

Investigating Diversity in Responses to Artwork through Latent Representations of Language Captions

Ananya Devarakonda
MAHE, Manipal

Samuel Showalter
UC, Irvine

October 2021

Abstract

The field of natural language processing (NLP) has seen immense growth in applications ranging from machine translation to text generation. There has been a rise in the use of generative models in NLP applications. The variational autoencoder (VAE) is a powerful generative model that learns probabilistic representations of latent space encodings. There are many versions of the VAE, including the β -VAE, which introduces a hyperparameter β to encourage disentangled representations in the latent space. This paper examines the latent space encodings of an LSTM β -VAE trained on the ArtEmis dataset to generate emotion-driven text and draw conclusions about specific works of art and genres in the dataset. The ArtEmis dataset is a collection of subjective captions of artwork based on the WikiArt dataset. We quantitatively evaluate the diversity of responses elicited by a work of art by calculating the Euclidean distance between points in the learned latent space representation of captions in the dataset and their respective centroid. We use the average distance as a metric to evaluate the diversity of responses elicited by a particular work of art. In addition, we qualitatively examine the latent representations of captions in the dataset using PCA. We found that analytical cubism, synthetic cubism, pointillism, contemporary realism, and action painting were the genres of art that elicited the most diverse responses. The source code for the project can be found here: <https://github.com/ananya/Artemis-VAE.git>

1 Introduction

Natural language processing (NLP), a subfield of artificial intelligence, deals with computational linguistics with applications in areas such as document classification Adhikari et al., 2019; Yang et al., 2016, machine translation Lample et al., 2018; Wang et al., 2019, and sentiment analysis Agarwal et al., 2011; C. Lin and He, 2009. Typically, applications in NLP use autoregressive models such as recurrent neural networks (RNNs) and their variants- gated recurrent units (GRUs) and long short-term memory (LSTM) models. An essential and well-studied application of NLP is text generation, which deals with the automatic production of readable and coherent text. Recently, the use of generative models has shown promising results in the application of text generation. Among the generative methods used, the variational autoencoder (VAE), as introduced by Kingma and Welling Kingma and Welling, 2014, has proven to be a powerful self-supervised model for text generation Bowman et al., 2016. Variational autoencoders (VAEs) are distinct from standard autoencoders as they take a Bayesian approach to learn continuous latent representations of features in the dataset. One version of the VAE, known as the β -VAE Higgins

et al., 2016, enforces the learning of disentangled representations in the latent space of the VAE by adding a hyperparameter β . It is valuable to learn disentangled or statistically independent representations to coherently infer properties of a particular dataset to draw conclusions about fundamental relationships between features in the dataset.

This paper outlines an emotion-driven text generation model using a β -VAE trained on the ArtEmis dataset Achlioptas et al., 2021. The ArtEmis dataset is an image captioning dataset that contains around 455,000 descriptions of 80,000 works of art. Unlike traditional image captioning datasets, such as the Flickr 8k Rashtchian et al., 2010, Flickr 30k Rashtchian et al., 2010, and MS-COCO datasets T.-Y. Lin et al., 2014 that focus on entirely objective descriptions of images, the ArtEmis dataset contains captions of artwork that highlight personal interpretation and emotional response.

This paper attempts to quantify a topic considered subjective: ranking artwork based on the diversity of the responses elicited. Quantification is done by examining the similarities in the latent encodings learned by the β -VAE. The paper also draws conclusions about artwork mentioned in the ArtEmis dataset and comments on the diversity of responses elicited by a particular work of art. Thus, our contributions are as follows:

1. Introducing an emotion-driven text generation method using an LSTM β -VAE model.
2. Exploring latent space representations of the captions in the ArtEmis dataset to quantitatively analyze features in the dataset.
3. Qualitatively analyzing diversity in response to artwork using PCA.

2 Literature Review

Several subfields of artificial intelligence, including computer vision and natural language processing, use the variational autoencoder (VAE) for various applications. In computer vision, the VAE is popularly used in image generation M.-Y. Liu et al., 2017; Pu et al., 2016; Vahdat and Kautz, 2021; Yan et al., 2016 by interpolating between the learned disentangled latent representations of features in the dataset. In addition to image generation, Pu et al. Pu et al., 2016 developed a novel VAE to generate accurate labels and captions along with images. The versatility of the VAE as a generative model is further highlighted by Yan et al. Yan et al., 2016, who introduced a VAE model to generate images from high-level textual descriptions.

The VAE has shown promising results in the field of natural language processing (NLP) as well. Hayashi et al. Hayashi and Watanabe, 2020 used a VAE for an end-to-end speech-to-text model and stated that they obtained state-of-the-art results. Sheng et al. Sheng et al., 2020 proposed a novel VAE model for accurate neural machine translation (NMT). Other applications of VAEs in NLP include document modelling Miao et al., 2016, dialog generation Wen et al., 2017, and text summarization Miao and Blunsom, 2016. Text generation, an essential application of natural language generation (NLG), is a subfield of NLP that involves the automatic generation of grammatically correct and coherent text. Probabilistic models used for sentence generation trained to predict words based on sequential input data are dubbed language models (LMs). Traditionally, language models used autoregressive models including recurrent neural networks (RNNs), long short-term memory (LSTM) models, and gated recurrent units (GRUs). Such models are called recurrent neural network language models, or RNNLMs Mikolov and Zweig, 2012. Vinyals et al. Vinyals et al., 2015 used an LSTM-based model for image caption generation and obtained realistic results. More recently, Liu et al. B. Liu et al., 2018 used an RNN decoder for poem generation in an image-to-poem model.

Although RNNLMs are shown to generate coherent sentences, they do not learn interpretable representations of dataset attributes; in fact, the behavior of an RNNLM remains black-box.

Unlike RNNLMs, models based on the variational autoencoder (VAE) introduce an element of interpretability by regularizing for a smooth latent representation of features in a dataset. Thus, VAEs have been extensively used in text generation models to not only generate accurate sentences but also to reveal properties of the dataset by analyzing learned latent representations. In particular, the β -VAE, encourages the learning of disentangled representations of features in the latent space. This quality is conducive to learning continuous hidden representations of captions from the ArtEmis dataset. These latent representations can then be analyzed quantitatively to help sort the works of art mentioned in the ArtEmis dataset by the diversity of response and interpretation- a topic that is inherently subjective.

Bowman et al. Bowman et al., 2016 introduced an LSTM-based VAE model to generate sentences by interpolating within the learned latent space. They found that the VAE model, without optimization, behaved like an RNNLM. This problem is dubbed “posterior collapse.” Pelsmaeker et al. Pelsmaeker and Aziz, 2020 describe and compare methods of overcoming the problem of posterior collapse. Bowman et al. Bowman et al., 2016 tackled this problem by optimizing the model using KL-cost annealing as well as word dropout and historyless decoding. Post-optimization, the model produced coherent and diverse sentences by interpolating between learned latent representations. Yang et al. Yang et al., 2017 proposed a variational autoencoder model with an LSTM encoder and a dilated convolutional neural network (CNN) decoder for conditional text generation.

In our implementation, we give priority to interpretability over the general text generation model and place more emphasis on the posterior, as our objective is to rank the diversity in the response of artwork using captions. The papers mentioned primarily focus on generating sentences as accurately as possible and do not provide much information on inferring properties of the datasets using learned latent representations. This paper provides information on the model used for emotion-driven text generation and analyses the features in the ArtEmis dataset by examining the learned latent encodings to sort works of art based on the diversity of their interpretations.

3 Experimental Design

3.1 Dataset

The ArtEmis dataset Achlioptas et al., 2021 is a large-scale image captioning dataset that contains over 455,000 emotion attributes and subjective captions of over 80,000 works of art. Image captioning datasets such as the Flickr 8k Rashtchian et al., 2010, Flickr 30k Rashtchian et al., 2010, and MS-COCO datasets T.-Y. Lin et al., 2014 contain only objective descriptions of real-world images. Unlike real-world image data, art is intrinsically subjective. Thus, for effective image captioning of art, it is essential to retain elements of personal interpretation. Recognizing that subjectivity is fundamental to artwork, the ArtEmis dataset contains captions highlighting the emotions elicited by a particular work of art. Much like real-world text data, the captions from the ArtEmis dataset contain grammatical errors, spelling errors, and show human bias. Captions are categorized by art style (“art_style”) and name of the painting (“painting”). The authors also specify a column to assign a general emotional attribute to each artwork (“emotion”). An example from the dataset, along with the WikiArt image, is given in Figure 1 Achlioptas et al., 2021.

It is essential to recognize that this emotional attribute is chosen based on eight emotion options (amusement, anger, awe, fear, sadness, contentment, disgust, excitement) and one default option (something else). Around 53,000 records categorize artwork as “something else.” Furthermore, much like real-world data, there are instances in the dataset where the chosen

emotion attribute does not align with the sentence. Thus, using only emotional attributes to judge the diversity of emotional responses to an artwork would be inaccurate. Keeping this in mind, we use raw captions from the ArtEmis dataset so that the model learns appropriate latent representations for further analysis.

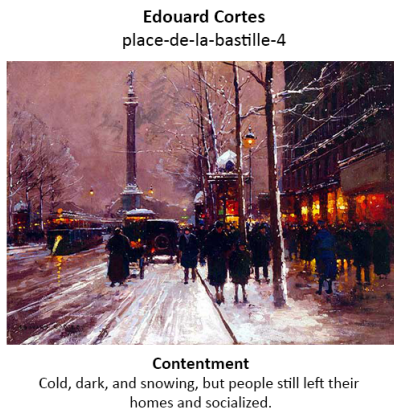


Figure 1: An example from the ArtEmis dataset

3.2 Preprocessing

The dataset used to train the model consists of all the captions in the ArtEmis dataset. First, we tokenize the captions. Tokenization refers to the splitting of each sentence in the dataset into its respective words and symbols. Next, we split the dataset containing all captions into a training (99%) and validation (1%) dataset. The validation dataset ensures that the model has not overfit and performs well with unseen data. We then build a vocabulary of all words and symbols in the tokenized training dataset using pre-trained GloVe vectors Pennington et al., 2014. While building the vocabulary, we add special tokens such as the $\langle \text{sos} \rangle$, $\langle \text{eos} \rangle$, and $\langle \text{pad} \rangle$ tokens that refer to start of sentence, end of sentence, and padding respectively. A few sentences in the ArtEmis dataset have no spaces between words. For such instances, we have added spaces wherever required.

3.3 Model

3.3.1 Variational Autoencoder (VAE)

Autoencoders are self-supervised deep generative neural network models that feature an encoder-decoder architecture with a bottleneck layer having reduced dimensions to learn latent representations of given data. They are trained to reconstruct a given input. The primary difference between a variational autoencoder (VAE) Kingma and Welling, 2014 and a standard autoencoder is that the VAE encodes features as a continuous distribution over the latent space. Thus, instead of encoding singular points for each feature, the variational autoencoder learns a “soft” representation of the data in the latent space.

Fundamentally, for an input x the VAE learns a latent representation z and generates a reconstruction of x , say, \hat{x} . In such a case, $p(z|x)$ is given as:

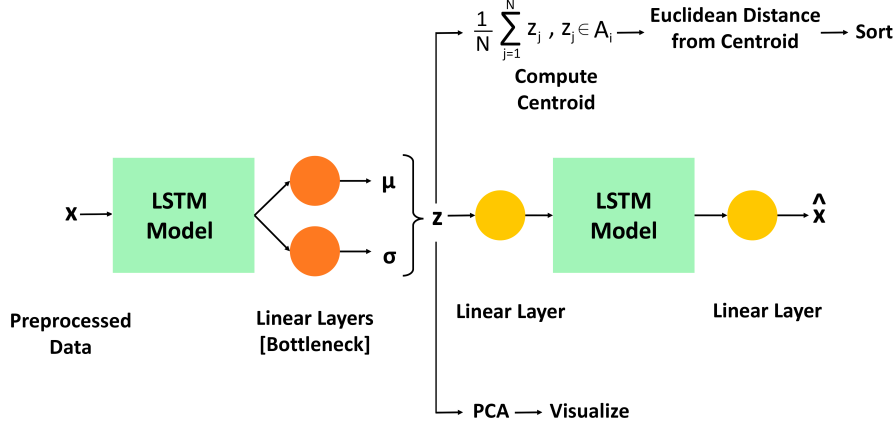


Figure 2: An illustration of the experimental design

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)} \quad (1)$$

Usually, $p(x)$ is an intractable distribution, and thus, variational inference is used to estimate the value of $p(z|x)$. The process involves selecting another distribution, say, $q(z|x)$ that we know is tractable (such as the Gaussian distribution) and making it as similar to $p(z|x)$ as possible. This is done by minimizing the KL divergence between the two distributions. Thus, the loss function used by the VAE is given in equation (2):

$$L(x, \hat{x}) = l(x, \hat{x}) + KL(q(z|x) || p(z|x)) \quad (2)$$

Where $l(x, \hat{x})$ represents the reconstruction loss. Reconstruction loss represents a metric such as mean squared error (MSE) or categorical cross-entropy and is chosen according to the assumed prior. Data generation using the VAE is done by sampling from the learned distribution in the latent space. As the output from sampling is discrete, backpropagation is not directly possible. Thus, Kingma et al. Kingma and Welling, 2014 proposed the reparameterization trick that made backpropagation possible. Assuming a Gaussian prior, sampling can be represented by:

$$z \sim q(z|x) = \mathcal{N}(z; \mu, \sigma)$$

Then, the reparameterization trick is as follows:

$$z = \mu + \epsilon \cdot \sigma, \text{ where } \epsilon \sim \mathcal{N}(0, I)$$

This trick helps to train the VAE model using gradient descent and backpropagation.

The β -VAE Higgins et al., 2016 is a version of the original VAE model that encourages the learning of continuous disentangled representations by adding a Lagrangian multiplier β to equation (2), as represented by:

$$L(x, \hat{x}) = l(x, \hat{x}) + \beta KL(q(z|x) || p(z|x))$$

When $\beta = 1$ the model behaves equivalently to the original VAE. For values of $\beta > 1$, the model encourages the learning of disentangled feature representations in the latent space.

3.3.2 LSTM β -VAE

This paper uses a β -variational autoencoder (β -VAE) Higgins et al., 2016 with an LSTM encoder and decoder as introduced by Bowmann et al. Bowman et al., 2016. The model architecture resembles that of the sequence-to-sequence (Seq2Seq) model but has fully connected layers to encode latent representations after the LSTM encoder. Figure 2 provides an illustration of the VAE model with an LSTM encoder and decoder.

The main advantage of using the β -VAE is that it forces the model to learn continuous disentangled representations of the features in the dataset. The model we define uses categorical cross-entropy as the reconstruction loss and Adam optimization. After training over 60 epochs, the model achieved a final training loss of 0.046. The LSTM β -VAE model was optimized using KL cost annealing by gradually increasing the hyperparameter β during training. The process uses sigmoid annealing where β was increased from a value of one to eight. KL annealing helps tackle the problem of “posterior collapse” and ensures that the value of the KL term does not diminish. Thus, KL annealing encourages the LSTM β -VAE model to learn continuous disentangled latent representations of captions in the ArtEmis dataset.

3.4 Latent Space Analysis

To analyze the latent space of captions from the ArtEmis dataset, after sorting by artwork, we compute the centroid by calculating the mean of the latent representation z learned by the LSTM β -VAE for a particular work of art or art style. We then compute the Euclidean distance of all latent encodings of captions from the centroid and save the average distance. We use the average distance as the metric to sort artwork and art styles mentioned in the ArtEmis dataset based on the diversity of the response it elicits. Thus, the higher the average distance from the centroid, the greater the diversity of response and vice-versa.

Finally, for qualitative analysis, we use principal component analysis (PCA) to linearly project the latent space from a 100-dimensional representation to a two-dimensional representation of captions ($\mathbb{R}^{100} \rightarrow \mathbb{R}^2$) for each work of art. After PCA, we plot the points onto a two-dimensional plane and qualitatively infer properties of the latent representations.

*this young lady appears to be very
pleased to sit for her portrait .*
*this young lady appears to be very pleased
to sit for her portrait .*
*the person walking in the woods look
peaceful and at ease .*
*the person walking in the woods look
peaceful and at ease .*

Table 1: Examples of reconstructed sentences from training data

4 Results

4.1 Text Generation

The LSTM β -VAE is trained to reconstruct captions from the ArtEmis dataset to accurately learn their latent space representations. To measure the accuracy of reconstruction, we use cross-

<i>beautiful painting of a tree with bright yellow and red leaves .</i>
<i>beautiful painting of a tree with bright yellow and red leaves .</i>
<i>it looks like a beautiful evening near the sea .</i>
<i>it looks like a beautiful evening behind the sea .</i>

Table 2: Examples of reconstructed sentences from validation data

entropy and found that the validation data had a final cross-entropy of **0.068**. Tables 1 and 2 present a few examples of the reconstructed captions from the train and validation datasets respectively. In the both tables, the sentence in bold represents the source sentence and the sentence below it represents the reconstructed sentence. As illustrated in tables 1 and 2, the model accurately learns to reconstruct the input data. The results also show that the model replicates grammatical errors present in the dataset. Hence, there is scope for future research in building a model that is capable of generating grammatically correct text even with errors in the dataset.

4.2 Sentence Interpolation

As the LSTM β -VAE model is trained to learn continuous latent representations of the captions in the ArtEmis dataset, it is possible to generate sentences by sampling from this continuous latent space. This property of the model makes it possible to choose two captions from the dataset and generate intermediate sentences by sampling from the smooth learned latent space. This process is called sentence interpolation. Table 3 provides an example of interpolation from a positive sentence expressing awe to a negative sentence expressing disgust. The two sentences selected from the dataset for interpolation are in bold. Although these sentences contain grammatical errors, they demonstrate a gradual shift in the overall emotion and meaning of each sentence. This gradual change in emotion through each step in interpolation is representative of the continuous learned latent space.

<i>the scene is quite beautiful and reminds me of fairy fairy tales .</i>
<i>the scene is majestic it and reminds me of magic queen .</i>
<i>the scene looks looks and seems exciting of and evil .</i>
<i>the scene made look and is serious and very detailed .</i>
<i>the man looks to and looks mean very heroic .</i>
<i>the man on the right looks very mean looking.</i>

Table 3: Sentence interpolation

Painting	Average Distance
alphonse-mucha holy-mount-athos-1926	3.130717
andrea-mantegna madonna-with-saints-st-john- thebaptist-st-gregory-i-the-great- st-benedict-1506	3.106543
brice-marden suicide-note-1973	3.051492
el-greco portrait-of-a-man-2	3.033126

Table 4: Top five artwork with the greatest average distances

Painting	Average Distance
thomas-eakins photograph-1910-8	0.001139
edwin-henry-landseer a-distinguished-member-of- the-humane-society	0.000937
louay-kayyali motherhood-1974	0.000897
jacob-jordaens bust-of-satyr-1621	0.000661
nikolay-bogdanov-belsky the-former-defender-of-the -homeland	0.000575

Table 5: Artwork with the least average distances

4.3 Ranking Based on Diversity of Response

4.3.1 Artwork

Tables 4 and 5 contain the results of the latent space analysis done on captions from each work of art from the ArtEmis dataset. Table 4 lists the top five most diverse works of art after analysis of the learned latent representations of the LSTM β -VAE. Similarly, table 5 presents the five works of art with the most similar (least diverse) responses to art from the ArtEmis dataset. Alphonse Mucha’s “Holy Mount Athos” was recognized as the work of art with the most diverse responses while Nikolay Bogdanov-Belsky’s “The former Defender of the Homeland” was recognized as the artwork with the most similar interpretations.

To qualitatively analyze the diversity of the responses, figure 3 provides a plot of the latent space of a work of art after PCA. Clearly, captions labelled with a particular emotion in the ArtEmis dataset seem to form distinct clusters. This is however, not always the case as in a few instances there exists variation in emotion within clusters. This variation can be attributed to the inherent subjectivity in the topic of emotions. It is well known that there is no quantitative distinction between one emotion and another. By nature, emotions are defined by personal

experience and this subjectivity can often lead to a blur between one emotion and another.

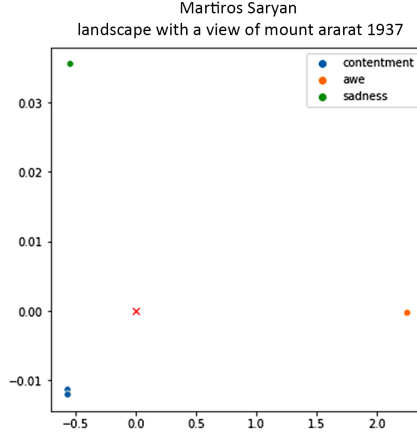


Figure 3: A scatter plot of latent representations after PCA

Art Style	Average Distance
Analytical Cubism	1.632765
Synthetic Cubism	1.562646
Pointillism	1.557388
Contemporary Realism	1.521055
Action painting	1.47832

Table 6: Art styles with the greatest average distances

4.3.2 Art Style

Along with individual works of art, we used the LSTM β -VAE model to rank art styles based on the general diversity of response to works of art within that particular art style. As humans beings easily relate to typical shapes and objects, we hypothesized that nonrepresentational designs would elicit more diverse responses. The results validated this hypothesis as we found that art in the styles of analytical and synthetic cubism had the most diversity in response while art in the style of expressionism had the most similar responses (least variation in response). Table 6 shows the top five styles of art that elicit the most diverse responses along with their average Euclidean distance. Table 7 shows the five art styles with the most similar responses. Our hypothesis is further supported by the result that action painting, a style of art that is by definition abstract and open to interpretation, is among the top five art styles that elicit the most diverse responses.

Art Style	Average Distance
Impressionism	1.254654
Art Nouveau Modern	1.254276
Baroque	1.251281
Post Impressionism	1.250021
Expressionism	1.247783

Table 7: Art styles with the least average distances

5 Discussion

The results obtained from the LSTM β -VAE fall into three main categories- text generation, sentence interpolation, and ranking of art based on the diversity of response elicited by that work of art. This paper primarily focuses on the subjective topic of ranking artwork based on the different responses people have by viewing a particular work of art or art style. The result that art styles such as analytical cubism, synthetic cubism, pointillism, contemporary realism, and action paintings elicit the most diverse responses in people, aligns with the historical context and type of art.

It is important to note that while evaluating our model, we have made the assumption that the data contains grammatically correct and accurate captions for artwork. However, as the dataset represents real-world data, this is not always the case. Thus, there is scope for future work in building a model that captures the essence of real-world data without its limitations such as grammatical errors and human bias.

6 Conclusion

While the generation of artwork using artificial intelligence is a popular topic today in the domain of computer science, analyzing the inherent subjectivity and variation in the interpretation of works of art remains unexplored. In this paper, we use an LSTM β -VAE model to learn latent representations of the captions in the ArtEmis dataset and calculate the average Euclidean distance between points in these latent representations for each work of art and art style. We use the average distance as a metric to sort individual works of art and art styles based on how differently human beings respond to them.

We found that our model produced accurate and coherent text with unseen data. Further, we demonstrated that it is possible to generate sentences by sentence interpolation and as illustrated in the results, there is scope for building a model with grammatically correct interpolations and higher accuracy. Finally, we found that certain works of art had a much greater average distance as compared to others and they primarily fall in the style of analytical cubism. In contrast, artwork with the least average distances mostly were in the style of expressionism.

Possible directions for future work include testing variants of the LSTM β -VAE model and accurately representing the captions from the ArtEmis dataset by considering them as characteristic of real-world data while taking into account inaccuracies, human bias, and grammatical errors.

References

- Achlioptas, P., Ovsjanikov, M., Haydarov, K., Elhoseiny, M., & Guibas, L. J. (2021). Artemis: Affective language for visual art. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11569–11579.
- Adhikari, A., Ram, A., Tang, R., & Lin, J. (2019). Docbert: Bert for document classification.
- Agarwal, A., Xie, B., Vovsha, I., Rambow, O., & Passonneau, R. J. (2011). Sentiment analysis of twitter data. *Proceedings of the workshop on language in social media (LSM 2011)*, 30–38.
- Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A., Jozefowicz, R., & Bengio, S. (2016). Generating sentences from a continuous space. *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, 10–21. <https://doi.org/10.18653/v1/K16-1002>
- Hayashi, T., & Watanabe, S. (2020). Discretalk: Text-to-speech as a machine translation problem. *arXiv e-prints*, arXiv–2005.
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., & Lerchner, A. (2016). Beta-vae: Learning basic visual concepts with a constrained variational framework.
- Kingma, D. P., & Welling, M. (2014). Auto-encoding variational bayes.
- Lample, G., Ott, M., Conneau, A., Denoyer, L., & Ranzato, M. (2018). Phrase-based & neural unsupervised machine translation. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 5039–5049.
- Lin, C., & He, Y. (2009). Joint sentiment/topic model for sentiment analysis. *Proceedings of the 18th ACM conference on Information and knowledge management*, 375–384.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollar, P., & Zitnick, L. (2014). Microsoft coco: Common objects in context (ECCV). *ECCV*. <https://www.microsoft.com/en-us/research/publication/microsoft-coco-common-objects-in-context/>
- Liu, B., Fu, J., Kato, M. P., & Yoshikawa, M. (2018). Beyond narrative description: Generating poetry from images by multi-adversarial training. *Proceedings of the 26th ACM international conference on Multimedia*, 783–791.
- Liu, M.-Y., Breuel, T., & Kautz, J. (2017). Unsupervised image-to-image translation networks. *Advances in neural information processing systems*, 700–708.
- Miao, Y., & Blunsom, P. (2016). Language as a latent variable: Discrete generative models for sentence compression. *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 319–328. <https://doi.org/10.18653/v1/D16-1031>
- Miao, Y., Yu, L., & Blunsom, P. (2016). Neural variational inference for text processing. *International conference on machine learning*, 1727–1736.
- Mikolov, T., & Zweig, G. (2012). Context dependent recurrent neural network language model. *2012 IEEE Spoken Language Technology Workshop (SLT)*, 234–239.
- Pelsmaeker, T., & Aziz, W. (2020). Effective estimation of deep generative language models. *ACL*.
- Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global vectors for word representation. *Empirical Methods in Natural Language Processing (EMNLP)*, 1532–1543. <http://www.aclweb.org/anthology/D14-1162>
- Pu, Y., Gan, Z., Hénao, R., Yuan, X., Li, C., Stevens, A., & Carin, L. (2016). Variational autoencoder for deep learning of images, labels and captions. *Advances in neural information processing systems*, 29, 2352–2360.

- Rashtchian, C., Young, P., Hodosh, M., & Hockenmaier, J. (2010). Collecting image annotations using amazon’s mechanical turk. *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk*, 139–147.
- Sheng, X., Xu, L., Guo, J., Liu, J., Zhao, R., & Xu, Y. (2020). Introvnmmt: An introspective model for variational neural machine translation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05), 8830–8837.
- Vahdat, A., & Kautz, J. (2021). Nvae: A deep hierarchical variational autoencoder. *stat*, 1050, 8.
- Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3156–3164.
- Wang, Q., Li, B., Xiao, T., Zhu, J., Li, C., Wong, D. F., & Chao, L. S. (2019). Learning deep transformer models for machine translation. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 1810–1822.
- Wen, T.-H., Miao, Y., Blunsom, P., & Young, S. (2017). Latent intention dialogue models. *International Conference on Machine Learning*, 3732–3741.
- Yan, X., Yang, J., Sohn, K., & Lee, H. (2016). Attribute2image: Conditional image generation from visual attributes. *European Conference on Computer Vision*, 776–791.
- Yang, Z., Hu, Z., Salakhutdinov, R., & Berg-Kirkpatrick, T. (2017). Improved variational autoencoders for text modeling using dilated convolutions. *International conference on machine learning*, 3881–3890.
- Yang, Z., Yang, D., Dyer, C., He, X., Smola, A., & Hovy, E. (2016). Hierarchical attention networks for document classification. *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, 1480–1489.