

Image Generation from Scene Graphs

Tae Gil Ha

Seoul National University

Main Reference : JOHNSON, Justin; GUPTA, Agrim; FEI-FEI, Li. **Image generation from scene graphs**. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018. p. 1219-1228.

Introduction

Theme : Image Generation from Scene Graphs

What I cannot create, I do not understand

-Richard Feynman

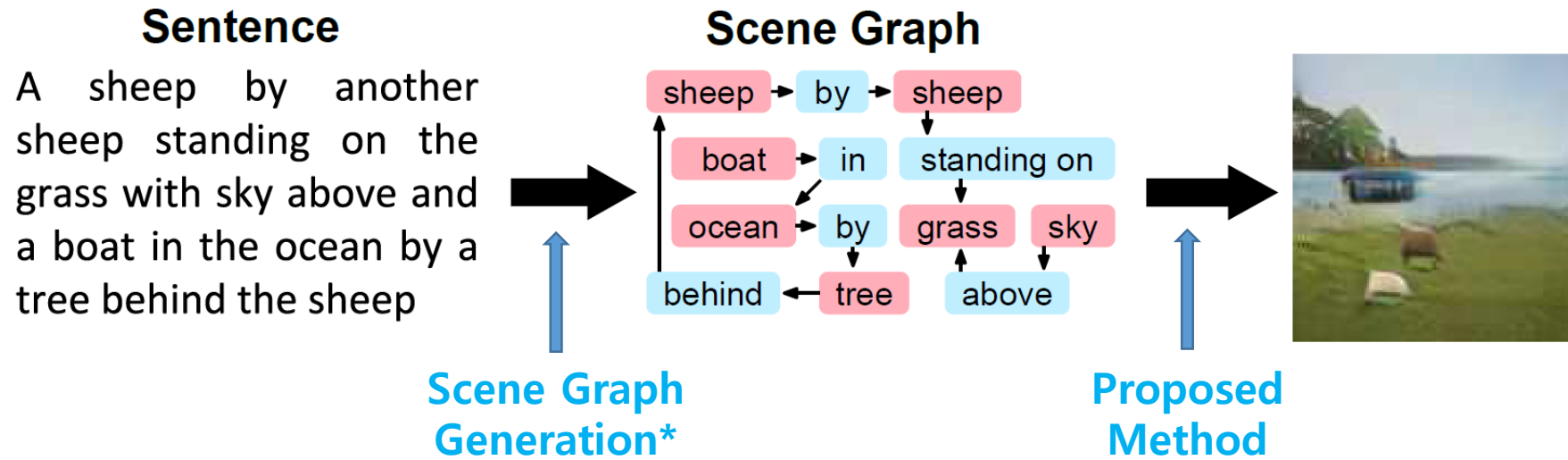
Sentence

A sheep by another
sheep standing on the
grass with sky above and
a boat in the ocean by a
tree behind the sheep



Introduction

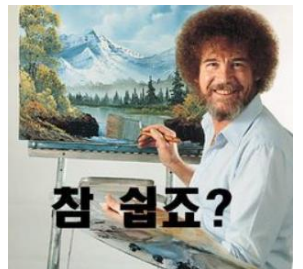
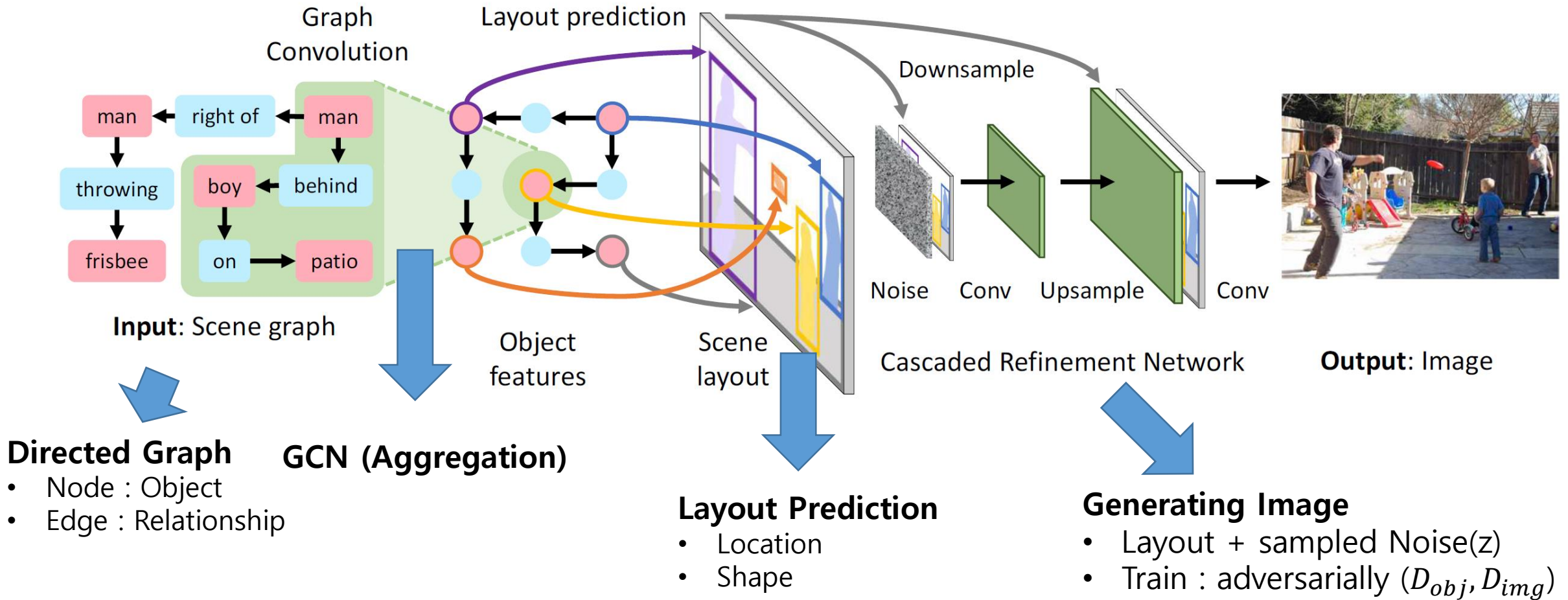
Theme : Image Generation from Scene Graphs



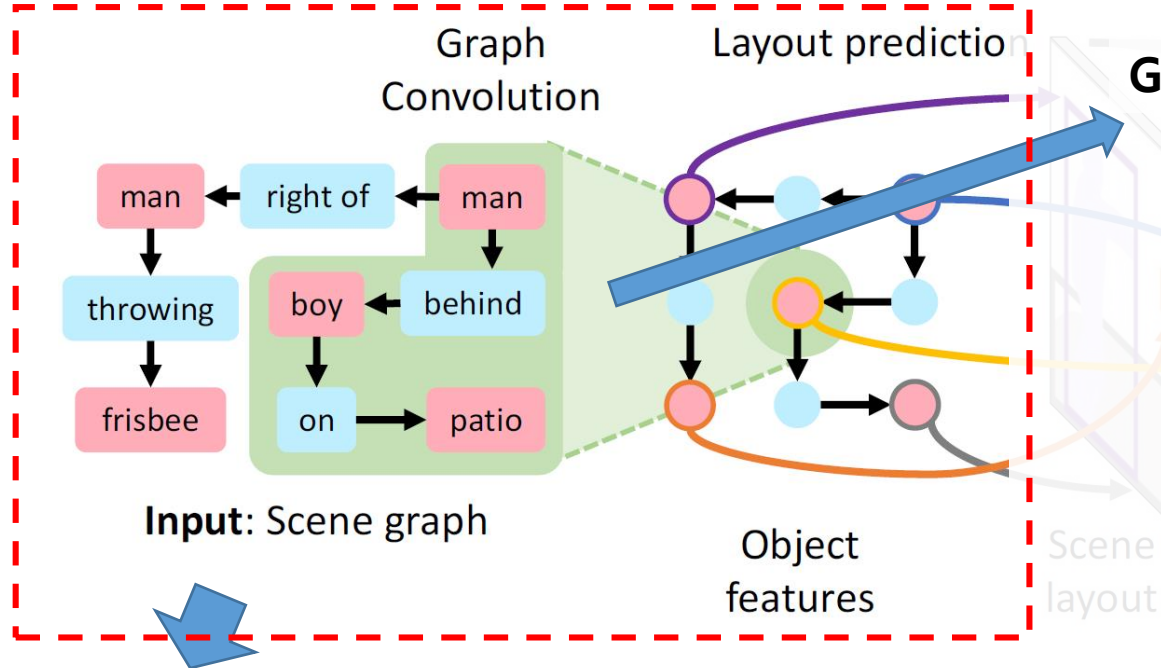
**Q. Chen and V. Koltun. Photographic image synthesis with cascaded refinement networks. In ICCV, 2017*

Overview

Theme : Image Generation from Scene Graphs



Scene Graph → Features



Directed Graph

- Node : Object
- Edge : Relationship

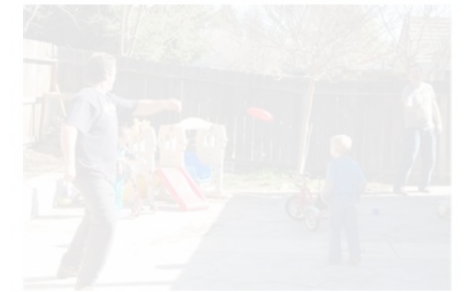
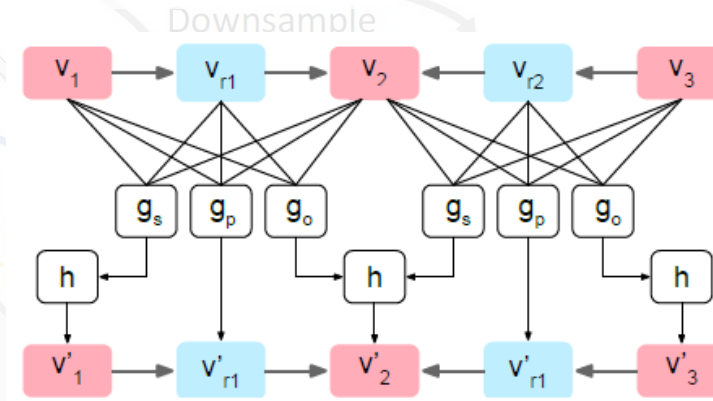
Scene Graph = (O, E)

where

- O = element from object category
- $E \subseteq O \times R \times O$, where R is relation

Theme : Image Generation from Scene Graphs

GCN (Aggregation)



Output: Image

- v_1, v_2, v_3 : object vector(node)
- v_{r1}, v_{r2} : relation vector(edge)
- g_s, g_p, g_o : MLP aggregates information
- h : function that matches output vector dimension
- $v'_1, v'_2, v'_3, v'_{r1}, v'_{r2}$: Output vectors

- ✓ relation vector : $v_r = g_p(v_i, v_r, v_j)$
- ✓ **subject** vectors $V_i^s = \{g_s(v_i, v_r, v_j); (v_i, v_r, v_j) \in E\}$
- ✓ **object** vectors $V_i^o = \{g_o(v_j, v_r, v_i); (v_j, v_r, v_i) \in E\}$
- ✓ $\Rightarrow v'_i = h(V_i^o \cup V_i^s)$

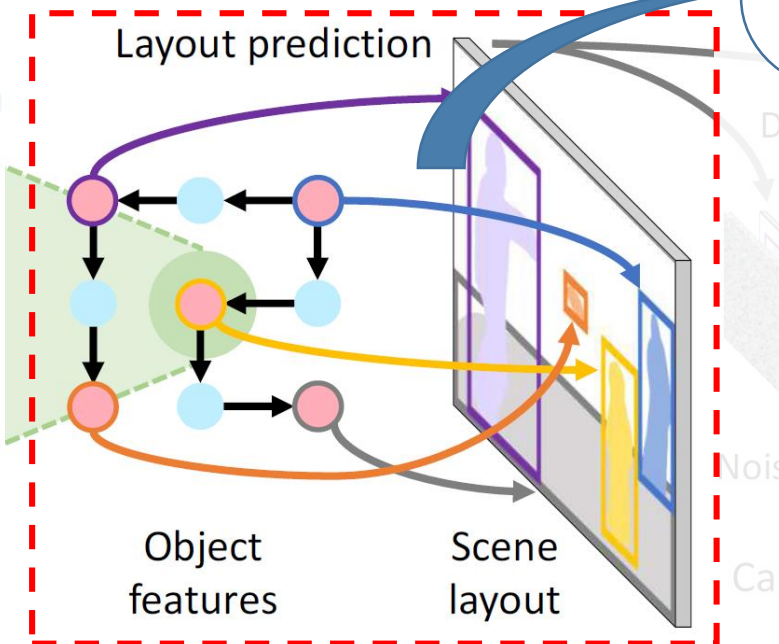
Features → Layout

Layout Prediction

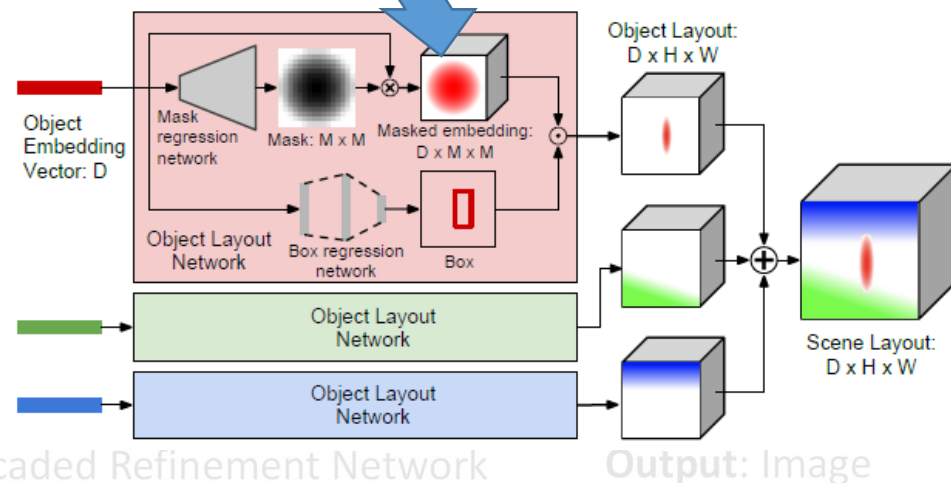
- Segmentation Mask
- Bounding Box



Input: Scene graph



Theme : Image Generation from Scene Graphs



Mask Regression Network

- Transpose Conv Net with last sigmoid
 - Soft binary $M \times M$ mask
- elementwise multiplication with input vector $v_i \in R^D$

Box Regression Network

- Use GT at training

Input

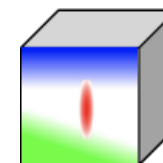
Object Feature $v_i \in R^D$

Object i 's layout

Object j 's layout

Object k 's layout

Σ



Scene Layout: $D \times H \times W$

Layout → Fine Image

Theme : Image Generation from Scene Graphs

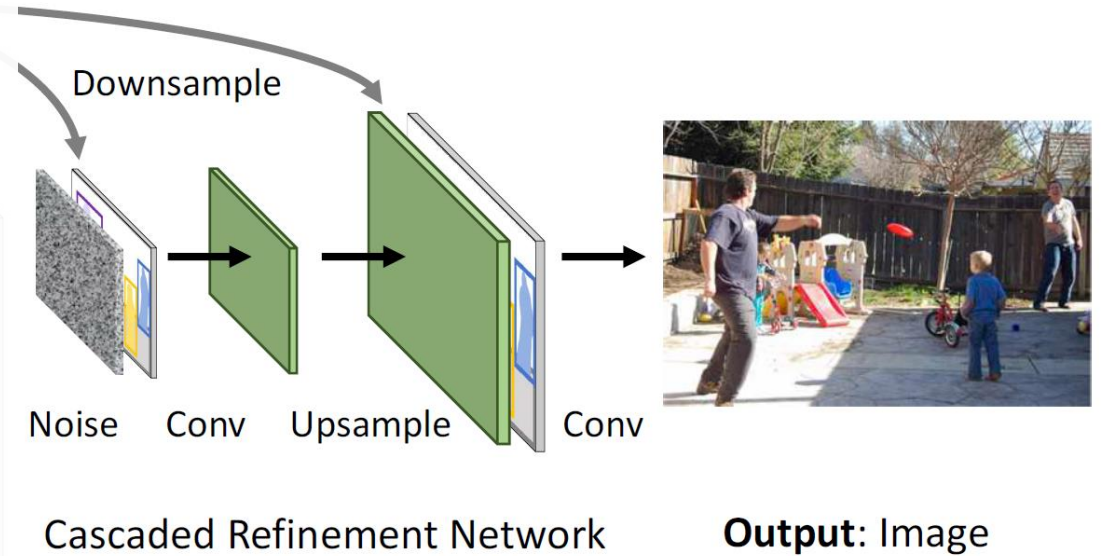
Cascaded Refinement Network

- resolution doubling (nearest neighbor interpolation)
- Previous output + layout prediction → output
- First module take Gaussian noise z

Discriminator

- Discriminator D_{img}, D_{obj}
- D_{obj} determines object is real or fake
- D_{obj} also use auxiliary classifier*, make the object recognizable

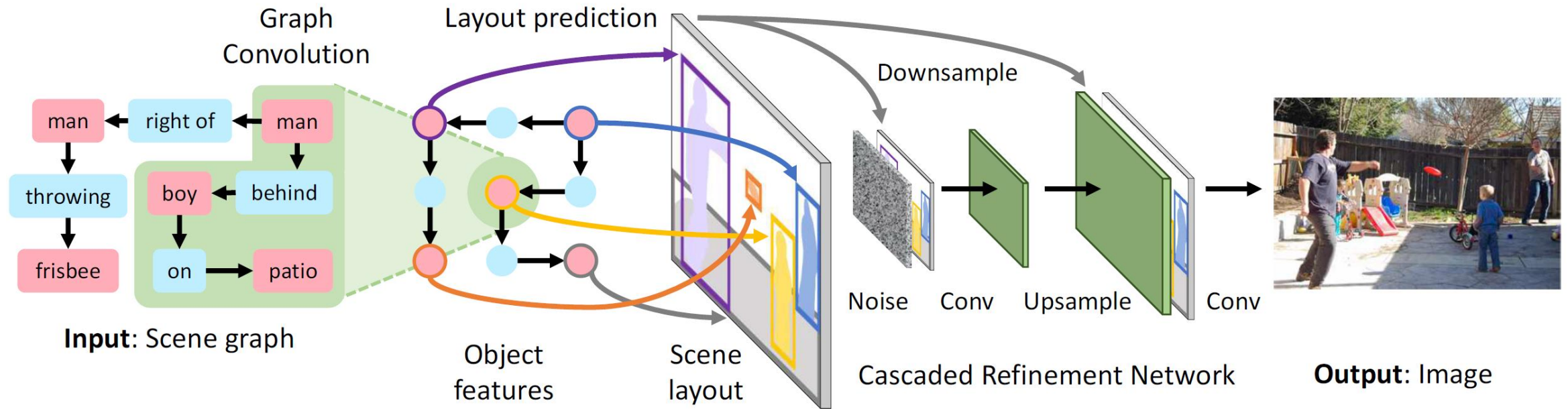
$$\mathcal{L}_{GAN} = \mathbb{E}_{x \sim p_{\text{real}}} \log D(x) + \mathbb{E}_{x \sim p_{\text{fake}}} \log(1 - D(x))$$



*A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier gans. In ICML, 2017

Training Losses

Theme : Image Generation from Scene Graphs



To Be Trained :

- **Generator**
 - **GNN**
 - **Layout prediction**
 - **Refinement Net**
- **Discriminator**



- All Jointly Trained
- Update generator first, then discriminator

Losses

- Box loss : $L_{box} = \sum_{i=1}^n ||b_i - \hat{b}_i||_1$, $b_i = (x_1, y_1, x_2, y_2)$
- Mask loss : $L_{mask} = \text{Pixel} - \text{Wise} - \text{CrossEntropy}$ of object mask
- Pixel loss : $L_{pix} = ||I - \hat{I}||_1$, pixel wise difference of images
- Adversarial loss : $L_{GAN}(D_{img})$, $L_{GAN}(D_{obj})$, $L_{Auxiliary Classifier}(D_{obj})$

Implementation of Source Code*

Theme : Image Generation from Scene Graphs

Code Structure

sg2im : master directory

- └ images : results and corresponding scene graph that author provided
- └ outputs : results that user generated will be saved here
- └ scene_graphs : input scene graph data, file format : .json
- └ **scripts** : source codes that author implemented
 - └ sg2im-models : trained parameters are saved here
 - └ download_OOOO.sh : download OOOO models that author trained
 - └ preprocess_vg.py : create word dictionaries, encode graphs
 - └ **run_model.py** : Test the trained model
 - └ **train.py** : Train the model (explained on next page)
- └ **sg2im** : common utilities –building layer of network, crop the image, and calculate losses

*Code Link : <https://github.com/google/sg2im>

Implementation of Source Code

Theme : Image Generation from Scene Graphs

Core Modules in Source Code : train.py

1. Sg2ImModel : Main backbone of entire code

- Graph Neural Network + Mask / Box Prediction
- 5-layer MLP with batchnormalization
- Train network in Training phase// test with this model

2. build_obj_discriminator : trained to discriminate an object in image

3. build_img_discriminator : trained to discriminate entire image

4. build_vg_dsets :

- Draw scenegraph from dataset
- Data sampler, collator, batch-formation

5. calculate_model_losses

- total loss : bbox, pixel, mask, discriminator, auxiliary classifier

Experiments & Result

Theme : Image Generation from Scene Graphs

conducted by [Author](#)



Which image matches the caption better?

User choice	332 / 1024 (32.4%)	692 / 1024 (67.6%)
-------------	-----------------------	-------------------------------



Which objects are present? motorcycle, person, clouds

Thing recall	470 / 1650 (28.5%)	772 / 1650 (46.8%)
Stuff Recall	1285 / 3556 (36.1%)	2071 / 3556 (58.2%)

Method	Inception	
	COCO	VG
Real Images (64×64)	16.3 ± 0.4	13.9 ± 0.5
Ours (No gconv)	4.6 ± 0.1	4.2 ± 0.1
Ours (No relationships)	3.7 ± 0.1	4.9 ± 0.1
Ours (No discriminators)	4.8 ± 0.1	3.6 ± 0.1
Ours (No D_{obj})	5.6 ± 0.1	5.0 ± 0.2
Ours (No D_{img})	5.6 ± 0.1	5.7 ± 0.3
Ours (Full model)	6.7 ± 0.1	5.5 ± 0.1
Ours (GT Layout, no gconv)	7.0 ± 0.2	6.0 ± 0.2
Ours (GT Layout)	7.3 ± 0.1	6.3 ± 0.2
StackGAN [59] (64×64)	8.4 ± 0.2	-

Experiments & Result

Theme : Image Generation from Scene Graphs

Algorithm run by **student(TaeGil Ha)**

car on street
line on street
sky above street



bus on street
line on street
sky above street



car on street
bus on street
line on street
sky above street



car on street
bus on street
line on street
sky above street
kite in sky



car on street
bus on street
line on street
sky above street
kite in sky
car below kite



car on street
bus on street
line on street
sky above street
building behind street



car on street
bus on street
line on street
sky above street
building behind street
window on building



sky above grass
zebra standing on grass



sky above grass
sheep standing on grass



sky above grass
sheep standing on grass
sheep' by sheep



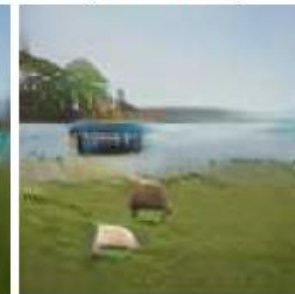
sky above grass
sheep standing on grass
sheep' by sheep
tree behind sheep



sky above grass
sheep standing on grass
tree behind sheep
sheep' by sheep
ocean by tree



sky above grass
sheep standing on grass
tree behind sheep
sheep' by sheep
ocean by tree
boat in ocean



sky above grass
sheep standing on grass
tree behind sheep
sheep' by sheep
ocean by tree
boat on grass

