

Vector Learning for Cross Domain Representations

Shagan Sah, Chi Zhang, Thang Nguyen, Dheeraj Kumar Peri, Ameya Shringi, Raymond Ptucha
Rochester Institute of Technology, Rochester, NY 14623, USA

Abstract—Recently, generative adversarial networks have gained a lot of popularity for image generation tasks. However, such models are associated with complex learning mechanisms and demand very large relevant datasets. This work borrows concepts from image and video captioning models to form an image generative framework. The model is trained in a similar fashion as recurrent captioning model and uses the learned weights for image generation. This is done in an inverse direction, where the input is a caption and the output is an image. The vector representation of the sentence and frames are extracted from an encoder-decoder model which is initially trained on similar sentence and image pairs. Our model conditions image generation on a natural language caption. We leverage a sequence-to-sequence model to generate synthetic captions that have the same meaning for having a robust image generation. One key advantage of our method is that the traditional image captioning datasets can be used for synthetic sentence paraphrases. Results indicate that images generated through multiple captions are better at capturing the semantic meaning of the family of captions.

I. INTRODUCTION

As few as five years ago, the automatic annotation of image and videos with natural language descriptions seemed quite distant. Recent discoveries in convolutional and recurrent neural networks have led to unprecedented vision and language understanding such that automatic captioning is now commonplace. These discoveries have fueled the growth of new capabilities such as improved description of visual stimulus for the blind, advanced image and video search, and video summarization.

An important objective for computer vision and natural language processing research is to be able to represent both modalities of data in a common latent vector representation. Concepts that are similar, lie close together in this space, while dissimilar concepts lie far apart. This allows, for example, a keyword search of “boat” to not only return images of boats, but similar words, similar sentences, videos, and audio clips. Much like the International Color Consortium’s device independent profile connection space for color management [1], [2], a source independent vector connection space requires each new modality to only define a single transformation into and out of this reference space, rather than define a transformation to and from all other modalities.

Our work borrows from recent advances in vector representations [3], [4], [5], generative models [6], image/video captioning [7], [8], [9] and machine translation [10] frameworks to form a multi-modal common vector connection space. We demonstrate the concept for both images and sentences,

and show the extension to other modalities such as words, paragraphs, video, and audio is possible.

The main contributions of this paper are: 1) Developing a robust source independent vector connection space between vision and language by conditioning image generation on multiple sentences; and 2) Integrating a language translation model for synthesizing artificial sentences with image generation.

The rest of this paper is organized as follows: Section II reviews relevant techniques. Section III presents the multiple caption conditioned image generation framework. Section IV discusses the experimental results. Concluding remarks are presented in Section V.

II. RELATED WORK

Image and video understanding has recently gained a lot of attention in deep learning research. Image classification [11], [12], [13], object detection [14], [15], semantic segmentation [16], [17], image captioning [7], [8], and localized image description [18] tasks have witnessed tremendous progress in the last few years.

Most machine learning algorithms require inputs to be represented by fixed-length feature vectors. This is a challenging task when the inputs are variable length sentences and paragraphs. Many studies have addressed this problem. For example, [3] presented a sentence vector representation while [4] created a paragraph vector representation. An application of such representations is shown by [19] that has used individual sentence embeddings from a paragraph to search for relevant video segments. An alternate approach uses an encoder-decoder [10] framework that encodes the inputs, one at a time to the RNN-based architecture. Such an approach is shown for video captioning tasks by S2VT [9] that encodes the entire video, then decodes one word at a time. [5] encodes sentences with common semantic information to similar vector representations. This latent representation of sentences has been shown useful for sentence paraphrasing and document summarization.

Deep convolutional networks have shown remarkable capability to generate images. Goodfellow *et al.* [20] proposed an easy and effective framework of generative models based on an adversarial process. This process involves two models: the generator that captures the data distribution to generate a fake image; and the discriminator that estimates the probability of a sample being fake or real. Radford [21] proposed a image generative network that was trained in an unsupervised manner and presented a set of guidelines for training the generative networks. Nguyen *et al.* [22] deals with the concept of activation maximization of a neuron and explored how each

layer captures varied information using encoder, generator and visualizing networks. Other recent works expand this idea into different types of multimedia. Reed *et al.*[23] explored the conversion between text to image space. Specifically, the text is encoded into a vector and consequently fed as input into the generator. Introducing an additional prior on the latent code, Plug and Play Generative Networks (PPGN) [6] drew a wide range of image types and a conditioner that tells the generator what to draw. Instead of conditioning on classes, the generated images were conditioned on text by attaching a recurrent, image-captioning network to the output layer of the generator, and performing similar iterative sampling. Our work borrows concepts from such captioning and generative models to form a common vector connection space.

The notion of a space where similar points are close to each other is a key principle of metric learning. The representations obtained from this formulation need to generalize well when the test data has unseen labels. Thus models based on metric learning have been used extensively in the domain of face verification [24], image retrieval [25], person-re-identification [26] and zero shot learning [27]. [28] used an auto-encoder model to learn cross modal representations and show results with audio and video datasets. Recently, [29] leveraged this concept to associate data from different modalities. Our work can be seen as an extension of this as we extend it to visual data while linking image and caption spaces to improve image generation.

III. METHODOLOGY

Inspired by Plug and Play Generative Networks [6], which reconstructs an image from a high-level feature space extracted from a pretrained encoder, we propose a similar architecture that conditions image generation on captions as shown in Figure 1. This model is comprised of three pretrained modules: the generator G , the CNN encoder and a image captioner model. Given an image x , a CNN encoder is used to extract the feature vector \hat{h} which is subsequently decoded by the captioner to generate a text description of the image. The reconstructed vector h lies in the vector connection space, which can be used to generate another image through the pretrained generator G . Representation h is updated based on the loss between a ground truth and generated sentence from the captioner. The loss from the captioner is back-propagated all the way back to h , forcing it to be representative of the semantics of the sentence. In an iterative fashion, the updated vector h is passed into the generative model and the process repeats. The parameters of the encoder CNN and the generator (G) are fixed and not updated during the iterative update process.

Equation (1) describes the update rule for h during inference. The γ_1 term represents the BLEU metric [30] loss associated with the generated and ground truth text. Specifically, we compute word level loss and scale it with the BLEU-1 score between the generated and ground truth text. This loss is back-propagated through the captioner and generator path to h . The γ_2 term is an image reconstruction loss

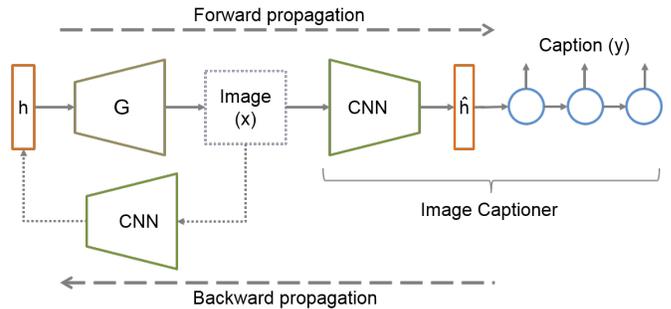


Fig. 1. Caption conditioned image generation model. The image generator, G and the CNN are pretrained and are fixed. During the forward pass, we start with a random vector h which generates an image (x) . The image is passed through the captioner to generate a caption. The generated caption is compared with the ground truth caption (y) and the error is back propagated to update h without updating the parameters in the CNN or G . After the update iterations, the image (x) is the resulting image generated using the final vector h .

calculated as a Euclidean distance between the pixel values of the generated and ground truth image. The γ_3 term is a vector h reconstruction loss calculated as a Euclidean distance between h and encoded h as shown with dotted line in Figure 1. The γ_4 term adds noise to the update rule. The update rule is generalizable to both *image-to-vector* and *sentence-to-vector* cases since the loss terms from captioner and the image reconstruction can be used in cases when ground truth text, images or both are available.

$$h_t = h_{t-1} + \gamma_1 \frac{\partial \mathcal{W}(C_{pred}, C_{gt})}{\partial h_{t-1}} + \gamma_2 \mathcal{L}(\hat{x}, x) + \gamma_3 \mathcal{L}(h, CNN(x)) + \mathcal{N}(0, \gamma_4^2) \quad (1)$$

Where, C_{pred} and C_{gt} are predicted and ground truth text and n is the length of the ground truth text, \mathcal{L} is Euclidean distance and \mathcal{W} is the word level caption loss, $\gamma_1 = BLEU(C_{pred}, C_{gt})/n$ is the scaling factor for the word loss. γ_2 , γ_3 and γ_4 are hyper-parameters.

To have a better generalization for caption conditioning, we use multiple copies of the ground truth caption. This is done by using a pretrained sentence-to-sentence model by inputting the reference caption and synthesizing paraphrase sentences. The generator is conditioned on this collection of synthesized sentences. We train a sequence-to-sequence model [10] for this task as described in the next section.

A. Synthesizing Artificial Sentences

Given a reference sentence, the objective is to produce a semantically related sentence as shown in Figure 2. Recent advances at vectorizing sentences represent exact sentences faithfully [31], [4], [32], or pair a current sentence with prior and next sentence [3]. Just like word2vec and GloVe map words of similar meaning close to one another, we desire a method to map sentences of similar meaning close to one another. For synthesizing sentences, we consider the sentence paraphrasing framework as an encoder-decoder model. Given a sentence, the encoder maps the sentence into a vector

(sent2vec) and this vector is fed into the decoder to produce a paraphrase sentence.

For the paraphrasing model, we represent the paraphrase sentence pairs as (S_m, S_n) . Let s_m denote the word embedding for sentence S_m ; and s_n denote the word embedding for sentence S_n . $S_m \in \{s_1 \dots s_M\}$, $S_n \in \{s_1 \dots s_N\}$ where M and N are the length of the paraphrase sentences. As shown in Figure 2, the input sentence y generates sentence y_1 , y_1 generates y_2 , and so on. In our model, we use an RNN encoder with LSTM cells. Specifically, the words in s_i are converted into token IDs and then embedded using GloVe [33]. To encode a sentence, the embedded words are iteratively processed by the LSTM cell [10].

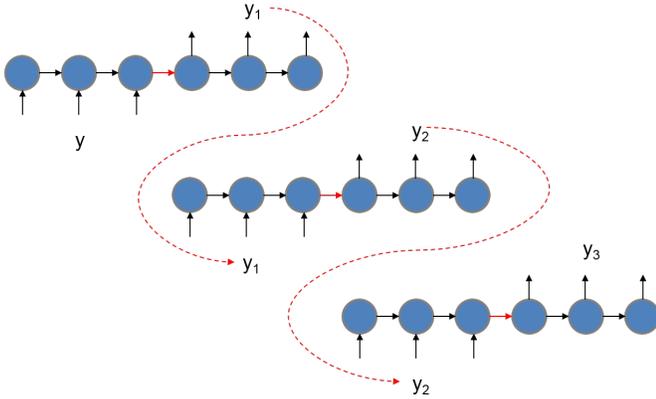


Fig. 2. Synthesizing multiple sentences using a sent2sent model. Given an input sentence y , a pre-trained sequence to sequence model trained on caption paraphrases is used to synthesize multiple copies of y with the same meaning. The input sentence y generates sentence y_1 , y_1 generates y_2 , and so on.

The paraphrasing model was trained on visual caption datasets. There are numerous datasets with multiple captions for images or videos. For example, MSR-VTT dataset [34] is comprised of 10,000 videos with 20 sentences each describing the videos. The 20 sentences are paraphrases since all the sentences are describing the same visual input. We form pairs of these sentences to create input-target samples. Likewise, MSVD [35], MS-COCO [36], and Flickr-30k [37] are used. Table I lists the statistics of datasets used.

TABLE I
SENTENCE PAIRS STATISTICS IN CAPTIONING DATASETS.

	MSVD	MSRVTT	MS-COCO	Flickr
#sent	80K	200K	123K	158K
#sent/samp.	~42	20	5	5
# sent pairs	3.2 M	3.8 M	2.4 M	600 K

IV. RESULTS AND DISCUSSION

A. MS-COCO sentences

Each image in the MS-COCO dataset has five human generated captions. We test the effect of using > 1 captions on image generation by randomly selecting three captions among the five captions from the dataset. Punctuations are removed



- A long fully baked pizza sitting on a table
- A long fully baked pizza sitting on a table
- A large pizza sitting on top of a cutting board
- A wood fired pizza containing mushrooms and pepperoni

Fig. 3. Image generation example conditioned on multiple captions. Top row are images generated and the corresponding input caption and bottom row are images and the corresponding three input captions.

and the sentences are lemmatized in order to eliminate out of vocabulary words. Word gradients are averaged across all three sentences and back-propagated back to the latent vector. This combined update was observed to have stronger impact on the image quality than using a single sentence. The results are shown in Figure 3. Using multiple sentences, the generator was able to generate an image with more details as shown in the pizza image example. One explanation is that multiple captions have more objects and attributes which the generator could attend to, hence improving the resulting image quality. Figures 4 and 5 show sample results of using three and five sentences respectively. With respect to Figure 5, we note the generator was trained on the ImageNet dataset. the ImageNet dataset does not have a “person” category. This makes it highly challenging for generating images with such categories as also noted in [6]. We observed that while conditioning on multiple captions, the generator was able to generate relevant images with “person” category. It can be seen from the example that the fine details were not clearly visible but the overall structure captures the semantic meaning of the captions.



- The kitchen and breakfast area of a modern house
- A kitchen filled with wooden floors and a stove top oven.
- A dining room and kitchen of a home is on display.

Fig. 4. Image generation example conditioned on three captions from the MS-COCO image captioning dataset.



- Some people sitting on a bench
- Four young people are sitting on a bench
- Some teens are sitting on a bench together
- Four people sit on a bench next to each other
- All of the people are sitting together on a bench

Fig. 5. Image generation example conditioned on five captions from the MS-COCO image captioning dataset.

B. Synthetic Sentences

The sentences used in the above experiment were sampled from the test set of MS-COCO and the image captioning model used in the experiment was trained on MS-COCO. In order to further validate our hypothesis, we generated synthetic sentences using a sequence-to-sequence model which produces paraphrases of an input sentence. This encoder-decoder model described in Figure 2 generates multiple paraphrase copies of an input sentence. These sentences were passed into the model for generating an image and an example is shown in Figure 6.

C. Image-to-Image Analysis

The essence of the image-to-image model can be amplified by generating a similar image given any input image. Given an input image, we pass it to an image captioning model. Multiple copies of the captions can be obtained either by using beam search in the captioner or by passing the caption through a sequence-to-sequence paraphrasing model as described earlier. Since this collection of captions capture the semantic meaning of the original image, conditioning on the captions generates images which have similar information as the original image. Different captioning models can be plugged in to compare the effect of varying generator and captioner architectures. The results are shown in Fig 6. The generator was able to attend to the white color mentioned in the captions and the overall bus structure in the original bus image.



- a table with chairs and a table on it
- a living room with couches and a table
- a dining room with a table and chairs
- a table with a table and a table
- a large white room with a table and a table
- a bus is driving down a street next to a building
- two buses are parked in a parking lot
- two buses are lined up on the street
- two buses are parked in a parking lot
- a white bus is driving down a street



Fig. 6. Image-to-image transition example. Top image is input image, center row sentences are five captions generated using the input image and bottom row images are generated images. The image generation is conditioned on multiple synthetic captions.

D. Generalization

Although back-propagating multiple word gradients did improve the search for a good quality image in the latent vector space, there were examples where we found images not representative of the captions. The effect is shown in Figure 7. Therefore, the generalization of this hypothesis is still under exploration and needs further investigation.



- Two sandwiches sitting on top of papers on a table
- Two sausages sit in buns on a counter
- Two hot dogs with mustard and sauerkraut
- Two hot dogs on gourmet bread with sauerkraut
- Two large hot dogs with sour kraut and mustard
- There are many zebras in the wild together
- Three zebras standing in a copse of bare trees
- A zebra eating in the middle of the dessert
- Several zebras can be seen milling within the tall trees
- Three zebras standing among some thicket and trees

Fig. 7. Some examples for which multiple captions do not help in improving the image quality.

V. CONCLUSION

This work advances the area of caption conditioned image generation by allowing image generation being conditioned on multiple sentences. The resulting vector also helps achieve a common vector space to be shared between vision and language representations. We conclude that iteratively sampling over multiple sentences indeed helps in improving the quality of the generated images. The model shows the robustness in performing cross-modal captioning. An extension of this work would be observing each layer of the generator during back-propagation. This would give more insight into the shapes and colors that are drawn by specific layers of the generator. Moreover, evaluation techniques of the generated images is still in nascent stage despite the area of image generative models having seen significant progress recently.

ACKNOWLEDGMENT

The authors would like to thank NVIDIA for some of the GPUs used in this work.

REFERENCES

- [1] I. C. Consortium *et al.*, “Specification icc. 1: 2010 (profile version 4.3.0.0) image technology colour management–architecture, profile format, and data structure,” 2010.
- [2] G. B. Pawle and L. Borg, “Evolution of the icc profile connection space,” in *9th Congress of the International Colour Association*, vol. 4421. International Society for Optics and Photonics, 2002, pp. 446–451.
- [3] R. Kiros *et al.*, “Skip-thought vectors,” in *NIPS*, 2015, pp. 3294–3302.
- [4] Q. V. Le and T. Mikolov, “Distributed representations of sentences and documents,” in *ICML*, vol. 14, 2014, pp. 1188–1196.
- [5] C. Zhang *et al.*, “Semantic sentence embeddings for paraphrasing and text summarization,” in *GlobalSIP*, 2017.
- [6] A. Nguyen, J. Yosinski, Y. Bengio, A. Dosovitskiy, and J. Clune, “Plug & play generative networks: Conditional iterative generation of images in latent space,” *arXiv preprint arXiv:1612.00005*, 2016.
- [7] J. D. *et al.*, “Long-term recurrent convolutional networks for visual recognition and description,” in *CVPR*, 2015, pp. 2625–2634.
- [8] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” *arXiv preprint arXiv:1502.03044*, vol. 2, no. 3, p. 5, 2015.

- [9] S. Venugopalan *et al.*, “Sequence to sequence-video to text,” in *ICCV*, 2015, pp. 4534–4542.
- [10] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *Advances in neural information processing systems*, 2014, pp. 3104–3112.
- [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [12] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [13] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *arXiv preprint arXiv:1512.03385*, 2015.
- [14] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” in *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [15] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” *arXiv preprint arXiv:1506.02640*, 2015.
- [16] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [17] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,” *arXiv preprint arXiv:1606.00915*, 2016.
- [18] J. Johnson, A. Karpathy, and L. Fei-Fei, “Densecap: Fully convolutional localization networks for dense captioning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [19] J. Choi *et al.*, “Textually customized video summaries,” *arXiv preprint arXiv:1702.01528*, 2017.
- [20] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [21] A. Radford, L. Metz, and S. Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” *arXiv preprint arXiv:1511.06434*, 2015.
- [22] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune, “Synthesizing the preferred inputs for neurons in neural networks via deep generator networks,” in *Advances in Neural Information Processing Systems*, 2016, pp. 3387–3395.
- [23] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, “Generative adversarial text to image synthesis,” *arXiv preprint arXiv:1605.05396*, 2016.
- [24] F. Schroff, D. Kalenichenko, and J. Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 815–823.
- [25] A. Gordo, J. Almazán, J. Revaud, and D. Larlus, “Deep image retrieval: Learning global representations for image search,” in *European Conference on Computer Vision*. Springer, 2016, pp. 241–257.
- [26] A. Hermans, L. Beyer, and B. Leibe, “In defense of the triplet loss for person re-identification,” *arXiv preprint arXiv:1703.07737*, 2017.
- [27] R. Socher, M. Ganjoo, C. D. Manning, and A. Ng, “Zero-shot learning through cross-modal transfer,” in *Advances in neural information processing systems*, 2013, pp. 935–943.
- [28] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, “Multimodal deep learning,” in *Proceedings of the 28th international conference on machine learning (ICML-11)*, 2011, pp. 689–696.
- [29] L. Wu, A. Fisch, S. Chopra, K. Adams, A. Bordes, and J. Weston, “Starspace: Embed all the things!” *arXiv preprint arXiv:1709.03856*, 2017.
- [30] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics, 2002, pp. 311–318.
- [31] N. Kalchbrenner, E. Grefenstette, and P. Blunsom, “A convolutional neural network for modelling sentences,” *arXiv preprint arXiv:1404.2188*, 2014.
- [32] H. Zhao, Z. Lu, and P. Poupart, “Self-adaptive hierarchical sentence model,” *arXiv preprint arXiv:1504.05070*, 2015.
- [33] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *EMNLP*, 2014, pp. 1532–1543. [Online]. Available: <http://www.aclweb.org/anthology/D14-1162>
- [34] J. Xu *et al.*, “Msr-vtt: A large video description dataset for bridging video and language,” in *CVPR*, 2016.
- [35] D. L. Chen and W. B. Dolan, “Collecting highly parallel data for paraphrase evaluation,” in *Proceedings of the 49th Association for Computational Linguistics: Human Language Technologies-Volume 1*. ACL, 2011, pp. 190–200.
- [36] T.-Y. Lin *et al.*, “Microsoft coco: Common objects in context,” in *European Conference on Computer Vision*. Springer, 2014, pp. 740–755.
- [37] P. Young *et al.*, “From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions,” *Transactions of the Association for Computational Linguistics*, 2014.