



# Captioning Images with Diverse Objects.

Subhashini Venugopalan<sup>1</sup>, Lisa Anne Hendricks<sup>2</sup>, Marcus Rohrbach<sup>2,3</sup>,

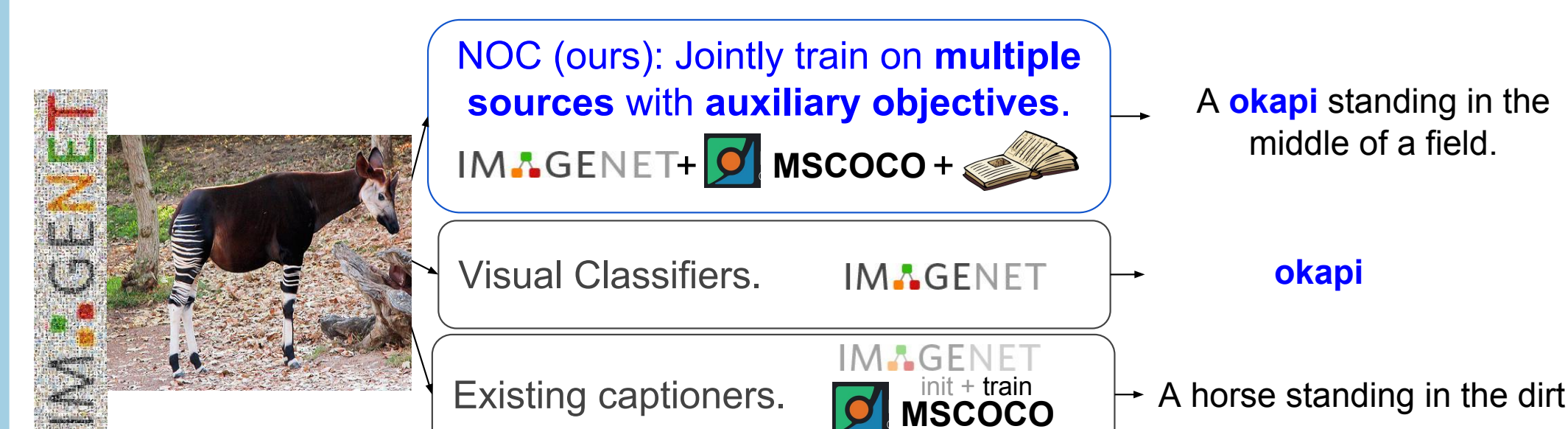
Raymond Mooney<sup>1</sup>, Trevor Darrell<sup>2</sup>, Kate Saenko<sup>4</sup>

<sup>1</sup> UT-Austin <sup>2</sup> UC-Berkeley <sup>3</sup> Facebook AI Research <sup>4</sup> Boston University



## GOALS

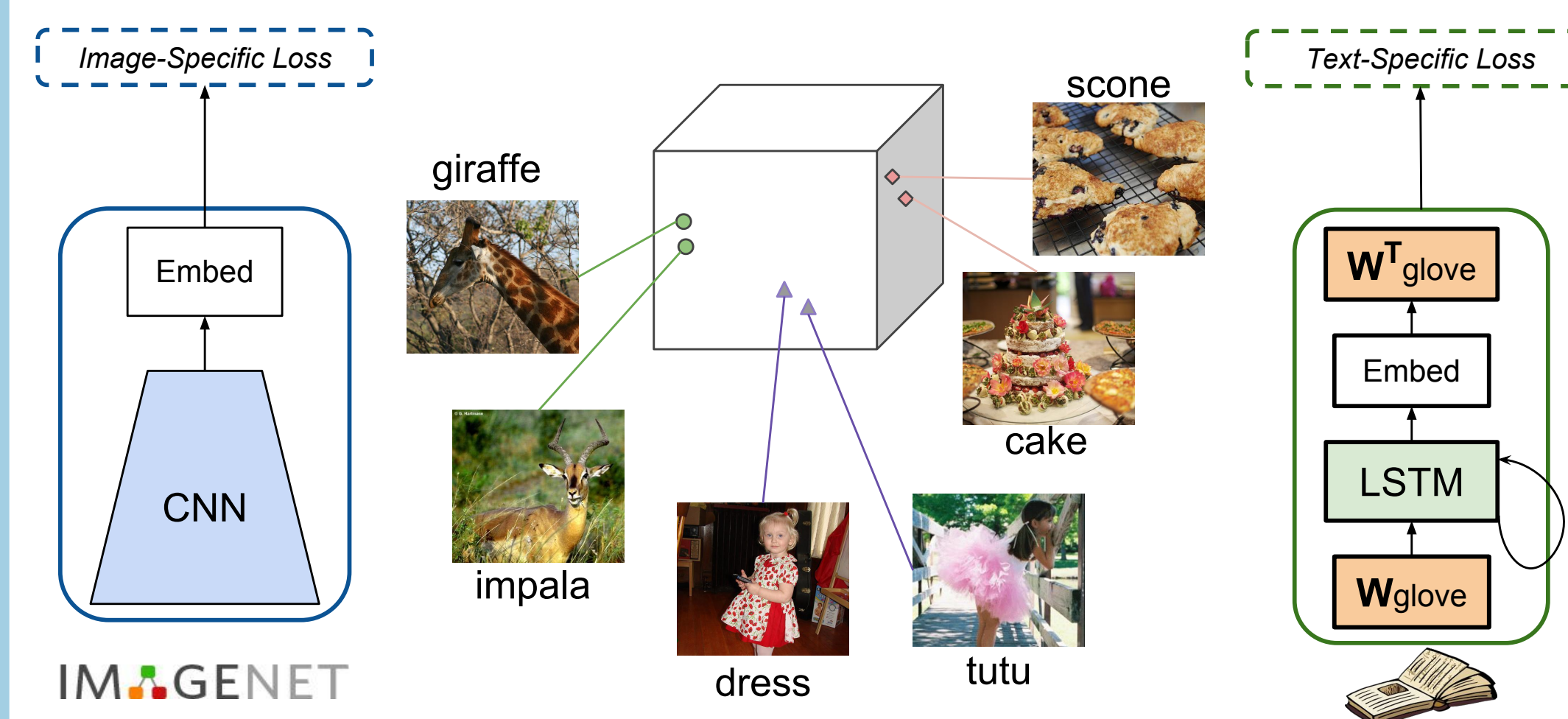
Existing visual classifiers can recognize hundreds of categories of objects. Can we describe these objects in context without paired image-caption training data?



We propose Novel Object Captioner which can describe objects unseen in paired image-caption data.

## NOC KEY INSIGHTS

Train jointly on multiple data sources.



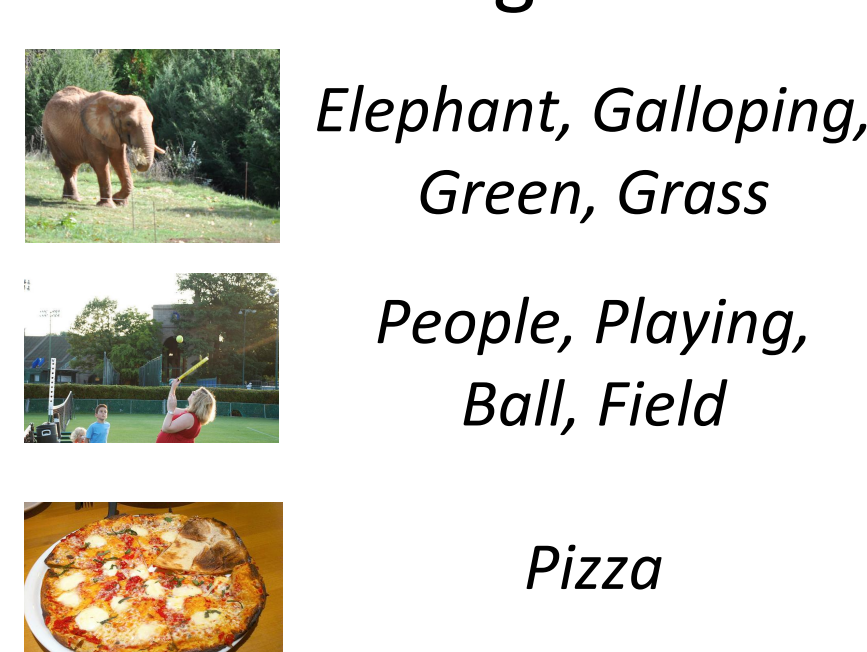
1. **Learn from unpaired data.** Train visual CNN on unpaired image data, and an LSTM Language Model on unannotated text data.
2. **Capture semantic similarity** of words in the language model using dense word embeddings.
3. **Train jointly to describe novel objects.** A visual recognition CNN, a language model, and an image-caption model [1] are trained jointly on different data sources with shared parameters.

## EVALUATION

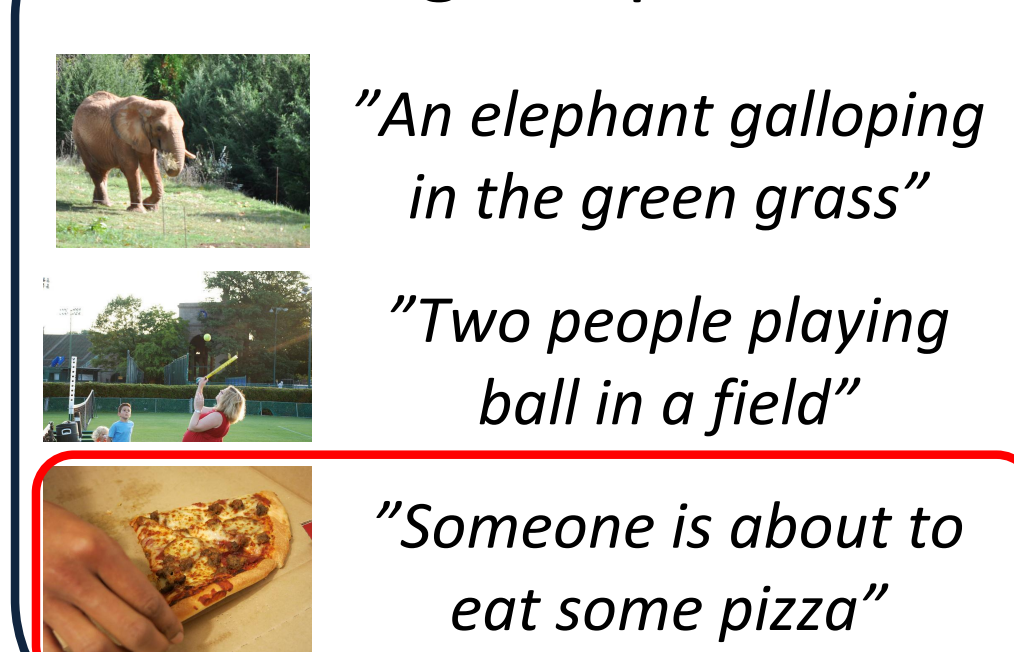
We hold out a subset of data from COCO [2].

1. COCO **Held-out** dataset

### COCO Image Data



### COCO Image-Caption Data

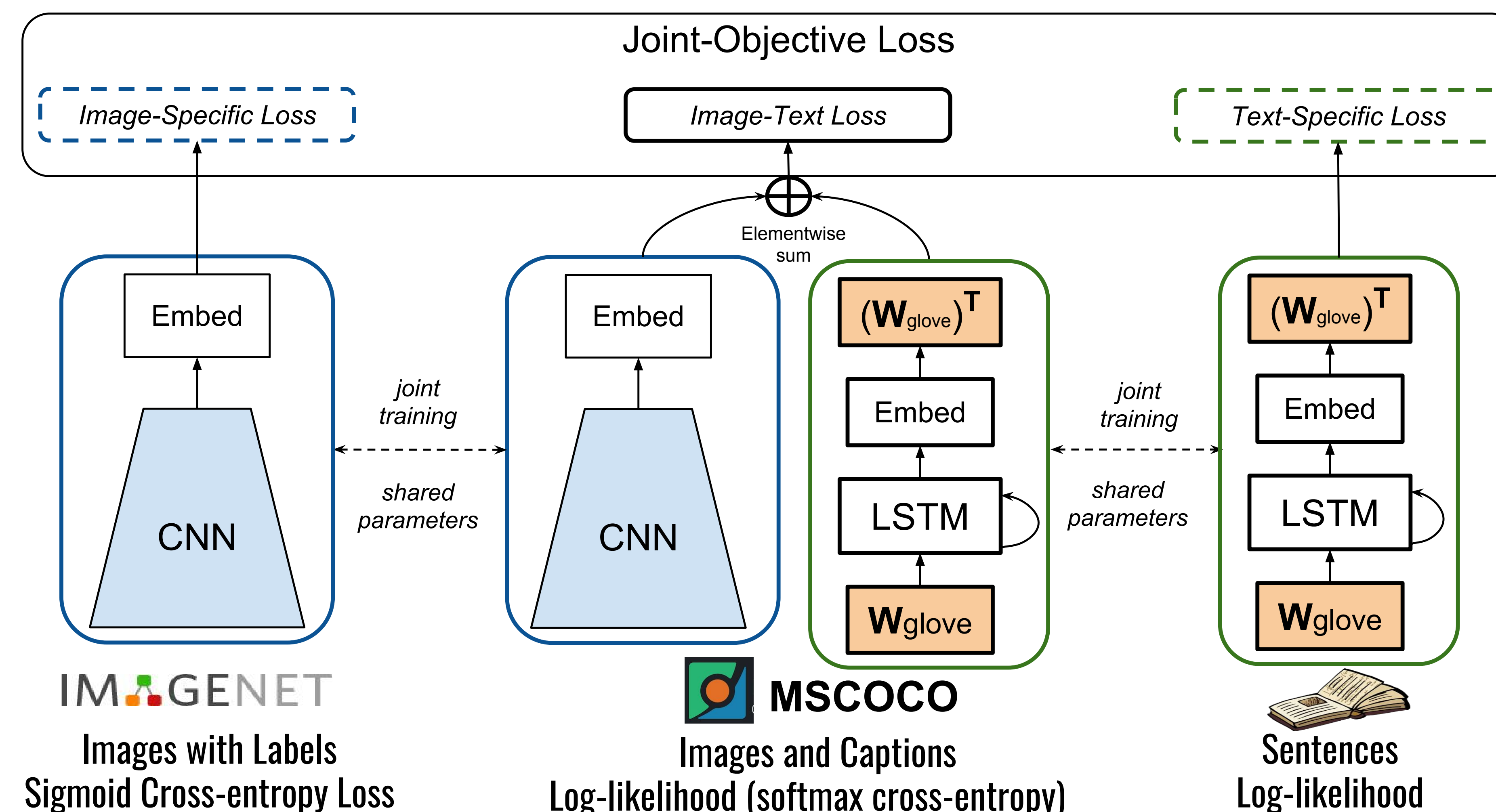


## 2. IMAGENET

638 categories from ImageNet not mentioned in COCO. 52 classes with rare mentions (med~5 images) in COCO.

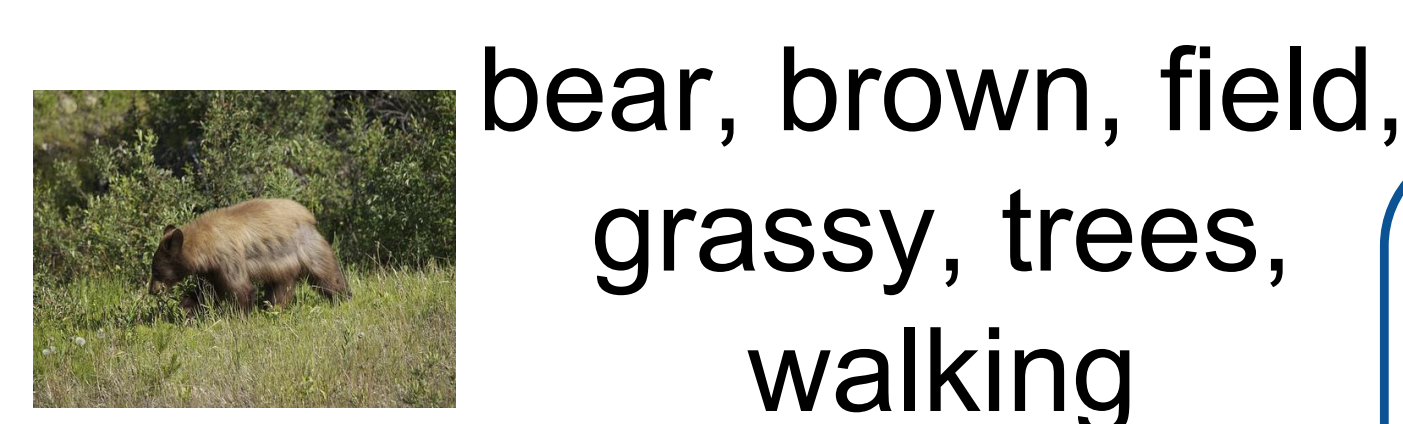
## MODEL

Share network parameters and train jointly on multiple data sources and with different objectives.

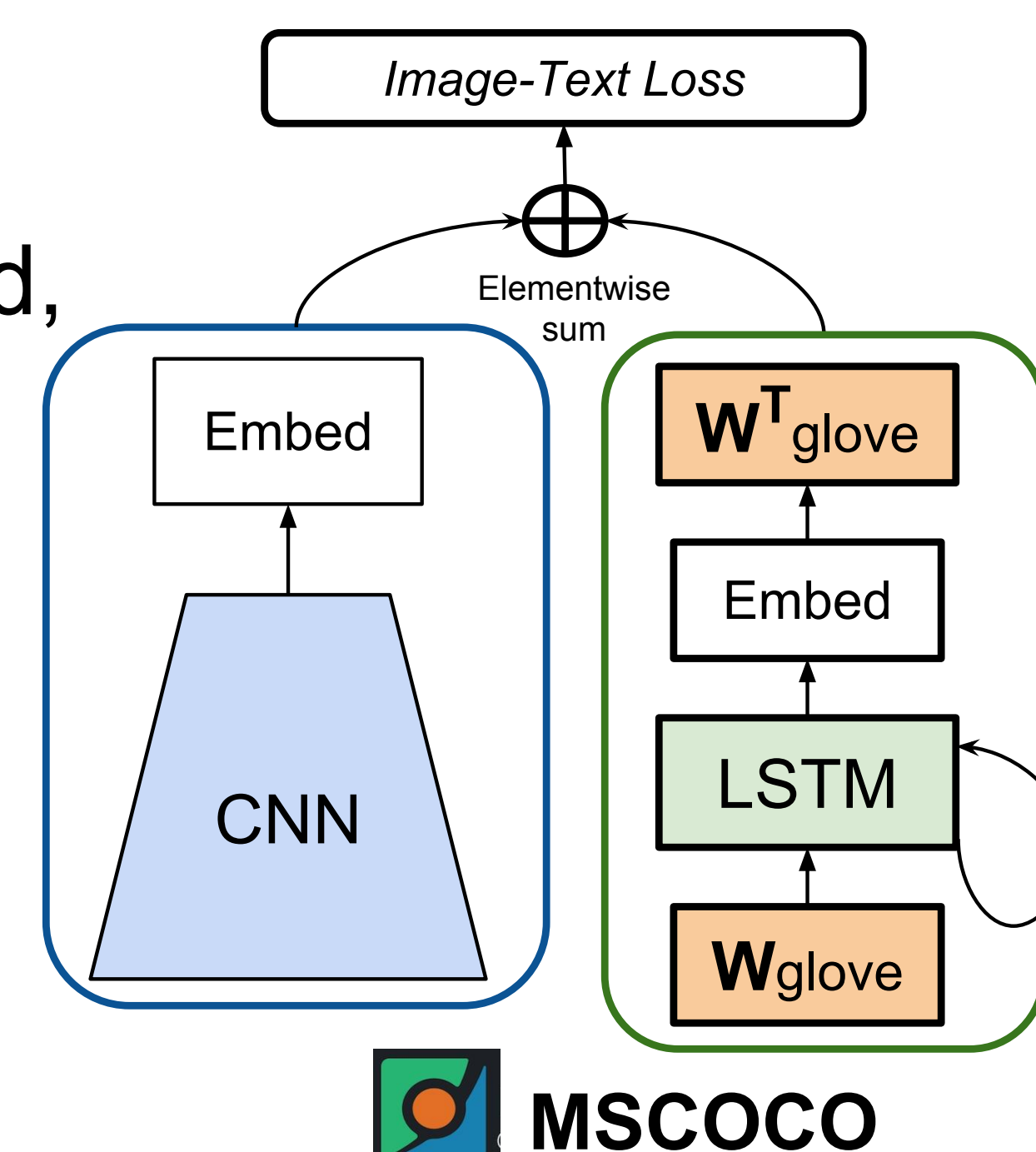


## TRAINING DATA

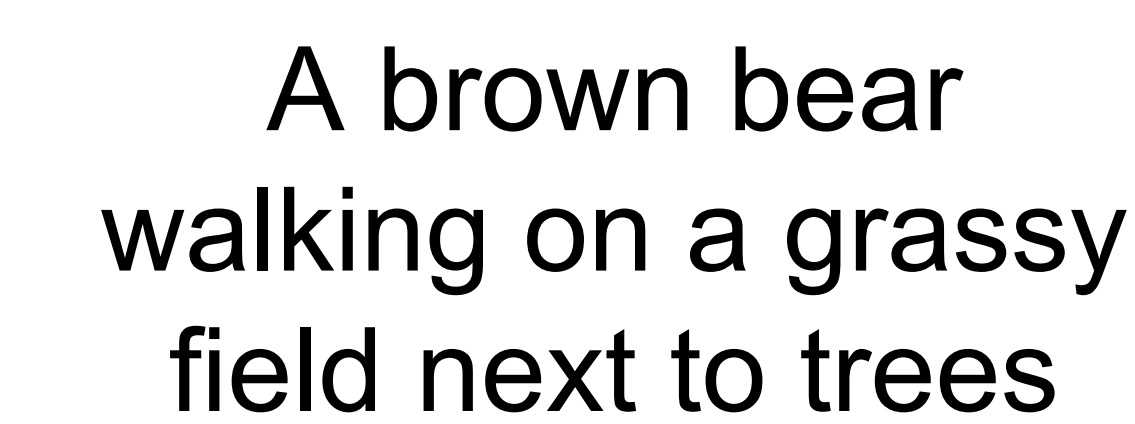
### In-Domain: COCO



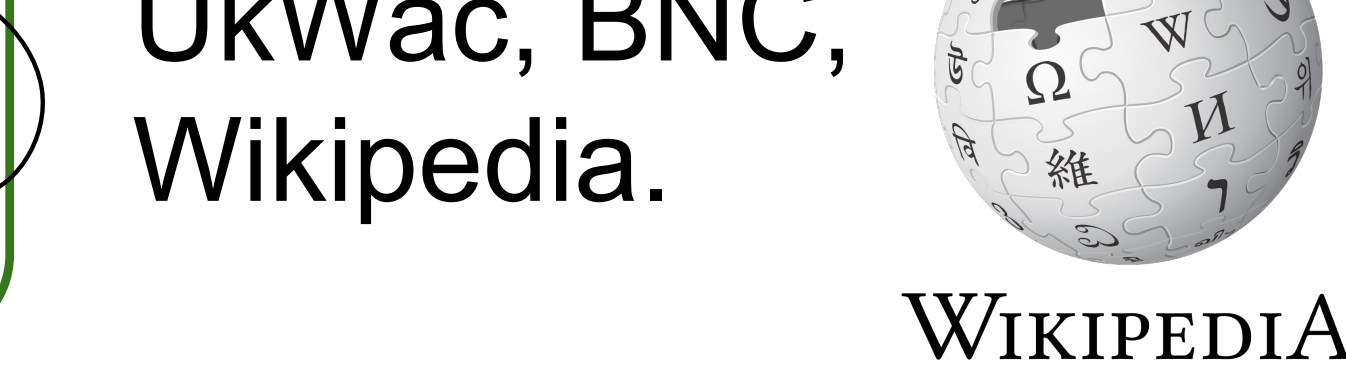
### Out of domain: IMAGENET



### In-Domain: COCO

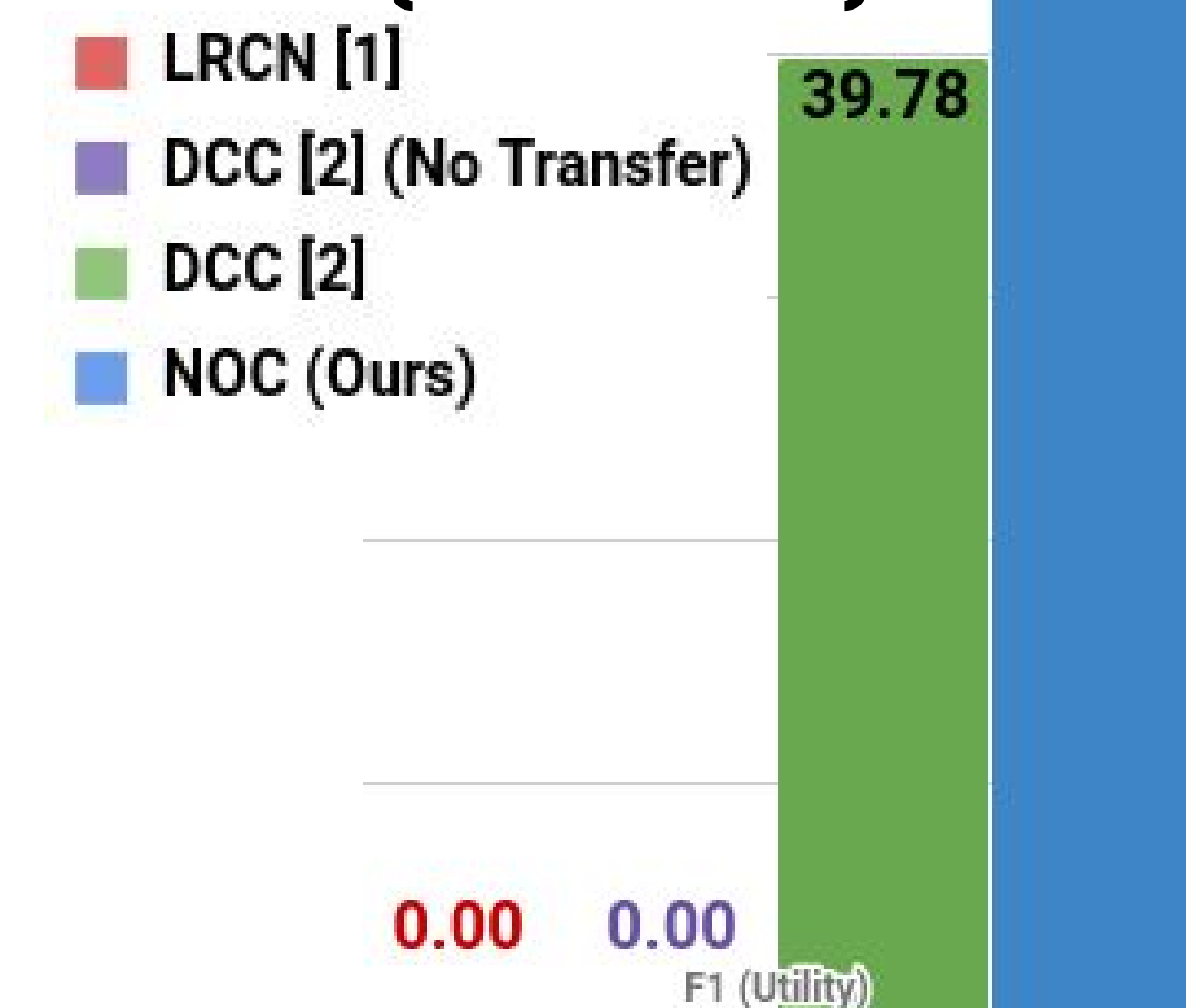


### Out of domain: Ukwac, BNC, Wikipedia.

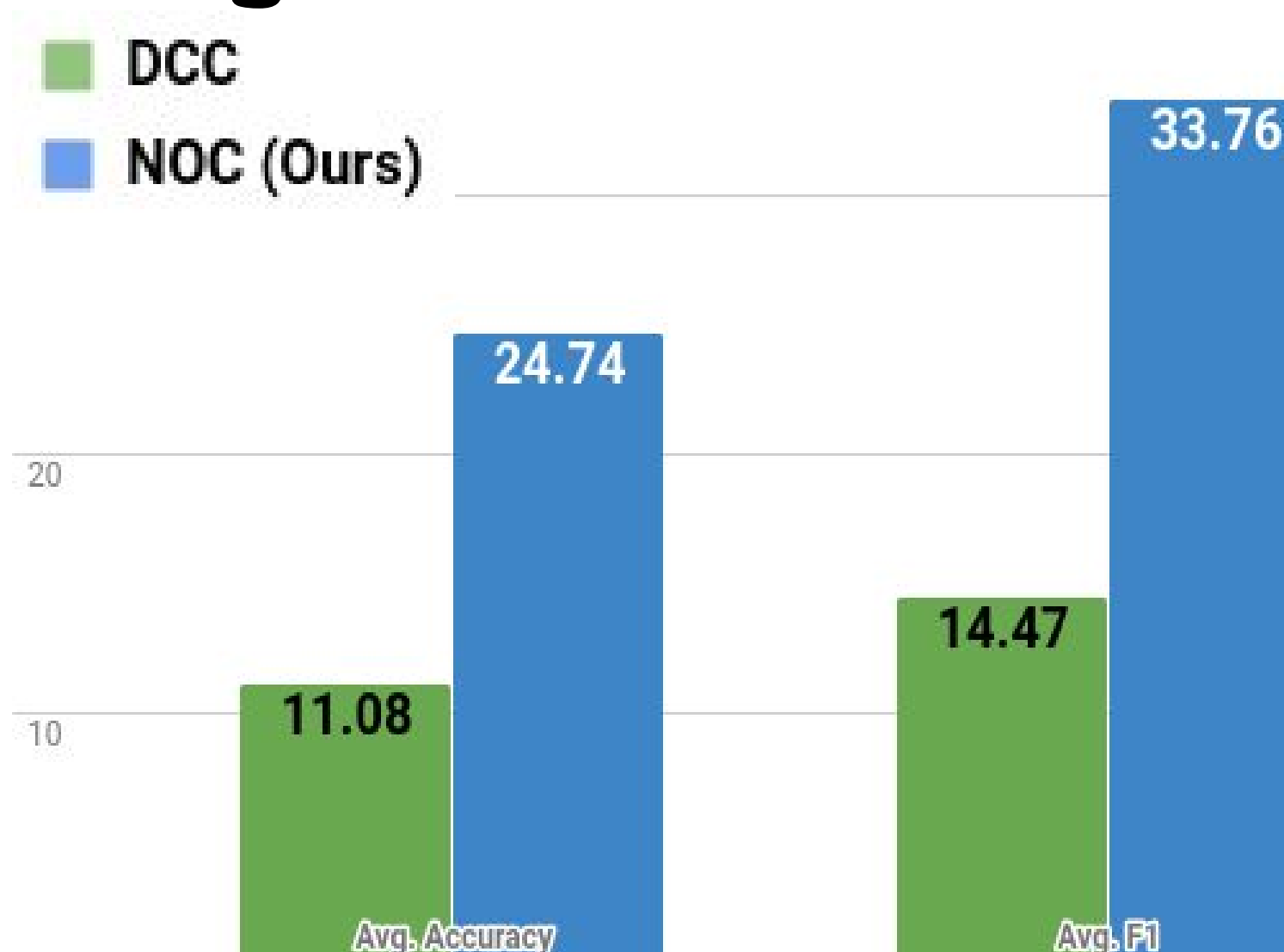


## RESULTS

### COCO (held-out)



### ImageNet



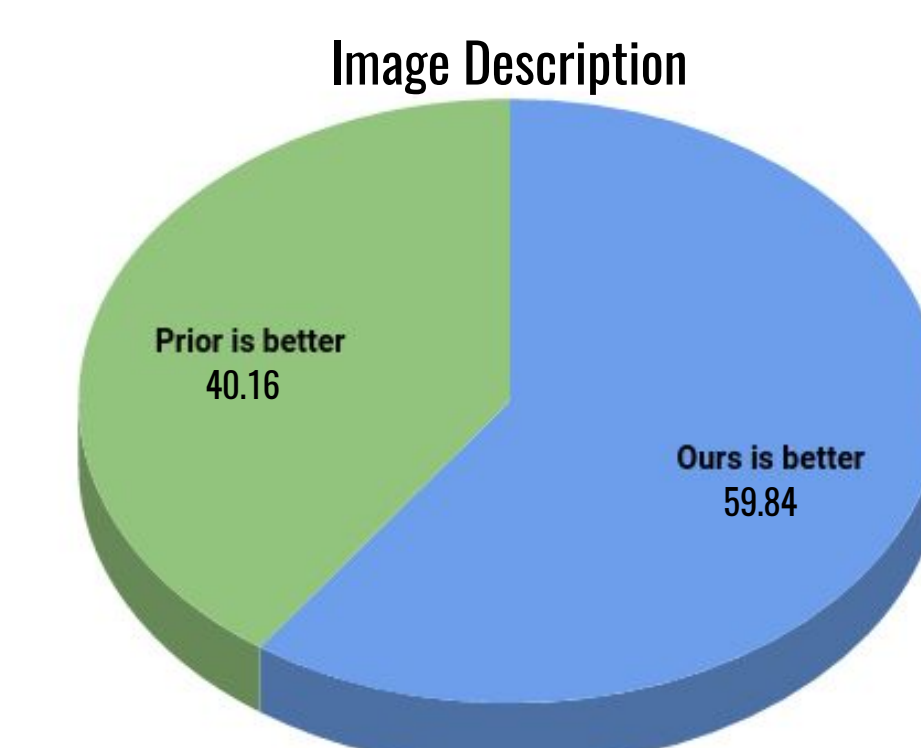
F1 (Utility): Ability to recognize and incorporate new words.

METEOR: Fluency and sentence quality.

## IMAGENET HUMAN EVAL.

**Word Incorporation:** Which model incorporates the word (name of the object) in the sentence better?

**Image Description:** Which describes the image better?



Intersection (both DCC and NOC can caption): NOC maintains descriptive quality but captions more objects.

## EXAMPLES



## CODE AND REFERENCES



Project Page

<http://vsubhashini.github.io/noc.html>

- [1] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In CVPR, 2015.
- [2] L. A. Hendricks, S. Venugopalan, M. Rohrbach, R. Mooney, K. Saenko, and T. Darrell. Deep compositional captioning: Describing novel object categories without paired training data. In CVPR, 2016.