

# Graph R-CNN for Scene Graph Generation

Keegan Lensink, Xiang Liu, Zicong (Alex) Fan



THE UNIVERSITY  
OF BRITISH COLUMBIA

# Table of Contents

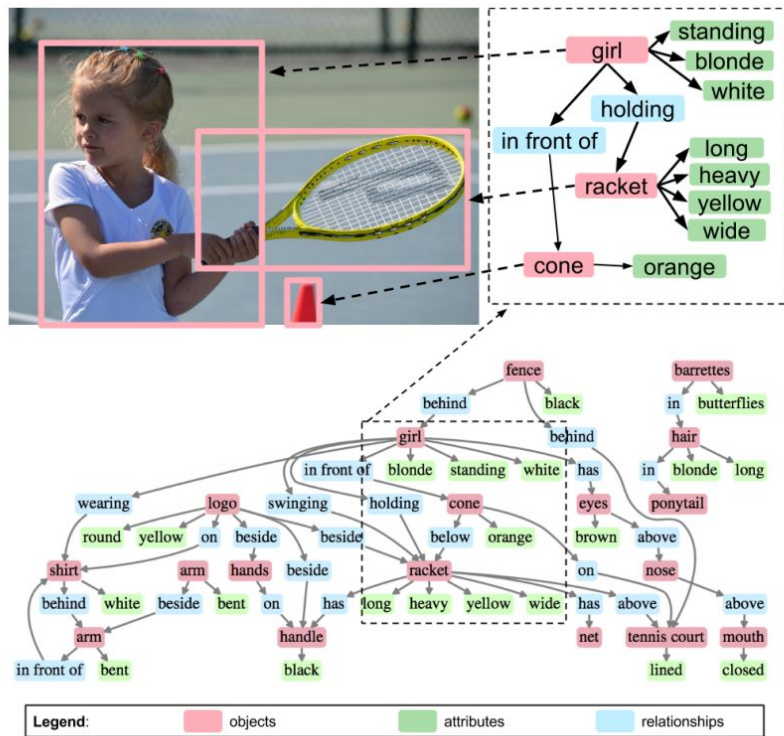
- Introduction
- Contributions
- Architecture:
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- Metric SGGen+
- Discussion

# Table of Contents

- Introduction
- Contributions
- Architecture:
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- Metric SGGen+
- Discussions

# Introduction

- Scene graphs form a structured representation of the image.
- Difficult to generate from a given image.

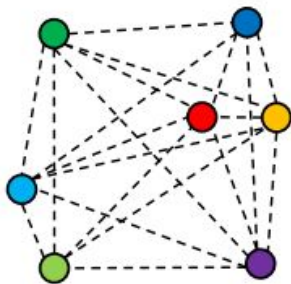


Johnson, J., Krishna, R., Stark, M., Li, L.J., Shamma, D.A., Bernstein, M., Fei-Fei, L.: Image retrieval using scene graphs. In: CVPR (2015)

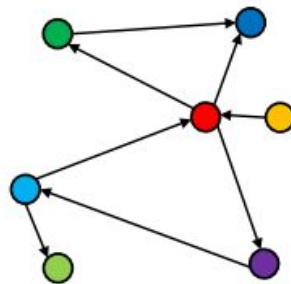
# Introduction



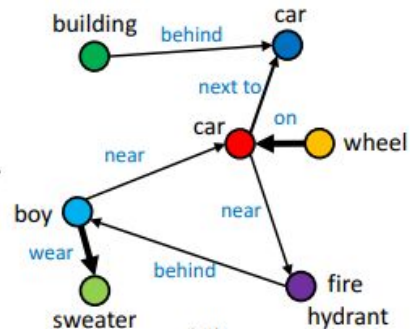
(a)



(b)



(c)



(d)

- Given  $n$  detected objects, a fully connected graph has  $n^2$  edges.
- Many of these edges are probably between unrelated objects.
- How can we prune this graph?

# Table of Contents

- Introduction
- Contributions
- Architecture:
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- Metric SGGen+
- Discussions

# Contributions

- RePN
  - Relation Proposal Network
- aGCN
  - attentional Graph Convolutional Networks.
- SGGEN+
  - Modified evaluation metric for scene graphs that gives more realistic results.
- Graph RCNN framework
  - Generates scene graph for a given image.
  - Provides node features for downstream tasks.

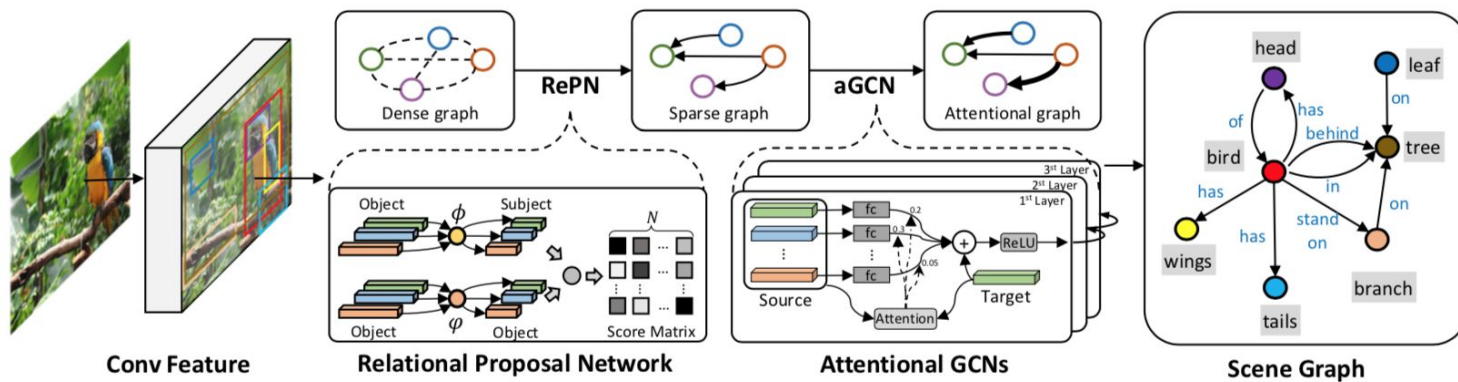
# Table of Contents

- Introduction
- Contributions
- **Architecture:**
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- Metric SGGen+
- Discussions



# Architecture

$$P(\mathcal{S}|\mathbf{I}) = \underbrace{P(\mathbf{V}|\mathbf{I})}_{\text{Object Region Proposal}} \underbrace{P(\mathbf{E}|\mathbf{V}, \mathbf{I})}_{\text{Relationship Proposal}} \underbrace{P(\mathbf{R}, \mathbf{O}|\mathbf{V}, \mathbf{E}, \mathbf{I})}_{\text{Graph Labeling}}$$



# Table of Contents

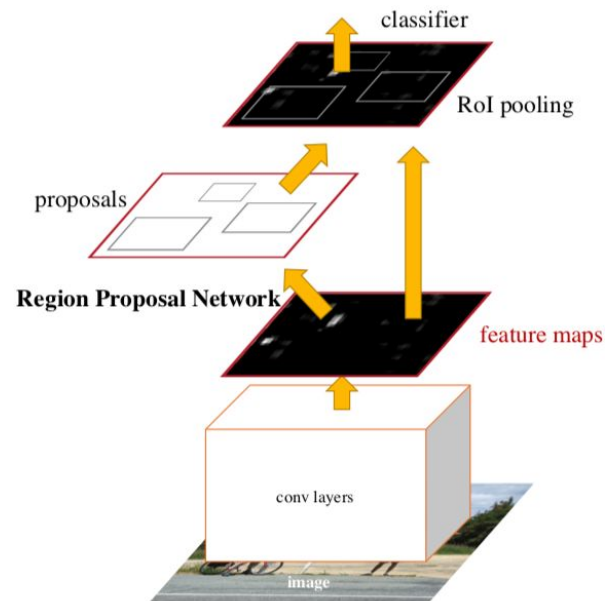
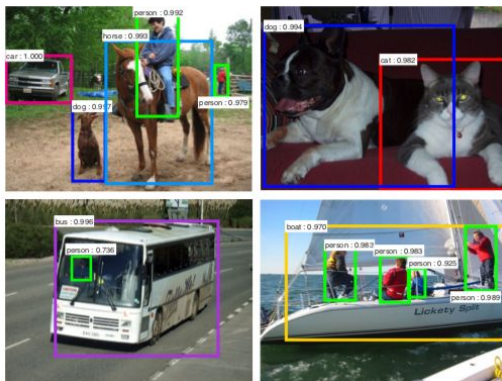
- Introduction
- Contributions
- **Architecture:**
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- Metric SGGen+
- Discussions

# Object Region Proposal

Given an image  $I$ , extract the following quantities with Faster-RCNN:

- Bounding boxes:
- Feature vectors:
- Label distributions:

$$R^o \in \mathbb{R}^{n \times 4}$$
$$X^o \in \mathbb{R}^{n \times d}$$
$$P^o \in \mathbb{R}^{n \times |C|}$$



# Table of Contents

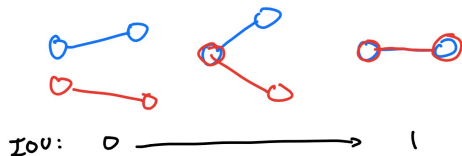
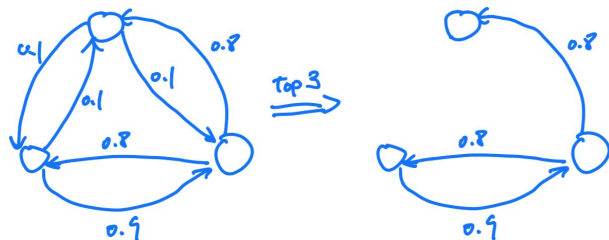
- Introduction
- Contributions
- **Architecture:**
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- Metric SGGen+
- Discussions

# Relationship Proposal (RePN)

Learn a relatedness function:

$$f(\mathbf{p}_i^o, \mathbf{p}_j^o) = \langle \Phi(\mathbf{p}_i^o), \Psi(\mathbf{p}_j^o) \rangle, i \neq j$$

- Binary classification:
  - Score of relatedness:  $[0, 1]$
  - $\{\text{edge}, \text{not\_edge}\}$
- Non-maximal suppression
  - $k$  top-scored edges
  - $m$  edges with least overlaps with others.



$$IoU(\{u, v\}, \{p, q\}) = \frac{I(r_u^o, r_p^o) + I(r_v^o, r_q^o)}{U(r_u^o, r_p^o) + U(r_v^o, r_q^o)}$$

# Table of Contents

- Introduction
- Contributions
- **Architecture:**
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- Metric SGGen+
- Discussions

# Conventional GCN

Recap on Graph Convolutional Network (GCN):

- Given a graph represented by:

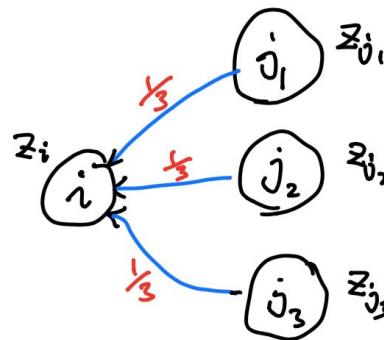
- Feature matrix:  $Z$  in  $N \times D_L$ .
- Adjacency matrix:  $\alpha$  in  $N \times N$ .
- Dimension map:  $W$

- Propagation rule:

$$z_i^{(l+1)} = \sigma \left( z_i^{(l)} + \sum_{j \in \mathcal{N}(i)} \alpha_{ij} W z_j^{(l)} \right)$$

- Interpretation:

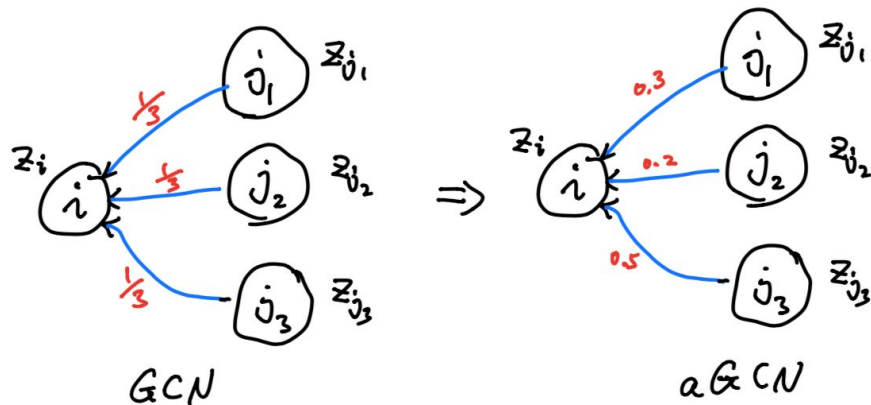
- Representation of each node at the next layer is aggregated by the equally-weighted average of its neighbours and itself.



$$z_i^{(l+1)} = \sigma \left( W Z^{(l)} \alpha_i \right)$$

# Attentional GCN (aGCN)

aGCN: “.. equally-weighted average of its neighbours”



$$z_i^{(l+1)} = \sigma \left( z_i^{(l)} + \sum_{j \in \mathcal{N}(i)} \alpha_{ij} W z_j^{(l)} \right)$$

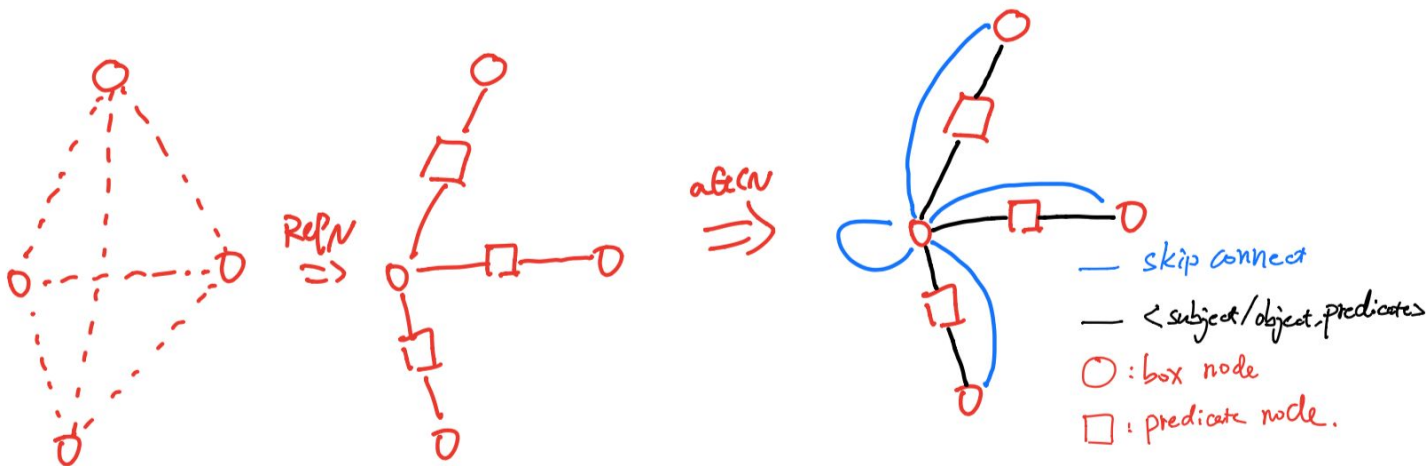
$$u_{ij} = w_h^T \sigma(W_a[z_i^{(l)}, z_j^{(l)}])$$

$$\alpha_i = \text{softmax}(\mathbf{u}_i),$$

$$\alpha_{ii} = 1 \text{ and } \alpha_{ij} = 0 \forall j \notin \mathcal{N}(i).$$



# Graph Formation

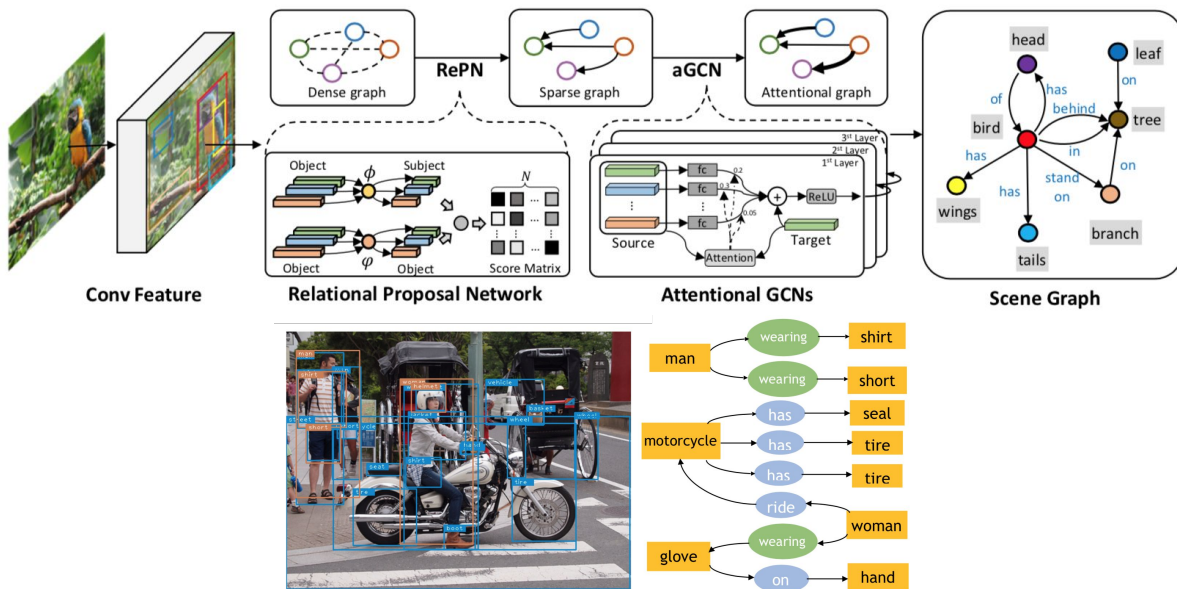


$$z_i^o = \sigma \left( \underbrace{W^{\text{skip}} Z^o \alpha^{\text{skip}}}_{\text{Message from Other Objects}} + \underbrace{W^{sr} Z^r \alpha^{sr} + W^{or} Z^r \alpha^{or}}_{\text{Messages from Neighboring Relationships}} \right)$$

$$z_i^r = \sigma \left( z_i^r + \underbrace{W^{rs} Z^o \alpha^{rs} + W^{ro} Z^o \alpha^{ro}}_{\text{Messages from Neighboring Objects}} \right).$$

# Graph Labeling

“Two multi-class cross entropy losses are used for object classification and predicate classification.”



# Table of Contents

- Introduction
- Contributions
- Architecture:
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- **Metric SGGen+**
- Discussions

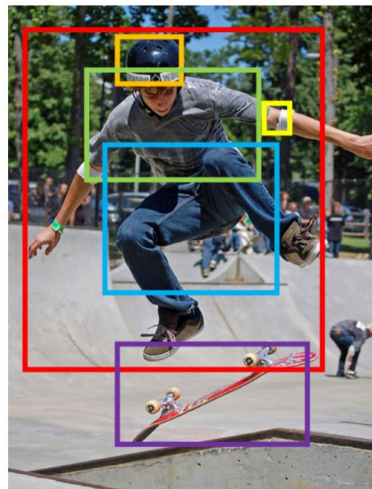
# Metric: SGGen

The ground truth scene graph: <subject, relationship, object>

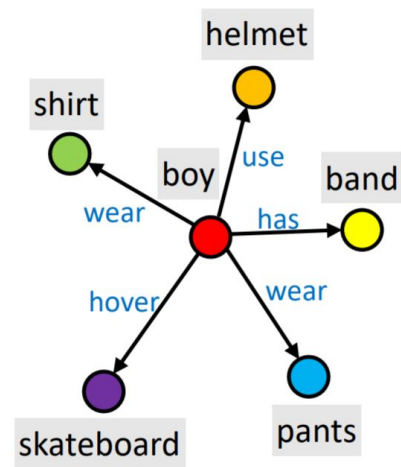
E.g.: <boy, wear, shirt>, <boy, use, helmet>, <boy, has, band> ...

SGGen counts one match when:

1. all three elements have been correctly labeled
2. both object and subject nodes have been properly localized (i.e., bounding box IoU > 0.5).



(a)



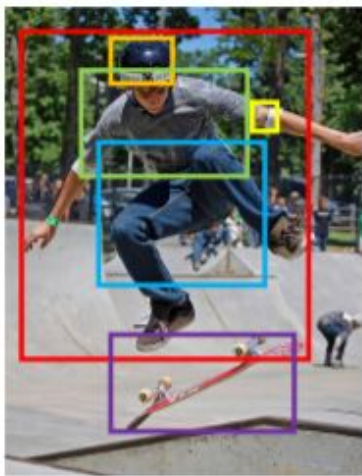
SGGen = 5

SGGen+ = 16

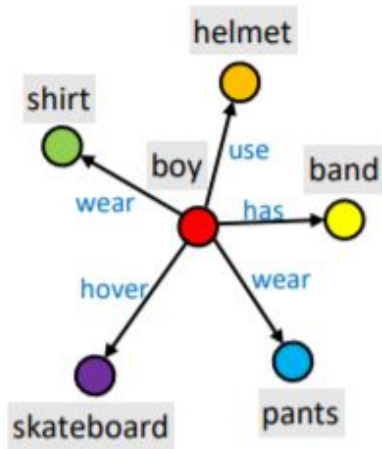
(b)

# Problem of SGGen

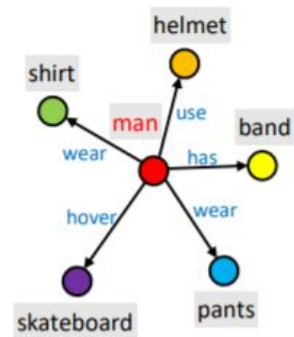
It only counts exact matches.



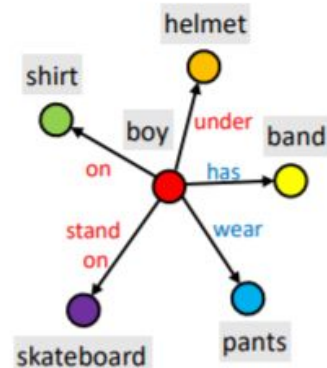
(a)



SGGen = 5      SGGen+ = 16  
(b)



SGGen = 0      SGGen+ = 10  
(d)



SGGen = 2      SGGen+ = 9  
(e)

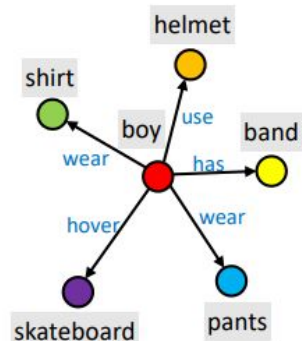
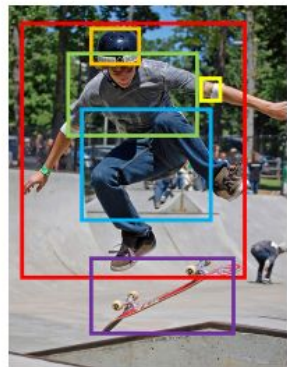
# “A More Comprehensive Metric”: SGGen+

$$\text{SGGen+} = C(O) + C(P) + C(T)$$

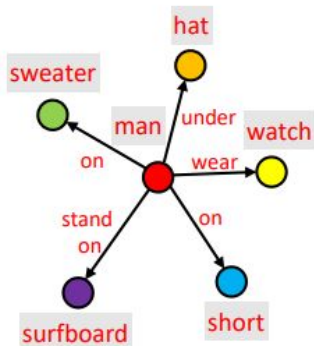
$C(O)$ : the number of object nodes correctly localized and recognized

$C(P)$ : the number of predicates correctly localized and recognized

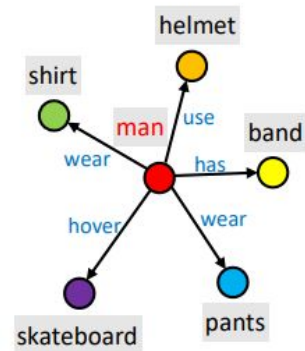
$C(T)$ : the number of matched triples, which is SGGen



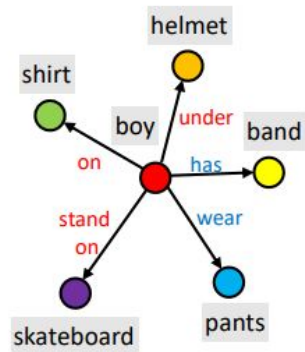
SGGen = 5      SGGen+ = 16



SGGen = 0      SGGen+ = 0



SGGen = 0      SGGen+ = 10



SGGen = 2      SGGen+ = 9

# Comparing of SGGen and SGGen+

Assign random incorrect labels to objects perturbing objects

Perturb Type	none	w/o relationship			w/ relationship			both		
Perturb Ratio	0%	20%	50%	100%	20%	50%	100%	20%	50%	100%
SGGen	100.0	100.0	100.0	100.0	54.1	22.1	0.0	62.2	24.2	0.0
SGGen+	100.0	94.5	89.1	76.8	84.3	69.6	47.9	80.1	56.6	22.8

**Table 1.** Comparisons between SGGen and SGGen+ under different perturbations.

SGGen is completely insensitive to the perturbation of objects without relationships.

SGGen is overly sensitive to label errors on objects with relationships.

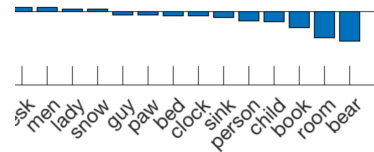
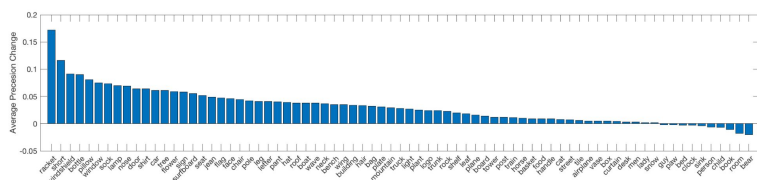
100% perturbation: the object localizations and relationships are still correct such that SGGen+ provides a non-zero score.

# Table of Contents

- Problem statement
- Contributions
- Architecture:
  - Object Proposal Network
  - Relationship Proposal Network (RePN)
  - Attentional GCN (aGCN)
- Metric SGGen+
- Discussion



# Discussion



- Less robust to objects that are used in varied contexts (person), trading performance for small objects that are used in relatively few contexts (racket).
- Use P to learn relatedness between two boxes.
  - Better solution:
    - learn feature representation of a class s.t. closer labels are closer in feature space.
    - e.g.  $D(\text{car}, \text{wheel}) < D(\text{car}, \text{noodles})$
    - Use this distance as a score instead.
- The choice of the new metric might be biased toward their model.
- Choosing  $m$  for remaining object pairs is not clear, and depends on the image.
- Only top 150 classes and top 50 relations; real world has more labels.
- Use similarity instead of count, e.g.  $\text{Similarity}(\text{boy}, \text{man}) = 0.8$

# References

- Yang, Jianwei, et al. "Graph r-cnn for scene graph generation." ECCV, 2018.
- Kipf, Thomas N., and Max Welling. "Semi-supervised classification with graph convolutional networks." ICLR, 2017.
- Ren, Shaoqing, et al. "Faster r-cnn: Towards real-time object detection with region proposal networks." NeurIPS, 2015.

