

Scene Graph Generation from Images

Satoshi Tsutsui
Manish Kumar

1 Introduction

Image understanding by computer is advancing exponentially these days due to the phenomenal success of deep learning, but there is still much work left for the computers to reach human level perception. Image classification (sometimes with localization) is one of the standard task, but this is far from the image understanding. The other tasks such as image caption generation or visual question answering have also reached to practical level of quality, but these are still far from the complete image understanding. Image caption generation, which is a task that generates a summary sentence from an image, cannot fully describe rich scenery in an image. Visual question answering also answers simple questions, but cannot answer questions that requires complex reasoning. In order to fully understand an image, we need to know; what objects are in the image, what are the characteristics of objects, and how these objects interact to each other.

In this project, we try two approaches to answer these questions. In the first approach, we use automatically generated caption as intermediate structure. In the second approach, we assume that objects with regions are given (or already detected by previous work), and mainly focus on two questions: what are the characteristics (attributes) of objects, and what are the relations between these objects. For example, in figure 1, four objects (man, women, bench, and river) are shown. A relation can be identified as something like "man sits on bench". Similarly, an attribute can be like "bench is wooden." We develop a method to classify the attribute given an object, and a method to classify the relation given two objects.

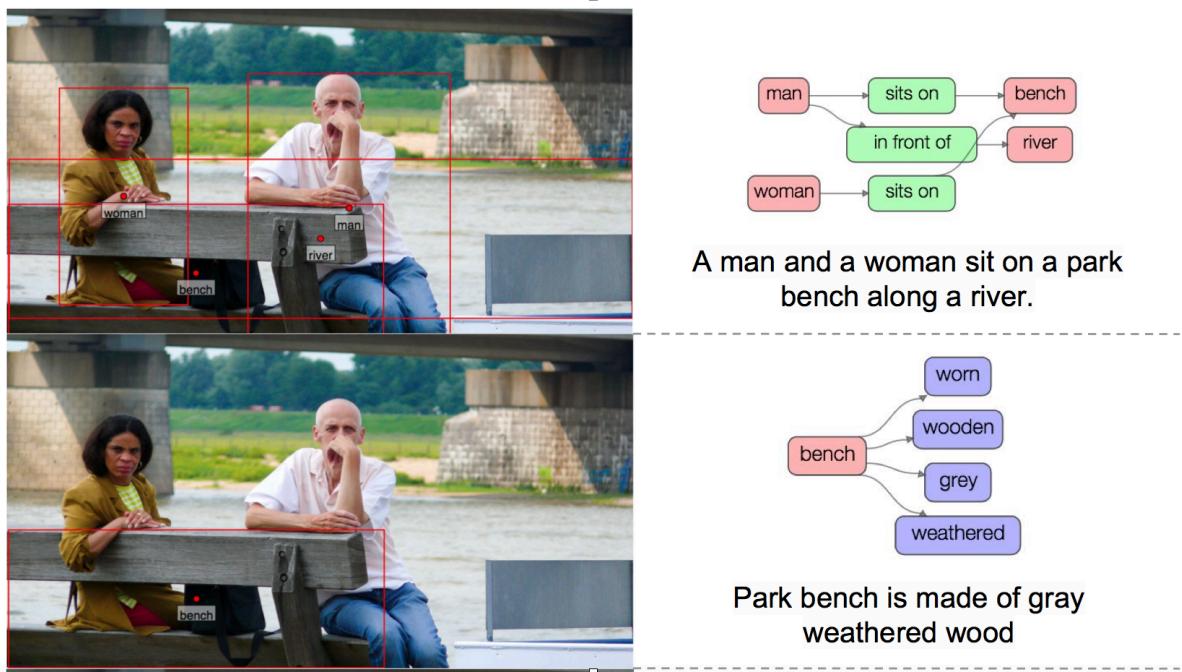


Figure 1 Annotation example. figure source: Krishna et al. (2016)

Our initial goal was to build a system to generate a *scene graph* from an arbitrary daily life image though we did not reach this level due to unexpected amount of time spent for environmental setup. Scene graph can be derived if we integrate objects, attributes of objects, and relations between objects. Part of scene graph is shown at the right side of Figure 1. The whole scene graph for the image in figure 1 is illustrated on figure 2. Scene graph is structured and more machine understandable than sentence descriptions or captions. It is actually shown to be more effective for image retrieval (Johnson, Krishna, et al., 2015).

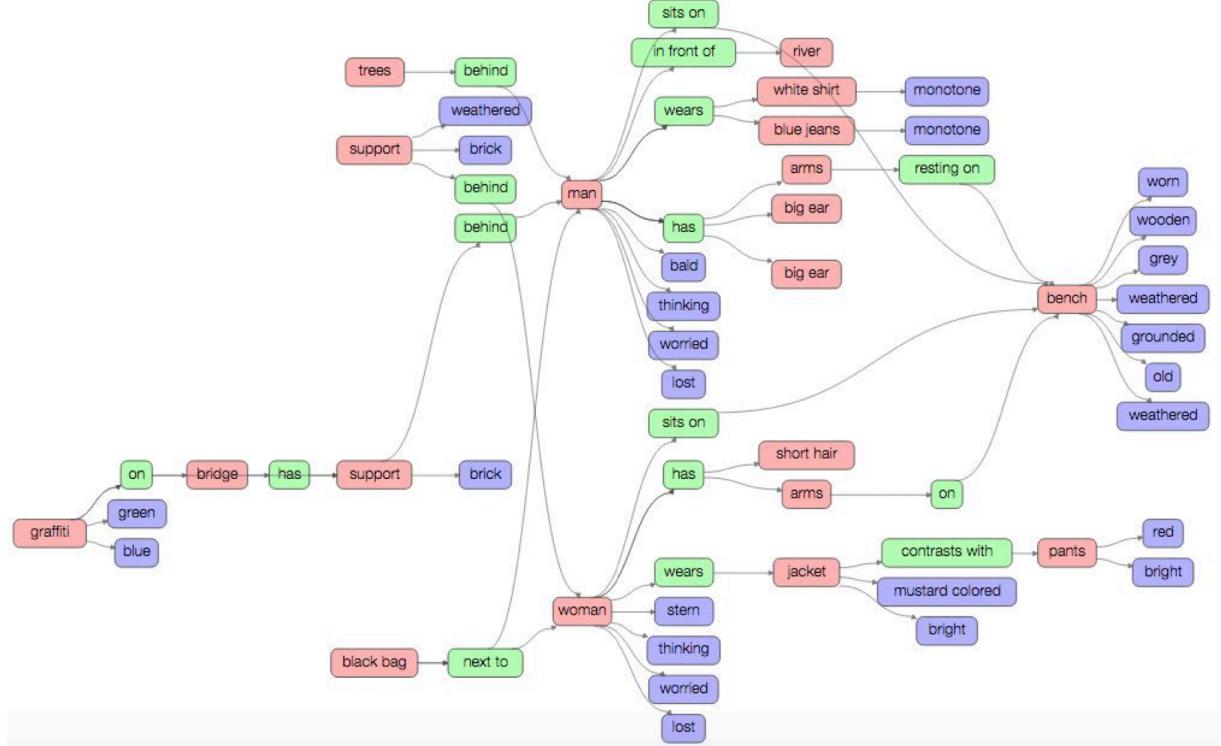


Figure 2 Scene graph example. figure source: Krishna et al. (2016)

2 Background and Related Work

2.1 Scene Graph

Scene graph is a graph structured image representation such as graph shown figure 2. This representation is easier to understand for computers than sentence descriptions. For example, it is shown that scene graph helps computer to retrieve images more efficiently than sentence description (Johnson, Krishna, et al., 2015). Aiming for better image retrieval, (Schuster, Krishna, Chang, Fei-Fei, & Manning, 2015) introduces sense graph generator from image descriptions. However, to our knowledge, no previous work tries to generate scene graph directly from images.

As a scene graph dataset, we used *Visual Genome*. Visual Genome¹ (Krishna et al., 2016) is the latest dataset aiming to accelerate research on image understanding. It has 100,000 images with regional descriptors, objects, attributes, relations, and question answer pairs. This is currently the largest and richest dataset available for image understanding. Actually, full dataset is note yet published, but we personally contacted the first author of the Visual Genome paper, and he was kind enough to give us preliminary dataset.

2.2 Caption Generation from Images

Automatic image captioning is one of the way to understand images, though caption is not machine friendly as scene graph, and one summary sentence is not enough for image understanding. Before the flood of deep learning, caption generation is modeled as template filling approach (Kulkarni et al., 2011) or sentence retrieval problem (Kuznetsova, Ordonez, Berg, Berg, & Choi, 2012). Now, the state of the art approaches uses deep learning to generate caption. They basically use pre-trained convolutional neural network (CNN) to extract image feature, and recurrent neural network to process sentences. (Kiros, Salakhutdinov, & Zemel, 2014) learns multimodal embedding space over images and sentences, and then generates captions from the embedded space. “Show and Tell” (Vinyals et al., 2014) generates captions directly with RNN from CNN image feature. (Xu et al., 2015) incorporates an attention model in that each word in a caption has regional attention over the original image. (Karpathy & Fei-Fei, 2015) and (Johnson, Karpathy, & Fei-Fei, 2016) generate captions with localized regions.

¹ <https://visualgenome.org>

Note that caption generation is often viewed within the framework of neural network machine translation (Sutskever, Vinyals, & Le, 2014) (Cho et al., 2014). Technically, caption generation with neural network can be regarded as translation from an image to a sentence.

3 Method Overview

Our final goal is to generate end to end scene graph generation system, so we propose two approaches for scene graph generation. First one is to generate scene graph via captions. Second one is to extract attributes and relations after object detection. Both approach is illustrated in figure 3. However, we did not complete to generate scene graph for the second approach, so we only report attribute extraction and relation extraction for the second one.

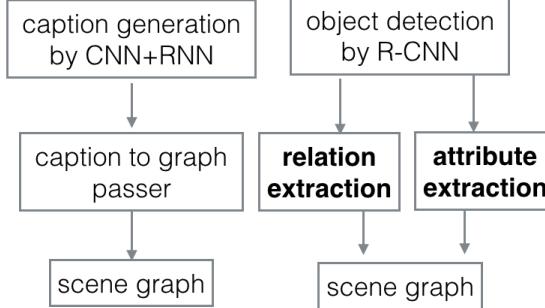


Figure 3 Two approaches. Left: generate caption, and parse caption to scene graph. Right: Detect object, infer attributes and relations, and connect them to scene graph. However, in this paper, we only report relation and attribute classification (bolded part).

3.1 Scene Graph via Caption

This approach generates a caption for an image first and then parses the generated sentence into a scene graph. It is originally intended to serve as a baseline implementation against the other approach. However, we did not complete actual graph generation for the other approach, so we did not compare. We use previous work for caption generations (Vinyals, Toshev, Bengio, & Erhan, 2015) and caption to scene graph (Schuster et al., 2015). We used publicly available codes for both.

3.2 Relation and Attribute Extraction after Object Detection

This approach was originally aimed to generate scene graph, but we did not complete the generation and evaluation. Here, we only report the relation and attribute extraction assuming that objects in an image are given with rectangular regions. For relations, we also assume that we know the objects that will be connected by a relation. In other words, we only consider pairs of objects that have at least one relation, and do not consider whether they have relations or not. Using these assumptions, we cast our problems into classification problems. We use a pre-trained convolutional neural network, specifically GoogleNet (Szegedy et al., 2015), to extract whole image features and localized image features. Then, we concatenate these feature vectors (whole image feature and an object feature for attributes, and whole image feature and two object features) into a vector, add two layers with ReLU activations, and finally a softmax activation to get probabilistic distribution over possible attributes/relations. These neural networks are illustrated in Figure 4 and 5.

We only train the weights after CNNs so we extracted all the CNN features (whole image and localized objects) in advance. This saves much time to run experiments. We tried several numbers of hidden units from 128 to 1024 for the intermediate layers. We also tried batch normalization (Ioffe & Szegedy, 2015) and dropout (Srivastava, Hinton, Krizhevsky, Sutskever, & Salakhutdinov, 2014). We initialized all the weights using uniform distribution of (-0.1, 0.1). We used Adam (Kingma & Ba, 2014) to train the networks.

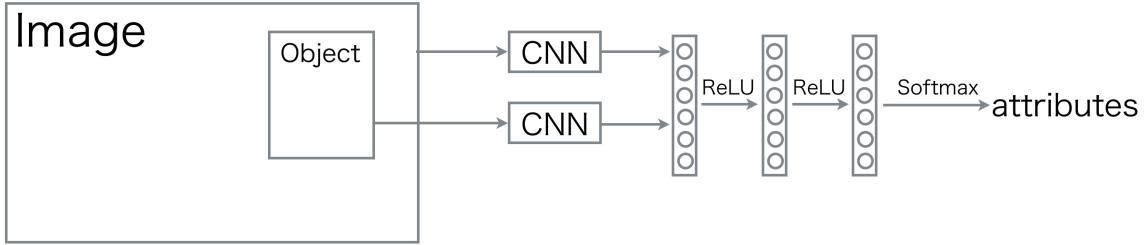


Figure 4 Neural Network Architecture for Attribute Classification

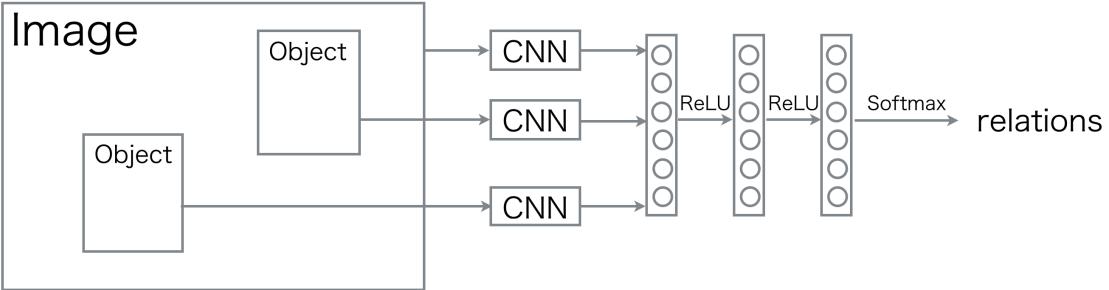


Figure 5 Neural Network Architecture for Relation Classification

3.3 Evaluation Procedure

We divided the dataset into three sets: train, validation, and test in the following way. Vocabulary size is the number of distinct relations/attributes whose frequency is over five.

Table 1 Dataset Split

	# of images			vocabulary size over five times appearance (original size)
	train	val	test	
relations	42728	4000	4000	4265 (30263)
attributes	38918	4000	4000	2834 (15912)

For the whole graph evaluation (Section 3.1), we evaluate based on elements of graph for each attributes and relations (Object-Attribute or Object-Relation-Object). We calculate Jaccard index ($|A \cap B| / |A \cup B|$) to evaluate the accuracy.

For the relations and attributes classification (Section 3.2), we evaluate just by accuracy.

4 Details for Implementation and Reproduce

In this section, we explain details including instruction on how to use the codes on IU git. First of all, in order to use codes on IU git, you need to set environment up by installing several dependencies. We use Python 2.7.x and Java 1.8. Also, in order to use GPU, you will be ready to make CUDA and CuDNN ready. As for the python environment, we suggest you to install Anaconda as the easiest way.

4.1 Commands to make python environment ready:

```

wget https://3230d63b5fc54e62148e-c95ac804525aac4b6dba79b00b39d1d3.ssl.cf1.rackcdn.com/Anaconda2-2.4.1-Linux-x86_64.sh
bash Anaconda2-2.4.1-Linux-x86_64.sh -b
echo 'export PATH=$HOME/anaconda/bin:$PATH' >> .bashrc
echo 'export PYTHONPATH=$HOME/anaconda/lib/python2.7/site-packages:$PYTHONPATH' >> .bashrc
source .bashrc
conda update conda -y
# install chainer

```

```
pip install chainer
```

*commands taken from: https://github.com/apple2373/chainer_caption_generation

4.2 Commands to make GPU environment ready

We do not cover how to install CUDA and CuDNN, but after you install them, you should configure the following path.

```
export PATH=/usr/local/cuda/bin:$PATH
export LD_LIBRARY_PATH=/usr/local/cuda/lib64:$LD_LIBRARY_PATH
export CPATH=/path/to/cudnn/include:$CPATH
export LIBRARY_PATH=/path/to/cudnn/lib:$LIBRARY_PATH
export LD_LIBRARY_PATH=/path/to/cudnn/lib:$LD_LIBRARY_PATH
```

*commands taken from <https://github.com/pfnet/chainer>

4.3 Commands to clone the IU git repository and download the related files

After cloning the main repository, you need to clone the nested repository. This repository is created by the first author of this paper. Also, you need to download several files.

```
git submodule update --init --recursive
cd chainer_caption_generation
bash download.sh
cd ..
bash data_download.sh
bash work_download.sh
cd graph_by_caption
bash scene_graph_parser_download.sh
```

4.4 Scene Graph via Caption

For captions, we use a publicly available implementation for caption generation. (https://github.com/apple2373/chainer_caption_generation). This is created by the first author of this paper to reproduce the work on caption generation (Vinyals et al., 2014). It uses MSCOC dataset to train, and it achieves CIDEr (Vedantam, Lawrence Zitnick, & Parikh, 2015) of 0.66 as an evaluation score. This repository is cloned inside the assignment repository.

For caption to scene graph, we use a publicly available caption to scene graph parser. (<http://nlp.stanford.edu/software/scenegraph-parser.shtml>). This is written in Java, so we wrote a program to call the java program from python using py4j (<https://www.py4j.org>). Be aware that, in order to generate a scene graph via caption, you need to set up java program in advance by the flowing command:

```
java -jar graph_by_caption/SceneGraphParserPython.jar
```

The sample codes to generate scene graph via caption is documented in ipython notebook in `/graph_by_caption/smaple.ipynb`

The code to generate scene graph for test images is `/codes/baseline_build.py` and its evaluation code is `/codes/baseline_eval.py`

4.5 Relation and Attribute Extraction after Object Detection

We use chainer (Tokui, Oono, Hido, & Clayton, 2015) to train neural networks. We only train the weights after CNNs, so we extracted all the CNN features (whole image and localized objects) in advance. This saves much time to run experiments. The codes to extract features are the followings:

```
/codes/pre_extract_googlenet_features_attr.py
/codes/pre_extract_googlenet_features_relation.py
/codes/pre_extract_googlenet_features_whole.py
```

We actually tried several neural network architectures:

```
/codes/train_attribute_model1.py
/codes/train_attribute_model2.py
/codes/train_relation_model1.py
/codes/train_relation_model2.py
```

```
/codes/train_relation_model3.py  
/codes/train_relation_model4.py  
/codes/train_relation_model5.py  
/codes/train_relation_model6.py
```

However, we only use the results on the followings:

```
/codes/train_attribute_model2.py  
/codes/train_relation_model1.py
```

To train the network, you can just execute:

```
cd codes  
python <file name>
```

This will generate a directory with log and models.

In order to evaluate performance, you need to execute corresponding evaluation file:

```
/codes/eval_attribute_model1.py  
/codes/eval_attribute_model2.py  
/codes/eval_relation_model1.py  
/codes/eval_relation_model2.py  
/codes/eval_relation_model3.py  
/codes/eval_relation_model4.py  
/codes/eval_relation_model5.py  
/codes/eval_relation_model6.py
```

The evaluated results are plotted on: */codes/Plot_restuls.ipynb*

A sample image is on */codes/Sample.ipynb*

5 Results

5.1 Samples

In order to effectively show what we have done, we use an example image to show the whole work flow. The example is the picture of IU campus shown below.



Figure 6 An example picture taken at IU. Image source: IU website: <http://www.indiana.edu/~stat/>

First, we generate caption for the image. The beam search (whose beam size is 3) generates the three sentences with probabilities:

- a group of people riding bikes down a street 0.00274640011112
- a group of people riding on the backs of horses 0.00139282029438
- a man riding a bike down a street 0.000736576015526

Taking the highest probability, the final caption will be “a group of people riding bikes down a street”. This will be parsed into the scene graph (* the number means the position of words in the sentence):

`[[u'group-2', u'of', u'people-4'], [u'people-4', u'ride', u'bike-6']]`

Next, we use the other approach using faster R-CNN. We do not claim this is our contribution, but we need to detect object in order to apply for new images other than Visualgenome. Hence for this sample, we used publicly available implementation (<https://github.com/rbgirshick/py-faster-rcnn>). We used a model trained on Pascal VOC dataset (not Image Net dataset), which means we can only detect 20 classes: person, bird, cat, cow, dog, horse, sheep, aeroplane, bicycle, boat, bus, car, motorbike, train, bottle, chair, dining table, potted plant, sofa, and tv/monitor. The results of R-CNN are shown below as red rectangular regions. We set the threshold as 0.7.

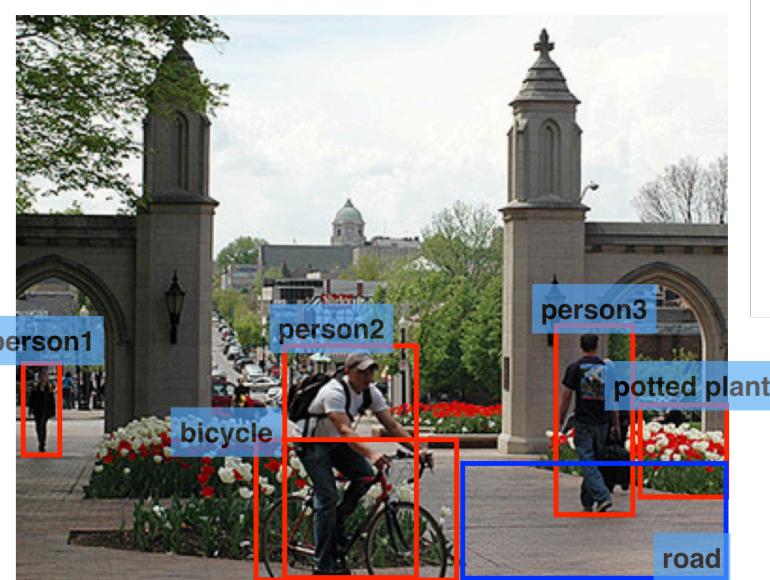


Figure 7 Objects in the sample picture. Red annotations are detected by faster R-CNN, blue one (road) is detected manually.

After we know the objects, we can use our trained model to infer relations and attributes. Four samples are shown below. We can confirm that relations are accurate, and attribute for potted plant is not accurate, but attribute for bicycle is not accurate.

Table 2 Top 5 softmax outputs for relation from person2 to bicycle.

1	0.224	on
2	0.117	riding
3	0.056	is on a
4	0.032	sitting on
5	0.020	on a

Table 3 Top 5 softmax outputs for relation from person3 to road.

1	0.630	on
2	0.071	on the
3	0.054	in

4	0.036	is on
5	0.031	walking on

Table 4 Top 5 softmax outputs for attributes of bicycle

1	0.006	white
2	0.005	black
3	0.004	blue
4	0.003	visible
5	0.003	green

Table 5 Top 5 softmax outputs of attributes of potted plant

1	0.006	white
2	0.005	black
3	0.004	blue
4	0.003	visible
5	0.003	green

In order to show an example of relation that is not accurate, we introduce another image from test set. In figure 8, a relation *trees - besides - house* is shown. Objects are give in the dataset. The relation prediction is shown in table 6, which is not accurate prediction.



Figure 8 Another example for relation: trees - besides – house

Table 6 Top 5 softmax outputs of relation from trees to house.

1	0.093	has
2	0.085	in
3	0.064	near
4	0.063	on

5	0.036	above
---	-------	-------

5.2 Scene Graph via Caption Results

The generated scene graph had zero Jaccard distance on test dataset on both relations and attributes. We infer two reasons for this. First, both graph parser and dataset are not normalized. “man” and “person” will be treated different. “on”, “on a”, and “is on a” will also be treated as different one, as shown in Table 2. Second, one summarizing sentence is not enough to capture the whole objects, attributes, and relations in an image. For example, in the sample picture, the generated caption “a group of people riding bikes down a street” does not capture the plants, flowers or gates. Therefore, we experimentally confirm that one summary sentence is never enough for full image understanding.

5.3 Relation and Attribute Classification Results

Table 7 Accuracies on relation and attribute classification

	train	val	test
relations	0.326	0.320	0.317
attributes	0.325	0.310	0.314

The accuracies on each dataset is shown in table 7. Overall, the accuracy is around 30%. We report the results on the network whose hidden units are 128, but the other networks also hold at least accuracy of 30%. Remember that the dataset is not normalized so predicting “on” where the dataset annotates “is on” will be counted as incorrect. However, even considering that aspect, we claim that these accuracies are never enough for practical applications such as autonomous driving. It is difficult to explain why the prediction worked or not because it uses neural networks, so we show some results with fail and success in the previous section.

Here we plot the accuracy on training data and validation data. These plots are used to determine the best model over the epochs. For relations, to our surprise, the model after the first epoch (just passing over the whole dataset once) is enough, and further training ends up with overfitting. For attributes, we use the model at epoch 2, which is the highest validation accuracy.

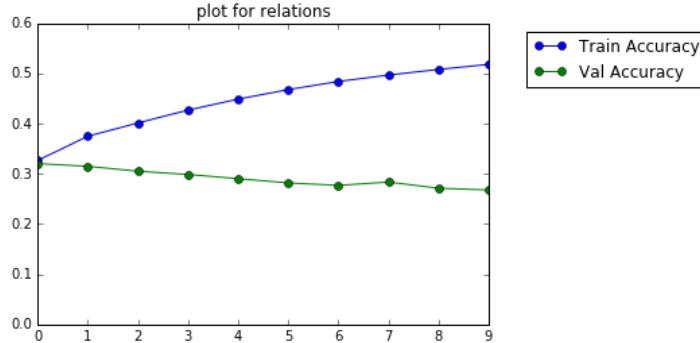


Figure 9 Accuracy plot for training for relations

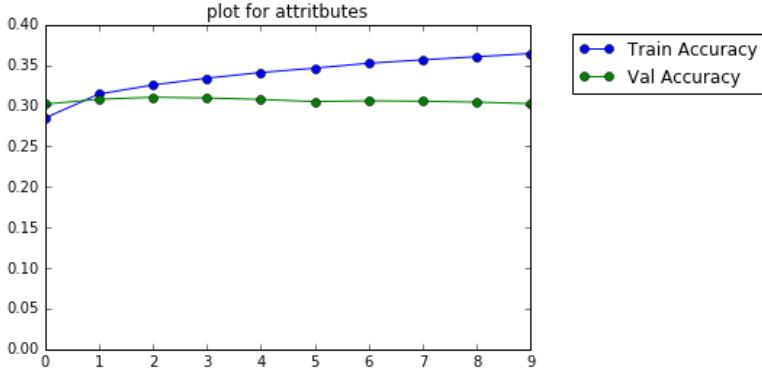


Figure 10 accuracy plot for training for attributes

6 Limitations and Future Work

6.1 Scene Graph via Caption

We confirmed that one summarizing sentence is not enough for image understanding. However, a recent work on image captioning generates multiple captions with localizations (Johnson, Karpathy, & Fei-Fei, 2015). This could be a future direction for generating scene graph via captions. To see the potential, we tried the recently open sourced implementation (<https://github.com/cjohnson/densecap>) to our sample image. The result is shown in Figure 11. This is more informative than one summary sentence. However, in order to construct a scene graph, it is challenging to disambiguate and identify objects in multiple captions. For example, there “man” on the right, middle and left are three distinct objects while all the “street” or “sidewalk” are the single object.

6.2 Relation and Attribute Classification

Not only the low accuracies, we have limitations for this approach. Our experiment assumes object pairs that have relations are given. This cause a problem when we finally want to build to scene graph. We do not know which two pairs of objects have relations or not in advance. For example, in figure 7, our system does not know that person3 and bicycle has no relations, so if we ask the relation, the system outputs some distribution that is shown in table 8. In order to solve this issue, we can think two approaches. First one is just add the label “no relations” to the training data by randomly selecting negative pairs of objects. Another approach is to train a binary classifier to judge whether two objects in an image has a relation or not. However, in both approaches, it is challenging to consider spatial information by CNN. For example, in figure 7, It might be challenging to distinguish features from person1 and person3 because CNN features are robust to spatial transformation.

Table 8 Top 5 softmax outputs for relation from person3 to bicycle.

1	0.169	on
2	0.047	riding
3	0.021	is on a
4	0.020	with
5	0.018	is on

Lastly, in the long run, we want to train end-to-end neural network to generate scene graph from an image without using R-CNN or relation/attribute classifier.



a man riding a bicycle. a man walking in the street. man walking with a backpack. a man walking on the sidewalk. flowers on the sidewalk. red flowers in a planter. tower on top of building. people riding a motorcycle. the man is wearing blue jeans. a blue sky with clouds. white clouds in the sky.

Figure 11 Dense caption for our sample image

7 Conclusion

We tackled the very fundamental yet difficult problem in computer vision: image understanding. Specifically, we aimed to generate scene graph, which is the structured and machine understandable representation of an image. We tried two possible approaches: scene graph by caption, and relation/attribute classification after detecting objects. For the first approach, we confirmed that one summary sentence is no sufficient to generate a scene graph, which has far more information than a summary sentence. For the second approach, we could not reach to scene graph generation due to environmental set up issues, but we are still able to learn that classifying attributes and relations based on CNN-features does not work very well. In short, we confirmed that much work is remained to generate scene graph of practical quality.

- Cho, K., van Merriënboer, B., Gürçehre, Ç., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. In *EMNLP*.
- Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv Preprint arXiv:1502.03167*.
- Johnson, J., Karpathy, A., & Fei-Fei, L. (2015). DenseCap: Fully Convolutional Localization Networks for Dense Captioning. *arXiv Preprint*. Retrieved from <http://arxiv.org/abs/1511.07571>
- Johnson, J., Karpathy, A., & Fei-Fei, L. (2016). DenseCap: Fully Convolutional Localization Networks for Dense

- Captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Johnson, J., Krishna, R., Stark, M., Li, L.-J., Shamma, D. A., Bernstein, M. S., & Fei-Fei, L. (2015). Image retrieval using scene graphs. In *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on* (pp. 3668–3678).
- Karpathy, A., & Fei-Fei, L. (2015). Deep Visual-Semantic Alignments for Generating Image Descriptions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Kingma, D. P., & Ba, J. (2014). Adam: A Method for Stochastic Optimization. *CoRR, abs/1412.6980*.
- Kiros, R., Salakhutdinov, R., & Zemel, R. S. (2014). Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models. *CoRR, abs/1411.2*. Retrieved from <http://arxiv.org/abs/1411.2539>
- Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., ... Fei-Fei, L. (2016). Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. Retrieved from <http://arxiv.org/abs/1602.07332>
- Kulkarni, G., Premraj, V., Dhar, S., Li, S., Choi, Y., Berg, A. C., & Berg, T. L. (2011). Baby talk: Understanding and generating simple image descriptions. In *CVPR*.
- Kuznetsova, P., Ordonez, V., Berg, A. C., Berg, T. L., & Choi, Y. (2012). Collective generation of natural image descriptions. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1* (pp. 359–368).
- Schuster, S., Krishna, R., Chang, A., Fei-Fei, L., & Manning, C. D. (2015). Generating semantically precise scene graphs from textual descriptions for improved image retrieval. In *Proceedings of the Fourth Workshop on Vision and Language* (pp. 70–80).
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research, 15*(1), 1929–1958.
- Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in Neural Information Processing Systems (NIPS)*, 3104–3112. Retrieved from <http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural>
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 1–9).
- Tokui, S., Oono, K., Hido, S., & Clayton, J. (2015). Chainer: a Next-Generation Open Source Framework for Deep Learning. In *Proceedings of Workshop on Machine Learning Systems (LearningSys) in The Twenty-ninth Annual Conference on Neural Information Processing Systems (NIPS)*. Retrieved from http://learningsys.org/papers/LearningSys_2015_paper_33.pdf
- Vedantam, R., Lawrence Zitnick, C., & Parikh, D. (2015). Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 4566–4575).
- Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3156–3164).
- Vinyals, O., Toshev, A., Bengio, S., Erhan, D., Google, O. V., Google, A. T., ... Google, D. E. (2014). Show and Tell: A Neural Image Caption Generator. Retrieved from <http://arxiv.org/abs/1411.4555>
- Xu, K., Ba, J., Kiros, R., Courville, A., Salakhutdinov, R., Zemel, R., & Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. *arXiv Preprint arXiv:1502.03044*.