# Imperial College London

# Phrase Localization Without Paired Training Examples

## Josiah Wang, Lucia Specia

**multi MT**

**Newton Fund**

## Phrase Localization

- Phrase Localization models are often trained with supervised paired training data.

- What if we don't have such training data?
- Can we still solve the problem?
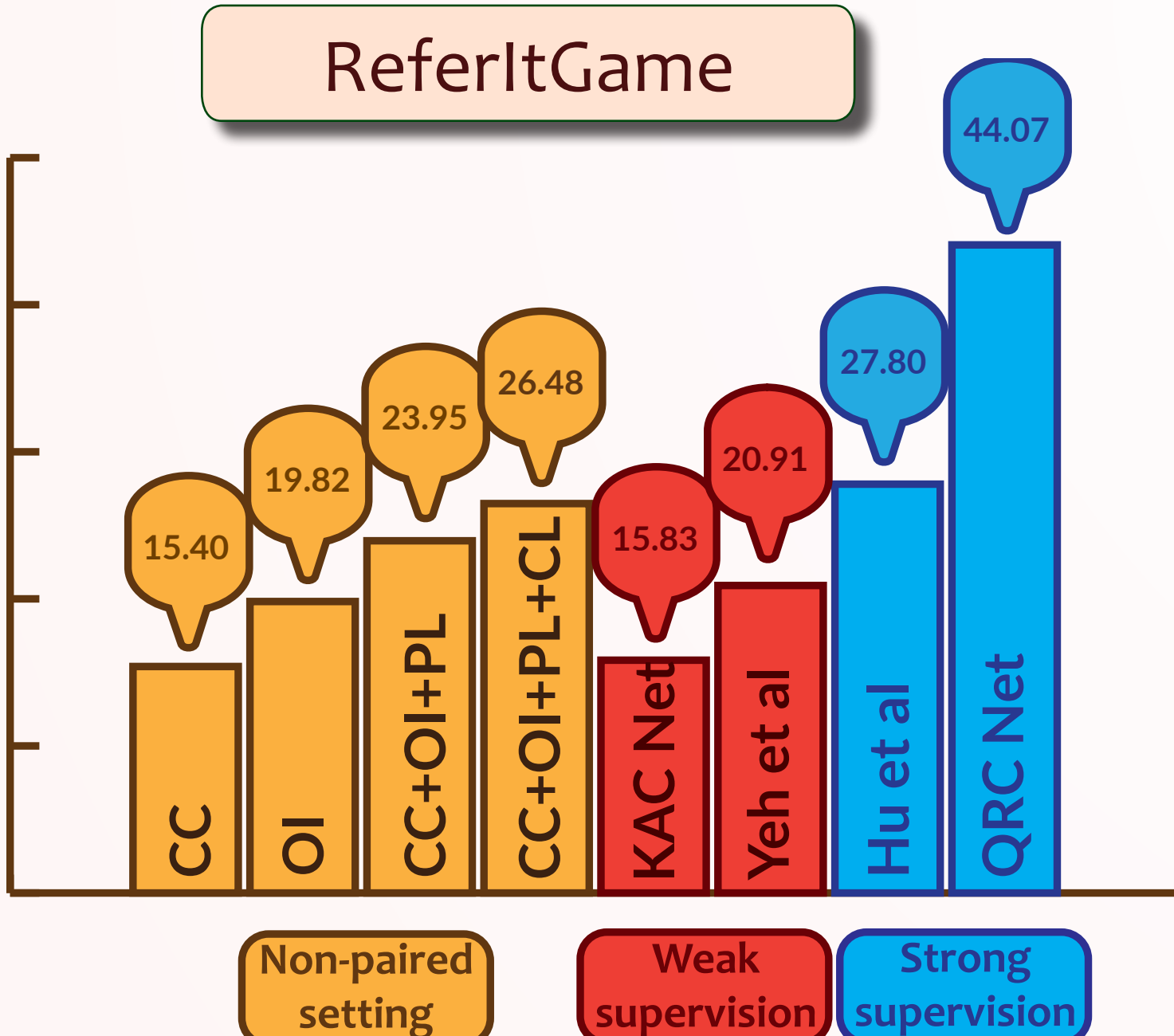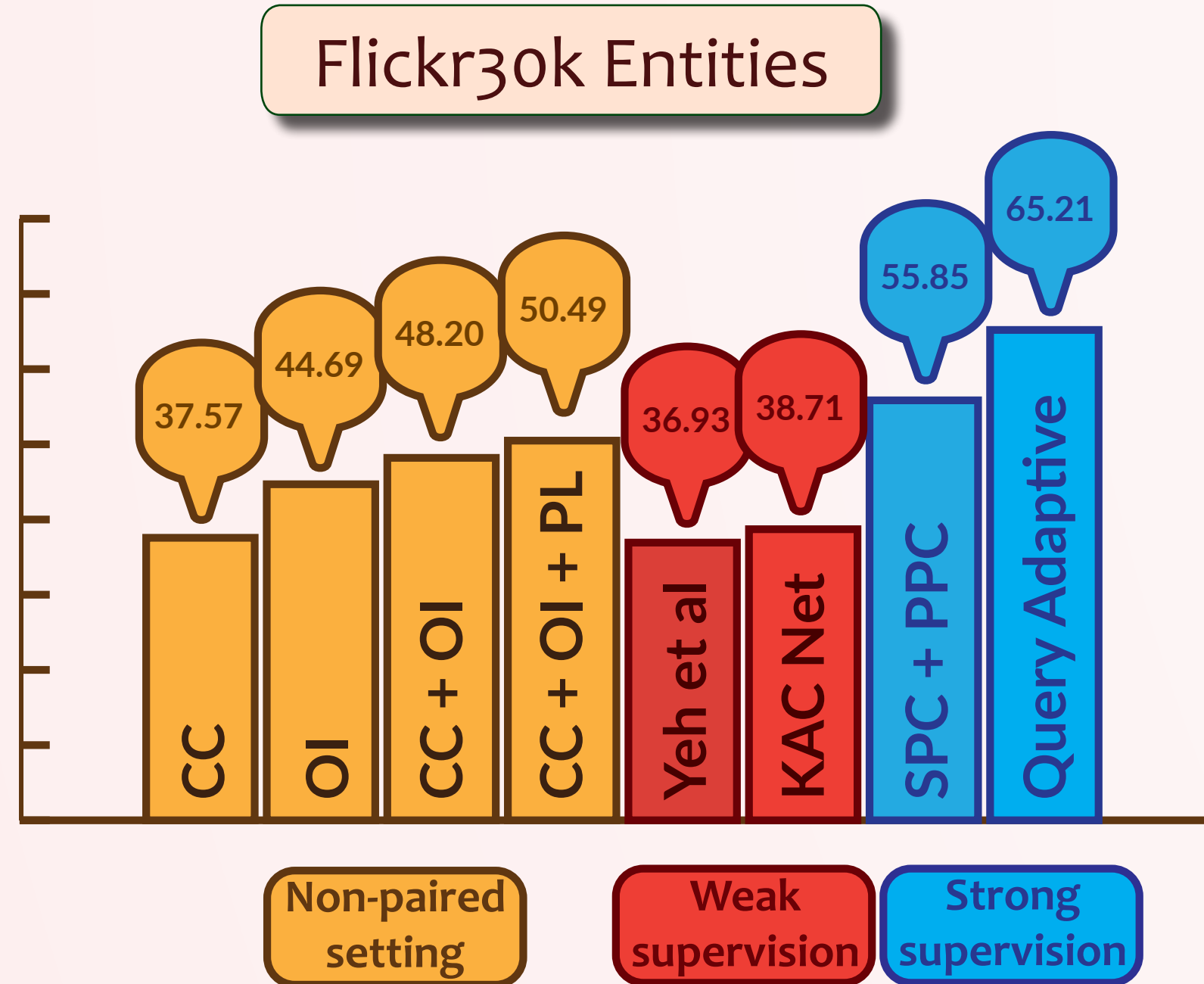- Maybe with off-the-shelf tools/model/resources?
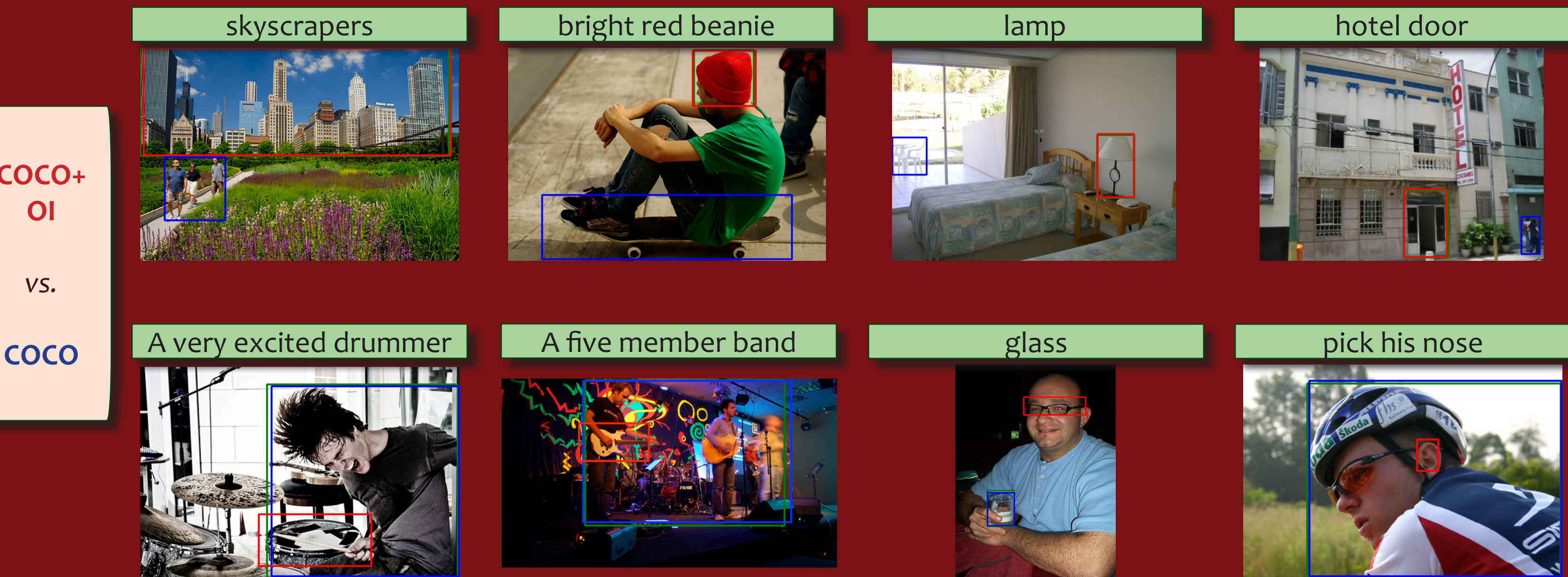
the handles of the slides

the handles of the slides
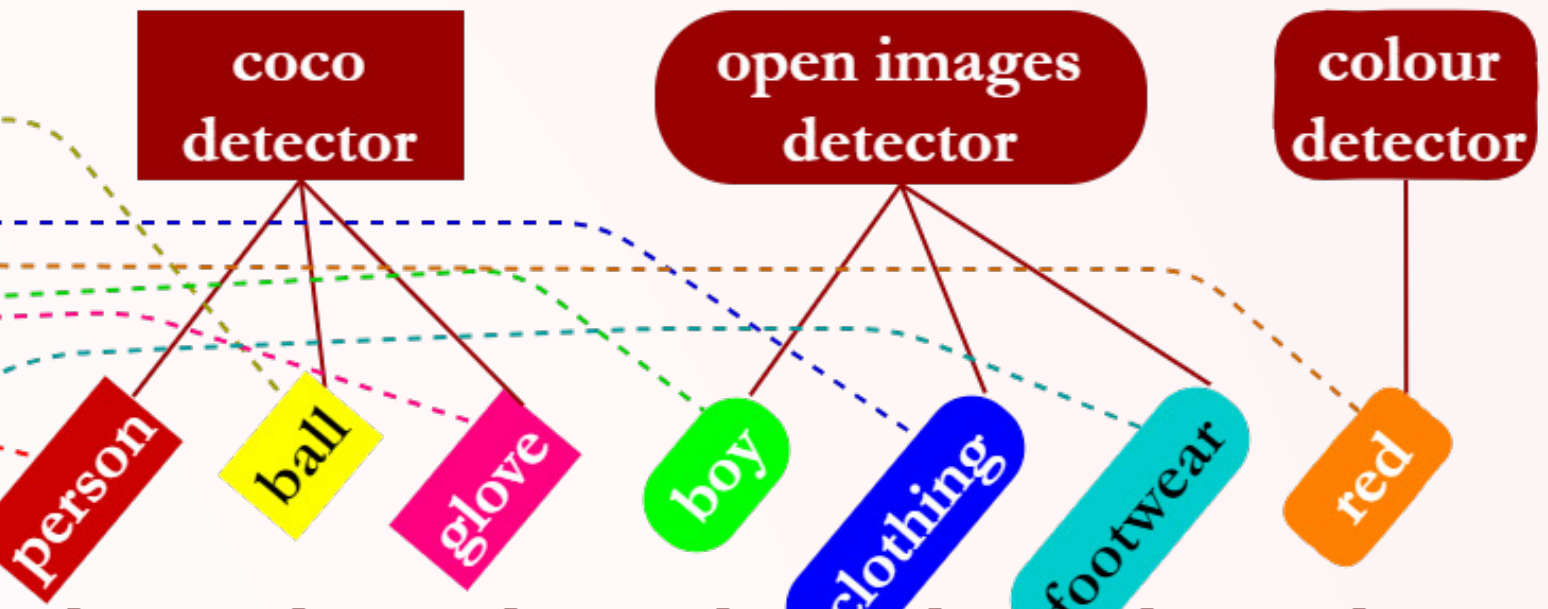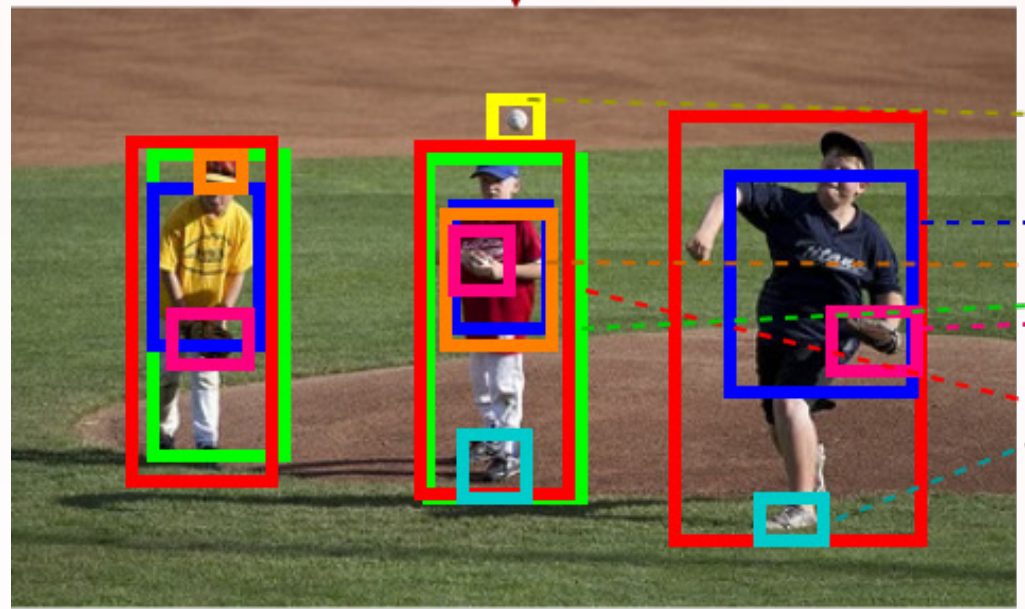
Strong supervision | Weak supervision

Non-paired setting
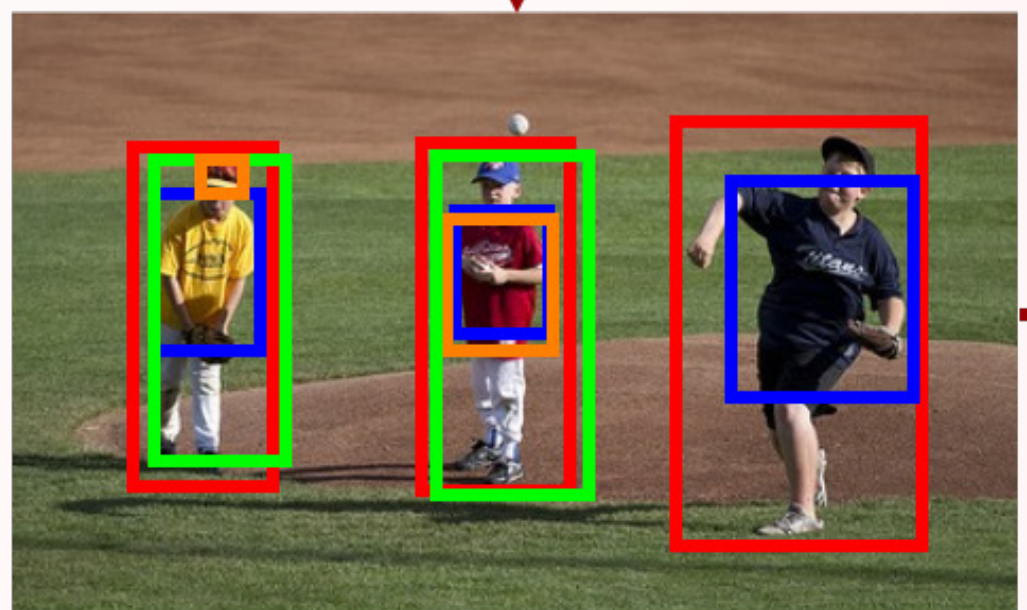
## Non-paired Setting Baseline

**1** Detect instances

Off-the-shelf pre-trained detectors:
- CC: COCO (80 categories)
- OI: Open Images (>500 categories)
- PL: Places365 (365 categories)
- YL: YOLO9000 (>9000 categories)
- CL: Colours (11 categories)
- Combinations of above

coco detector | open images detector | colour detector

person | ball | glove | boy | clothing | footwear | red

**boy in red shirt**

|  | 0.30 | 0.26 | 0.30 | 0.63 | 0.37 | 0.18 | 0.55 |
| boy | 0.34 | 0.21 | 0.20 | 1.00 | 0.16 | 0.16 | 0.16 |
| red | 0.09 | 0.15 | 0.16 | 0.16 | 0.16 | 0.16 | 1.00 |
| shirt | 0.16 | 0.24 | 0.29 | 0.22 | 0.42 | 0.29 | 0.26 |

- Compute semantic similarity between a query phrase & detector labels
  - cosine, word2vec
- Aggregate words in phrases:
  - average embeddings using one word from phrase
- Output:
  - Ranked/filtered list of bounding boxes

**2** Select relevant concepts

- Select/aggregate from instances most semantically similar to query phrase
- Tie-breakers:
  - Select random
  - Select largest
  - Select most confident
  - Union of bounding boxes
  - Consensus

**3** Localize phrase

## Experimental Results

### Flickr30k Entities

| CC | OI | CC + OI | CC+OI+PL | Yeh et al | KAC Net | SPC + PPC | Query Adaptive |
|----|----|---------|----------|-----------|---------|-----------|----------------|
| 37.57 | 44.69 | 48.20 | 50.49 | 36.93 | 38.71 | 55.85 | 65.21 |

Non-paired setting | Weak supervision | Strong supervision

### ReferItGame

| CC | OI | CC+OI+PL | CC+OI+PL+CL | KAC Net | Yeh et al | Hu et al | QRC Net |
|----|----|----------|-------------|---------|-----------|----------|---------|
| 15.40 | 19.82 | 23.95 | 26.48 | 15.83 | 20.91 | 27.80 | 44.07 |

Non-paired setting | Weak supervision | Strong supervision

## Example Ouput

COCO+OI vs. COCO

skyscrapers | bright red beanie | lamp | hotel door

A very excited drummer | A five member band | glass | pick his nose

WITH COLOUR vs. NO COLOUR

a yellow tennis suit | a long green shirt | pink blanket | trees

a red toy | Older male | guy in yellow shirt | bag below women in orange

## Discussion

- Non paired setting can be used as a strong baseline for phrase localization (or other V&L tasks)
- Paired data should be used more effectively, *on top of* what can be achieved with simpler methods without paired data
- Need to understand datasets better and not just blindly running complex models
- General, human-like AI: Better generalisation to different tasks