Reconstruction of Misclassified Samples

How do the softmax classification nodes and the style memory interact when a digit or letter is misclassified? The first column in Fig.5 shows an example where the digit ‘3’ was misclassified as a ‘5’ with 71% confidence. The resulting reconstruction in the middle row looks more like a ‘5’ (although there is a hint of a ‘3’). However, correcting the softmax neurons to the one-hot ground truth label for ‘3’ changed the reconstruction to look more like a ‘3’, as shown in the bottom row of Fig.5. Similar results were observed when we used letters from the EMNIST dataset, as shown in Fig.6. The network was able to correct the reconstruction of the letters simply by changing the predicted class to the corrected one-hot labels.

We believe that the generative abilities of these classifier networks enable them to identify misclassified inputs. Our results suggest that if the reconstruction does not closely match the input, then it is likely that the input was misclassified. This idea forms the crux of how these networks might defend

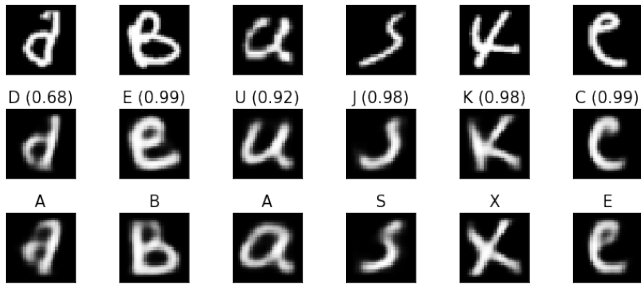


Fig. 6: Comparison of EMNIST letter reconstruction using the prediction from the network versus ground truth label. The format follows the format of Fig. 5.

against being fooled by adversarial or ambiguous inputs.

C. Style Memory Representation

To better understand what was being encoded in the style memory, we generated digits that were close together in the style memory space (16-dimensional) and compared them with digits that are close together in the image space (784-dimensional). The distance, in either space, was calculated using the Euclidean norm.

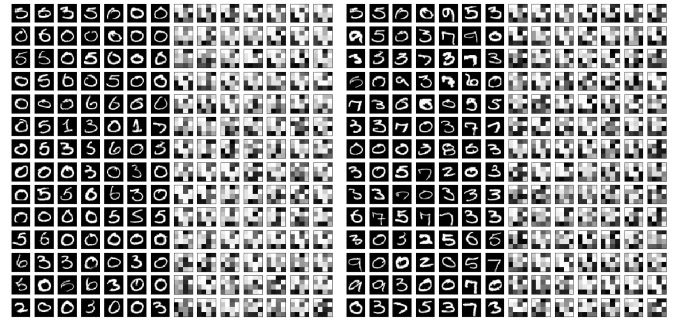
From Fig. 7 and Fig. 8, we can observe that proximity in the style-memory space has different semantic meaning than proximity in the image space. Figure 7a, showing the 97 images that are closest to the ‘5’ image in the top-left corner, displays many digits that share common pixels. However, Fig. 7b, which shows the 97 digits with the closest style memories, displays digits that come from various different classes. Similarly, Fig. 7c shows many digits of class ‘3’, while Fig. 7d is less dominated by ‘3’s.

There are 18 digits of ‘5’ in Fig. 7a, while there are only 13 digits of ‘5’ in Fig. 7b. However, Fig. 7a is actually dominated by ‘0’, even though the base digit is a ‘5’. There are 54 digits of ‘0’ in Fig. 7a, while there are only 25 digits of ‘0’ in Fig. 7b. Similarly, there are 76 digits of ‘3’ in Fig. 7c, while there are only 46 digits of ‘3’ in Fig. 7d. We also observed that the image distance between Fig. 7a and Fig. 7b increased from 8.6 to 9.3, while the style distance decreased from 1.2 to 0.98. The image distance between Fig. 7c and Fig. 7d also increased from 8.5 to 9.5, while the style distance decreased from 1.2 to 1.0.

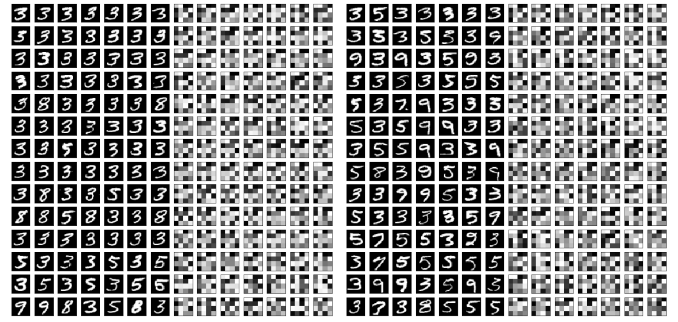
Similarly, there are 52 letters of ‘S’ in Fig. 8a, while there are only 6 letters of ‘S’ in Fig. 8b. Furthermore, there are 47 letters of ‘P’ in Fig. 8c, while there are only 17 letters of ‘P’ in Fig. 8d. The image distance between Fig. 8a and Fig. 8b increased from 9.1 to 10.5, while the style distance decreased from 1.3 to 0.91. Lastly, The image distance between Fig. 8c and Fig. 8d also increased from 8.5 to 9.8, while the style distance decreased from 1.3 to 0.99. These results show that style memory successfully separates some of the class information from the data, while also being class-agnostic.

D. Style Memory Interpolation

In this experiment, we attempted to reconstruct a continuum of images that illustrate a gradual transformation between two



(a) Image Dist=8.6, Style Dist=1.2 (b) Image Dist=9.3, Style Dist=0.98



(c) Image Dist=8.5, Style Dist=1.2 (d) Image Dist=9.5, Style Dist=1.0

Fig. 7: Nearest neighbours in image space and style-memory space. Figures (a) and (c) show the 97 digit images closest to the image in the top-left, as well as their corresponding style-memories. Figures (b) and (d) show the 97 style memories closest to the style memory in the top-left, as well as their corresponding digit images. The order of elements (across rows, then down) indicate increasing Euclidean distance. The subfigure captions give the average distance from the top-left element, both in image space, and style-memory space.

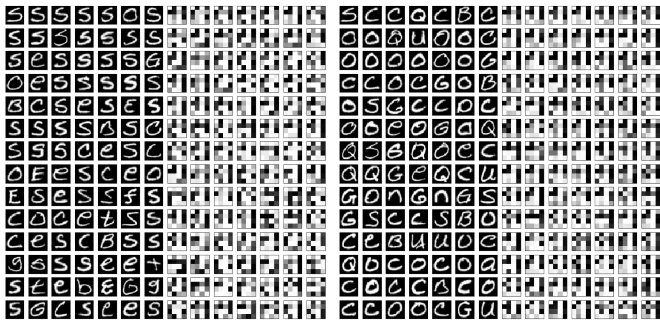
different styles of the same character class. For example, we encoded two different digits for each MNIST class, as shown in Fig. 9, as well as a selection of characters classes from EMNIST, shown in Fig. 10. We then generated a sequence of images that slowly evolve from one style to the other. We performed the interpolation by simply taking convex combinations of the two style memories, using

$$\hat{m}(\lambda) = \lambda m_1 + (1 - \lambda) m_2, \quad (4)$$

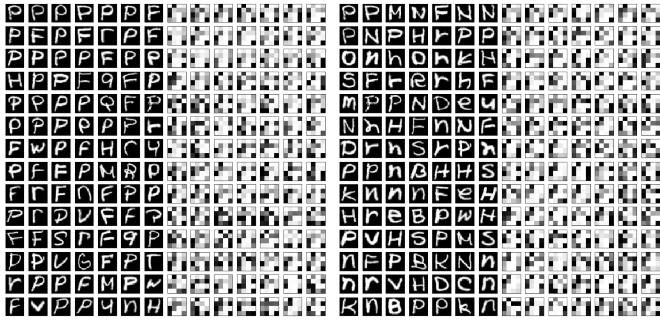
where $\lambda \in [0, 1]$, and m_1 and m_2 denote the style memories. The interpolated style memory is denoted by $\hat{m}(\lambda)$.

Figure 11 shows the interpolated digits and letters, illustrating that the generated images transform smoothly when the style memory is interpolated. The results of within-class interpolation suggest that style memory captures style information about how a digit was drawn. The figure also shows examples of attempted interpolations between incongruous letter forms (eg. ‘A’ to ‘a’, and ‘r’ to ‘R’). Not surprisingly, the interpolated characters are nonsensical in those cases.

An obvious experiment is to try transferring the style memory of one digit onto another digit class. Although not



(a) Image Dist=9.1, Style Dist=1.3 (b) Image Dist=10.5, Style Dist=0.91



(c) Image Dist=8.5, Style Dist=1.3 (d) Image Dist=9.8, Style Dist=0.99

Fig. 8: Nearest neighbours in image space and style-memory space of EMNIST dataset.



Fig. 9: Two different styles of digits form the endpoints for the style interpolation experiment.

shown here, we observed that the style memory of a digit can, in some cases, be transferred to some other classes. However, in general, the reconstructions did not look like characters.

V. CONCLUSION AND FUTURE WORK

Classification neural networks do not typically maintain enough information to reconstruct the input; they do not have to. Their goal is to map high-dimensional inputs to a small number of classes, typically using a lower-dimensional vector representation. Thus, simply adding feed-back connections to a classification network and training it as an autoencoder will not yield a network that simultaneously performs classification

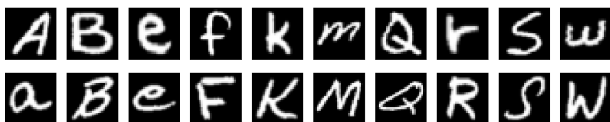
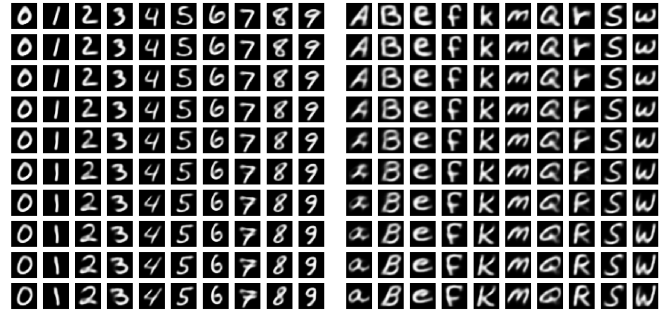


Fig. 10: Two different styles of letters form the endpoints for the style interpolation experiment.



(a) (b)

Fig. 11: Image reconstruction with style memory interpolation between digits and letters shown in Fig. 9 and Fig. 10, where λ was increasing from 0.1 to 1.0 with a step of 0.1 from top to bottom.

in the feed-forward direction, and specific sample generation in the feed-back direction.

In order for a classification network to be capable of generating samples, additional information needs to be maintained. In this paper, we proposed the addition of “style memory” to the top layer of a classification network. The top layer is trained using a multi-objective optimization, trying to simultaneously minimize both classification and reconstruction errors.

Our experiments suggest that the style memory encodes information that is largely disjoint from the classification vector. For example, proximity in image space yields digits that employ an overlapping set of pixels. However, proximity in style-memory space yielded a different set of digits. Further investigation into this phenomenon is needed. In fact, it would be interesting if we could articulate what the style memory is encoding. But we concede that words might fail us in this endeavour; the encoding of style could be very difficult to pinpoint with language.

For the style interpolation experiment, we generated images from a straight line in style-memory space. Each position on the line generates a sample in image space – an image; it would be interesting to see what shape that 1-dimensional manifold takes in image space, and how it differs from straight-line interpolation in image space. However, the fact that we were able to interpolate digits and letters within the same class using novel style-memory activation patterns suggests that the style memory successfully encodes additional, abstract information about the encoded input.

To our knowledge, existing defense mechanisms to combat adversarial inputs do not involve running the network in both feed-forward and feed-back directions during inference. The closest defense approaches that involve generative capacity of the network were proposed in PixelDefend [21] and DefenseGAN [22]. While PixelDefend only allows the network to run in a feed-forward manner, DefenseGAN allows the generative network to update its parameters by performing stochastic gradient descent during inference to achieve better reconstruction.

tions. However, Defense-GAN’s generative model (Wasserstein GAN [23]) was trained separately from the classifier network. The generative model in PixelDefend and Defense-GAN were only designed to “purify” or “denoise” any inputs, with the hope that the generative model creates versions of the inputs that have the same distribution as the training data before passing them to the classifier network. This approach makes the prediction and reconstruction functions of the network independent of each other, which is in contrast with what we believe: classification and reconstruction should be coupled together. Moreover, both of these defense tactics have already been successfully attacked [24].

We saw that the network has a property where the reconstruction generated was affected both by the classification neurons and style memory. Motivated by the results in Sec. IV-A and Sec. IV-B, preliminary experiments that we have done suggest that treating perception as a two-way process, including both classification and reconstruction as part of training and inference, is effective for guarding against being fooled by adversarial or ambiguous inputs.

Finally, inspired by how human perception is influenced by expectation [5], we believe that this work opens up opportunities to create a classifier network that takes advantage of its generative capability to detect misclassifications since the reconstruction of the input does not match the input when the classification is incorrect. For example, perception and inference could be the result of running the network in feed-forward and feed-back directions simultaneously, like in the wake-sleep approach [25]. Moreover, predictive estimator networks might be a natural implementation for such feed-back networks [5], [26], [27]. Continuing in this vein is left for future work and these experiments are ongoing.

REFERENCES

- [1] J. Bullier, M. McCourt, and G. Henry, “Physiological studies on the feedback connection to the striate cortex from cortical areas 18 and 19 of the cat,” *Experimental Brain Research*, vol. 70, no. 1, pp. 90–98, 1988.
- [2] J. Poort, F. Raudies, A. Wannig, V. A. F. Lamme, H. Neumann, and P. R. Roelfsema, “The role of attention in figure-ground segregation in areas V1 and V4 of the visual cortex,” *Neuron*, vol. 75, no. 1, pp. 143–156, 2012.
- [3] L. Reddy, N. Tsuchiya, and T. Serre, “Reading the mind’s eye: decoding category information during mental imagery,” *NeuroImage*, vol. 50, no. 2, pp. 818–825, 2011.
- [4] R. P. N. Rao and D. H. Ballard, “Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects,” *Nature Neuroscience*, vol. 2, pp. 79 EP –, Jan 1999, article. [Online]. Available: <http://dx.doi.org/10.1038/4580>
- [5] C. Summerfield and F. P. De Lange, “Expectation in perceptual decision making: Neural and computational mechanisms,” *Nature Reviews Neuroscience*, vol. 15, no. 11, pp. 745–756, oct 2014. [Online]. Available: <http://dx.doi.org/10.1038/nrn3838>
- [6] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” in *International Conference on Learning Representations*, 2014. [Online]. Available: <http://arxiv.org/abs/1312.6199>
- [7] I. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” in *International Conference on Learning Representations*, 2015. [Online]. Available: <http://arxiv.org/abs/1412.6572>
- [8] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik, “EMNIST: an extension of MNIST to handwritten letters,” *CoRR*, vol. abs/1702.05373, 2017.
- [9] H. Luo, J. Fu, and J. R. Glass, “Bidirectional backpropagation: Towards biologically plausible error signal transmission in neural networks,” *CoRR*, vol. abs/1702.07097, 2017.
- [10] G. E. Hinton, “To recognize shapes, first learn to generate images,” in *Computational Neuroscience: Theoretical Insights into Brain Function*, ser. Progress in Brain Research, P. Cisek, T. Drew, and J. F. Kalaska, Eds. Elsevier, 2007, vol. 165, no. Supplement C, pp. 535 – 547.
- [11] R. Salakhutdinov and G. Hinton, “Deep Boltzmann machines,” in *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, D. van Dyk and M. Welling, Eds., vol. 5. Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA: PMLR, 16–18 Apr 2009, pp. 448–455.
- [12] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, “Greedy layer-wise training of deep networks,” in *Proceedings of the 19th International Conference on Neural Information Processing Systems*, ser. NIPS’06. Cambridge, MA, USA: MIT Press, 2006, pp. 153–160.
- [13] G. E. Hinton, S. Osindero, and Y. W. Teh, “A fast learning algorithm for deep belief nets,” *Neural Computation*, vol. 18, pp. 1527–1554, 2006.
- [14] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, “Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion,” *J. Mach. Learn. Res.*, vol. 11, pp. 3371–3408, Dec. 2010.
- [15] R. Salakhutdinov and G. E. Hinton, “Learning a nonlinear embedding by preserving class neighbourhood structure,” in *International Conference on Artificial Intelligence and Statistics*, 2007, pp. 412–419.
- [16] G. F. Elsayed, S. Shankar, B. Cheung, N. Papernot, A. Kurakin, I. Goodfellow, and J. Sohl-Dickstein, “Adversarial Examples that Fool both Human and Computer Vision,” *ArXiv e-prints*, Feb. 2018.
- [17] M. C. Potter, B. Wyble, C. E. Hagmann, and E. S. McCourt, “Detecting meaning in RSVP at 13 ms per picture,” *Attention, Perception, & Psychophysics*, vol. 76, no. 2, pp. 270–279, Feb 2014. [Online]. Available: <https://doi.org/10.3758/s13414-013-0605-z>
- [18] A. Kurakin, I. J. Goodfellow, and S. Bengio, “Adversarial examples in the physical world,” *CoRR*, vol. abs/1607.02533, 2016. [Online]. Available: <http://arxiv.org/abs/1607.02533>
- [19] —, “Adversarial machine learning at scale,” 2017. [Online]. Available: <https://arxiv.org/abs/1611.01236>
- [20] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- [21] S. N. S. E. N. K. Yang Song, Taesup Kim, “Pixeldefend: Leveraging generative models to understand and defend against adversarial examples,” *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=rJUYGxbCW>
- [22] R. C. Pouya Samangouei, Maya Kabkab, “Defense-GAN: Protecting classifiers against adversarial attacks using generative models,” *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=BkJ3ibb0>
- [23] M. Arjovsky, S. Chintala, and L. Bottou, “Wasserstein generative adversarial networks,” in *Proceedings of the 34th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, D. Precup and Y. W. Teh, Eds., vol. 70. International Convention Centre, Sydney, Australia: PMLR, 06–11 Aug 2017, pp. 214–223. [Online]. Available: <http://proceedings.mlr.press/v70/arjovsky17a.html>
- [24] A. Athalye, N. Carlini, and D. Wagner, “Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples,” 2018. [Online]. Available: <https://arxiv.org/abs/1802.00420>
- [25] G. E. Hinton, P. Dayan, B. J. Frey, and R. M. Neal, “The “Wake-Sleep” Algorithm for Unsupervised Neural Networks,” *Science*, vol. 268, no. 5214, pp. 1158–1161, 1995.
- [26] D. Xu, A. Clappison, C. Seth, and J. Orchard, “Symmetric predictive estimator for biologically plausible neural learning,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. PP, no. 99, pp. 1–12, 2017.
- [27] J. Orchard and L. Castricato, “Combating Adversarial Inputs Using a Predictive-Estimator Network,” in *Proc. of the International Conference on Neural Information Processing*, vol. LNCS 10638, 2017, pp. 118–125.