# Consistent Scene Graph Generation by Constraint Optimization

**Boqi Chen**[1], Kristóf Marussy[2], Sebastian Pilarski[1], Oszkár Semeráth[2], Daniel Varro[1]

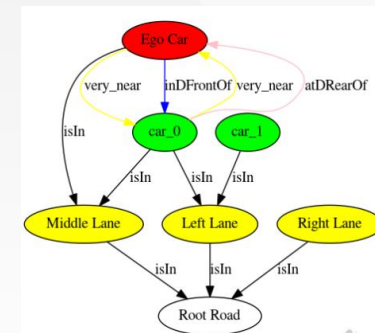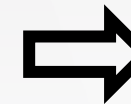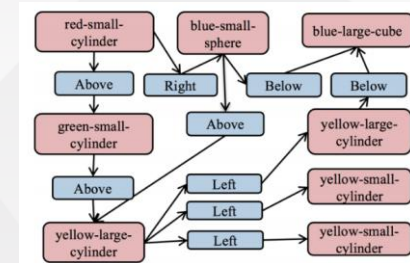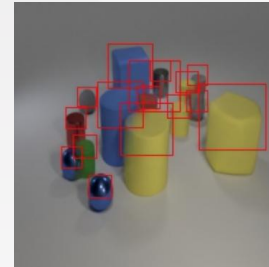[1]Department of Electrical & Computer Engineering, McGill University
[2]Department of Measurement and Information Systems, Budapest University of Technology and Economics

boqi.chen@mail.mcgill.ca, marussy@mit.bme.hu, sebastian.pilarski@mail.mcgill.ca, semerath@mit.bme.hu, daniel.varro@mcgill.ca
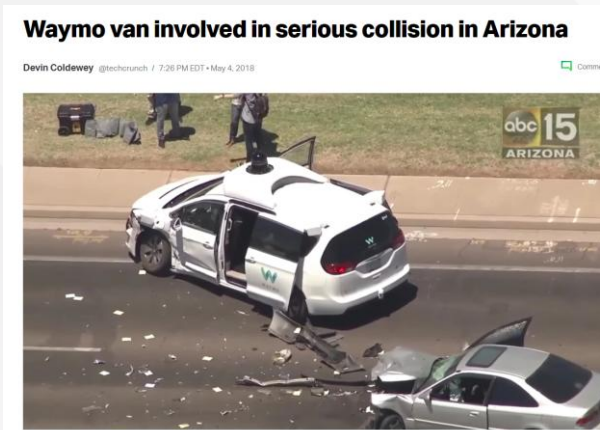
# Scene Graph Generation (SGG)

- **Scene graphs** represent objects and their relations in an image

- **Scene graph generation** produces a graph from a given image containing *focused objects* and their *relationships*

- Scene graph generation is a common challenge in computer vision

Examples taken from:
Herzig, Roei, et al. "Learning canonical representations for scene graph to image generation." *European Conference on Computer Vision*. Springer, Cham, 2020.
Yu, Shih-Yuan, et al. "Scene-graph augmented data-driven risk assessment of autonomous vehicle decisions." *IEEE Transactions on Intelligent Transportation Systems* (2021).

Introduction | Problem Formulation | Approach | Evaluation | Conclusion

- **Scene graphs** can be used in safety critical fields such as autonomous driving and robotics



- In such applications, it is important to provide safety guarantee on the produced scenes under consistent situations
  - Law of physics: Car cannot (yet) fly
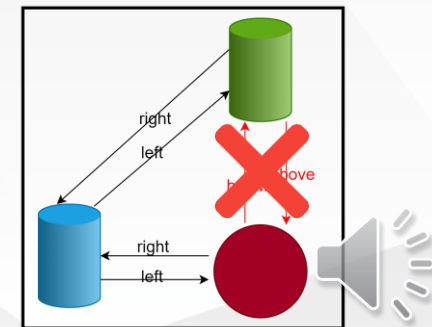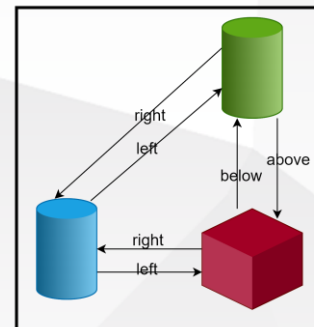  - Traffic rules: No contradictory traffic signs

# Motivation: Constraints Formulation

- **What is safety?** One aspect of safety is *consistency*: The system should comply with a set of *consistency constraints* $\Phi$

- *Consistency constraints* can be expressed with logic or constraint languages (**FOL**, OLC and VIATRA-Query etc.)

Rule: Nothing can be above a sphere

FOL →

$$\forall a, b: \text{Above}(a, b) \Rightarrow \neg\text{Shape}(b, \text{'Sphere'})$$

Assumption: ground truth scenes are consistent

# Problem Statement

- To *guarantee consistency* in such systems, we define and tackle the problem of *consistent scene graph generation*

- Given a set of constraints $\Phi$, an image $I$ with underlying ground truth scene graph $SG_{gt}$, find a model $\mathcal{M}$ for scene graph $SG$ such that

  1. The generated $SG$ is close to the ground truth $SG_{gt}$ (accurate):

  $$P(SG \cong SG_{gt} | \mathcal{M})$$

  2. The generated $SG$ satisfy all consistency constraints $\Phi$ (consistent)

  $$SG \vDash \Phi$$

accurate

consistent

$$P(SG \cong SG_{gt}|\mathcal{M}) \wedge SG \vDash \Phi$$

Many existing deep learning approaches

No explicit guarantees

Assumption: $SG_{gt}$ satisfies $\Phi$: $SG_{gt} \vDash \Phi$

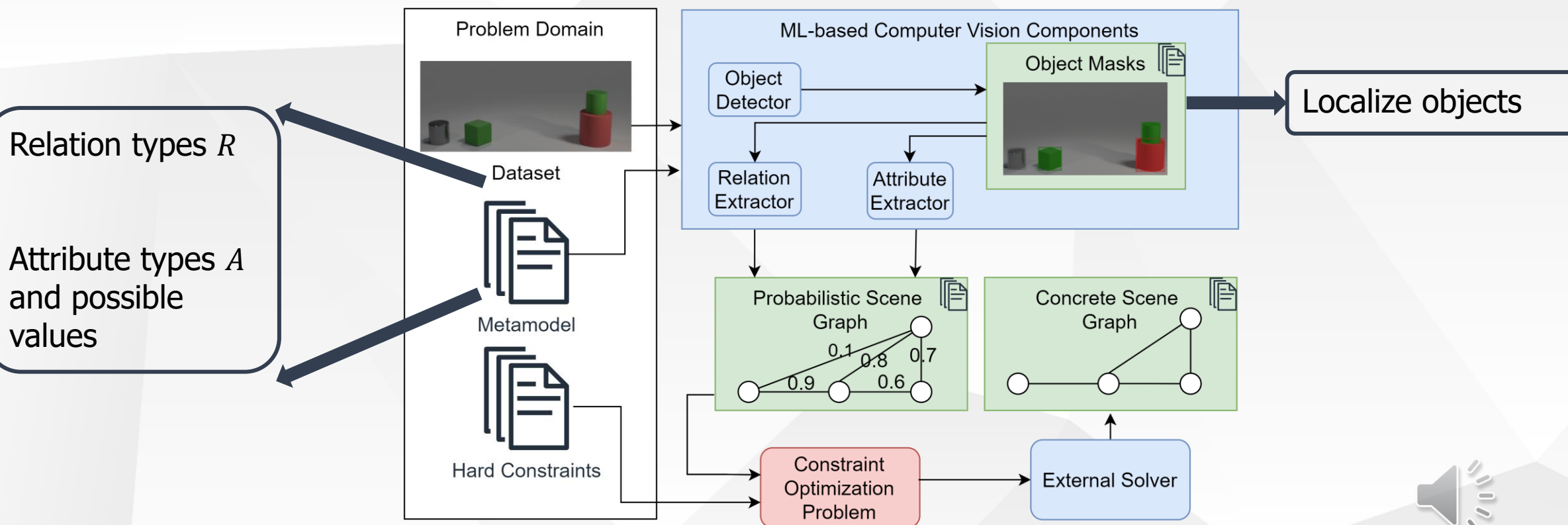How can we guarantee constraints are always satisfied?

Can help?

$$P(SG \cong SG_{gt}|\mathcal{M}) \qquad SG \vDash \Phi$$

help

Core Idea: use existing DL methods to optimize for $P(SG \cong SG_{gt}|\mathcal{M})$, and handle $SG \vDash \Phi$ later with *constraint optimization*
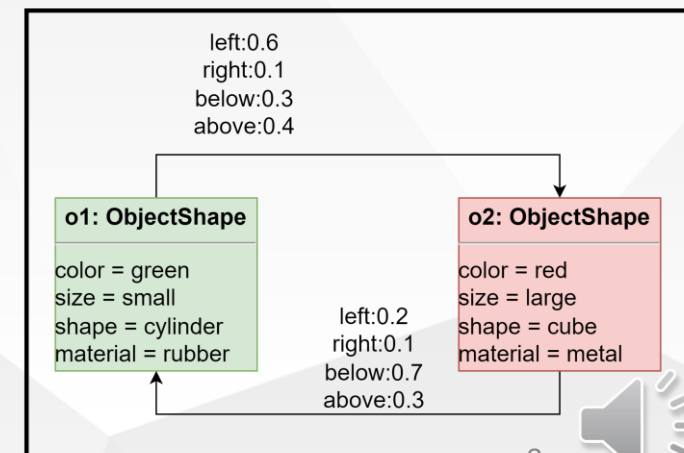
A ML-based vision model outputs two types of independent probabilities for an input image to form a *probabilistic graph*

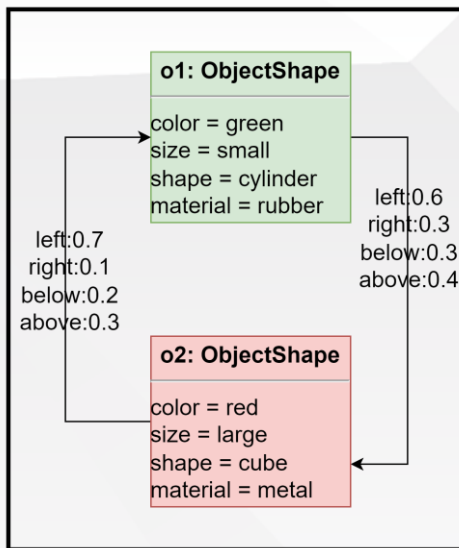1. The probability of a pair of objects and a relation type

$$P_R: N \times \mathbb{R} \times N \rightarrow [0,1], : P_R\left(\overrightarrow{n_1 n_2}^{\,r}\right)$$

2. The probability of an object, an attribute type and an attribute value

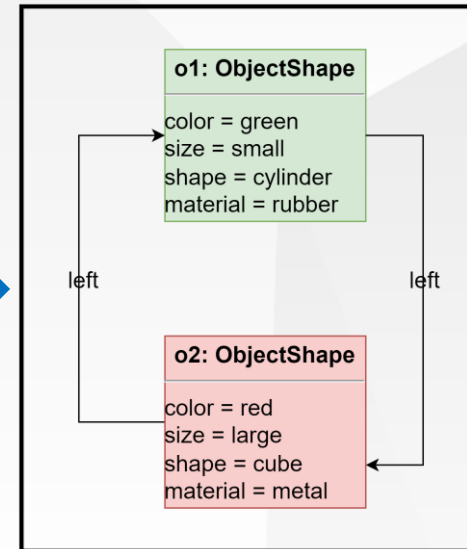$$P_A: N \times A \times V_a \rightarrow [0,1]: P_A\left(\overrightarrow{n\,v}^{\,a}\right)$$

Commonly, we can choose a *concrete scene graph* from the probabilistic SG by selecting the *most probable relations and attributes* individually
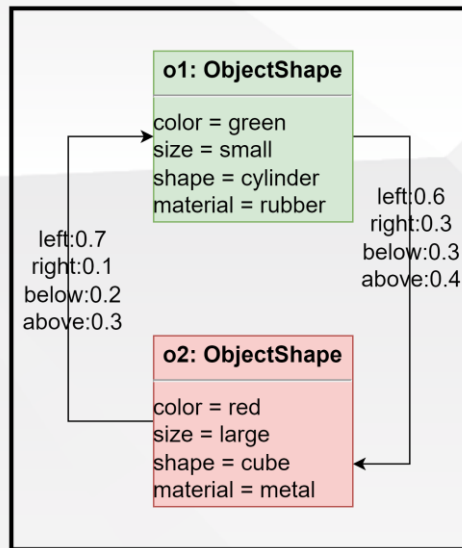


(a)
Probabilistic graph

(b)
Concrete graph

$\Phi$:
1. $left \leftrightarrow right$
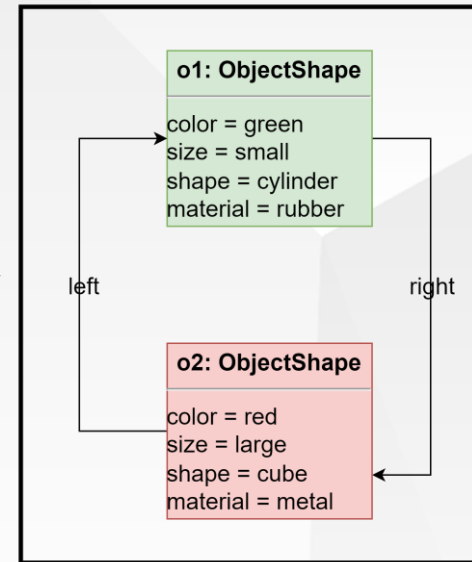2. below $\leftrightarrow above$

**No consideration of $\Phi$!**

Instead, we propose to select the *most probable scene* subject to Φ



(a)
Probabilistic graph

Φ

(c)
Concrete graph

Φ:
1. $left \leftrightarrow right$
2. $below \leftrightarrow above$

Constraint Optimization!

# Approach: MAXSAT

MAXSAT is an optimization problem aiming to find the maximum subset of clauses with weights (in CNF)

Many existing solvers: MaxSatz, WBO, SAT4J, **Gurobi**

Given a set of hard constraints $\Phi$, a set of clauses $C$ with weights $w$

$$\neg\phi: -\infty \text{ for each } \phi \in \Phi$$

$$C_i: w_i \text{ for } i \in [0, n]$$

$$max_x \sum_i w_i \cdot 1(x \vDash C_i) \text{ subject to } x \vDash \phi \text{ for each } \phi \text{ } in \text{ } \Phi$$

# Approach: MAXSAT

Given a probabilistic graph $\mathbb{G}$ with $P_R$ and $P_A$, our approach transform it into a MAXSAT problem by:

1. The hard constraints are respected

$$\neg\phi: -\infty \text{ for each } \phi \in \Phi$$

2. Edges and attributes are clauses with weights being the log probabilities

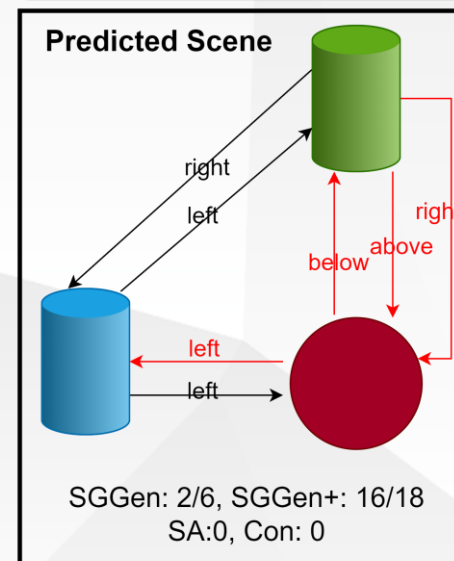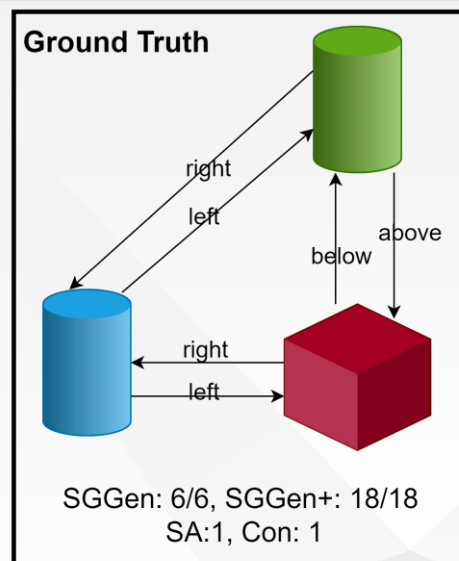$$x_{\overrightarrow{n_1 n_2}r}: \log P_R(\overrightarrow{n_1 n_2}^r) \qquad x_{\overrightarrow{n_1 v}a}: \log P_A(\overrightarrow{n_1 v}^a)$$

3. The optimization target is to maximize the sum of log probabilities

- **SGGen**: measures recall of relations if all attribute of an object is identified correctly

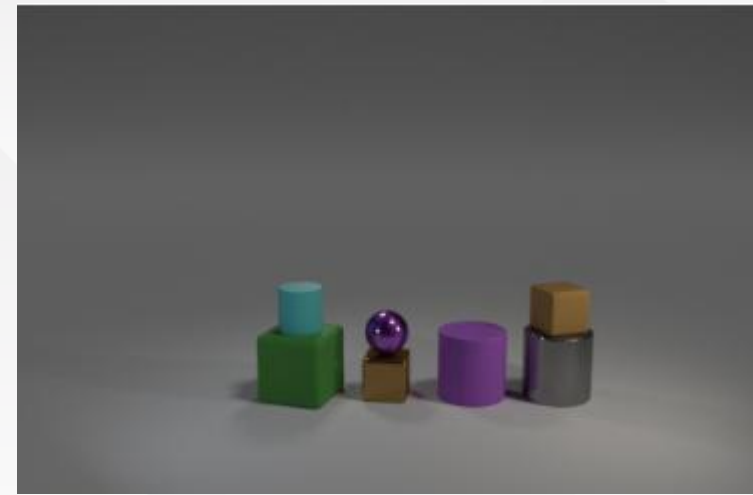- **SGGen+**: measures recall separately for relations and attributes

- **Con**: measures *consistency* of the scene

- **SA**: *scene accuracy* measures if the predicted scene is *isomorphic* to the ground truth scene



Ground Truth

right
left
above
below
right
left

SGGen: 6/6, SGGen+: 18/18
SA:1, Con: 1



Predicted Scene

right
left
right
above
below
left
left

SGGen: 2/6, SGGen+: 16/18
SA:0, Con: 0
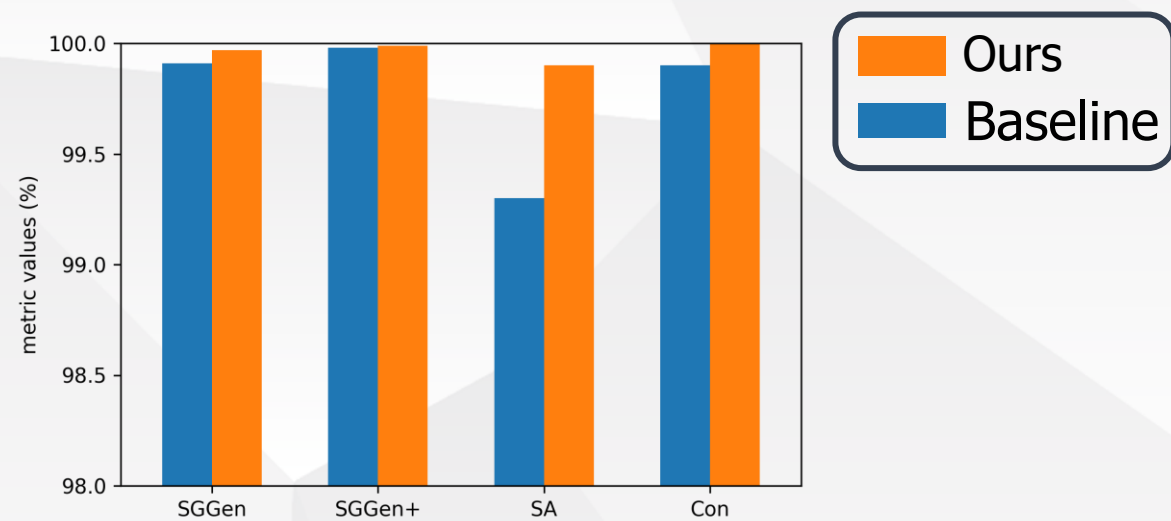
CLEVR:

BLOCKWORLD:





- 4 types of constraints with different complexity were created
- Scenes are generated to satisfy the constraints
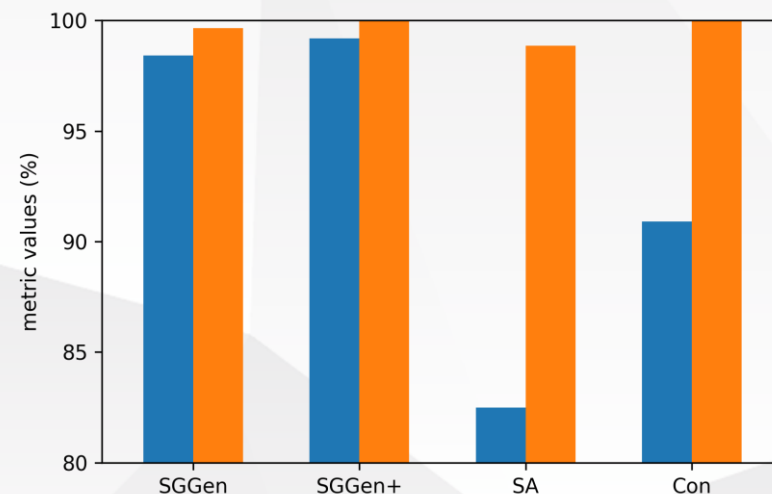- 4000 scenes for training and 2000 for testing

- Our approach is better than the baseline in all cases

- High values in relation recall (SGGen, SGGen+) does not mean high SA

- Our approach always improves SA by improving Con
  - In fact, we can prove SA is *always at least as good* as the baseline scenes



CLEVR

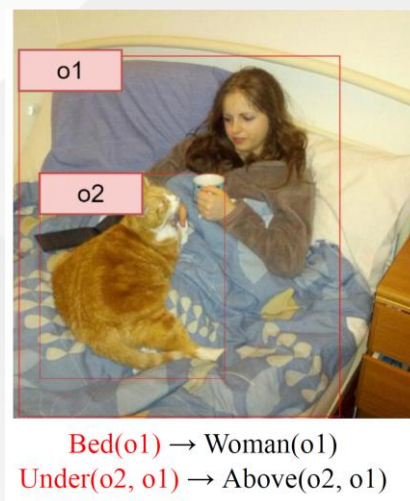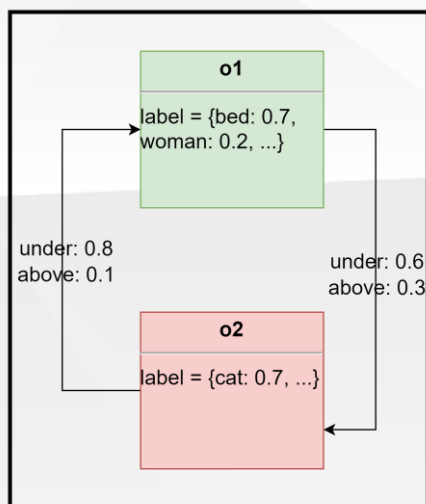

Block

# Evaluation: Real Dataset

- What about the performance on real-world images?

- We applied our approach on a subset of the Visual Genome dataset with two types of constraints:
  - There must be at least one person in the scene
  - There is no cycle on relations such as 'Above', 'Under'
  - We filtered the datasets with the first constraint
  - 99.85% ground truth satisfies the second type of constraints

- Probabilistic scenes are derived from a model pre-trained on *original VG dataset*



Bed(o1) → Woman(o1)
Under(o2, o1) → Above(o2, o1)

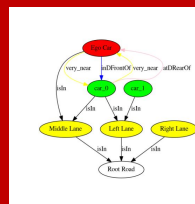| Metric | Improvement (%) |
|---|---|
| SGGen | 33.04 → **33.15** |
| SGGen+ | 63.44 → **63.48** |
| Con | 64.43 → **100** |

- SA is not measured because the labelled graph is not complete
- Our approach is still able to improve on all metrics while ensuring consistency

Image credit: https://www.flickr.com/photos/todoleo/8310164456/

17

# Conclusion

## Scene Graphs

Scene Graphs represents objects and their relations in an image.
Scene graphs can be used by safety critical systems in which consistency is a key.
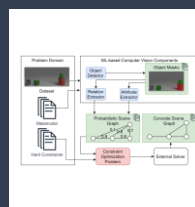


## Consistent Scene Graph Generation

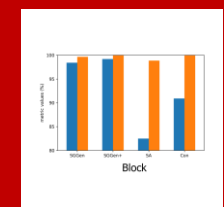Generate scene graphs from images that comply to the set of constraints $\Phi$.



## Constraint optimization

use existing DL methods to optimize for $P(SG \cong SG_{gt}|\mathcal{M})$, and handle $SG \vDash \Phi$ later with a *constraint optimization.*



## Neural Symbolic Reasoning

- Applications to autonomous vehicles
- Incorporate with neural MAXSAT solvers
- Certify consistency directly from deep learning component

# Thank you

Artifacts available at: https://github.com/20001LastOrder/Clevr-Relational