# Scene Understanding

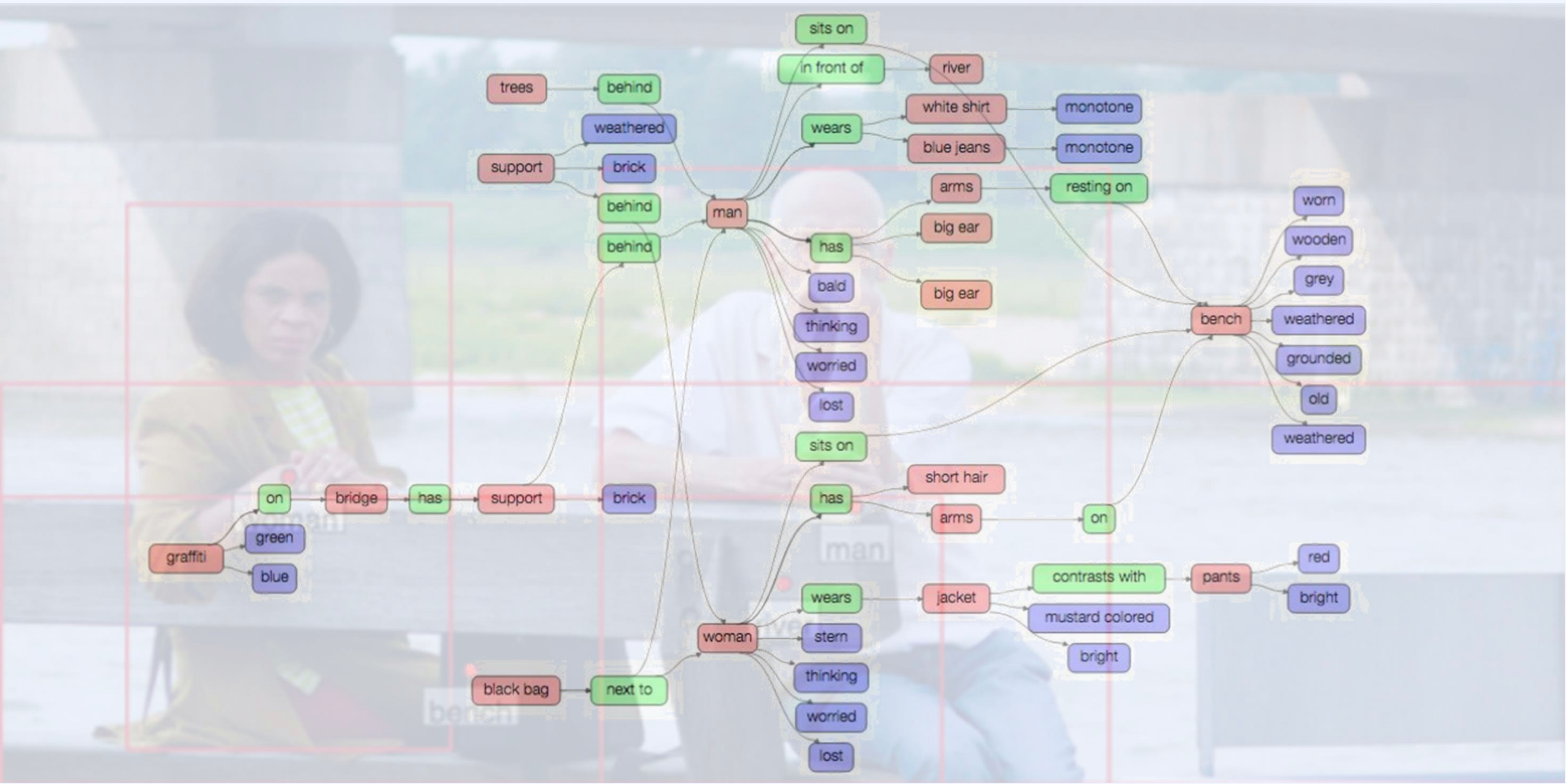**2** Scene Representation Learning (data)

# Scene Understanding
# Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv

# Scene Understanding
## Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv



- **Region descriptions**

- Objects

- Attributes

- Relationships
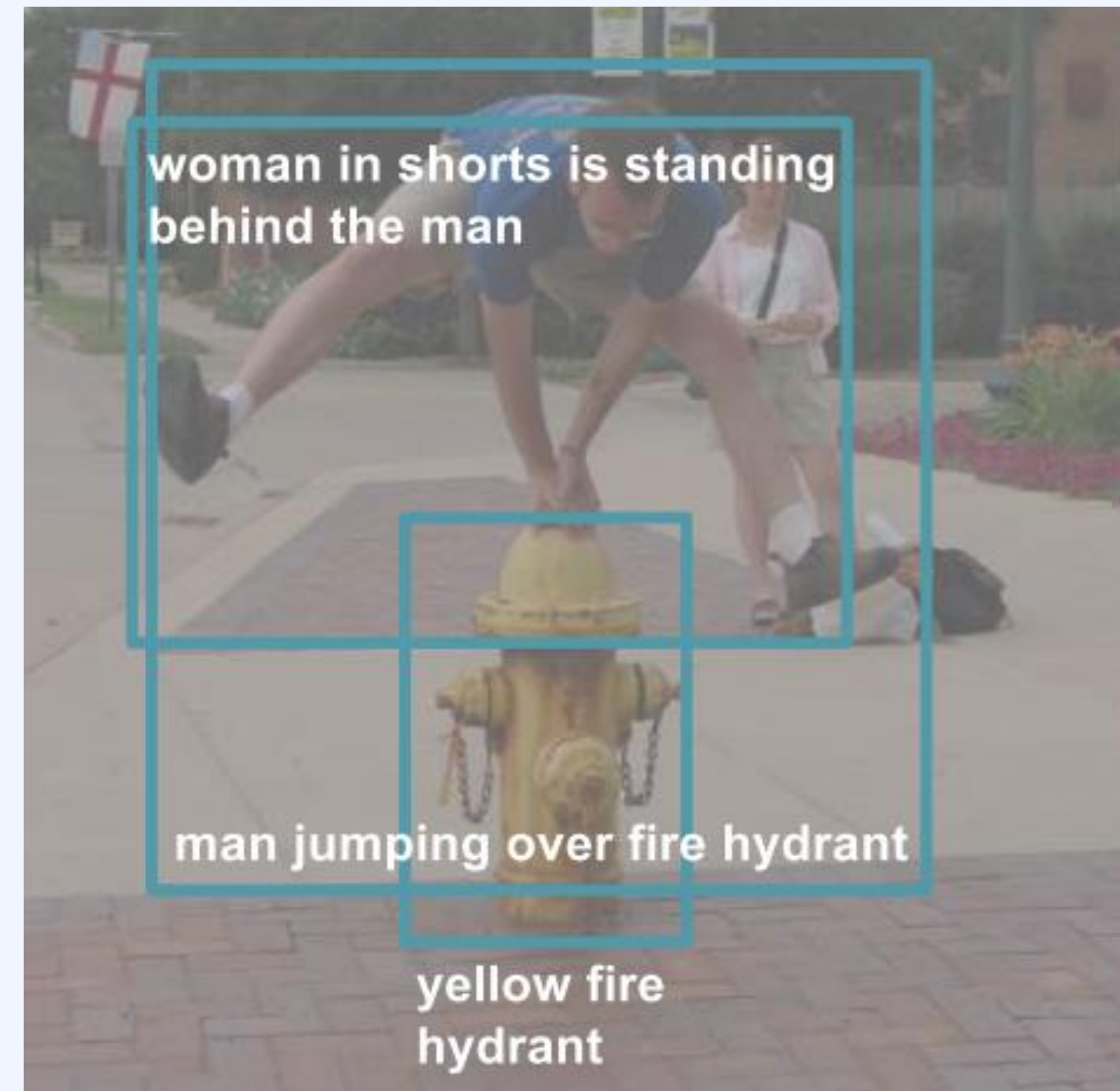
- Region graphs

- Scene graphs

- Question answer pairs

# Scene Understanding
## Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv



- Region descriptions
- **Objects**
- Attributes
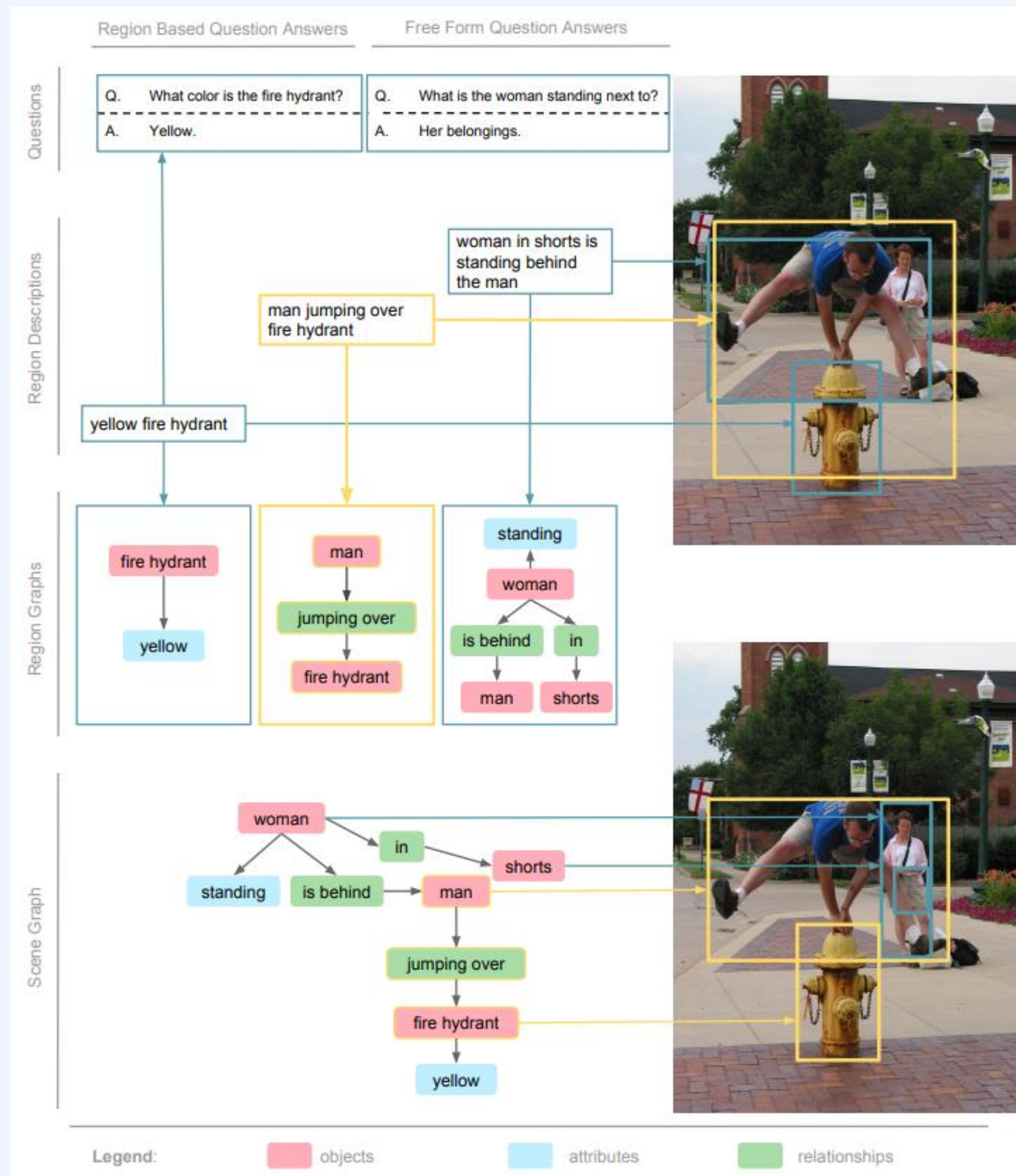- Relationships
- Region graphs
- Scene graphs
- Question answer pairs

# Scene Understanding
## Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv



- Region descriptions

- Objects

- **Attributes**

- Relationships

- Region graphs
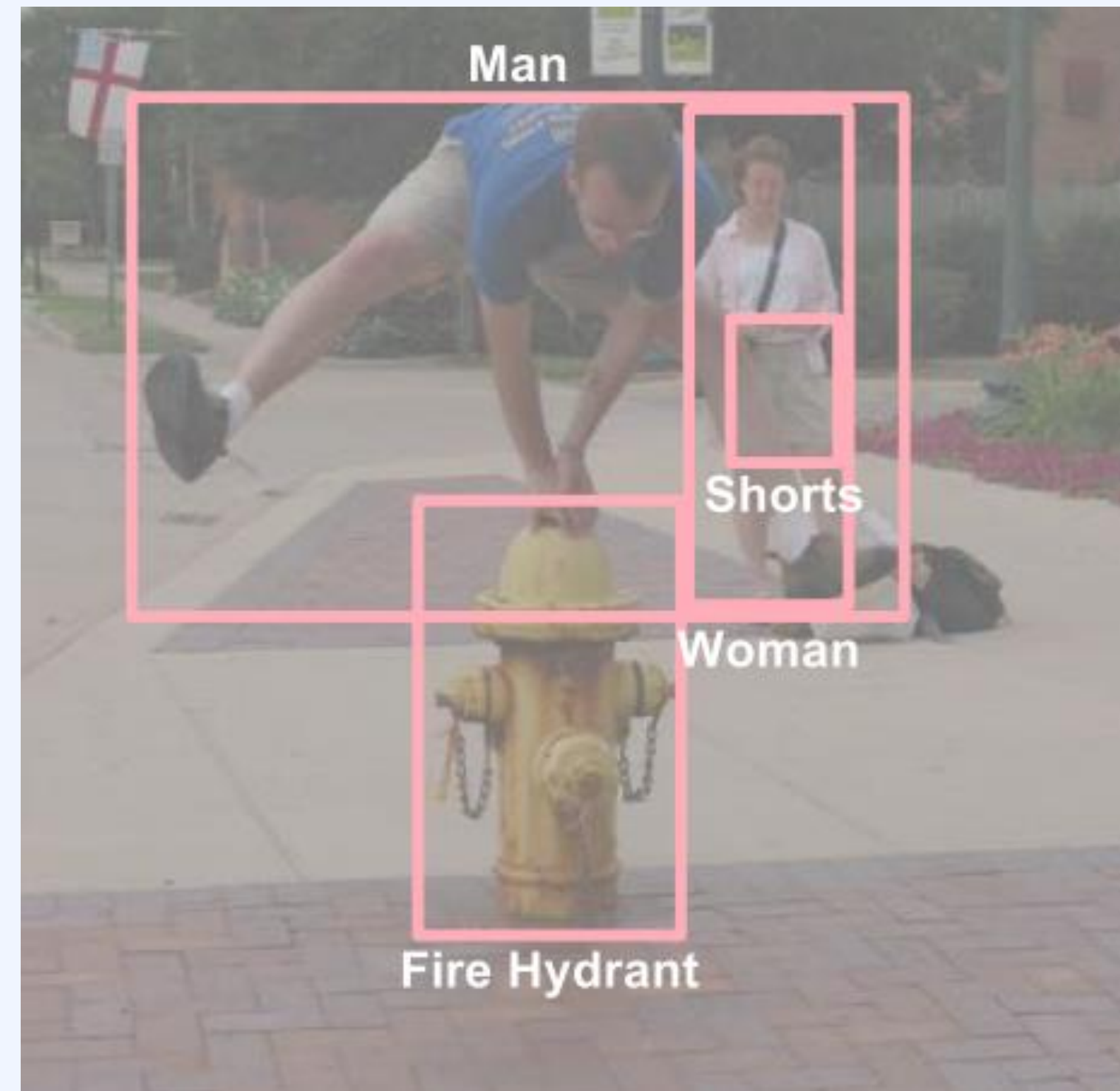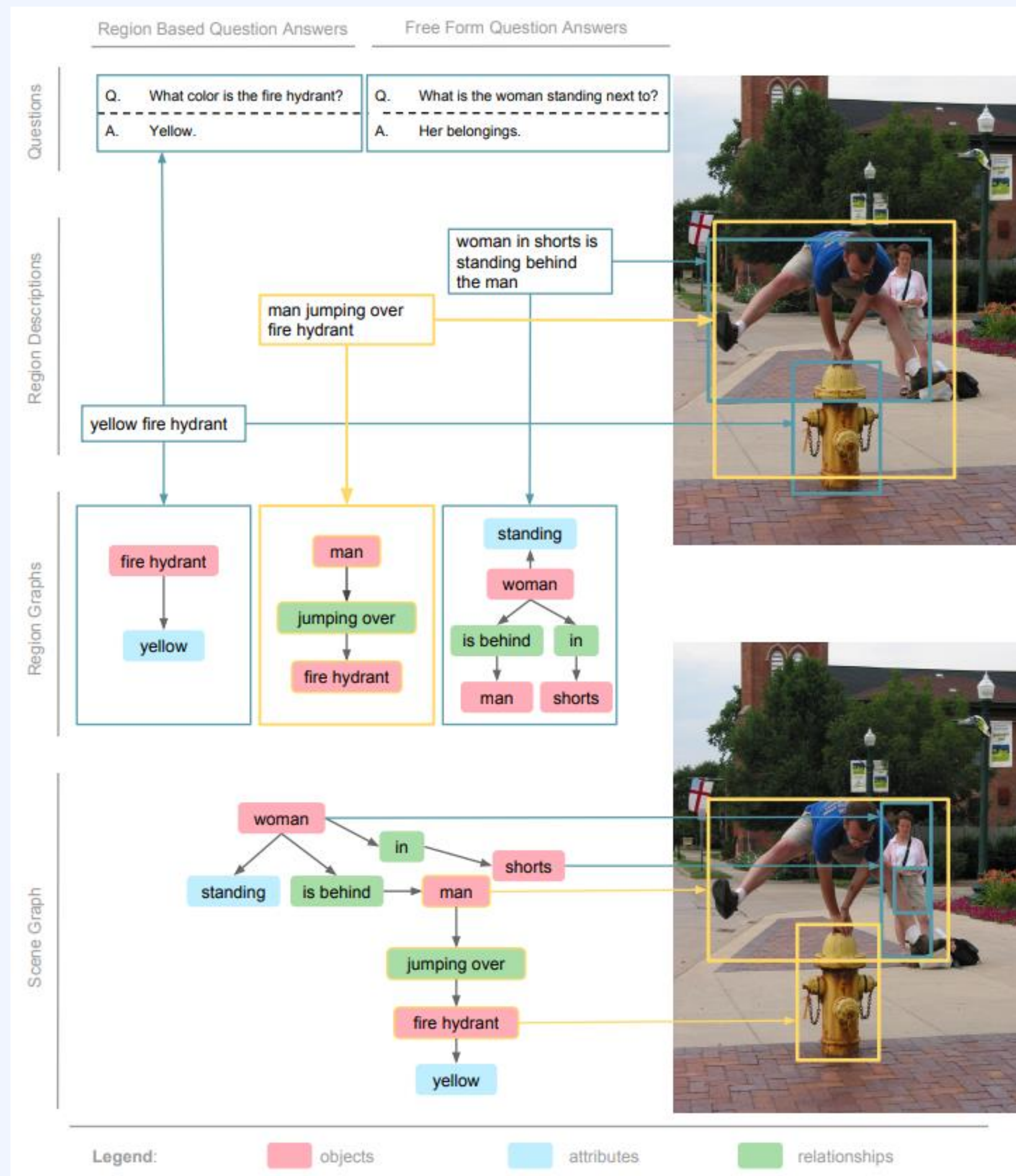
- Scene graphs

- Question answer pairs

# Scene Understanding
## Scene representation learning (data)

**2**.
Scene data

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv



- Region descriptions

- Objects

- Attributes

- **Relationships**

- Region graphs

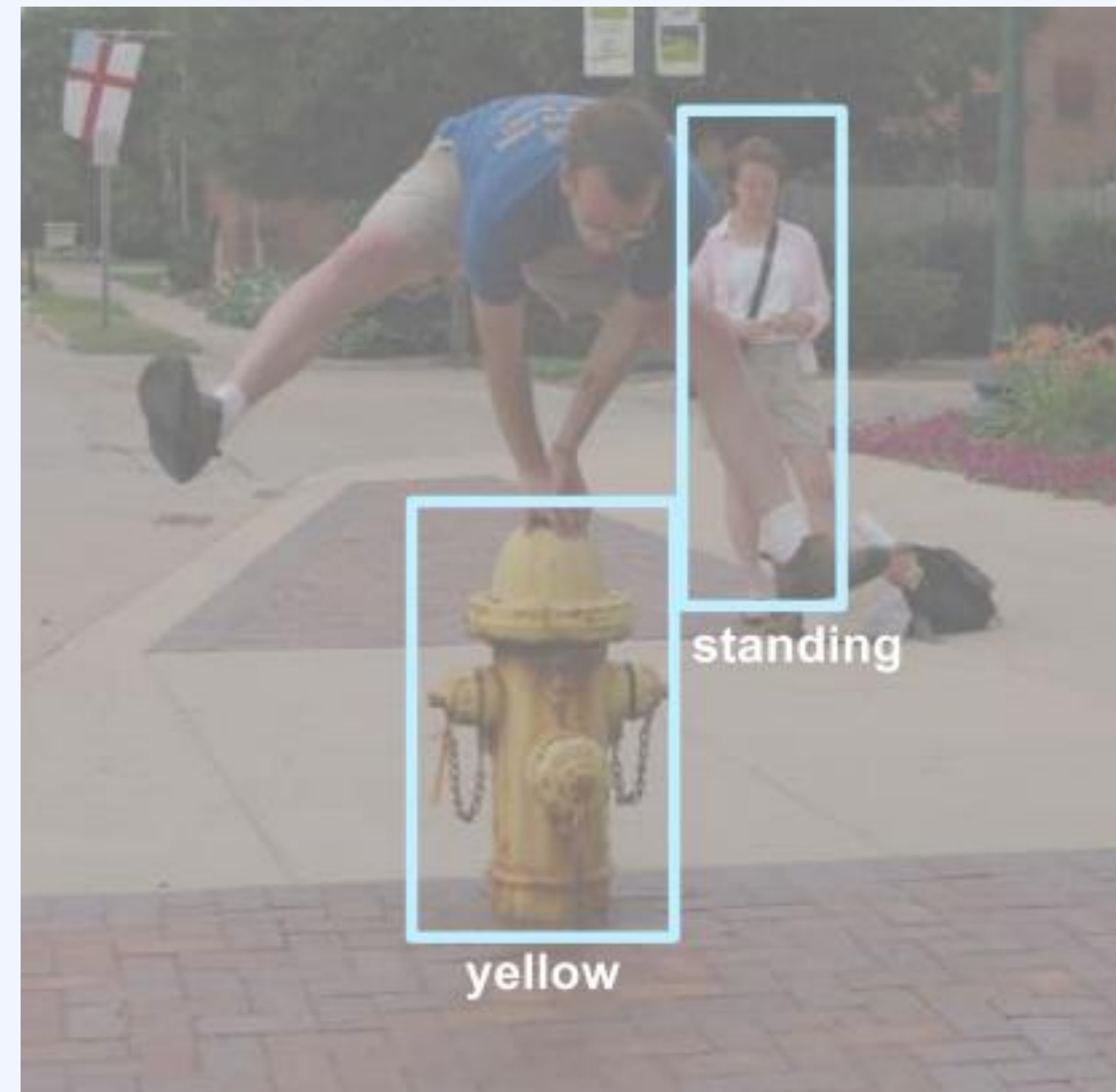- Scene graphs

- Question answer pairs

# Scene Understanding
## Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv



- Region descriptions

- Objects

- Attributes

- Relationships

- **Region graphs**
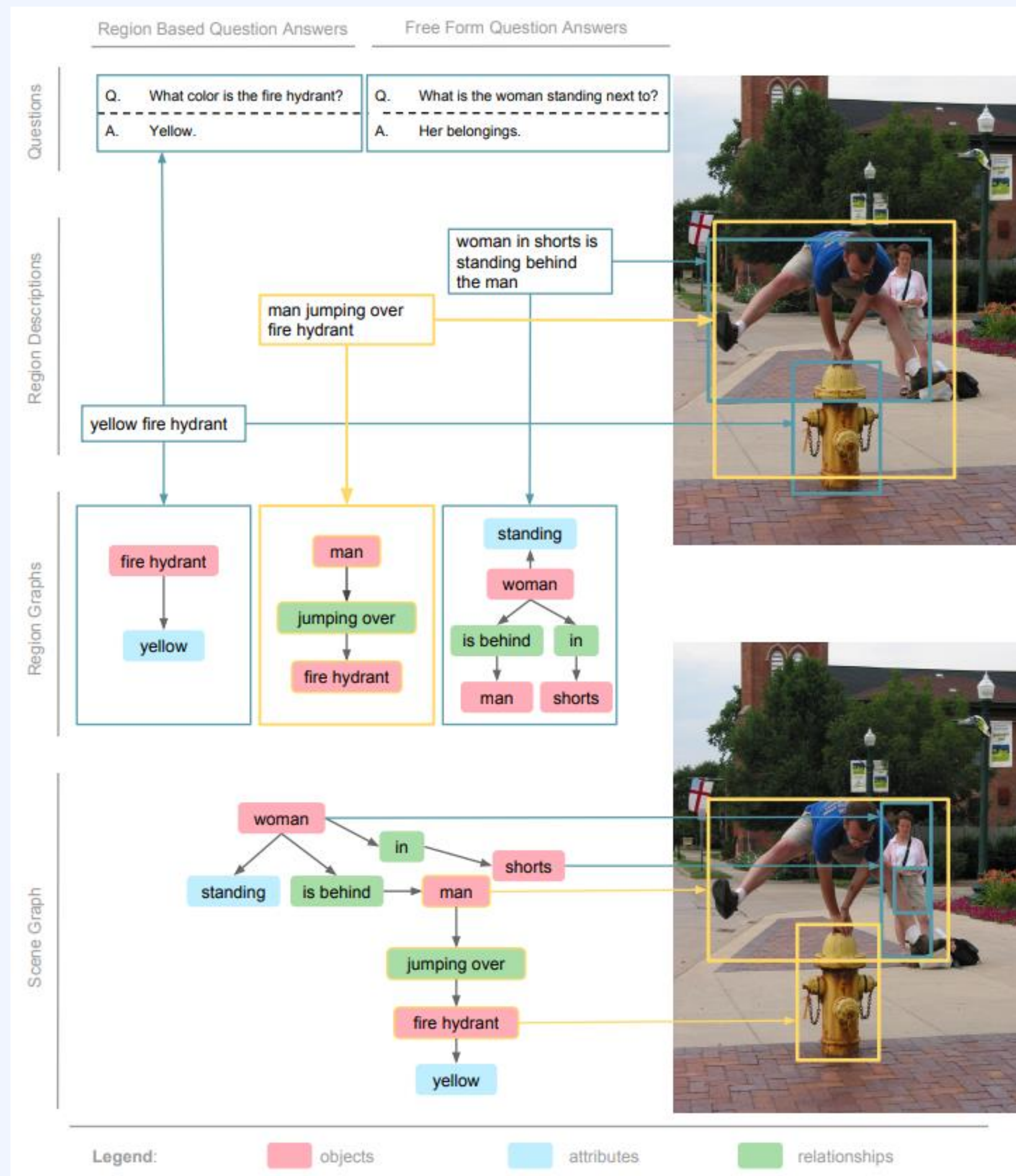
- Scene graphs

- Question answer pairs

# Scene Understanding
# Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv



- Region descriptions

- Objects

- Attributes

- Relationships

- Region graphs
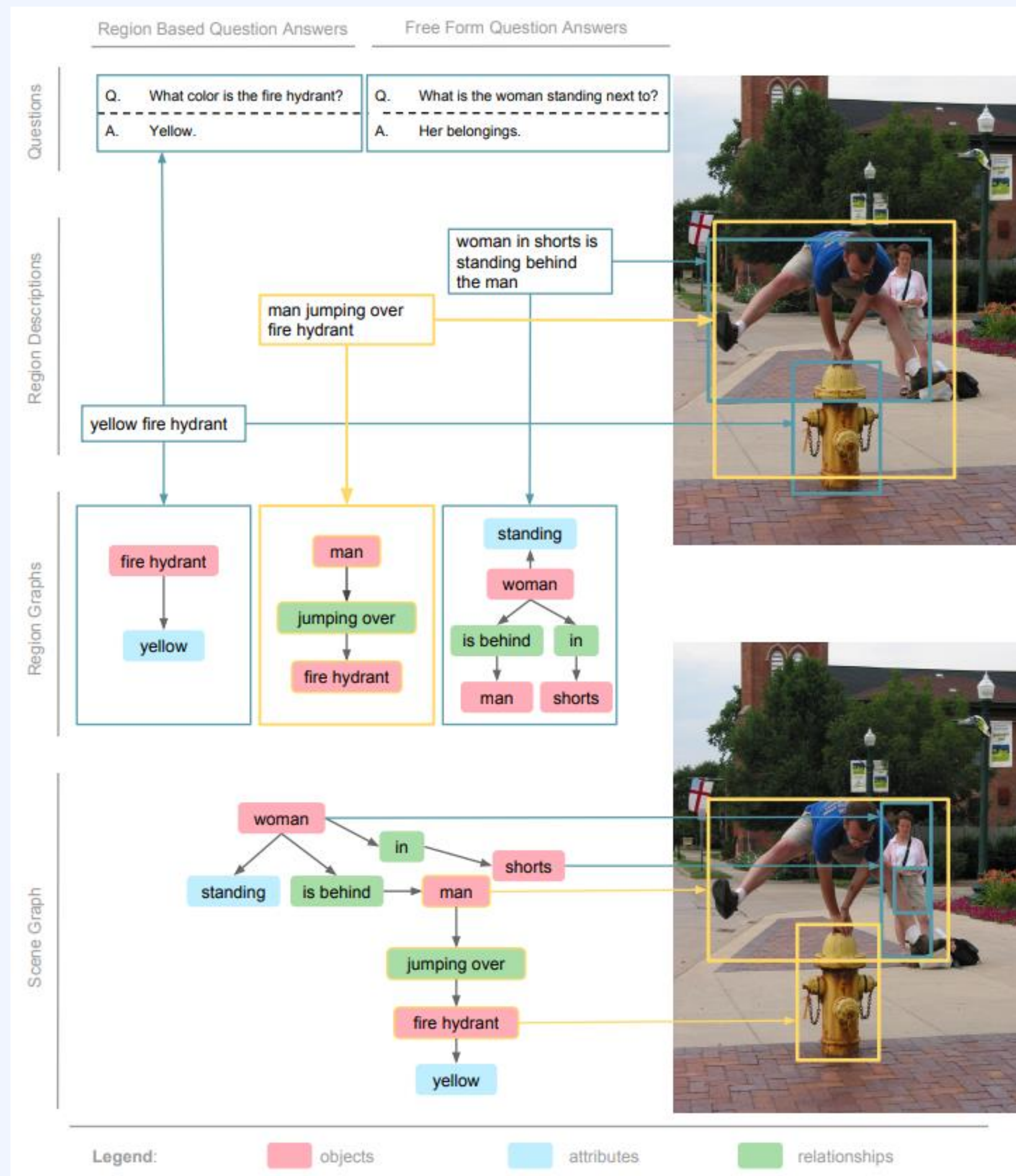
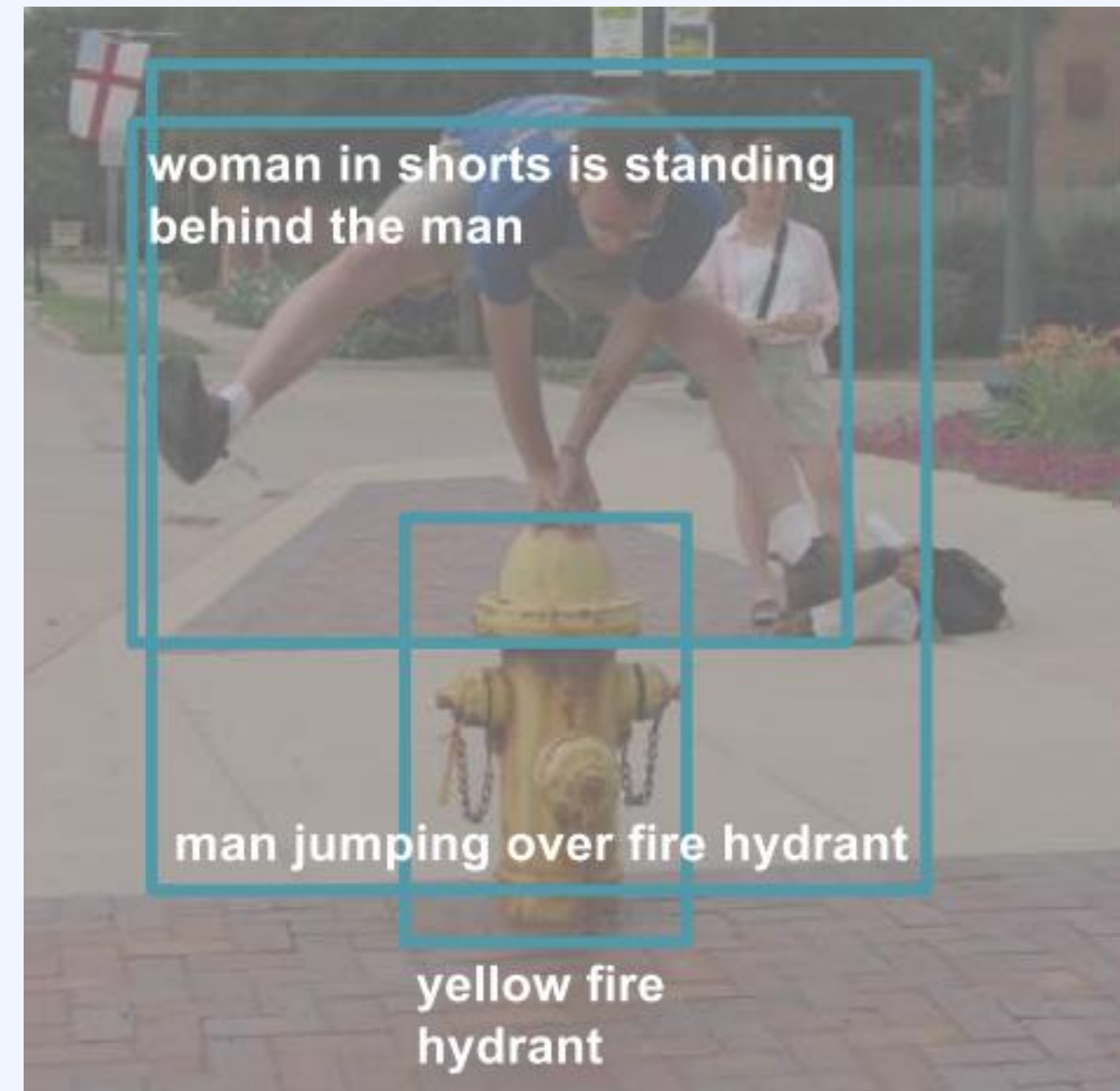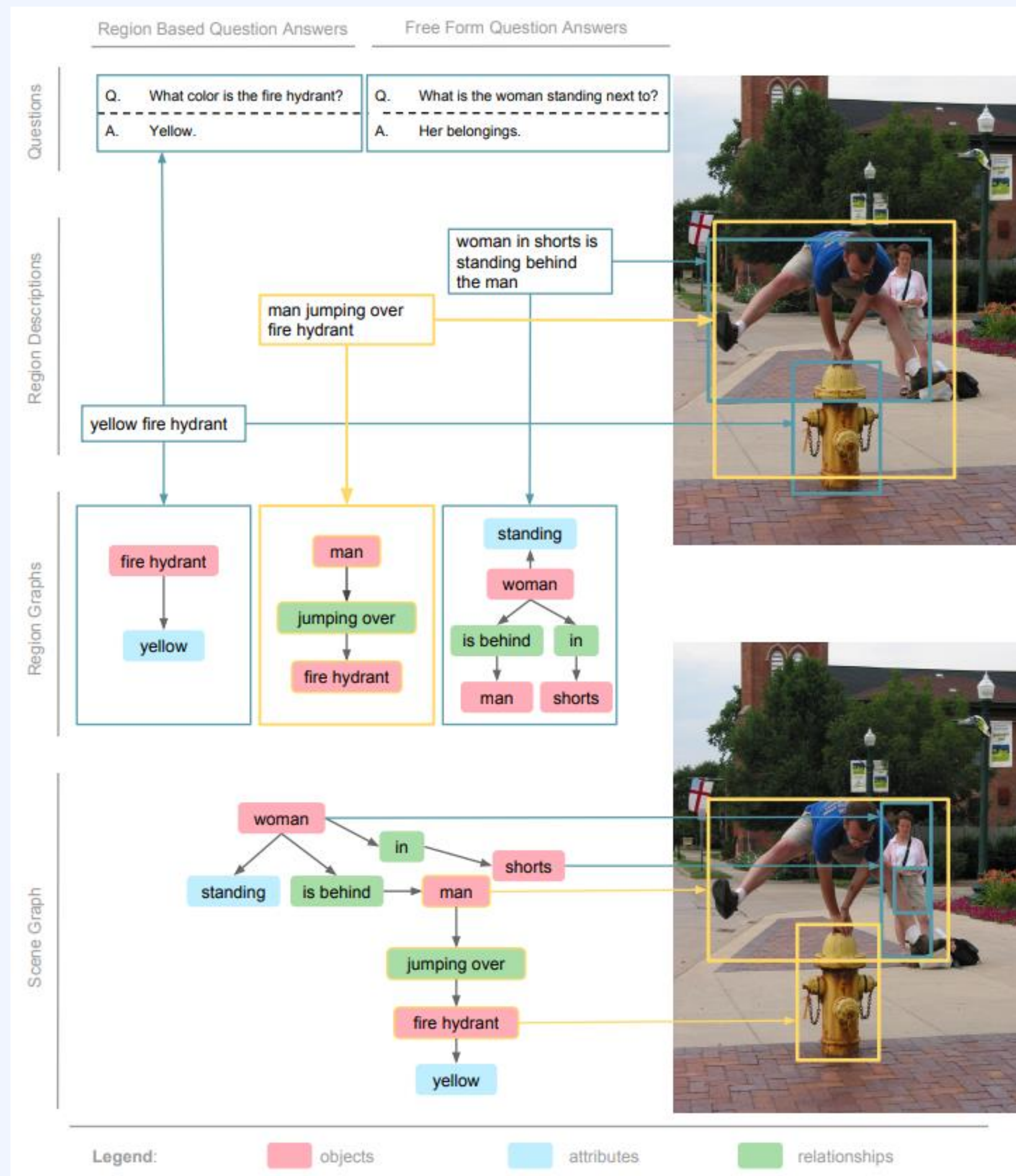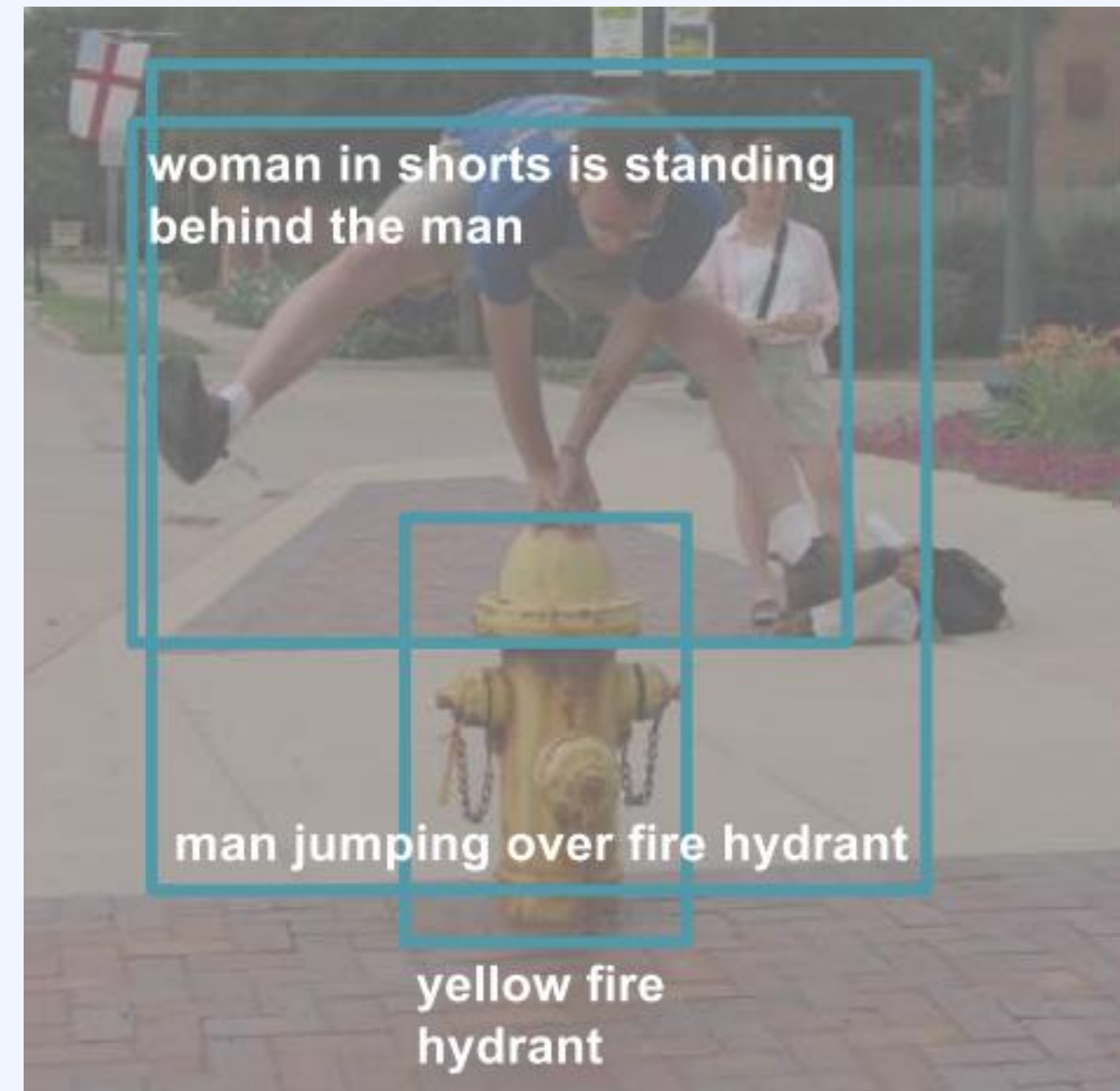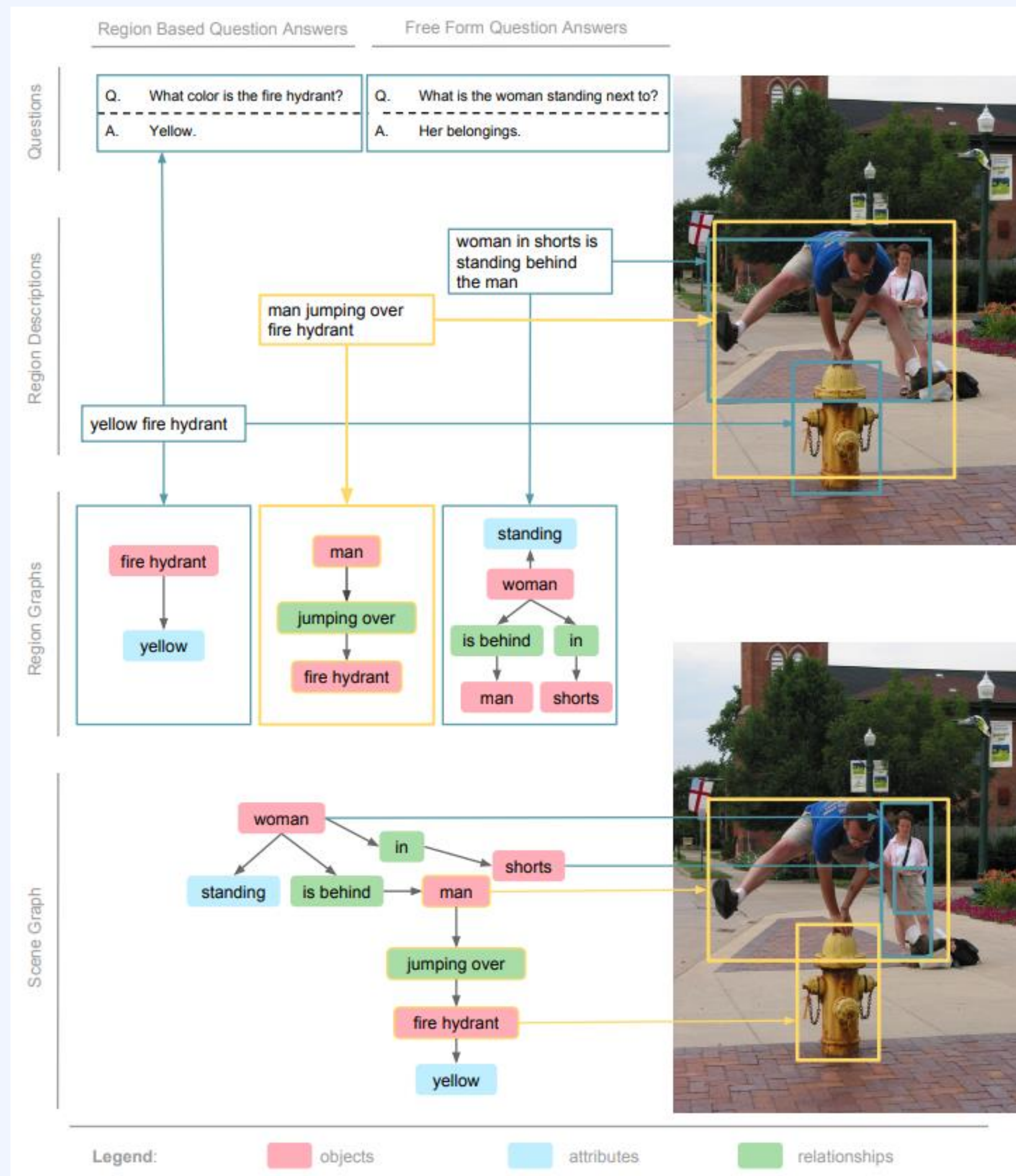- **Scene graphs**

- Question answer pairs

# Scene Understanding
# Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv



## Basic Statistics

- Total images: 108,077
- Average image size (px): 500.14
- Maximum image size (px): 1,280
- Minimum image size (px): 72

- Total region descriptions: 4,297,502
- Total image object instances: 1,366,673
- Unique image objects: 75,729

- Total object-object relationship instances: 1,531,448
- Unique relationships: 40,480
- Total attribute-object instances: 1,670,182
- Unique attributes: 40,513

- Total Scene Graphs: 108,249
- Total Region Graphs: 3,788,715

- Total Question Answers: 1,773,258

# Scene Understanding
## Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv

|  | Average number of objects | Average number of relationships | Average number of attributes |
|---|---|---|---|
| **Per region annotation** | 1.01 | 0.63 | 0.77 |
| **Per image** | 21.24 | 17.68 | 16.08 |

| Most common objects | Most common predicates | Most common attributes |
|---|---|---|
| •man<br>•person<br>•woman<br>•building<br>•sign<br>•table<br>•bus<br>•window<br>•sky<br>•tree | •ON<br>•has<br>•IN<br>•WEARING<br>•wears<br>•behind<br>•next to<br>•with<br>•near<br>•in front of | •white<br>•blue<br>•red<br>•black<br>•green<br>•yellow<br>•brown<br>•large<br>•on<br>•parked |

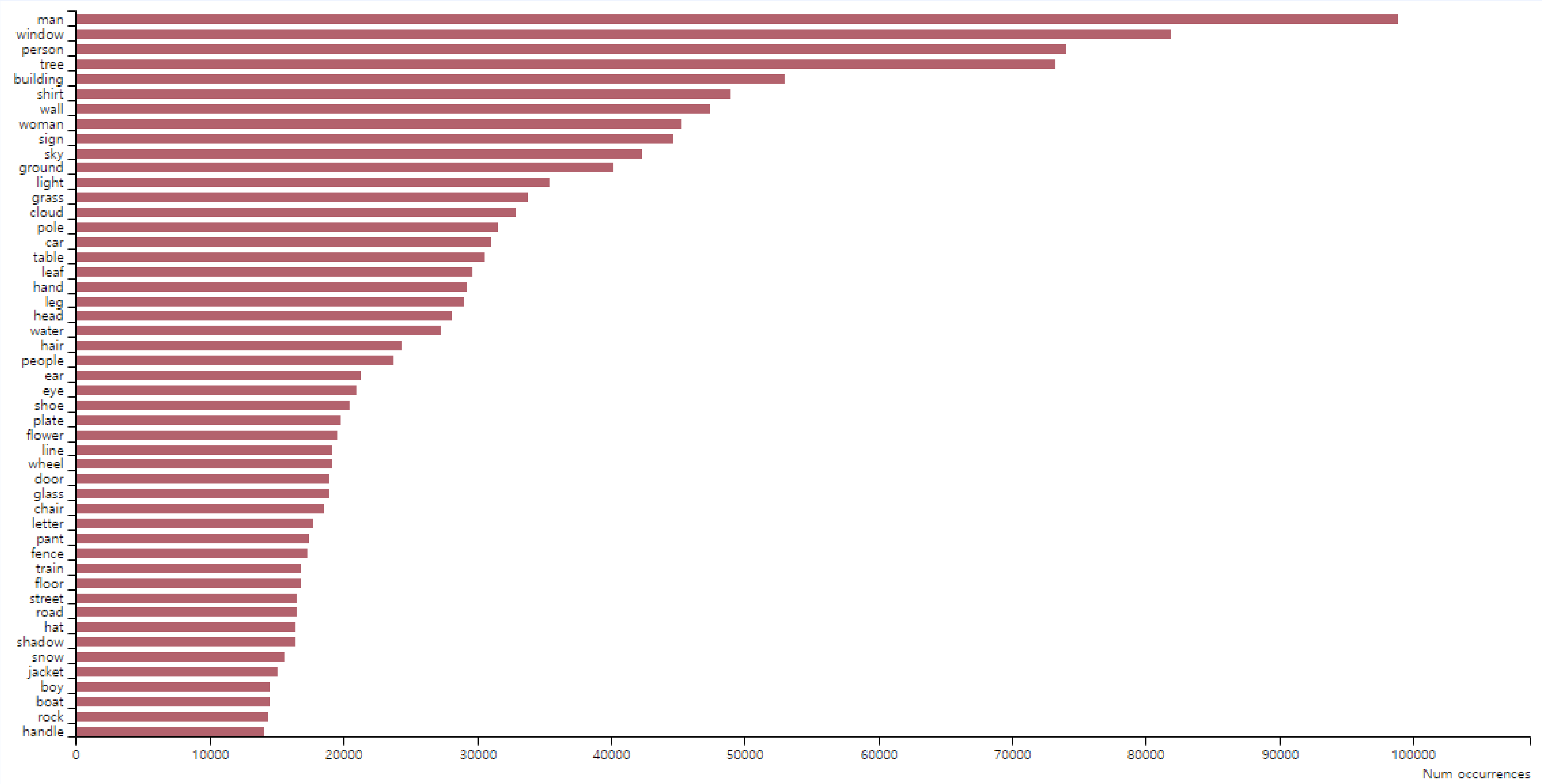# Scene Understanding
# Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv

## Top Object Names

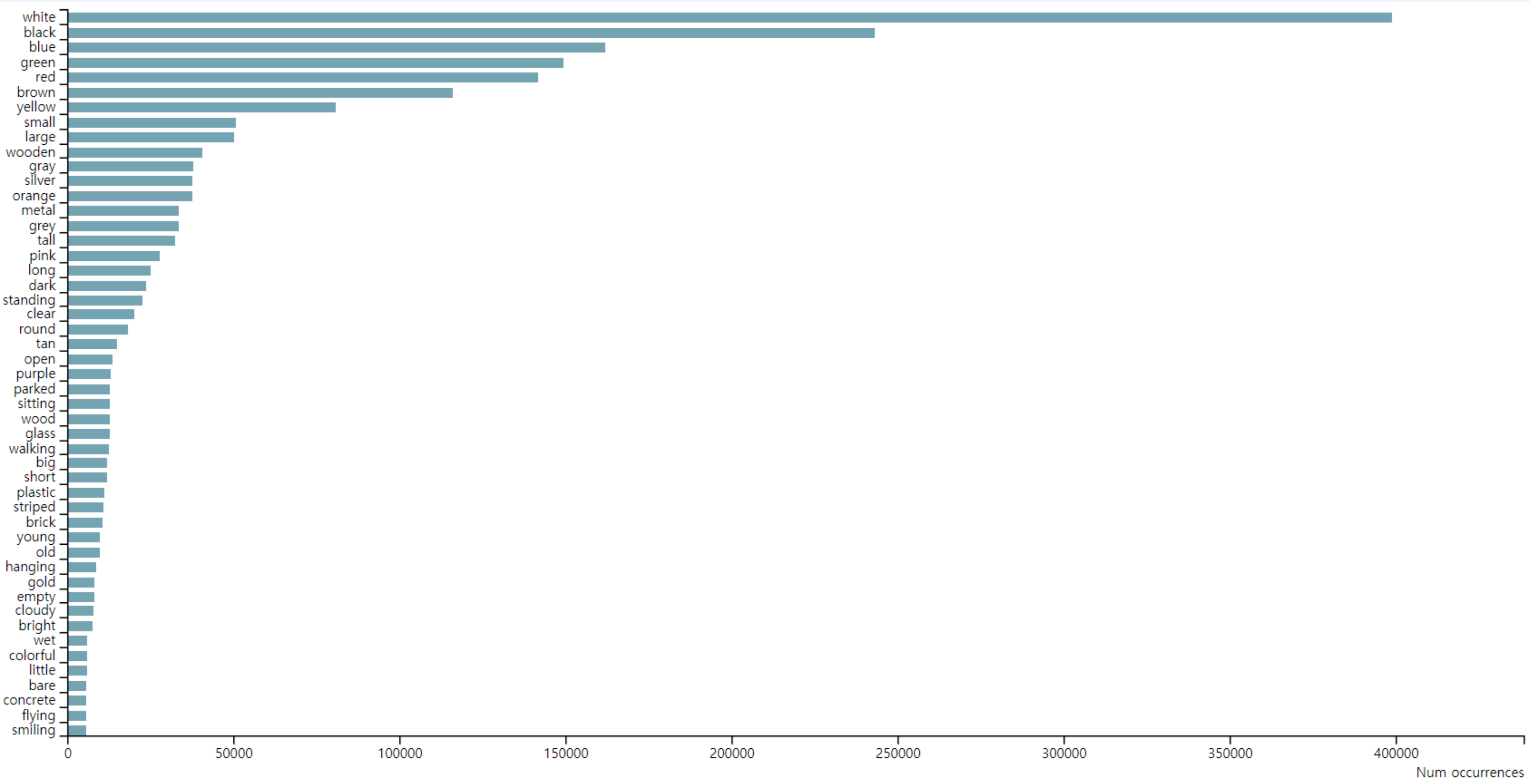# Scene Understanding
# Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv

## Top Attributes

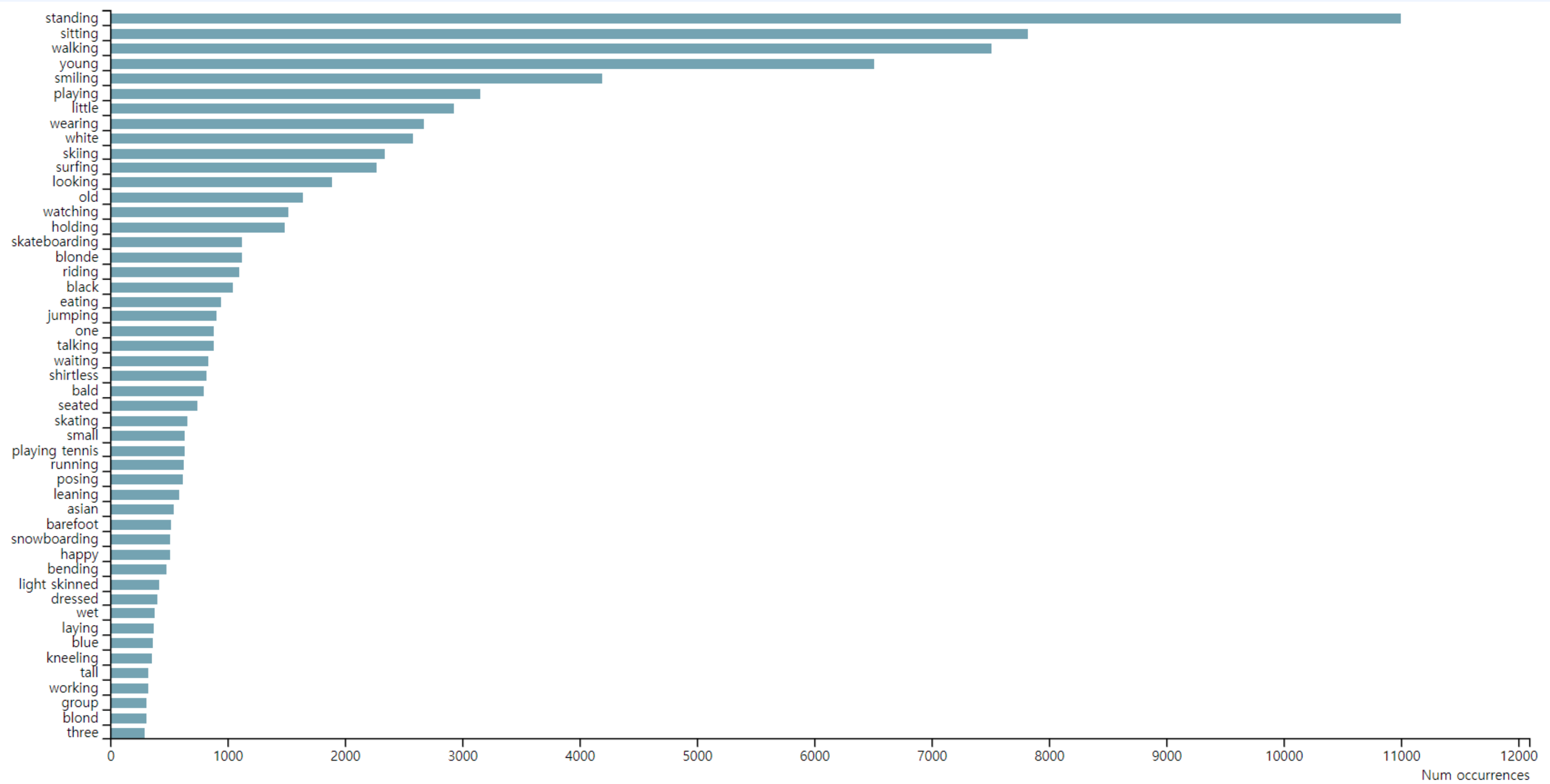# Scene Understanding
## Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv

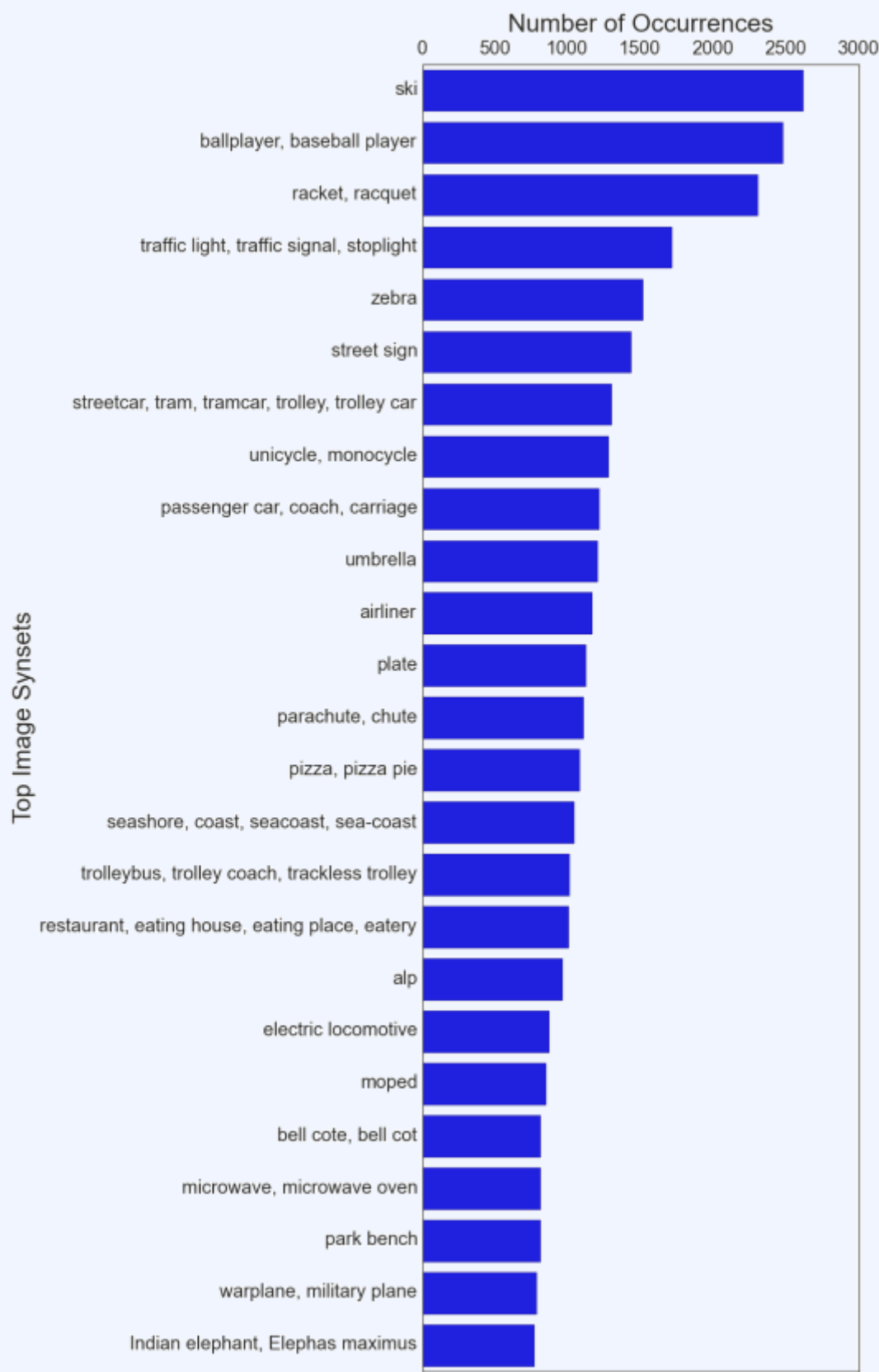### Top Attributes on People (man, women, person)

# Scene Understanding
## Scene representation learning (data)

R. Krishna et al. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations . arxiv

## Top Image Synsets



## Region Descriptions