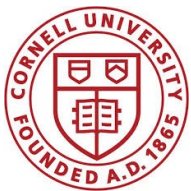# Natural Language for Visual Reasoning

Alane Suhr, Mike Lewis, James Yeh, Yoav Artzi

lic.nlp.cornell.edu/nlvr/

# Language and Vision



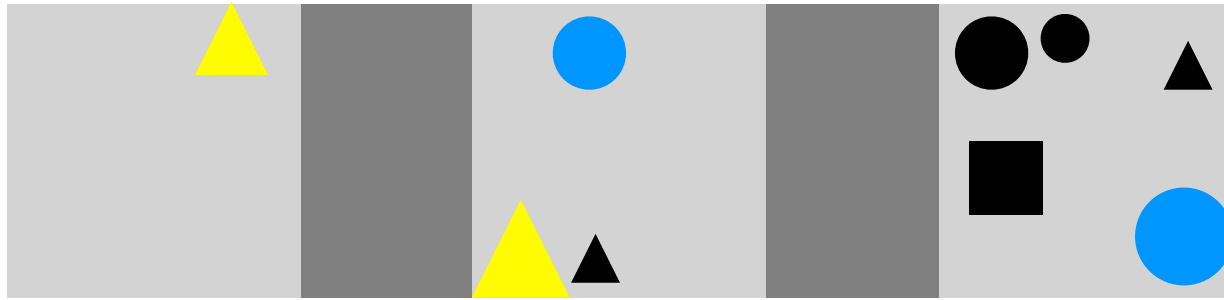A small herd of cows in a large grassy field.

(Chen et al 2015)



What is the dog carrying?

(Agrawal et al 2015)

**Our goal:** natural language with a diverse set of semantic and syntactic phenomenon
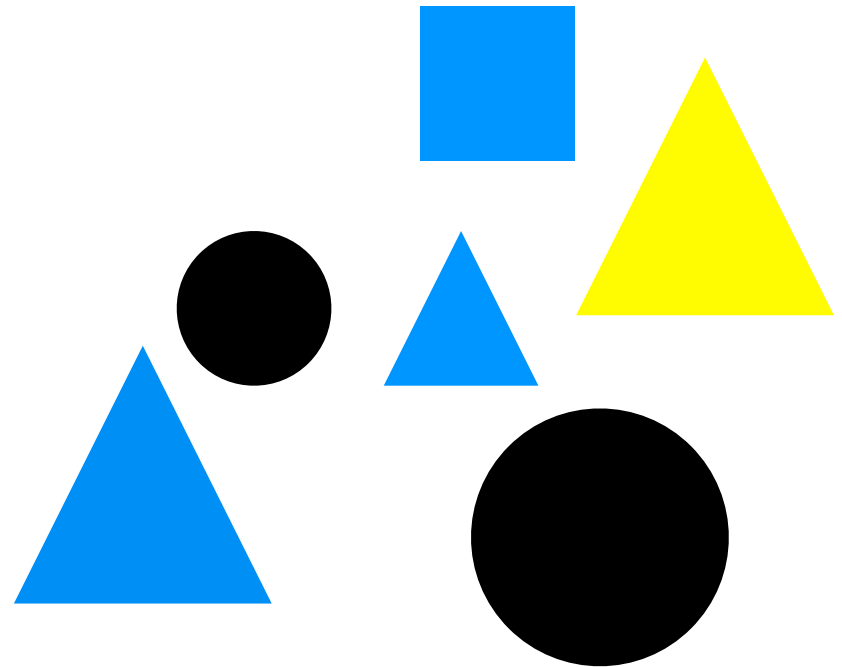
# Natural Language for Visual Reasoning



There is a box with 3 items of all 3 different colors.

**TRUE**

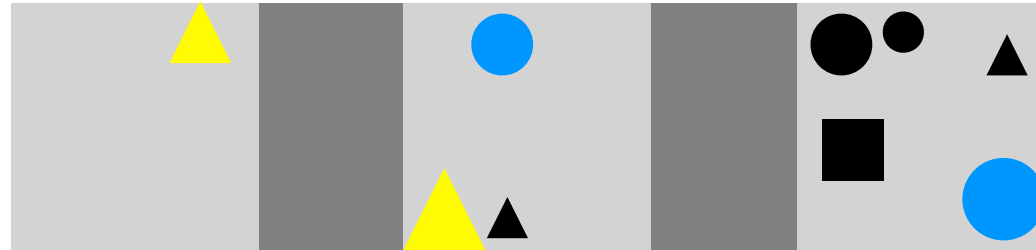**Task:** determine whether the statement is true or false for the image.

# Outline

- Task and environments

- Data collection
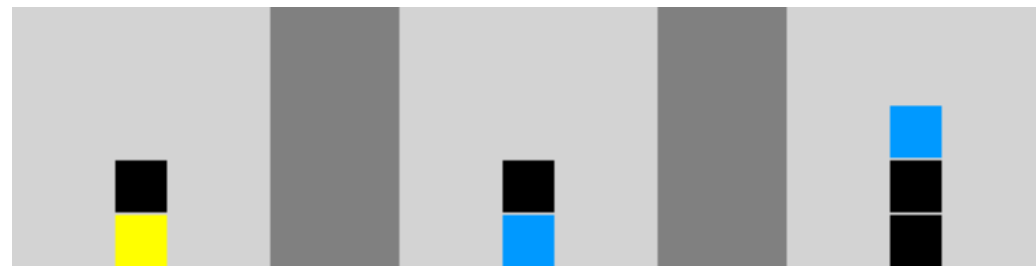
- Analysis

- Baselines

# Task and Environments

**Scatter**



There is a box with 3 items of all 3 different colors.

**TRUE**

**Tower**



There are only two towers which has the same base color.

**FALSE**

# Data collection

- **Goal:** collect natural language descriptions of images and true/false judgments

- Generate images

- Collect natural language sentences

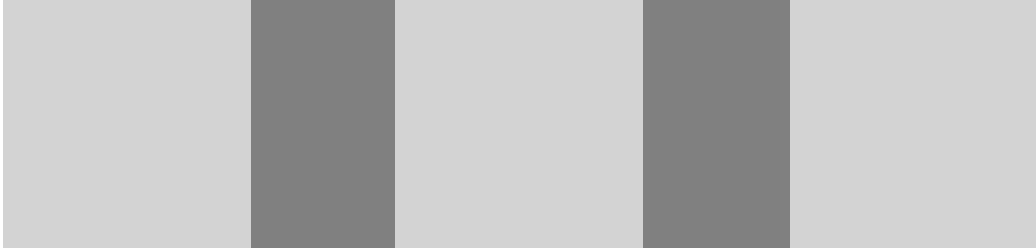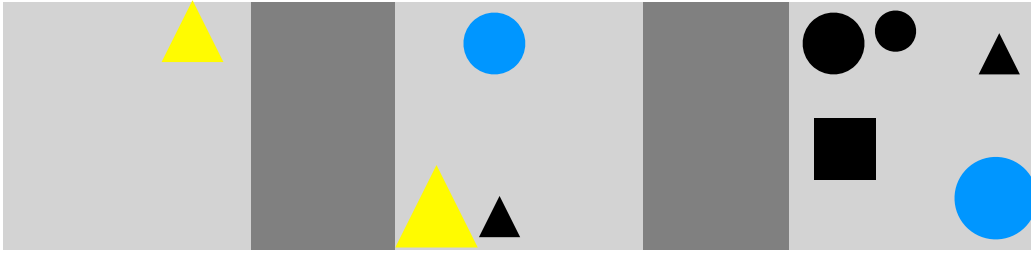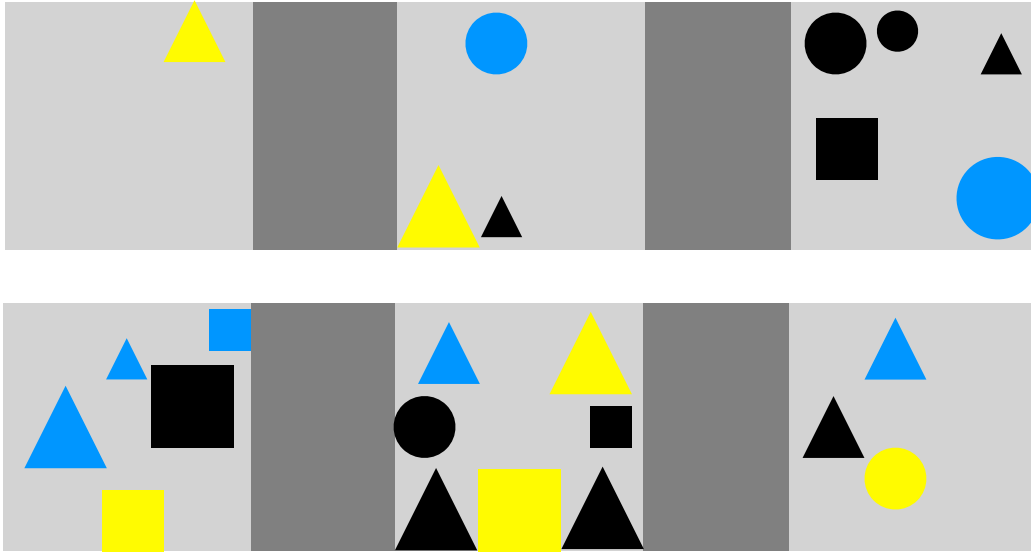- Validate image/sentence pairs

# Image Generation
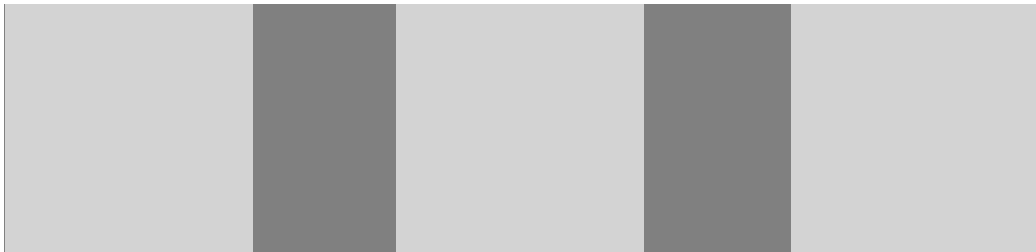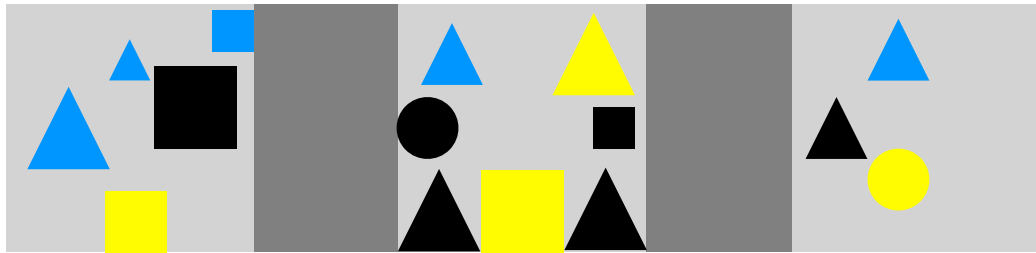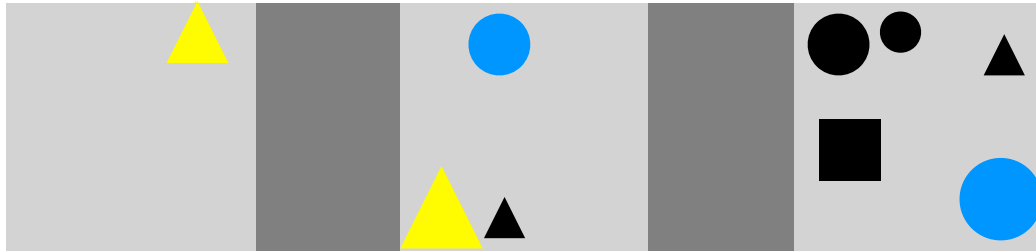
# Image Generation



- Randomly choose number of items per box and item shapes, colors, sizes, and positions (without overlap)

# Image Generation



- Randomly choose number of items per box and item shapes, colors, sizes, and positions (without overlap)
- Construct second image with the same type

# Image Generation


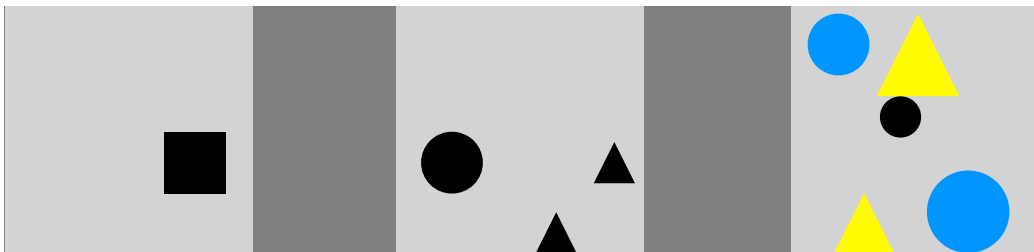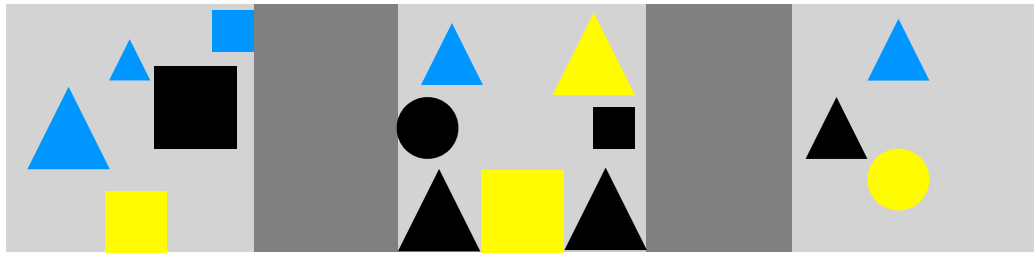
- Randomly choose number of items per box and item shapes, colors, sizes, and positions (without overlap)
- Construct second image with the same type

# Image Generation



- Randomly choose number of items per box and item shapes, colors, sizes, and positions (without overlap)
- Construct second image with the same type
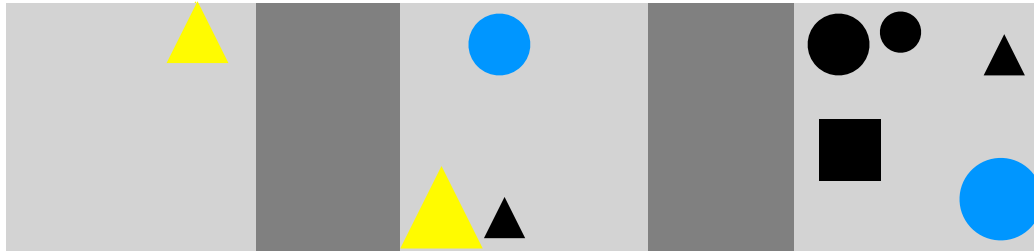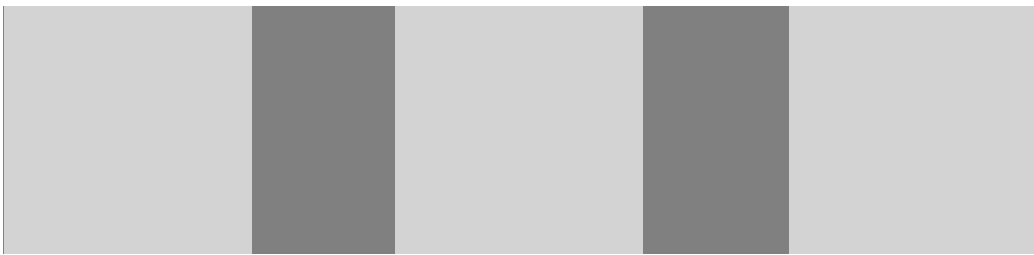- Construct third image by shuffling items in the first image
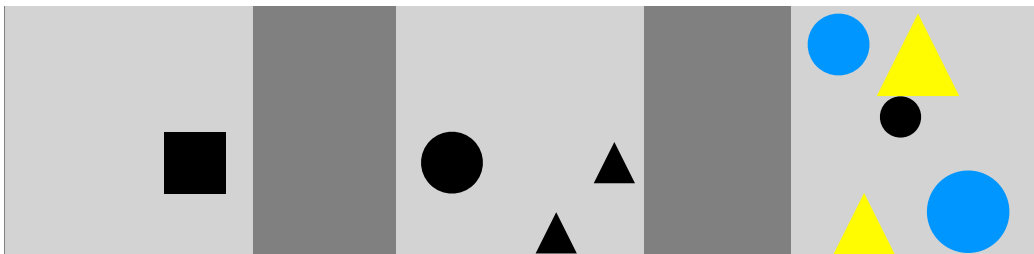
# Image Generation



- Randomly choose number of items per box and item shapes, colors, sizes, and positions (without overlap)
- Construct second image with the same type
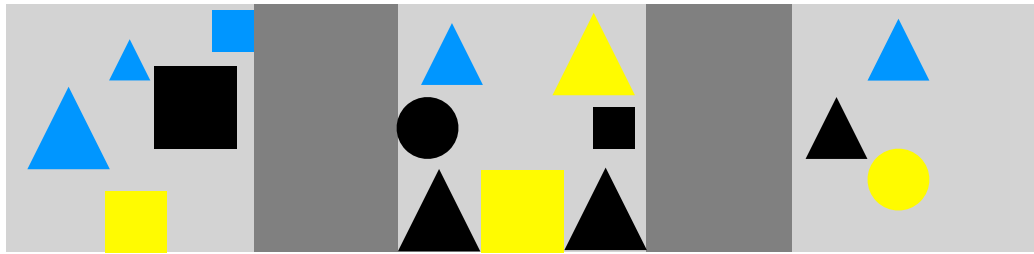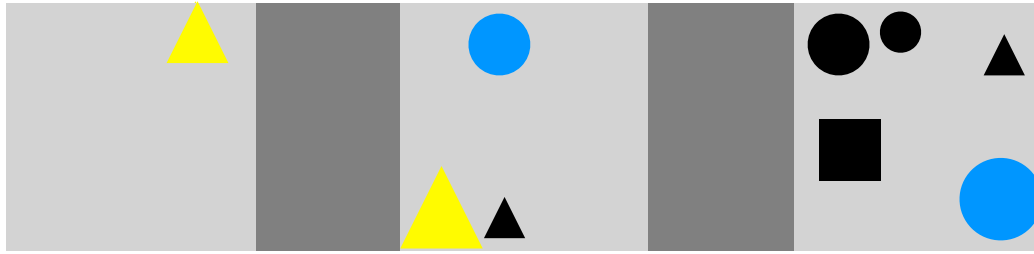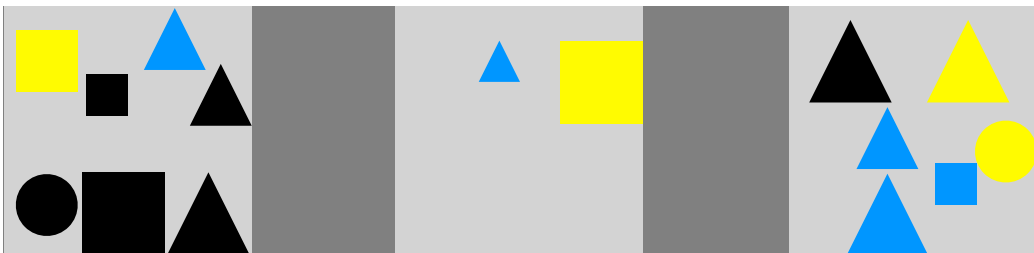- Construct third image by shuffling items in the first image
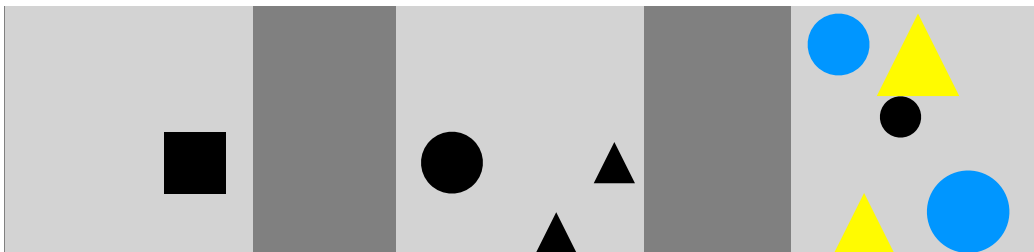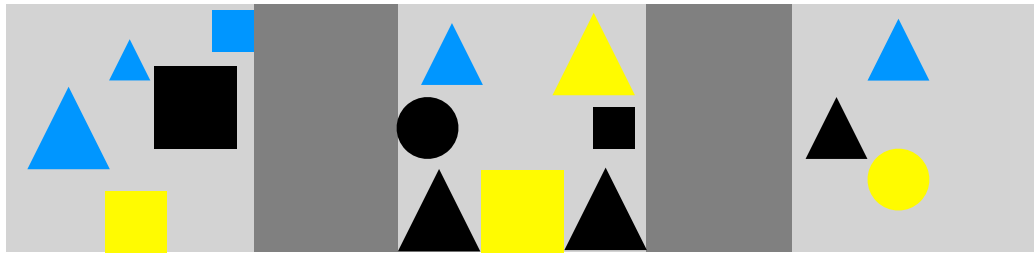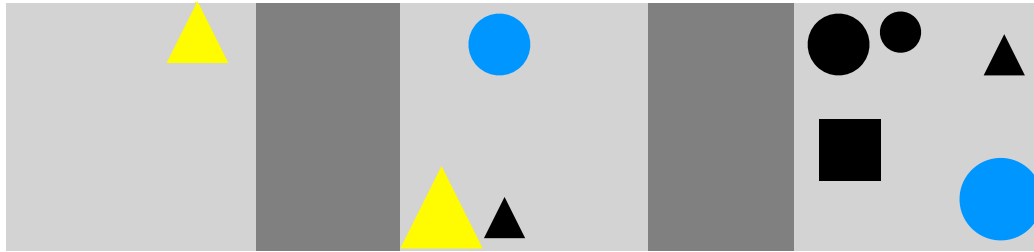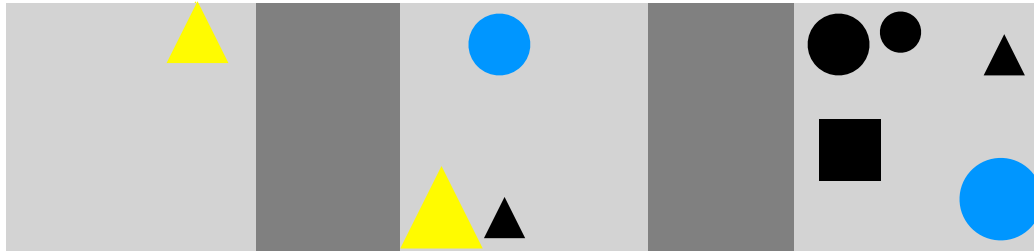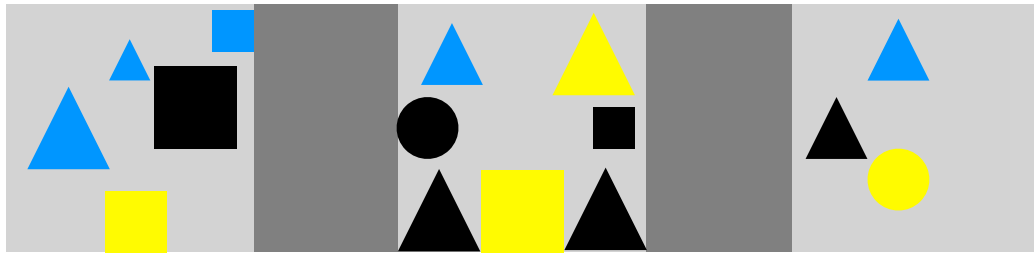
# Image Generation



- Randomly choose number of items per box and item shapes, colors, sizes, and positions (without overlap)
- Construct second image with the same type
- Construct third image by shuffling items in the first image
- Construct fourth image by shuffling items in the second image

**Generate two unique images and permute their items to create two other images**

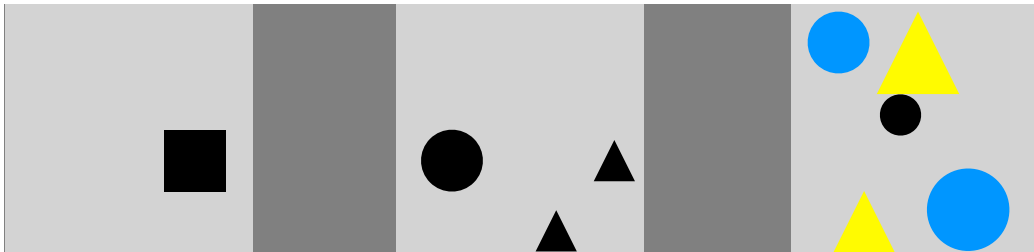# Sentence Writing

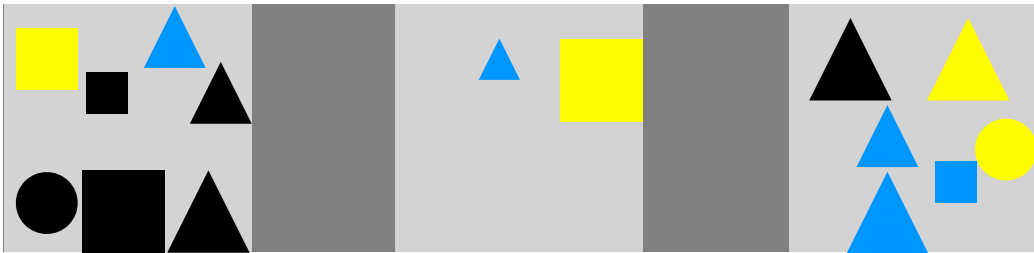**Write a sentence that is <span style="color:green">true</span> about the top two images and <span style="color:red">false</span> about the bottom two.**

- Don't refer to the order of the images.
- Don't refer to the order of the boxes.

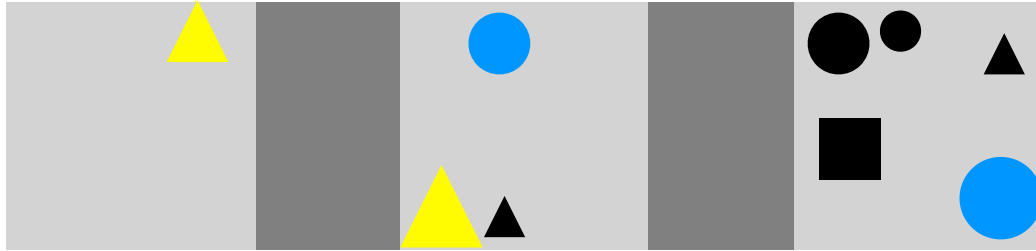There is a box with 3 items of all 3 different colors.

**Setup encourages set reasoning, counting, and comparisons**

# Sentence Writing



There is a box with 3 items of all 3 different colors.

**TRUE**

There is a box with 3 items of all 3 different colors.

**TRUE**

There is a box with 3 items of all 3 different colors.

**FALSE**

There is a box with 3 items of all 3 different colors.

**FALSE**

# Validation



There is a box with 3 items
of all 3 different colors.

- Higher-quality data
- Measure agreement
- Make sure sentences follow the guidelines

Fleiss' κ: **0.709 ➡ 0.808**

# Validation



There is a box with 3 items
of all 3 different colors.
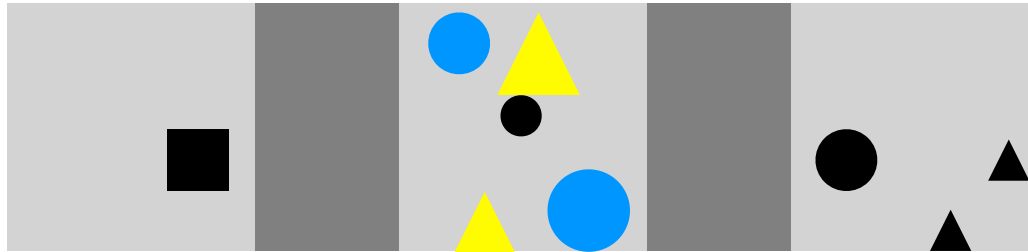
☐ **TRUE**

☑ **FALSE**

# Permutation



There is a box with 3 items
of all 3 different colors.

☐ **TRUE**

☑ **FALSE**

# Corpus Statistics

- 92,244 examples
- 3,962 unique sentences
- Krippendorff's α: 0.831
- Fleiss' κ: 0.808
  - (Landis and Koch, 1977)
- 262 words in the vocabulary
- Average sentence length of 11.2

- Four data splits
  - 80.7% training
  - 6.4% development
  - 6.4% public test
  - 6.4% unreleased test

lic.nlp.cornell.edu/nlvr

# Related Corpora

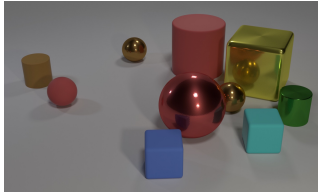| | Task | Examples | |
|---|---|---|---|
| **MSCOCO**<br>(Chen et al 2015) | **Caption generation** |  | A small herd of cows in a large grassy field. |
| **CLEVR**<br>(Johnson et al 2016) | **Question answering** |  | How many objects are either small cylinders or red things? |
| **VQA — real**<br>(Agrawal et al 2015) | **Question answering** |  | What is the dog carrying? |
| **VQA — abstract**<br>(Agrawal et al 2015) | **Question answering** |  | Is this a forest? |
| **NLVR**<br>(Suhr et al 2017) | **Binary classification** |  | there are exactly three blue objects not touching any edge |

# Related Corpora

| | Task | Real images? | Natural language? |
|---|---|---|---|
| **MSCOCO** (Chen et al 2015) | Caption generation | ✔ | ✔ |
| **CLEVR** (Johnson et al 2016) | Question answering | ✘ | ✘ |
| **VQA — real** (Agrawal et al 2015) | Question answering | ✔ | ✔ |
| **VQA — abstract** (Agrawal et al 2015) | Question answering | ✘ | ✔ |
| **NLVR** (Suhr et al 2017) | Binary classification | ✘ | ✔ |

# Lengths



Legend:
- **NLVR (ours)** (blue)
- VQA real images (green)
- VQA abstract images (yellow)
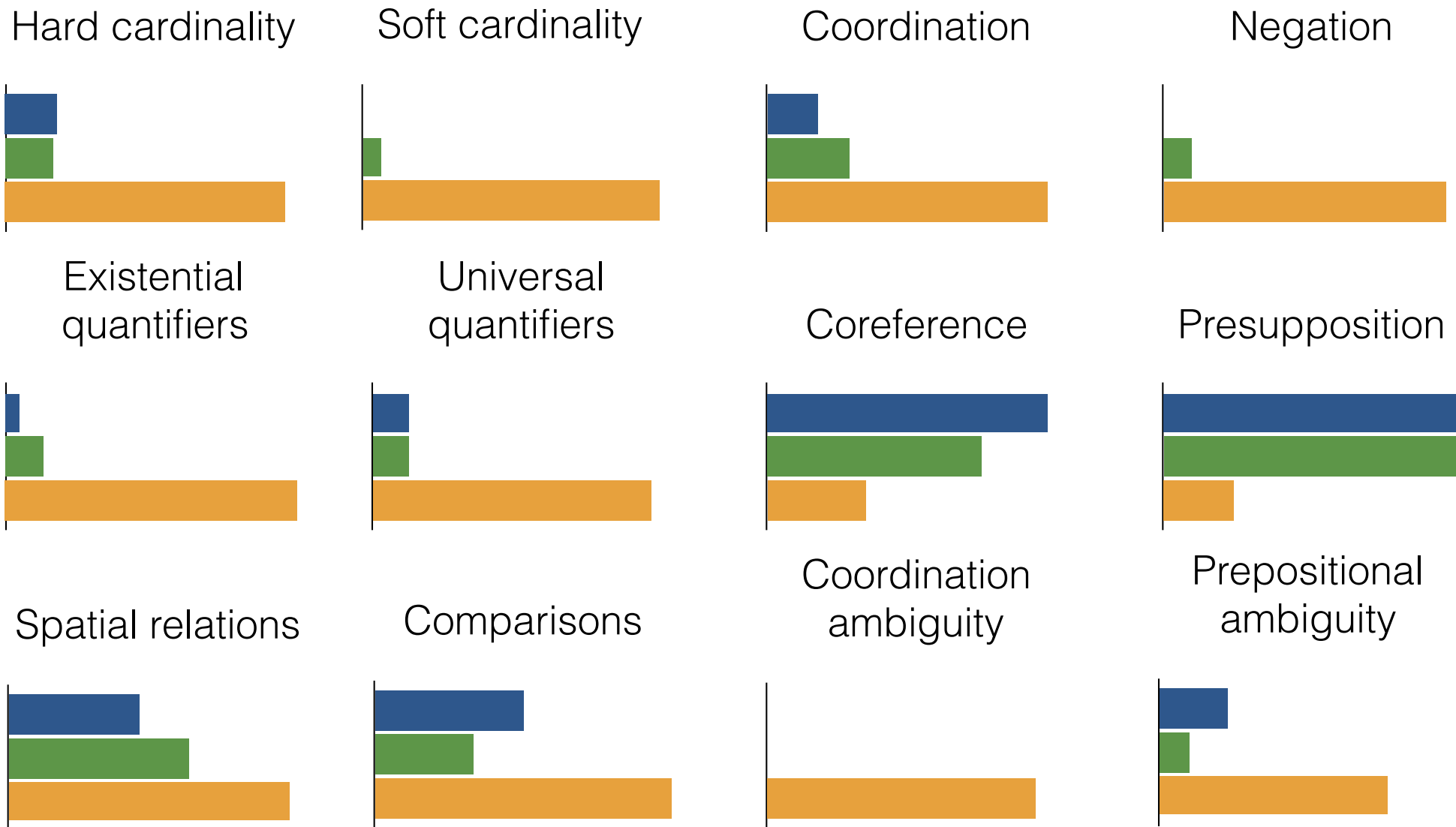- MSCOCO (orange)
- CLEVR (red)

**Longer than VQA
Similar to MS COCO**

# Linguistic Analysis

## Analyzed 200 random development sentences.



Legend: VQA (abstract), VQA (real), NLVR

Charts: Hard cardinality, Soft cardinality, Coordination, Negation, Existential quantifiers, Universal quantifiers, Coreference, Presupposition, Spatial relations, Comparisons, Coordination ambiguity, Prepositional ambiguity

# Numerical Expressions



Hard cardinality

66%

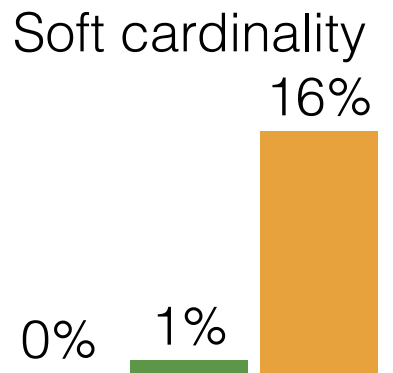12% 12%

Soft cardinality

16%

0% 1%

- VQA (abstract)
- VQA (real)
- NLVR

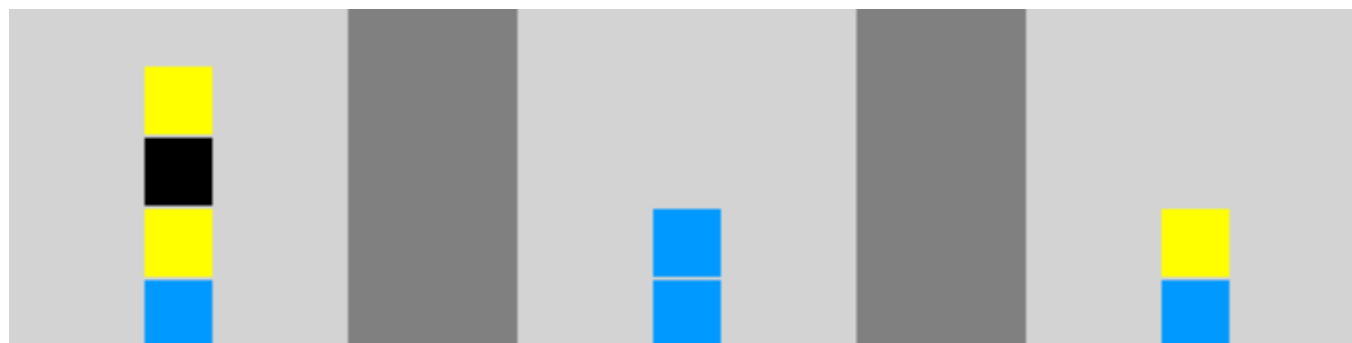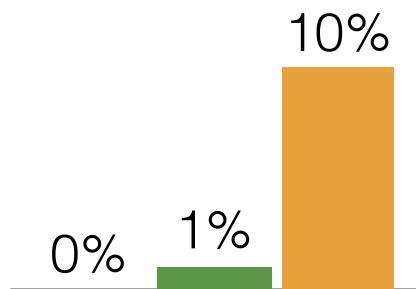There is a tower with **exactly three** blocks, and it has a yellow block and **two** blue blocks.

**TRUE**

there are **at least two** yellow squares not touching any edge

**TRUE**

# Negation and Coordination
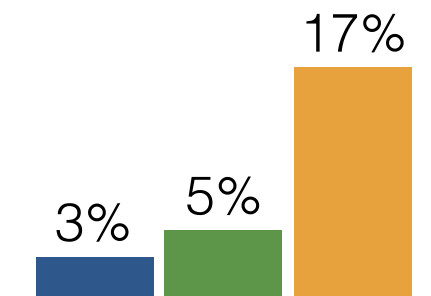


Negation

10%

1%

0%

There is a box with a black item between 2 items of the same color and **no item on top of that.**
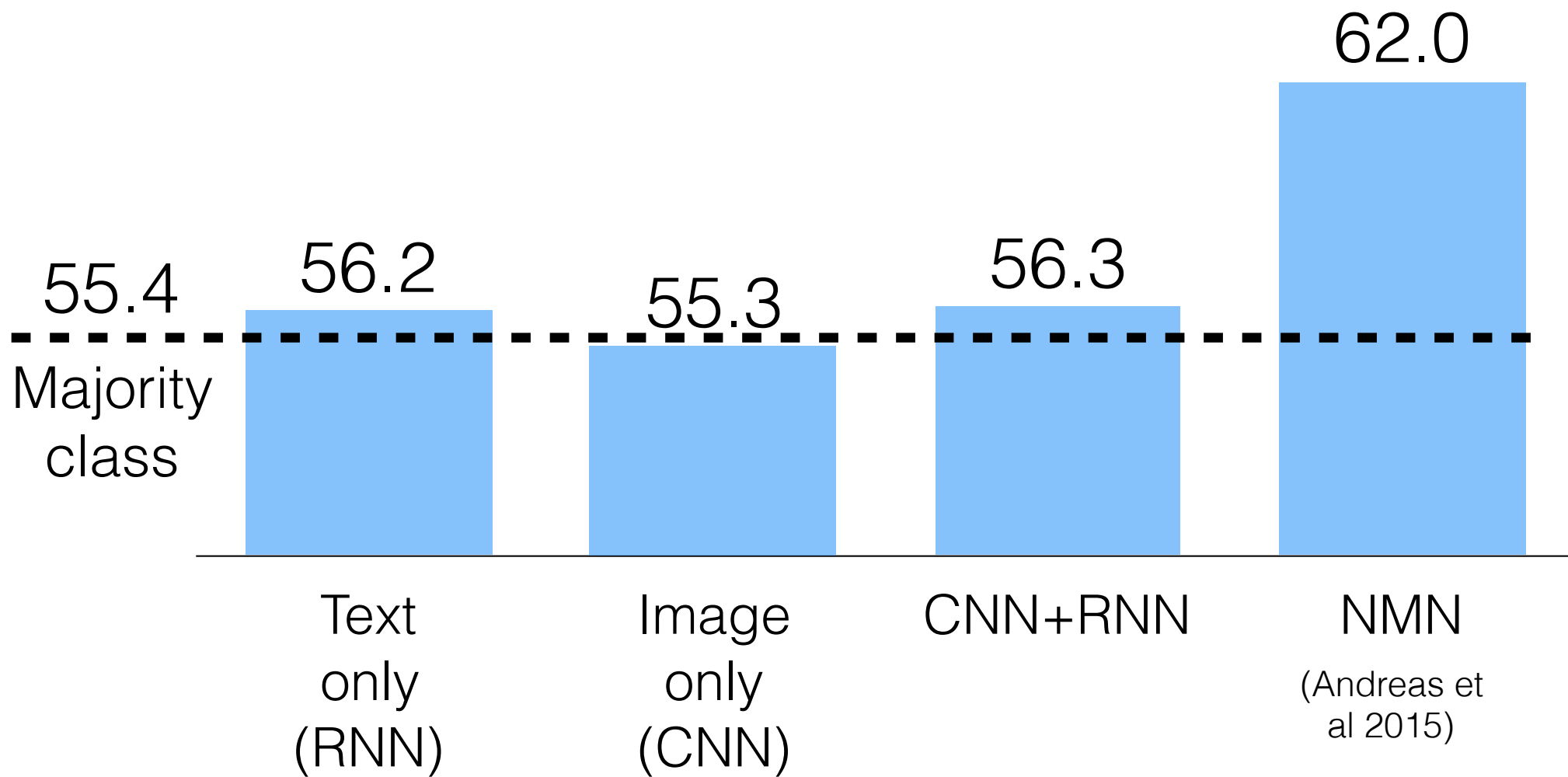
TRUE

Coordination

17%

5%

3%

There is a box with a yellow item **and** three black items.
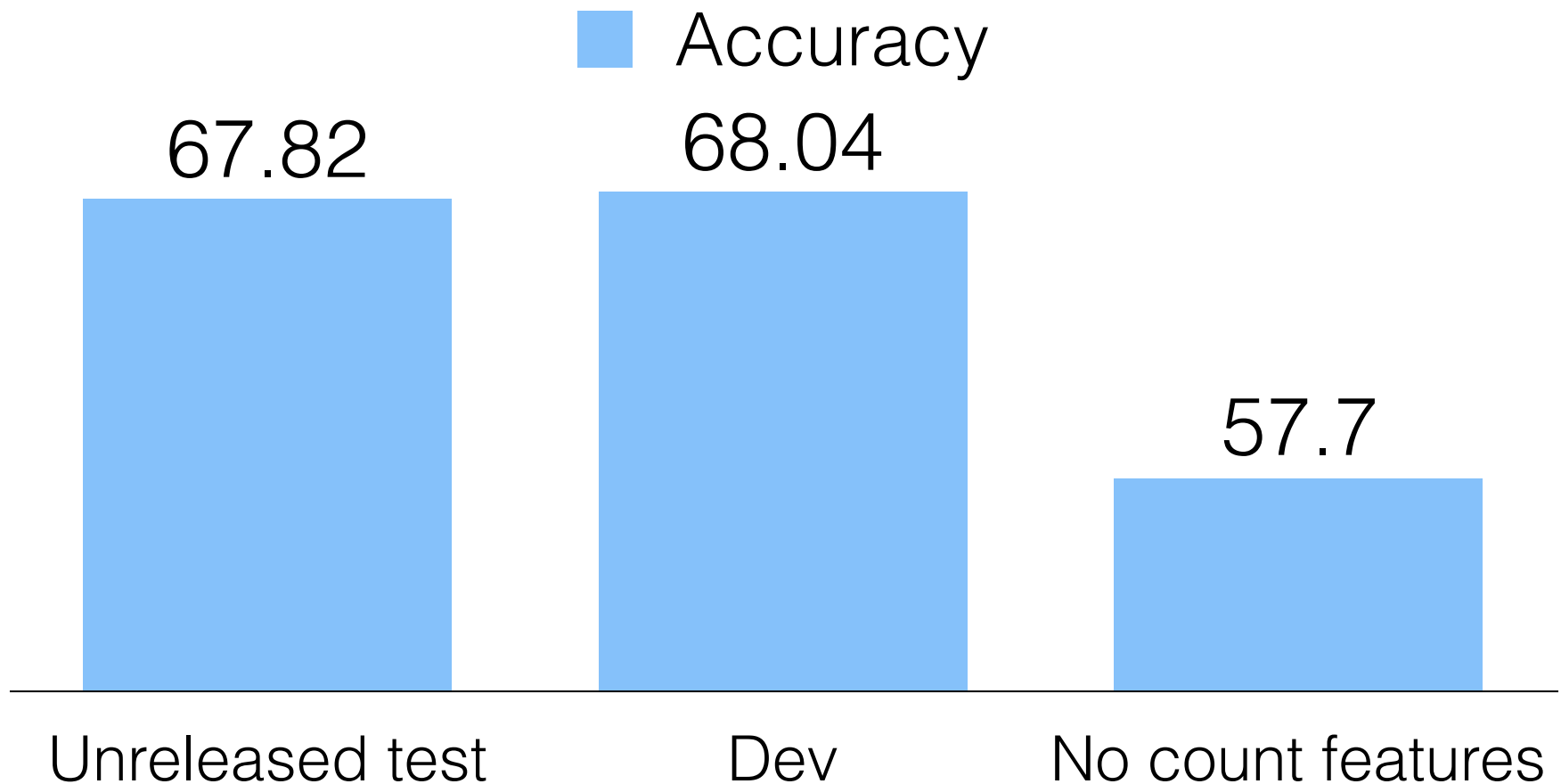
TRUE
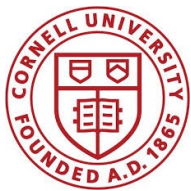
■ VQA (abstract)
■ VQA (real)
■ NLVR

# Baselines



Accuracy on unreleased test set

| | | | |
|---|---|---|---|
| 55.4 Majority class | 56.2 | 55.3 | 56.3 | 62.0 |

Text only (RNN)     Image only (CNN)     CNN+RNN     NMN (Andreas et al 2015)

# Feature-based Analysis

- Features text and structured representation
- Use maximum entropy model



Accuracy

67.82    68.04

57.7

Unreleased test    Dev    No count features

http://lic.nlp.cornell.edu/nlvr/

**Thank you!**