# Captioning Images with Diverse Objects



**Subhashini Venugopalan**

Lisa Anne Hendricks

Marcus Rohrbach

Raymond Mooney

Kate Saenko

Trevor Darrell

UT Austin  UC Berkeley  Boston Univ.

# Object Recognition

Can identify hundreds of categories of objects.

IMAGENET   14M images, 22K classes [Deng et al. CVPR'09]



mammal → placental → carnivore → canine → dog → working dog → husky

vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

# Visual Description



Berkeley LRCN [Donahue et al. CVPR'15]:
A brown bear standing on top of a lush green field.

MSR CaptionBot [http://captionbot.ai/]:
A large brown bear walking through a forest.

**MSCOCO**
**80 classes**

# Novel Object Captioner (NOC)

We present Novel Object Captioner which can compose descriptions of 100s of objects in context.
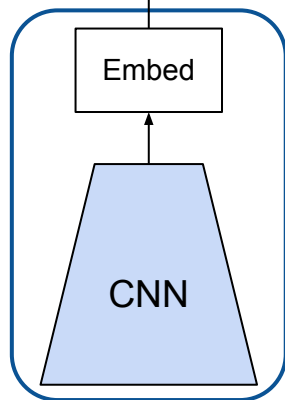
# Insights

1. Need to recognize and describe objects outside of image-caption datasets.

 ➡️ **okapi**

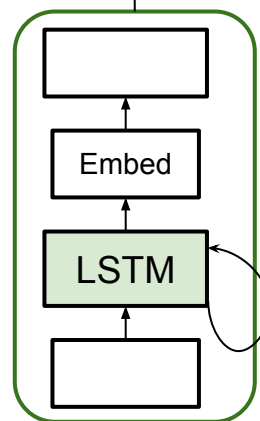# Insight 1: Train effectively on external sources

Image-Specific Loss

Embed

CNN

IMAGENET

Visual features from
unpaired image data

Language model from
unannotated text data

Text-Specific Loss

Embed

LSTM

# Insights

2. Describe unseen objects that are similar to objects seen in image-caption datasets.
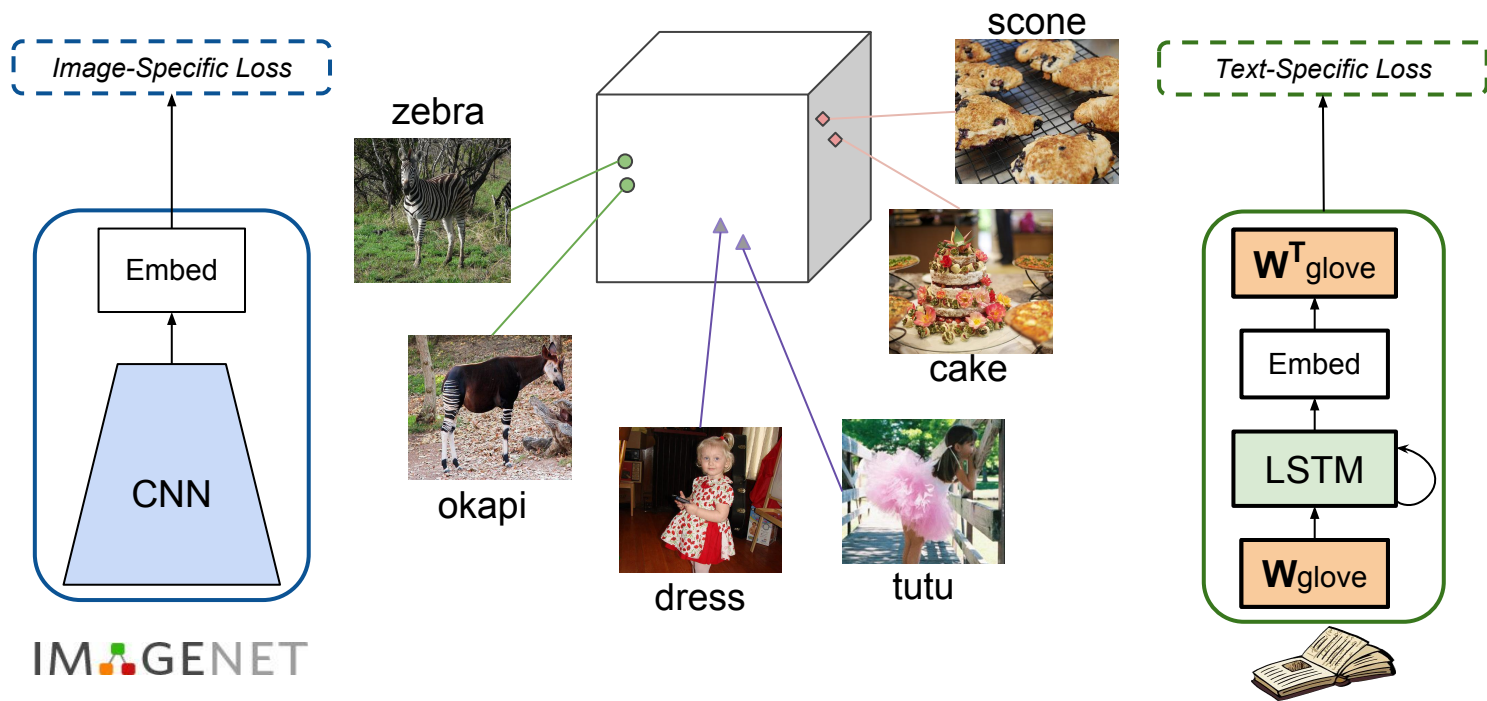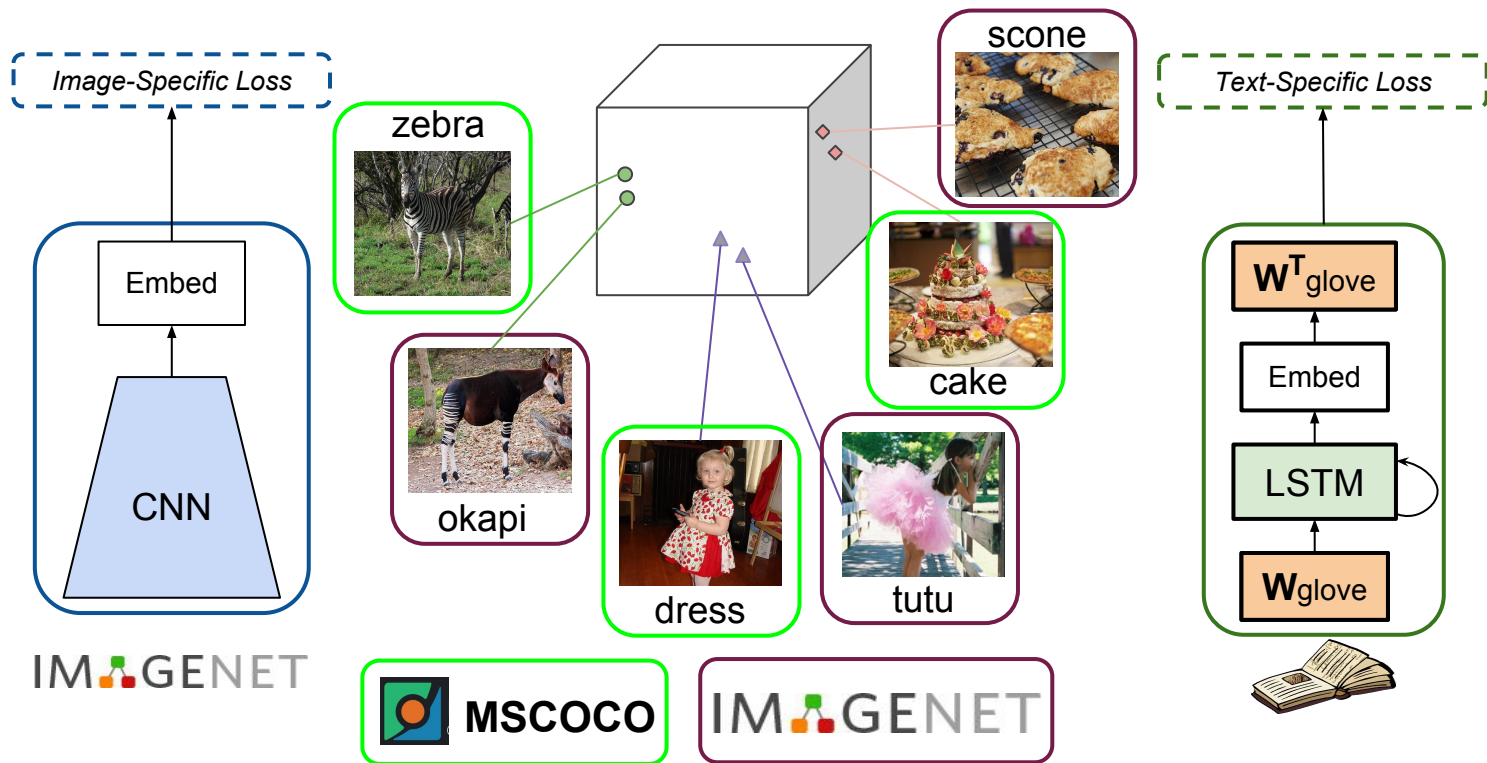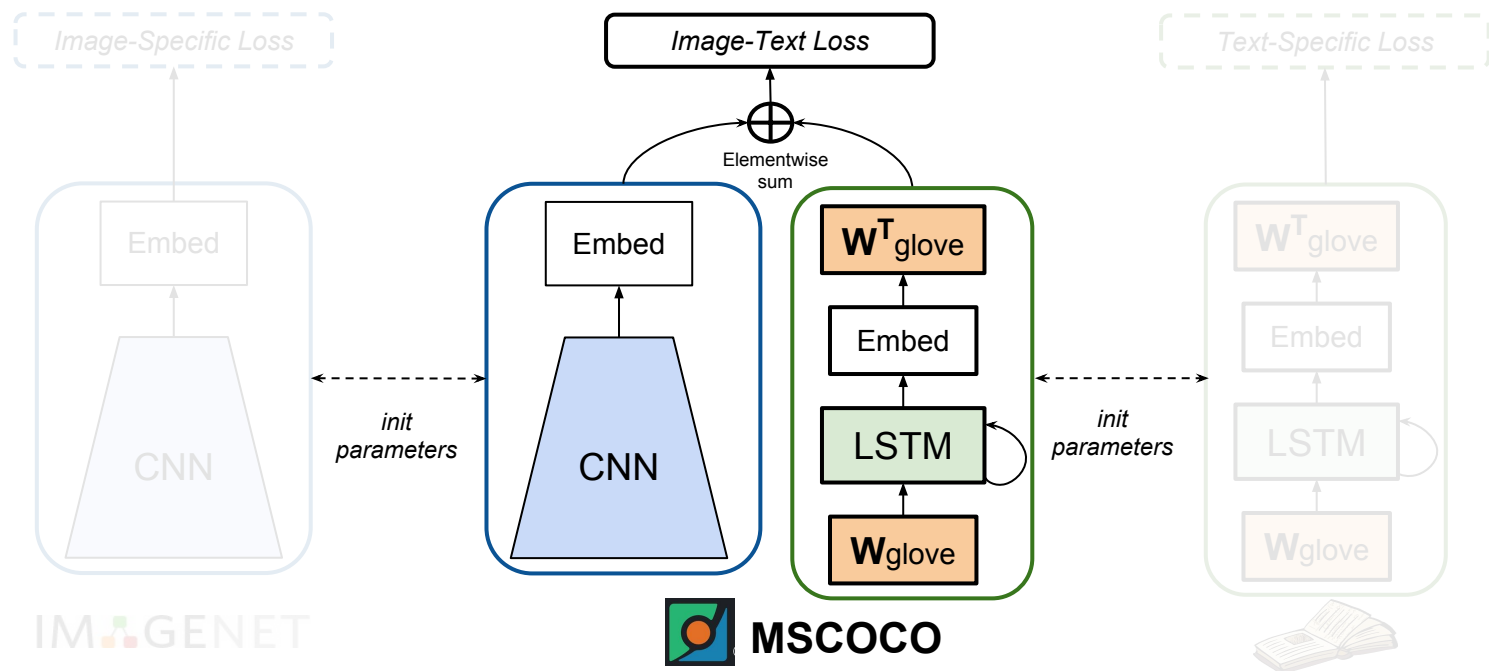


okapi ⟷ zebra

# Insight 2: Capture semantic similarity of words

# Insight 2: Capture semantic similarity of words

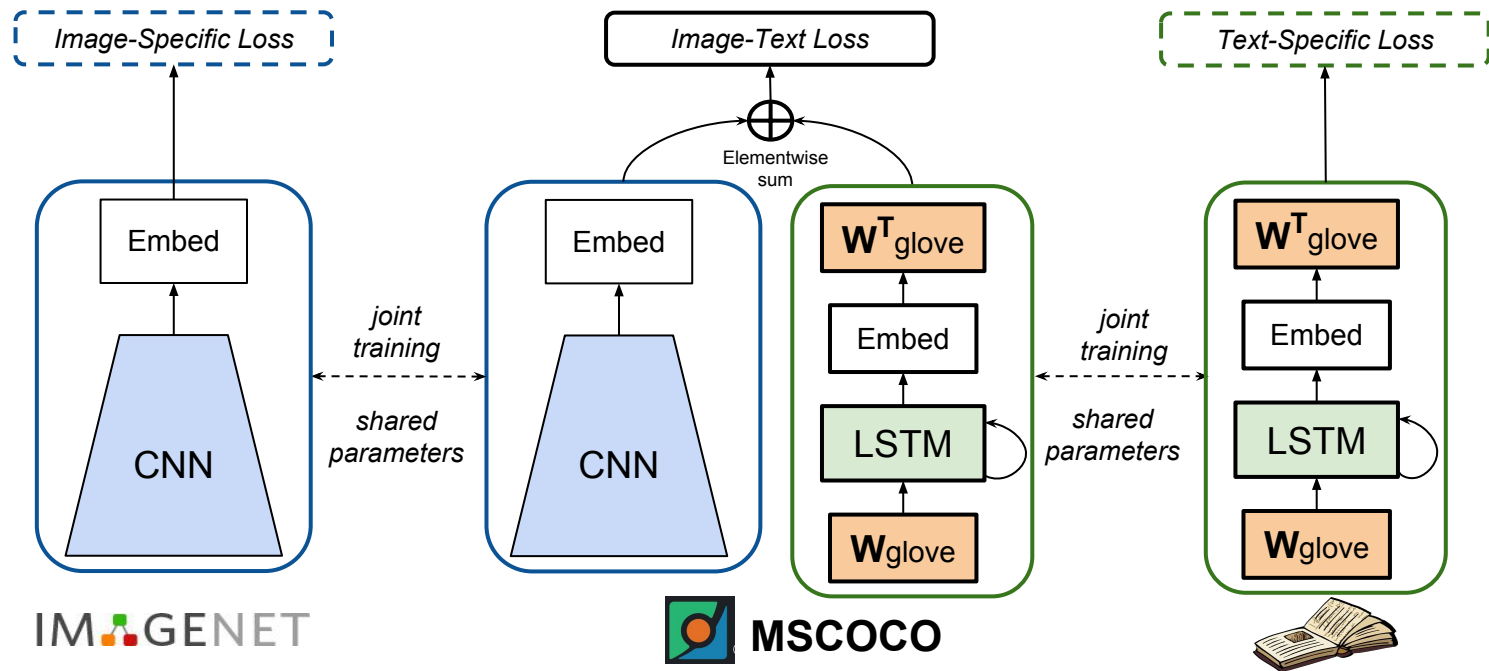# Combine to form a Caption Model



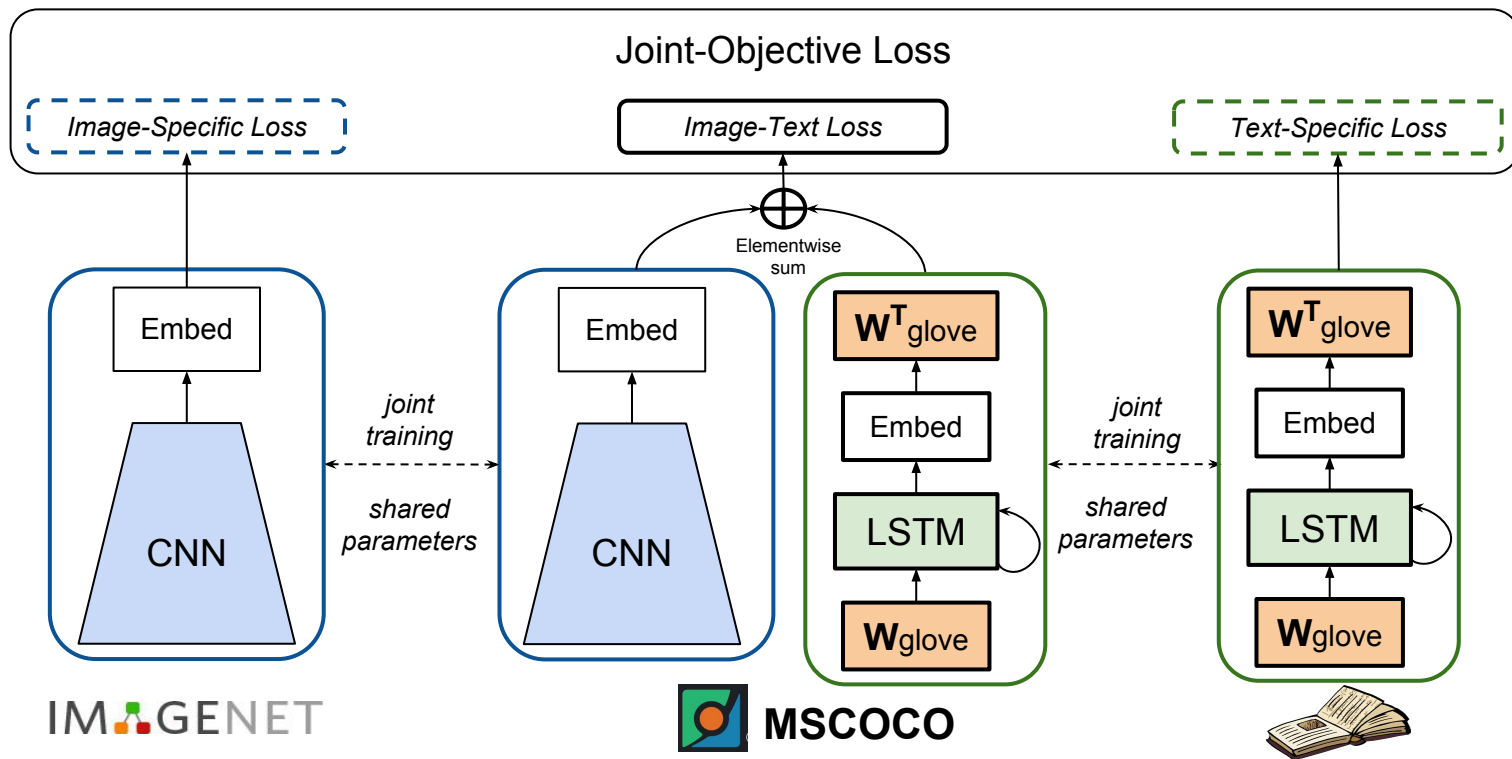Not different from existing caption models. Problem: Forgetting.

# Insights

3. Overcome "forgetting" since pre-training alone is not sufficient.

[Catastrophic Forgetting in Neural Networks. Kirkpatrick et al. PNAS 2017]

# Insight 3: Jointly train on multiple sources

# Novel Object Captioner (NOC) Model

# Visual Network



impala:0.86
green: 0.72
...
cut: 0.04

**Network:** VGG-16 with multi-label loss [sigmoid cross-entropy loss]

**Training Data:** Unpaired image data

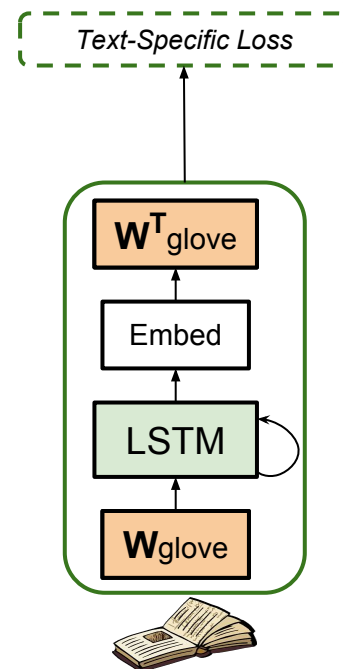**Output:** Vector with activations corresponding to scores for *words in the vocabulary*.

# Language Model

**Network:** Single LSTM layer. Predict next word $w_{t+1}$ given previous words $w_{0..t}$ $\quad p(w_{t+1} \mid w_{0..t})$

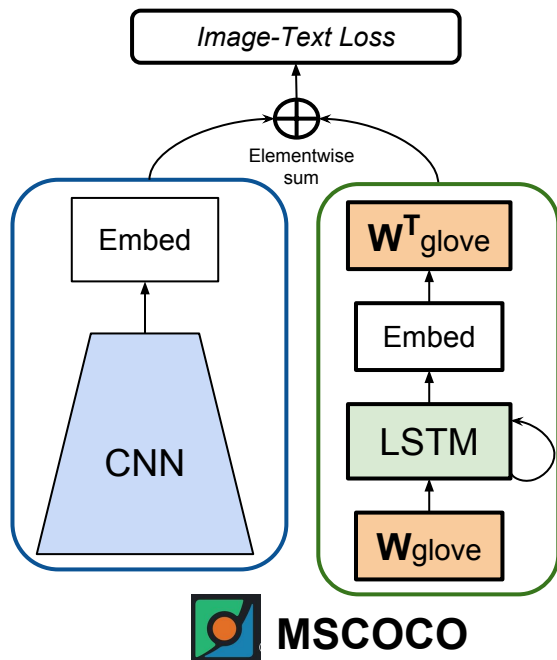$(\mathbf{W}_{glove})^{\mathbf{T}}$ : Shared weights with input embedding.

**Training Data:** Unannotated text data (BNC, ukWac, Wikipedia, Gigaword)

**Output:** Vector with activations corresponding to scores for *words in the vocabulary*.



*Text-Specific Loss*

$\mathbf{W}^{\mathbf{T}}_{glove}$

Embed

LSTM

$\mathbf{W}_{glove}$

# Caption Network

**Network:** Combine output of the visual and text networks. (softmax + cross-entropy loss)

# Caption Model

**Training Data:**
COCO images with
multiple labels

bear, brown, field,
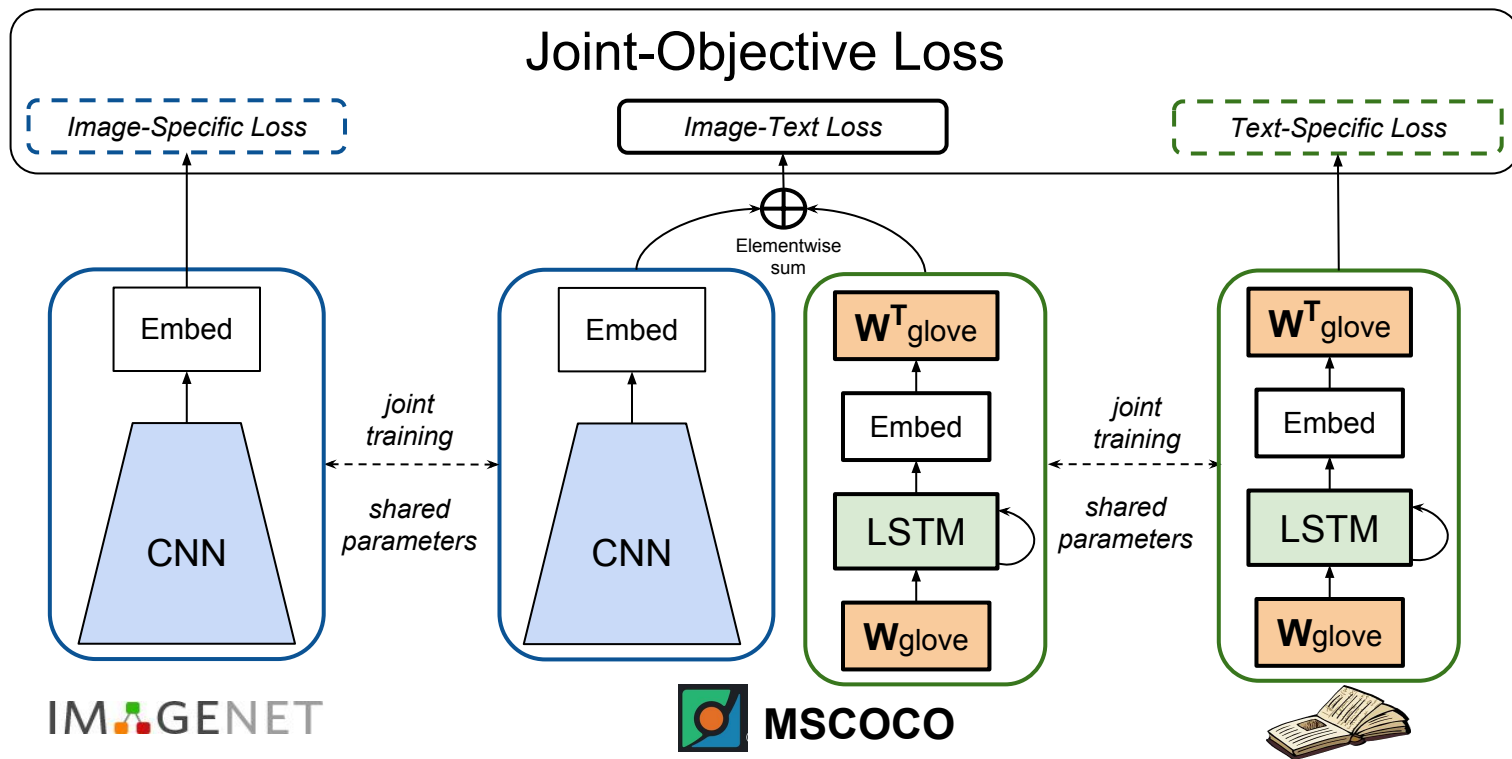grassy, trees,
walking

Image-Text Loss

$\oplus$ Elementwise sum

Embed

CNN

$\mathbf{W^T}_{glove}$

Embed

LSTM

$\mathbf{W}_{glove}$

**MSCOCO**

**Training Data:**
Captions from
COCO

A brown bear
walking on a grassy
field next to trees

# NOC Model: Train simultaneously



Joint-Objective Loss

Image-Specific Loss

Image-Text Loss

Text-Specific Loss

Elementwise sum

Embed

Embed

$W^T_{glove}$

$W^T_{glove}$

Embed

Embed

*joint training*

*joint training*

LSTM

LSTM

CNN

CNN

*shared parameters*

*shared parameters*

$W_{glove}$

$W_{glove}$

IMAGENET

MSCOCO

# Evaluation

- Empirical: COCO held-out objects
  - **In-domain [Use images from COCO]**
  - Out-of-domain  [Use imagenet images for held-out concepts]
- Ablations
  - Embedding & Joint training contribution
- ImageNet
  - Quantitative
  - **Human Evaluation - Objects not in COCO**
  - Rare objects in COCO

# Evaluation

- Empirical: COCO held-out objects
  - **In-domain [Use images from COCO]**
  - Out-of-domain  [Use imagenet images for held-out concepts]
- Ablations
  - Embedding & Joint training contribution
- ImageNet
  - Quantitative
  - **Human Evaluation - Objects not in COCO**
  - Rare objects in COCO

# Empirical Evaluation: COCO dataset In-Domain setting

## MSCOCO Unpaired Image Data



*Elephant, Galloping, Green, Grass*



*People, Playing, Ball, Field*



*Black, Train, Tracks*



*Eat, Pizza*



*Kitchen, Microwave*

## MSCOCO Paired Image-Sentence Data



*"An elephant galloping in the green grass"*



*"Two people playing ball in a field"*



*"A black train stopped on the tracks"*



*"Someone is about to eat some pizza"*



*"A kitchen counter with a microwave on it"*

## MSCOCO Unpaired Text Data

*"An elephant galloping in the green grass"*

*"Two people playing ball in a field"*

*"A black train stopped on the tracks"*

*"Someone is about to eat some pizza"*

*"A microwave is sitting on top of a kitchen counter "*

# Empirical Evaluation: COCO <span style="color:red">heldout</span> dataset

## MSCOCO Unpaired Image Data


*Elephant, Galloping, Green, Grass*


*People, Playing, Ball, Field*


*Black, Train, Tracks*


*Pizza*


*Microwave*

## MSCOCO Paired Image-Sentence Data


*"An elephant galloping in the green grass"*


*"Two people playing ball in a field"*


*"A black train stopped on the tracks"*


*"Someone is about to eat some pizza"*


*"A kitchen counter with a microwave on it"*

## MSCOCO Unpaired Text Data

*"An elephant galloping in the green grass"*

*"Two people playing ball in a field"*

*"A black train stopped on the tracks"*

*"A white plate topped with cheesy pizza and toppings."*

*"A white refrigerator, stove, oven dishwasher and microwave"*

<span style="color:red">Held-out</span>

22

# Empirical Evaluation: COCO



**MSCOCO Unpaired Image Data**

*Two, elephants, Path, walking*

*Baseball, batting, boy, swinging*

*Black, Train, Tracks*

*Pizza*

*Microwave*

**MSCOCO Paired Image-Sentence Data**

*"An elephant galloping in the green grass"*

*"Two people playing ball in a field"*

*"A black train stopped on the tracks"*

**MSCOCO Unpaired Text Data**

*"A small elephant standing on top of a dirt field"*

*"A hitter swinging his bat to hit the ball"*

*"A black train stopped on the tracks"*

*"A white plate topped with cheesy pizza and toppings."*

*"A white refrigerator, stove, oven dishwasher and microwave"*

- CNN is pre-trained on ImageNet

# Empirical Evaluation: Metrics

**F1** (Utility)**:** Ability to recognize and incorporate new words.
(Is the word/object mentioned in the caption?)

**METEOR:** Fluency and sentence quality.

# Empirical Evaluation: Baselines

LRCN [1]
DCC [2] (No Transfer)
DCC [2]
NOC (Ours)

**LRCN [1]:** Does not caption novel objects.

**DCC [2]** : Copies parameters for the novel
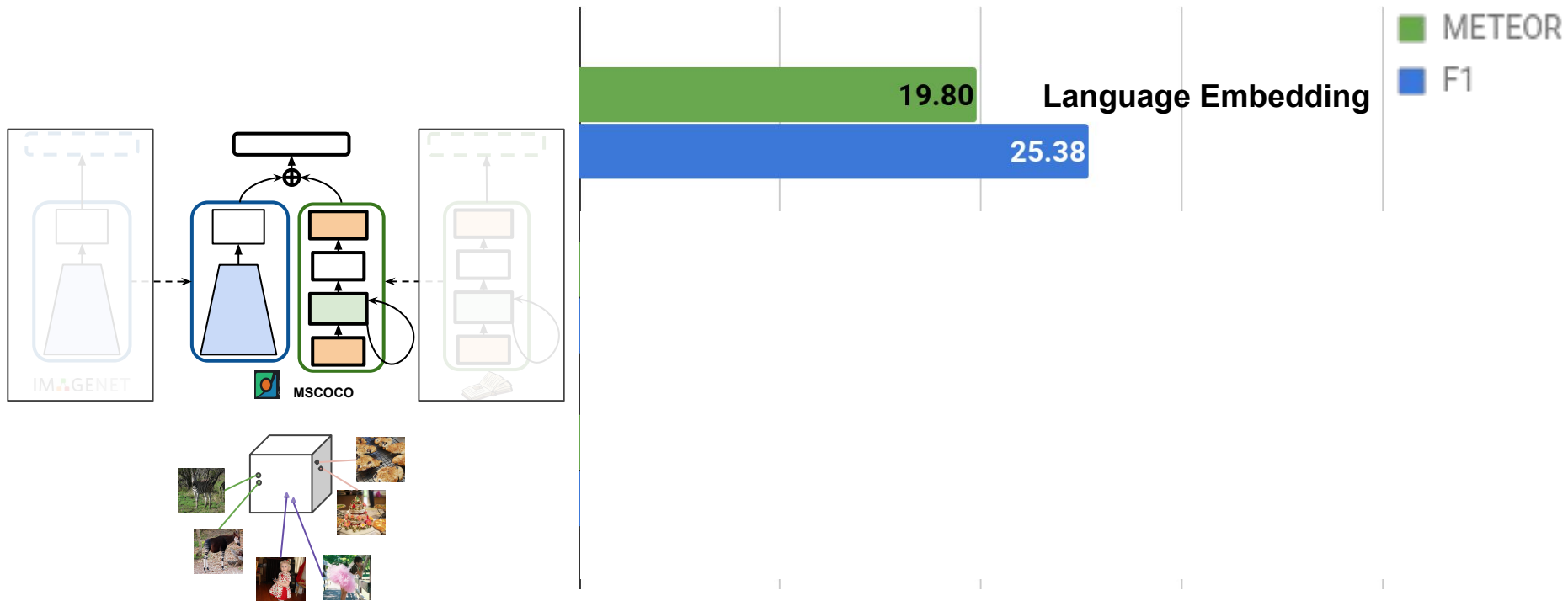object from a similar object seen
in training. (also not end-to-end)

[1] J. Donahue, L.A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, T. Darrell. CVPR'15
[2] L.A. Hendricks, S. Venugopalan, M. Rohrbach, R. Mooney, K. Saenko, T. Darrell CVPR'16

# Empirical Evaluation: Results



LRCN [1]
DCC [2] (No Transfer)
DCC [2]
NOC (Ours)

F1 (Utility): 0.00, 0.00, 39.78, 48.79

METEOR (Fluency): 19.33, 19.90, 21.00, 21.32

[1] J. Donahue, L.A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, T. Darrell. CVPR'15
[2] L.A. Hendricks, S. Venugopalan, M. Rohrbach, R. Mooney, K. Saenko, T. Darrell CVPR'16

# Ablations



IM🐾GENET

MSCOCO

Evaluated on held-out
COCO objects.

# Ablation: Language Embedding



METEOR

F1

Language Embedding

19.80

25.38

# Ablation: Freeze CNN after pre-training



METEOR
F1

Language Embedding
- 19.80
- 25.38

Frozen CNN
- 18.91
- 39.70

[Catastrophic forgetting in Neural Networks Kirkpatrick et al. PNAS 2017]

# Ablation: Joint Training



- Language Embedding: METEOR 19.80, F1 25.38
- Frozen CNN: METEOR 18.91, F1 39.70
- Joint Training: METEOR 21.32, F1 48.79

# ImageNet: Human Evaluations

- **ImageNet:** 638 object classes not mentioned in COCO

NOC can describe 582 object classes
(60% more objects than prior work)

# ImageNet: Human Evaluations

- **ImageNet:** 638 object classes not mentioned in COCO

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

- **Image Description:** Which sentence (model) describes the image better?

# ImageNet: Human Evaluations



**Word Incorporation**

**Image Description**

# Qualitative Evaluation: ImageNet



**Instruments**

A man holding a **banjo** in a park.

A large **chime** hanging on a metal pole

**Vehicles**

A **snowplow** truck driving down a snowy road.

A group of people standing around a large white **warship**.

**Land Animals**

A **okapi** is in the grass with a **okapi**.

A small brown and white **jackal** is standing in a field.

**Household**

A large metal **candelabra** next to a wall.

A black and white photo of a **corkscrew** and a **corkscrew**.

# Qualitative Evaluation: ImageNet

**Birds**



A small **pheasant** is standing in a field.



A **osprey** flying over a large grassy area.

**Outdoors**



A large **glacier** with a mountain in the background.



A group of people are sitting in a **baobab**.

**Water Animals**



A **humpback** is flying over a large body of water.



A man is standing on a beach holding a **snapper**.

**Misc**



A table with a **cauldron** in the dark.



A woman is posing for a picture with a **chiffon** dress.

# Qualitative Examples: Errors

*Balaclava (*n02776825*)*          Error: Repetition
NOC: A **balaclava** black and white photo of a man in a **balaclava**.

*Sunglass (n04355933)*          Error: Grammar
NOC: A **sunglass** mirror reflection of a mirror in a mirror.

*Gymnast (n10153594)*          Error: Gender, Hallucination
NOC: A man **gymnast** in a blue shirt doing a trick on a skateboard.

*Cougar (n02125311)*          Error: Description
NOC: A **cougar** with a **cougar** in its mouth.

# Novel Object Captioner - Take away

Semantic embeddings and joint training to caption 100s of objects.



A **okapi** standing in the middle of a field.

# Poster 11



Captioning Images with Diverse Objects.

Subhashini Venugopalan[1], Lisa Anne Hendricks[2], Marcus Rohrbach[2,3], Raymond Mooney[1], Trevor Darrell[2], Kate Saenko[4]

[1] UT-Austin [2] UC-Berkeley [3] Facebook AI Research [4] Boston University

38