

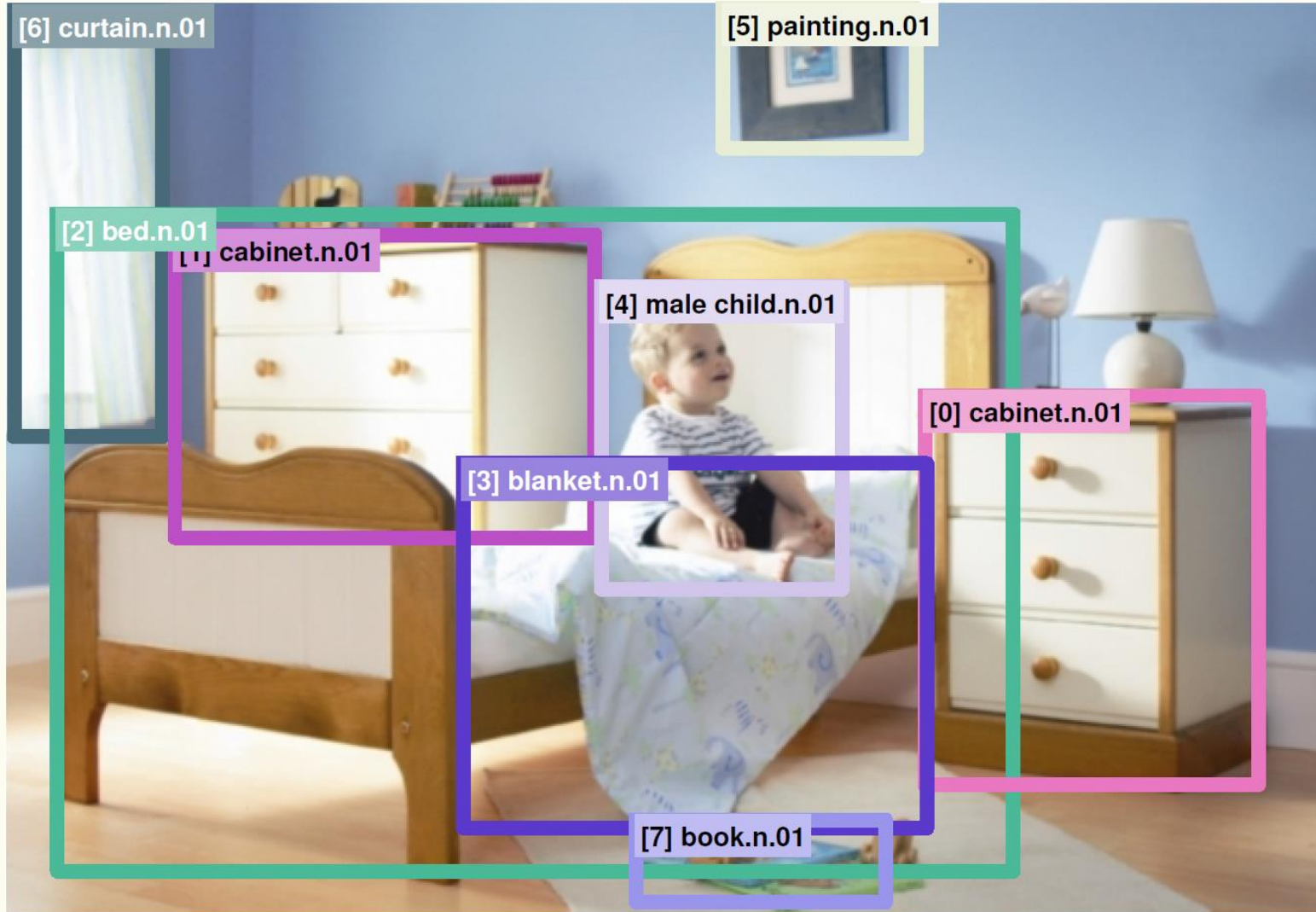
# Generating Image Descriptions with Gold Standard Visual Inputs: Motivation, Evaluation and Baselines

Josiah Wang & Robert Gaizauskas



# Task: Generate a description...

... given labelled bounding boxes



# This Presentation

- ImageCLEF2015 image description generation task
- Fine-grained evaluation metric for content selection
- Introduce and evaluate baselines

# Motivation

- Why gold standard visual input?
  - ‘Detect and generate’ methods – sensitive to noisy vision input
  - Assume ‘perfect’ input: Can we generate? How?
- Why fine-grained evaluation?
  - Human judgments are expensive, not scalable
  - Existing measures (BLEU, ROUGE, METEOR, CIDEr) are global
  - Better to evaluate tasks/phases of NLG pipeline explicitly: referring expression, verb/event/attribute/spatial relation expression, content ordering, **content selection**

# ImageCLEF

- Evaluation campaign since 2003
- Benchmark automatic image annotation & indexing for wide range of source images
- Part of CLEF initiative
  - Conference and Labs of the Evaluation Forum
  - Evaluation of information access systems since 2000

# ImageCLEF 2015

- Four main tasks in 2015
  - **Scalable image annotation**
  - Medical classification
  - Medical clustering
  - Liver CT annotation

# ImageCLEF 2015: Scalable Image Annotation

- Training set: 500k (image, webpage) pairs
- Subtask 1: Image annotation + localisation
  - For each 500k image, annotate + localise with 251 concepts
- Subtask 2: Image description generation
  - Noisy track: Generate descriptions for all 500k images
  - Clean track: Generate descriptions for 450 test images, given labelled bounding boxes ← **this talk**

# ImageCLEF 2015: Subtask 2 (Clean Track)

- Validation set (500 images from 500k)
  - Minimum 5 descriptions per image
    - Mean 9.5: Median: 8, Max: 51
  - Labelled bounding box annotations for 251 WordNet synsets
  - Correspondence annotation between text and bounding boxes





A [woman]<sup>2</sup> in a white [dress]<sup>0</sup> and gold [boots]<sup>5</sup> leaning on a [car]<sup>3</sup> .

A [woman]<sup>2</sup> poses along a [car]<sup>3</sup> .

[Woman]<sup>2</sup> dressed in white with gold [boots]<sup>5</sup> poses next to a police [car]<sup>3</sup> .

A [woman]<sup>2</sup> dressed in white leans against a white [car]<sup>3</sup> .

A [woman]<sup>2</sup> is leaning against a [car]<sup>3</sup> .

A [woman]<sup>2</sup> with gold [boots]<sup>5</sup> leans against an Indy pace [car]<sup>3</sup> .

A blonde [woman]<sup>2</sup> wearing gold shiny [boots]<sup>5</sup>, a white [top]<sup>0</sup> and short white skirt is leaning on a [car]<sup>3</sup> .

# Content Selection Evaluation Metric

**bbox instances referenced  
in gold standard**

$$P^I = \frac{1}{M} \sum_{m=1}^M \frac{|G_m^I \cap S^I|}{|S^I|}$$

**# gold standard descriptions  
for image  $I$**

**bbox instances referenced  
in generated sentence**

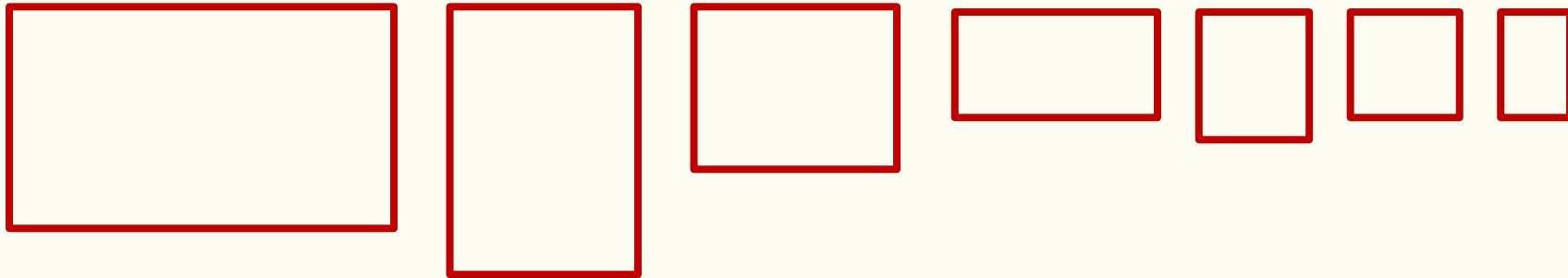
$$R^I = \frac{1}{M} \sum_{m=1}^M \frac{|G_m^I \cap S^I|}{|G_m^I|}$$

$$F^I = 2 \times \frac{P^I \times R^I}{P^I + R^I}$$

- Final score: Average over all test images

# Generating Descriptions: Baselines

- Visual Cues
  - Bounding box size (choose biggest first)
  - Bounding box position (choose most central object first)
- Threshold selected number of bounding boxes,  $k$



# Generating Descriptions: Baselines

- Textual Priors (from validation set)
  - Unigram
  - Bigram

A [woman] in a white [dress] and gold [boots] leaning on a [car] .

<start> → woman

woman → dress

dress → boots

boots → car

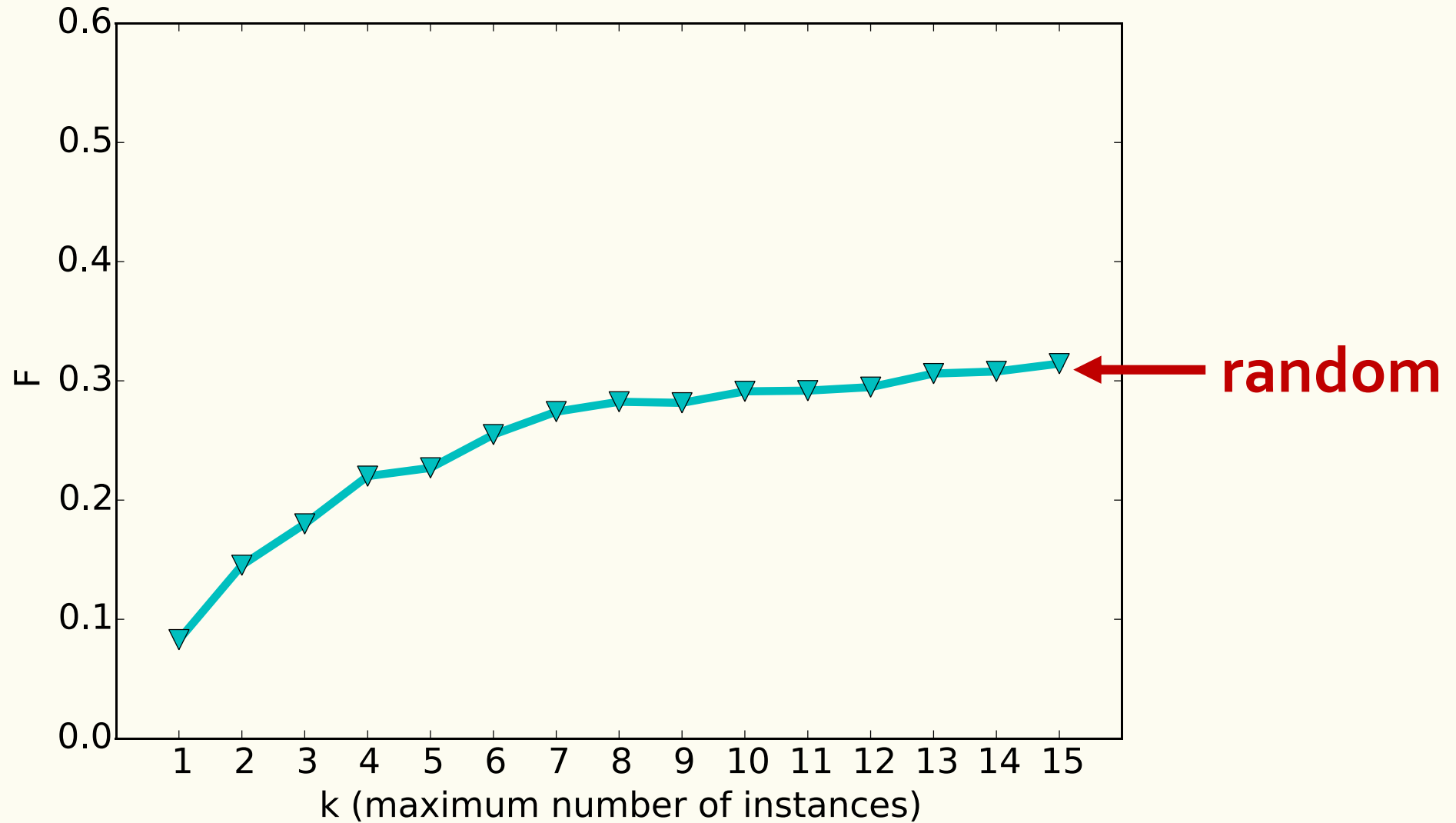
car → <end>

# Generating Descriptions: Baselines

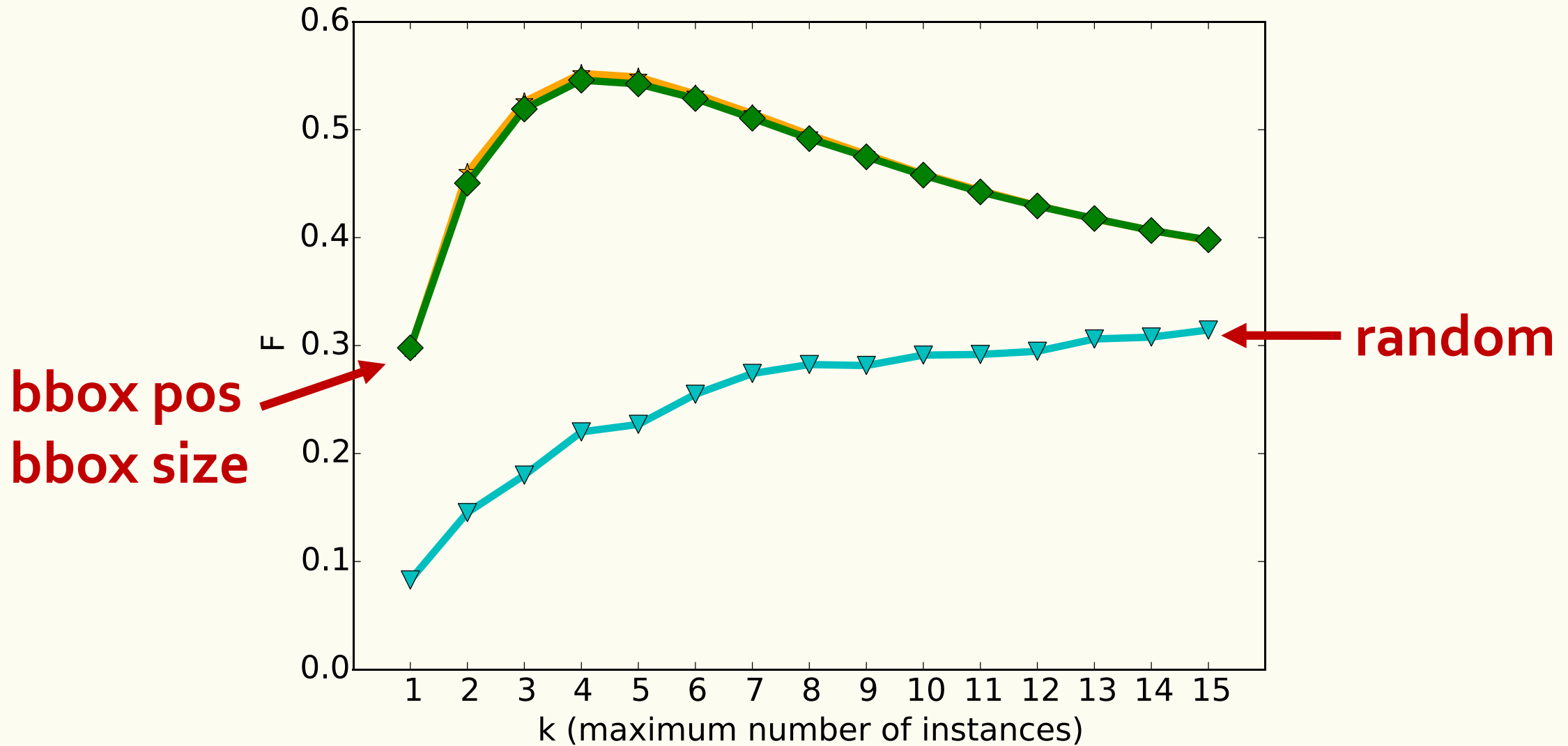
- Function Words
  - Random preposition/conjunction followed by optional *'the'*

[Woman]<sup>2</sup> with [dress]<sup>0</sup> and [boots]<sup>5</sup> on the [car]<sup>3</sup> .

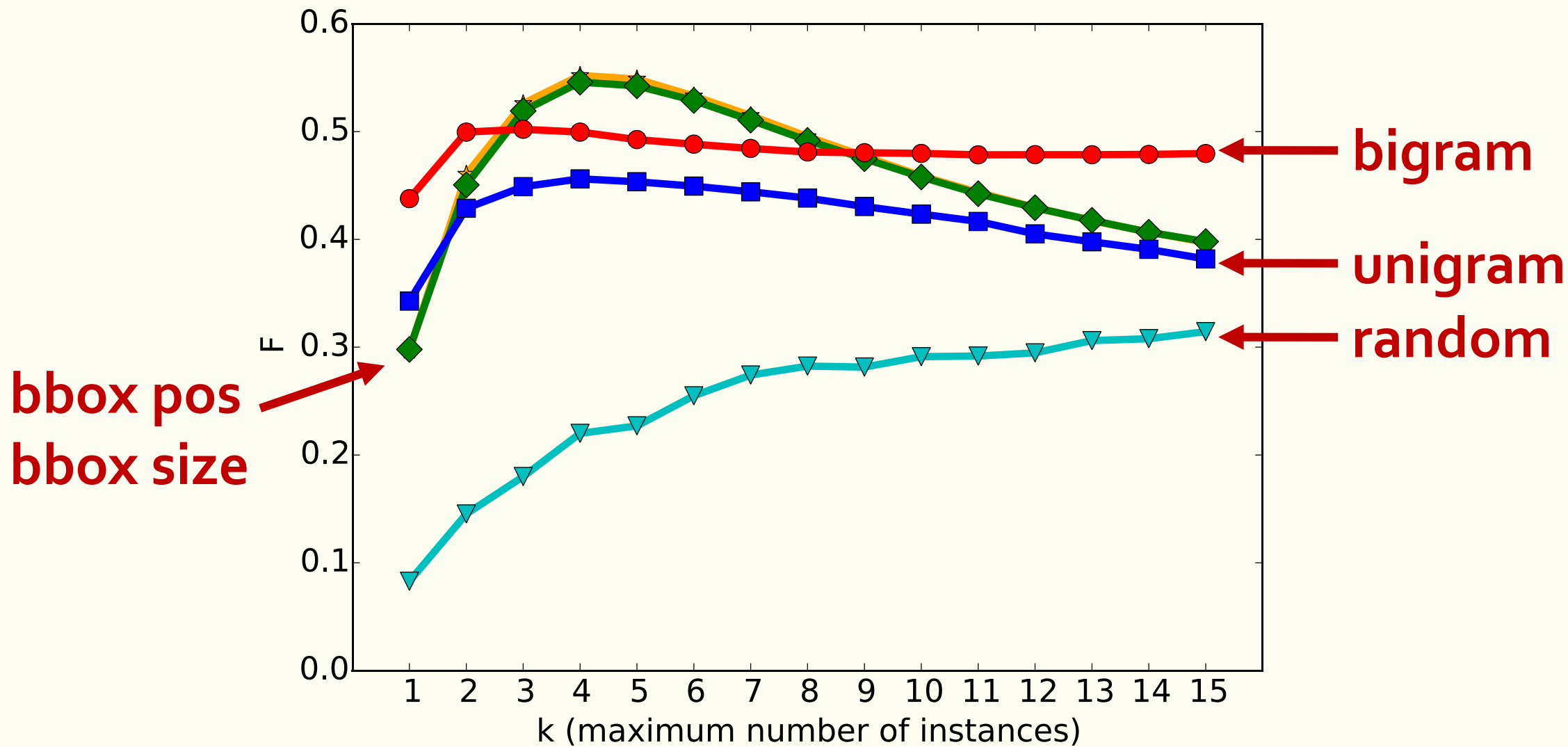
# Generating Descriptions: Results (F-score)



# Generating Descriptions: Results (F-score)

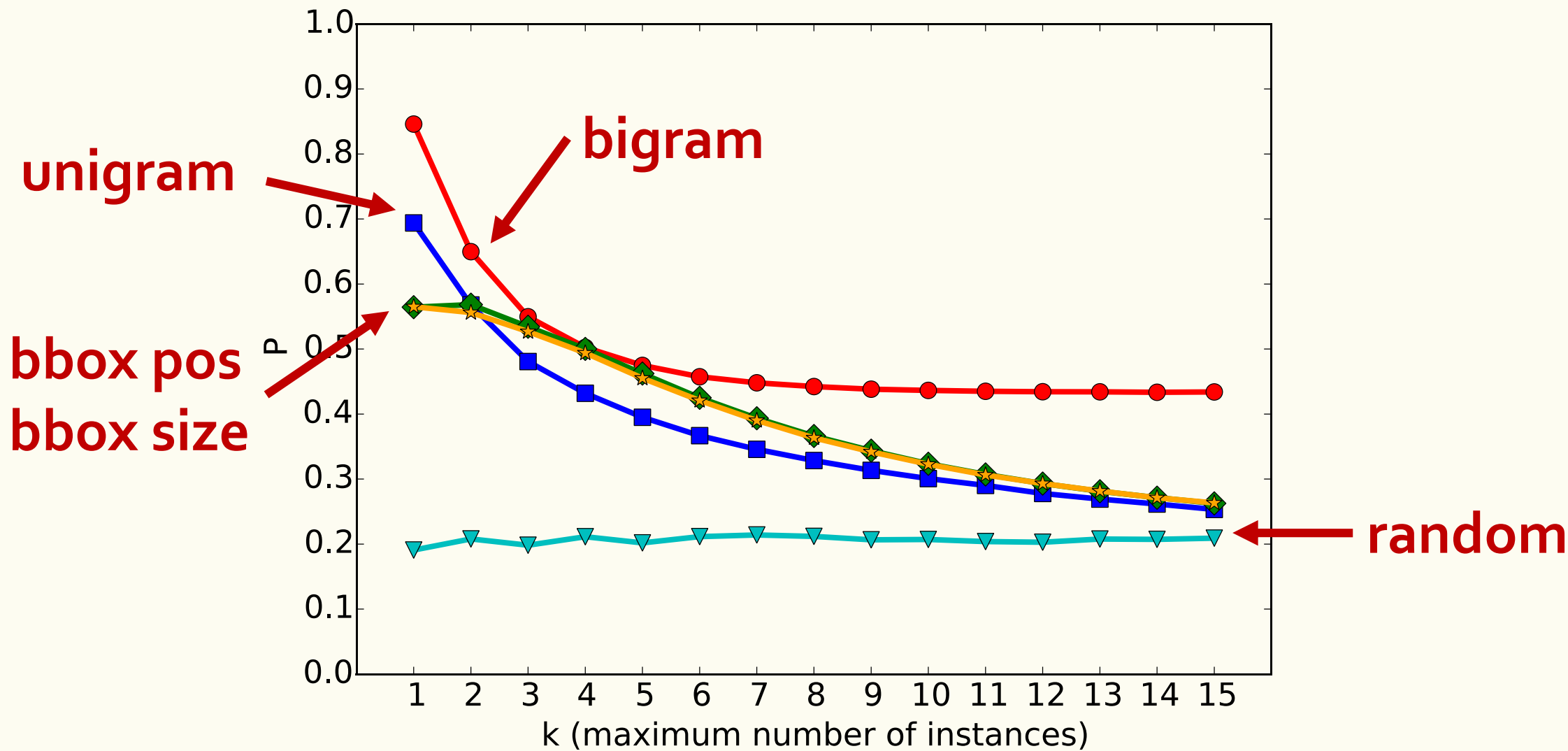


# Generating Descriptions: Results (F-score)

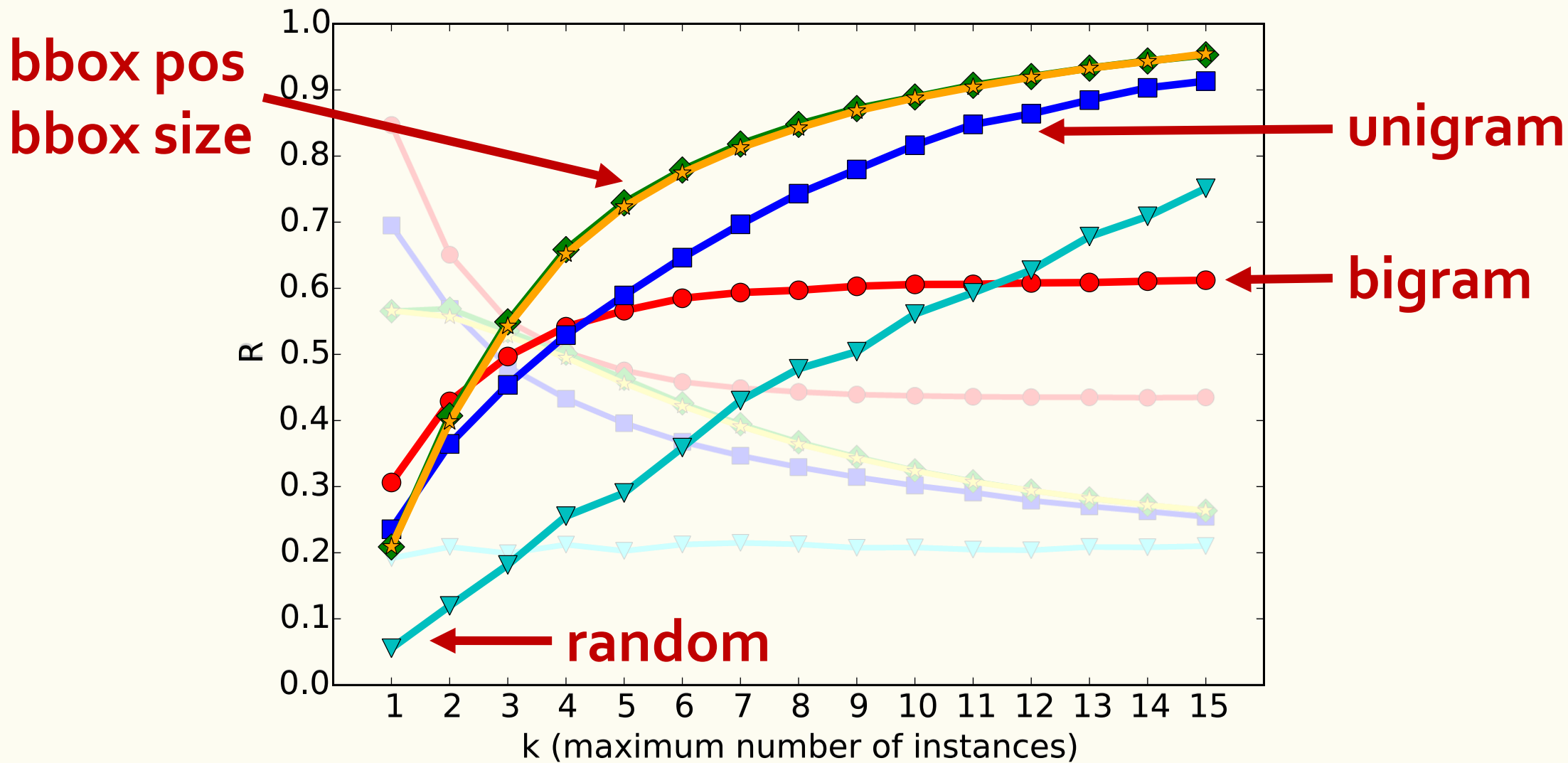




# Generating Descriptions: Results (Precision)

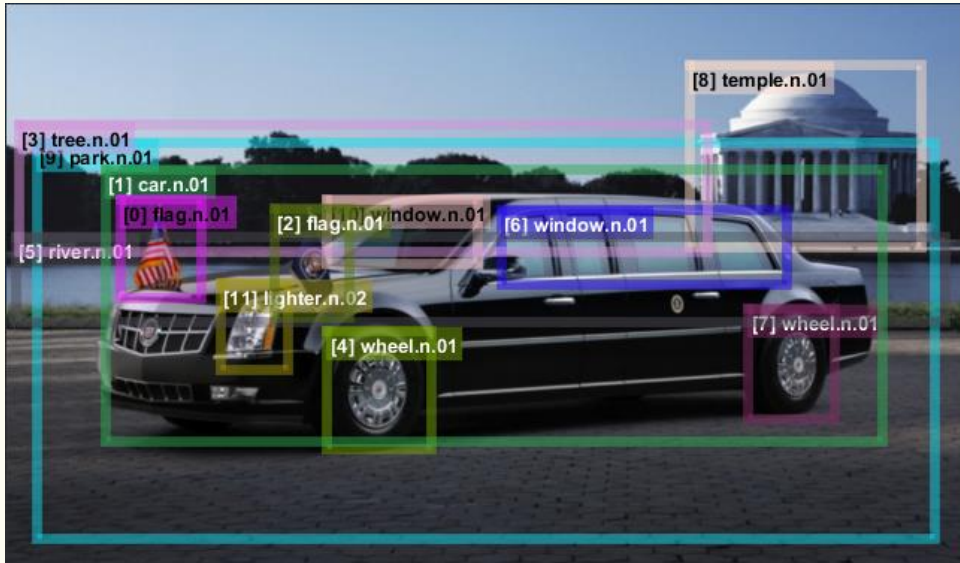


# Generating Descriptions: Results (Recall)

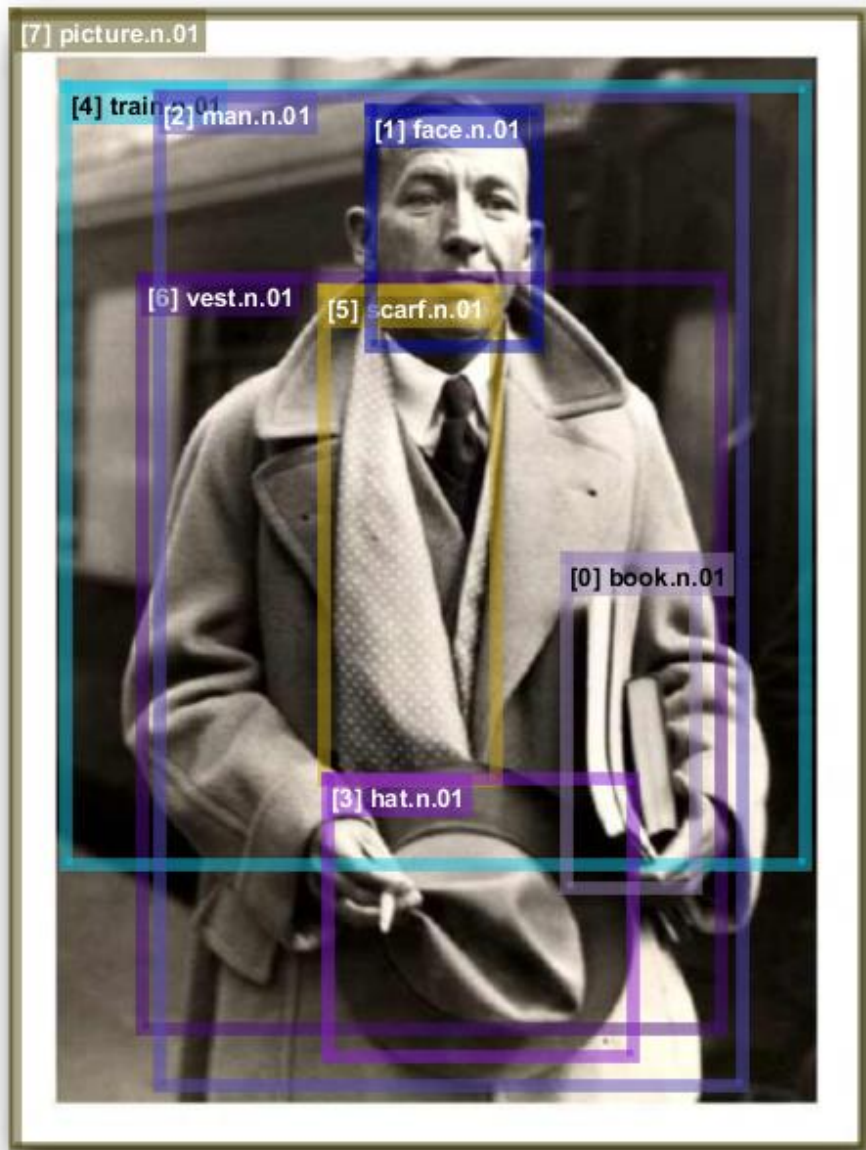




- random [F: 0.43]
  - [Wall]<sup>4</sup> inside [door]<sup>3</sup> around the [bicycle]<sup>0</sup> .
- bbox position [F: 0.79]
  - [Bicycle]<sup>0</sup> in [floor]<sup>1</sup> below [wall]<sup>2</sup> .
- bbox size [F: 0.79]
  - [Bicycle]<sup>0</sup> on [floor]<sup>1</sup> with [wall]<sup>2</sup> .
- unigram [F: 0.34]
  - [Table]<sup>7</sup> in the [wall]<sup>4</sup> around [wall]<sup>2</sup> .
- bigram [F: 0.03]
  - [Table]<sup>7</sup> near [door]<sup>3</sup> .



- random [F: 0.05]
  - [Park]<sup>9</sup> behind [wheel]<sup>7</sup> underneath the [window]<sup>6</sup> .
- bbox position [F: 0.59]
  - [Park]<sup>9</sup> on the [car]<sup>1</sup> below [river]<sup>5</sup> .
- bbox size [F: 0.44]
  - [Park]<sup>9</sup> behind the [car]<sup>1</sup> against the [tree]<sup>3</sup> .
- unigram [F: 0.42]
  - [Tree]<sup>3</sup> beneath [car]<sup>1</sup> by [window]<sup>6</sup> .
- bigram [F: 0.71]
  - [Car]<sup>1</sup> inside [flag]<sup>0</sup> underneath the [flag]<sup>2</sup> .



- random [F: 0.39]
  - [Vest]<sup>6</sup> at [hat]<sup>3</sup> behind the [picture]<sup>7</sup> .
- bbox position [F: 0.49]
  - [Picture]<sup>7</sup> on [man]<sup>2</sup> beside the [train]<sup>4</sup> .
- bbox size [F: 0.49]
  - [Picture]<sup>7</sup> among [man]<sup>2</sup> on the [train]<sup>4</sup> .
- unigram [F: 0.77]
  - [Man]<sup>2</sup> below the [hat]<sup>3</sup> at [book]<sup>0</sup> .
- bigram [F: 0.77]
  - [Man]<sup>2</sup> around the [hat]<sup>3</sup> along the [book]<sup>0</sup> .

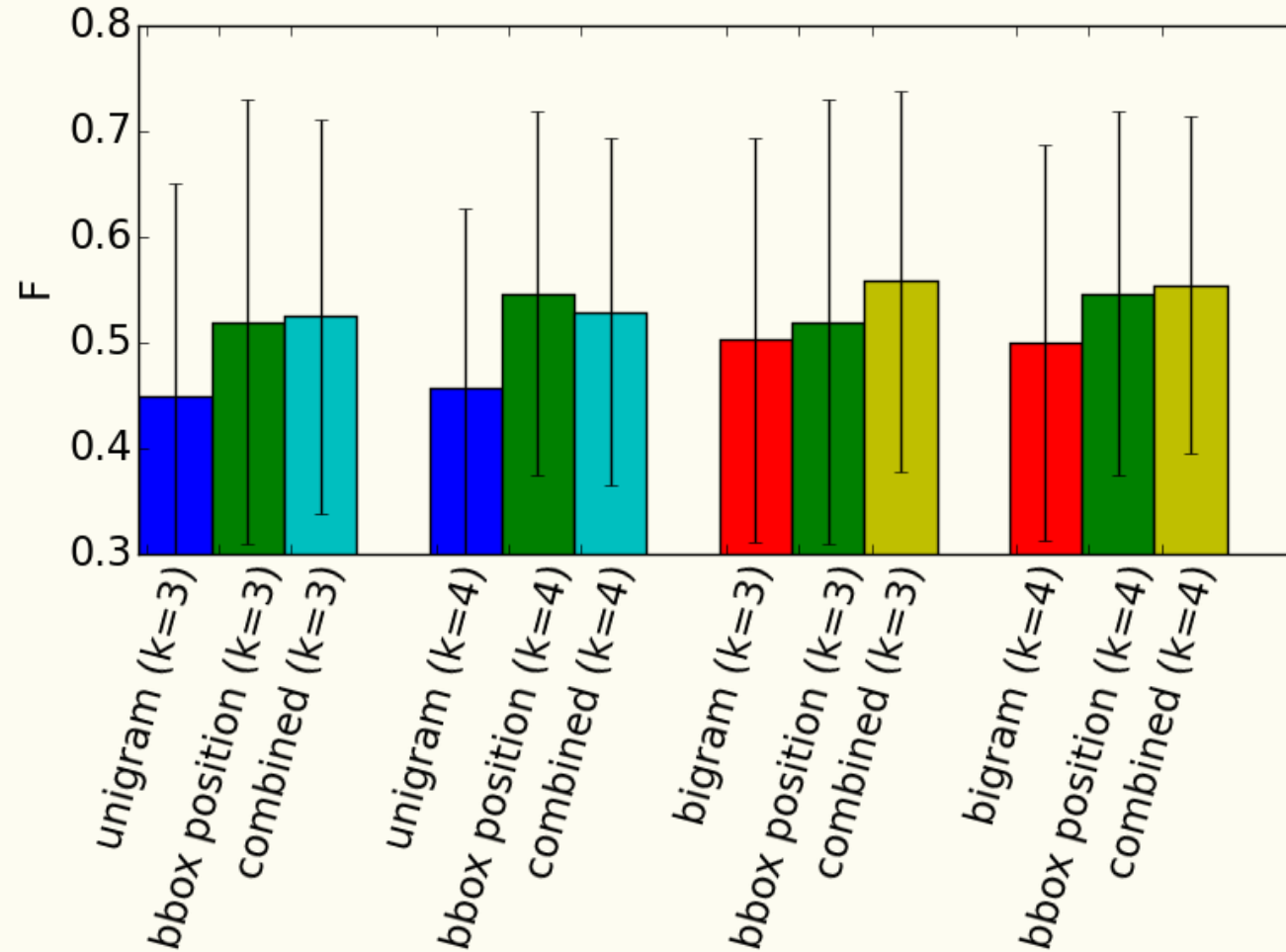
# Generating Descriptions: Combined

- Combining Text + Visual Priors
  - Re-rank by average rank from two systems

|           | Bigram | Bbox Position | Average Rank    | Combined |
|-----------|--------|---------------|-----------------|----------|
| [0] dress | 2      | 3             | $(2+3)/2 = 2.5$ | 2        |
| [1] wheel | -(6)   | 4             | $(6+4)/2 = 5$   | 5        |
| [2] woman | 1      | 1             | $(1+1)/2 = 1$   | 1        |
| [3] car   | 4      | 2             | $(4+2)/2 = 3$   | 3        |
| [4] hair  | -(6)   | 6             | $(6+6)/2 = 6$   | 7        |
| [5] boot  | 3      | 7             | $(3+7)/2 = 5$   | 4        |
| [6] sign  | -(6)   | 5             | $(6+5)/2 = 5.5$ | 6        |

[Woman]<sup>2</sup> in [dress]<sup>0</sup> by [car]<sup>3</sup> . (for  $k=3$ )

# Generating Descriptions: Combined



# Discussion

- We presented the sentence generation task of ImageCLEF
  - Proposed content selection evaluation metric and baselines
- Challenges
  - Bounding box annotations may not be informative enough
  - Suitability of fine-grained metrics
- Future work
  - Fine-grained metrics
    - Content ordering, referring expressions, verbs/predicates/prepositions
  - Generation of image descriptions
    - Stronger cues (co-occurrences, spatial relations)



# Generating Image Descriptions with Gold Standard Visual Inputs: Motivation, Evaluation and Baselines

Josiah Wang & Robert Gaizauskas

