INTERPRETABLE AND GLOBALLY OPTIMAL PREDICTION FOR TEXTUAL GROUNDING USING IMAGE CONCEPTS Raymond A. Yeh¹ Jinjun Xiong² Wen-mei W. Hwu¹ Minh N. Do¹ Alexander G. Schwing¹

LLINOS

MOTIVATION

Task: Phrase localization (*i.e.*, given a phrase and image find the corresponding bounding box described in the phrase)

- Challenging multimodal task
- Joint understanding of language and image
- Variety of objects and instance level descriptions

Contributions:

- A unified framework
- Does not rely on object proposals
- Globally optimal prediction
- Learned model parameters are interpretable

INTRODUCTION

Problem Formulation:

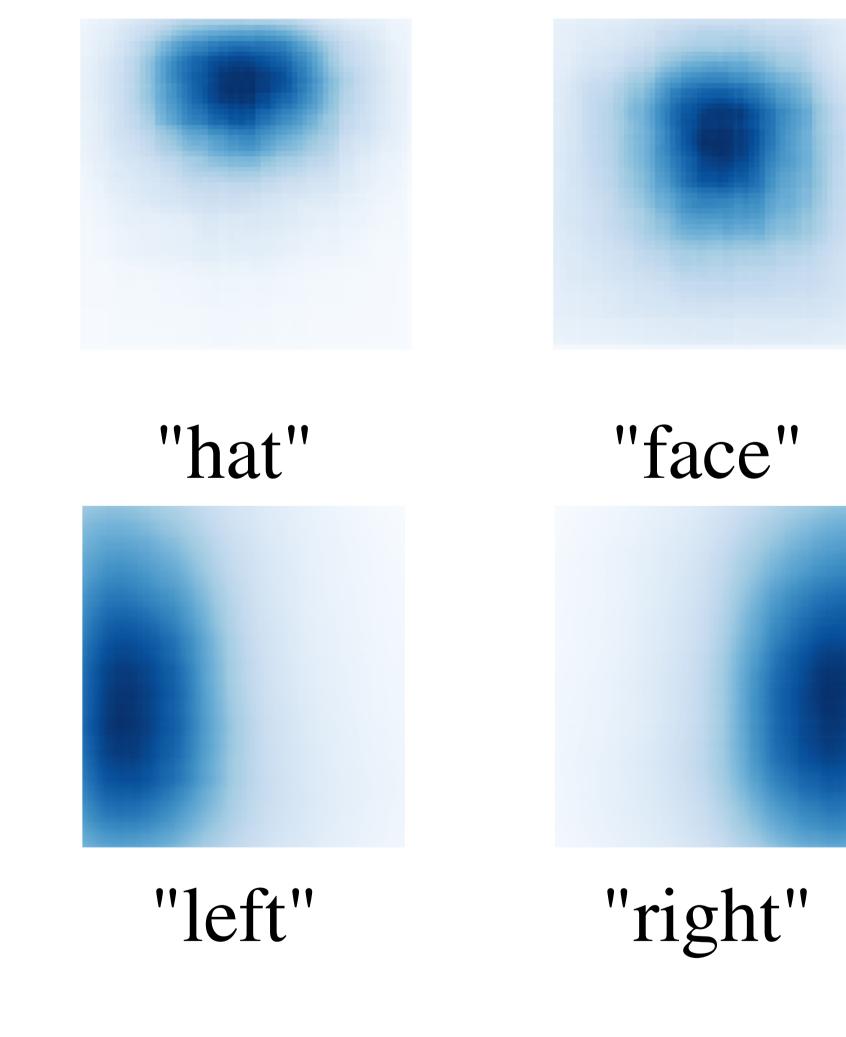
- Input: x = (Q, I)
- Output: $y = (y_1, ..., y_4)$
- Bounding box prediction: $\hat{y} = \arg\min_{y \in \mathcal{Y}} E(x, y, w)$

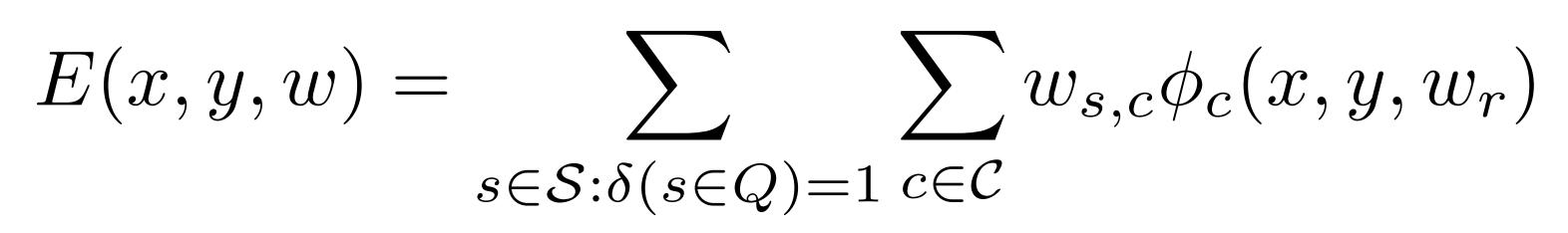
Intuition of Our Approach:

- Explicit word-concept modeling
- Combine "image concepts"

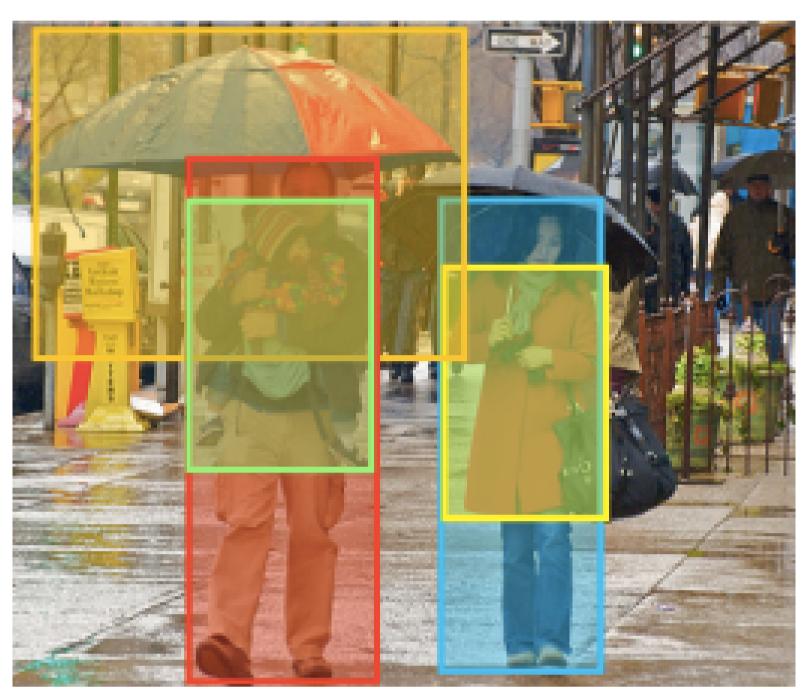
Energy Function Details:

- Image concept $c \in C$ is attached to a parametric score map $\hat{\phi}_c(x, w_r) \in \mathbb{R}^{W \times H}$
- $\phi_c(x, y, w_r) \in \mathbb{R}$ refers to the score accumulated within the bounding box y of score map $\phi_c(x, w_r)$
- Designed energy function:







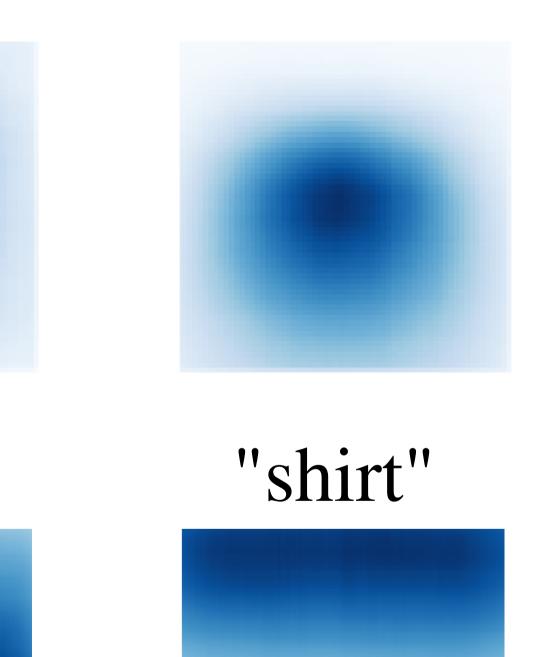


From Plummer et al. 2017

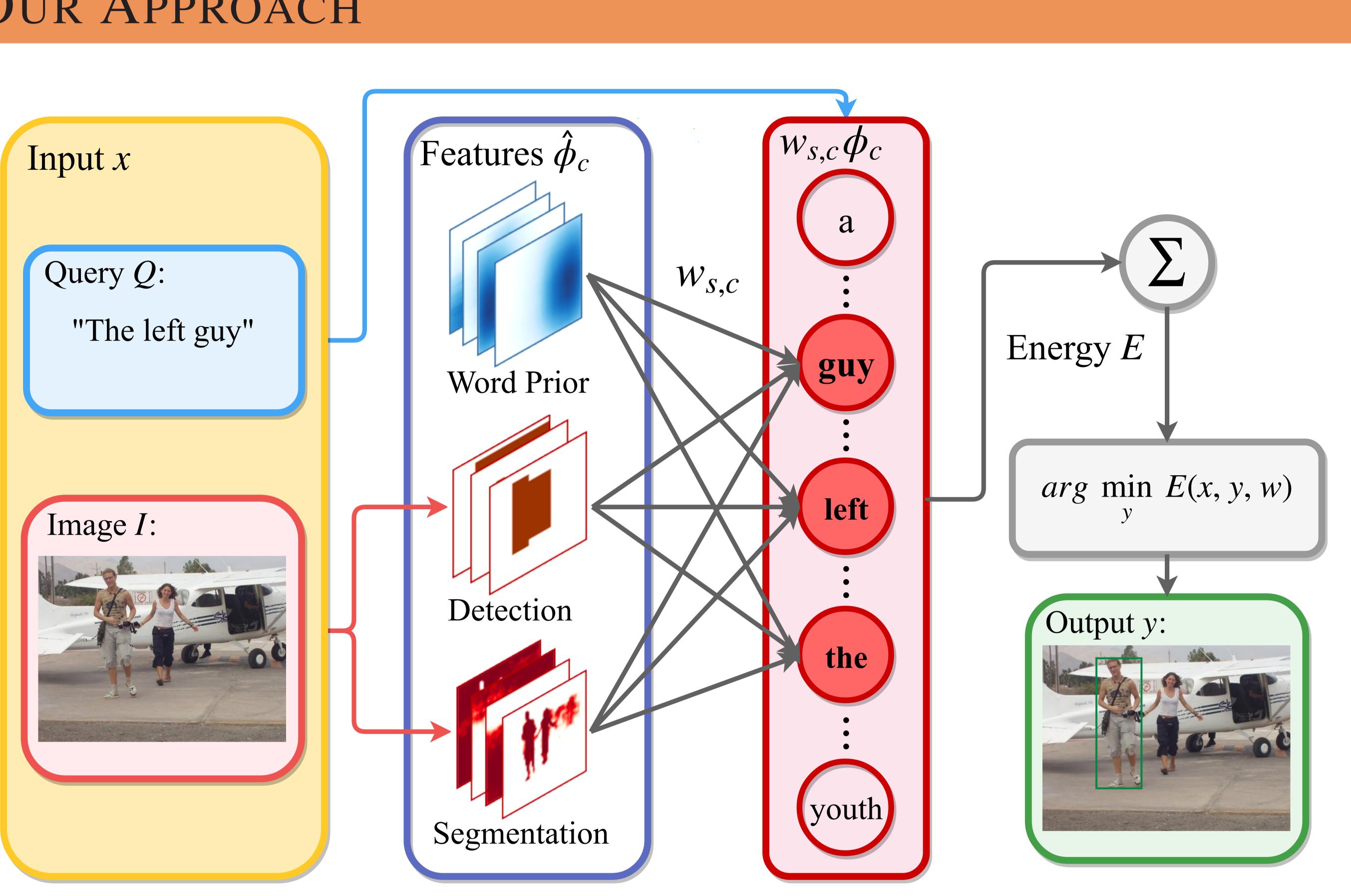
1. Dept. of ECE, University of Illinois at Urbana-Champaign 2. IBM Research

OUR APPROACH

A man carries a baby under a red ue umbrella next to a woman in a red jacket







Learning the Parameters:

• Structured support vector machine based surrogate loss minimization:

$$\min_{w} \quad \frac{C}{2} \|w\|_{2}^{2} + \sum_{(x,y)\in\mathcal{D}} \max_{\hat{y}\in\mathcal{Y}} (-E(x,\hat{y},w) + L(\hat{y},y)) + E(x,y,w)$$

- Task loss $L(\hat{y}, y)$ is Intersection over Union (IoU)
- Parameter $w_{s,c}$ connects word $s \in S$ and concept $c \in C$
- With fixed w_r , training is equivalent to training a structured SVM
- The cutting-plane algorithm of Joachims *et al.* 2009 works well
- by Lampert et al. 2009

• Parameters w includes score map parameters w_r and top layer parameters $w_{s,c}$ • To solve the inference problem effectively, we utilize efficient subwindow search

RESULTS

Method

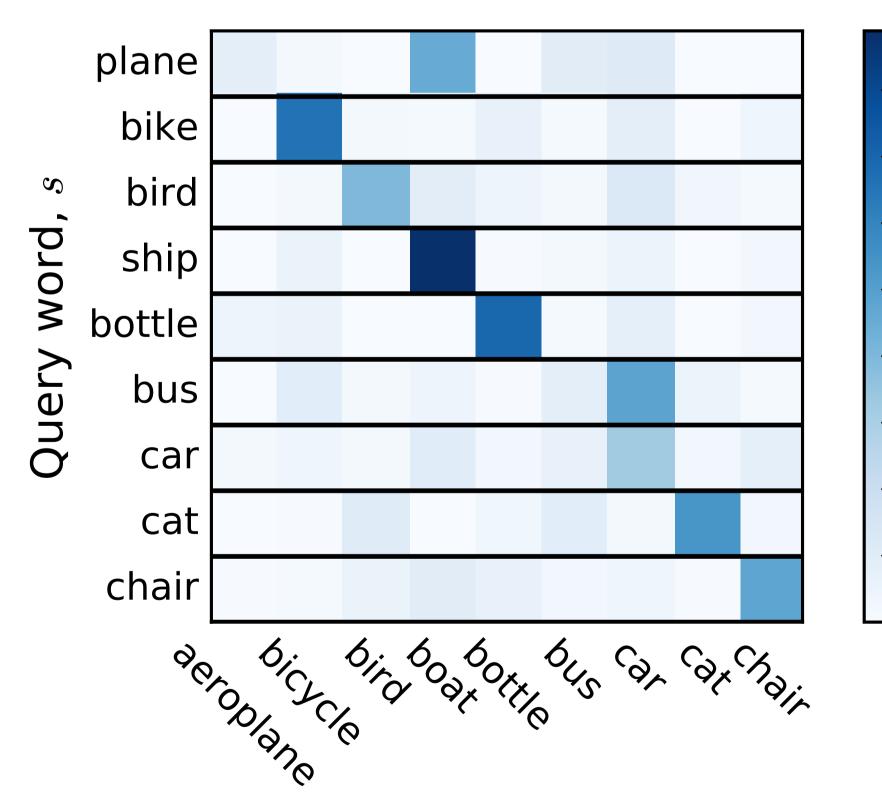
SCRC (Hu 2016 DSPE (Wang 201 GroundeR(Rohrbach CCA (Plummer 20 Ours

Results on Flickr30k dataset **Success Cases:**



second bike from right in front

Weights Interpretability:



Concept, a

Failure Cases:



her shoes

Ú		

	Acc. (%)
5)	27.80
16)	43.89
2016)	47.81
)17)	50.89
	53.97

Acc. (%) Method 17.93 SCRC (Hu 2016) 26.93 GroundeR (Rohrbach 2016) Ours [Prior + Geo] 25.56 33.36 Ours [Prior + Geo + Seg] Ours [Prior + Geo + Seg + Det] 34.70



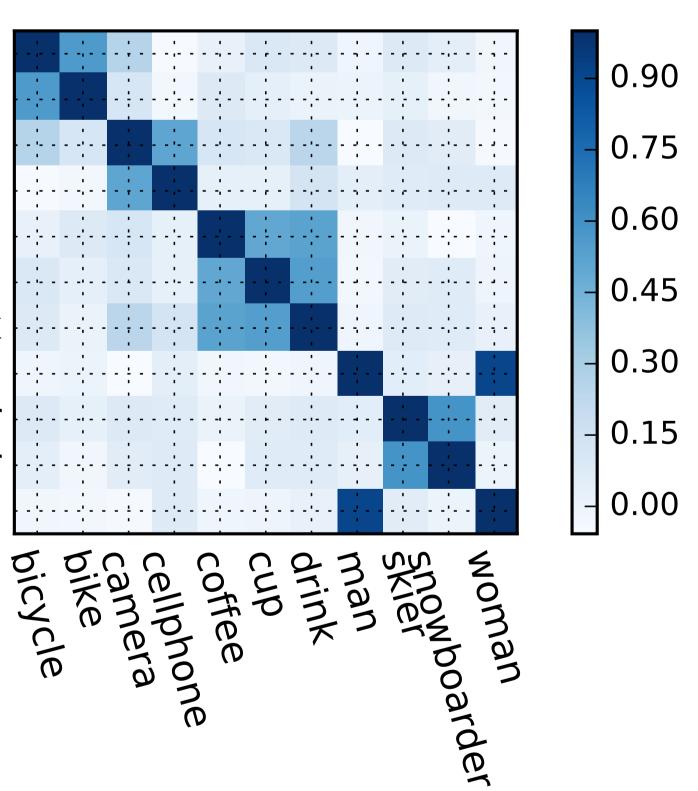
painting next to the two on the

Results on ReferItGame dataset



person all the way to the right

camera cellphone coffee drink man skier snowboarder woma



Query word, s'



a red shirt



a dirt bike