



Topics in AI (CPSC 532S): Multimodal Learning with Vision, Language and Sound

Lecture 14: Unsupervised Learning, Autoencoders [Part 3]

Logistics

- **Project pitches** next week (**November 1 & 3**)
9 groups per class (~8 minutes / group, 5-6 min presentation + questions)
- Project proposals are **NOT** due next week (due **November 15th**)
- **Assignment 4** — Remember you only need to do 1 PART

Final **Project** (40% of grade total)

- Group project (groups of 3 are encouraged, but fewer maybe possible)
- Groups are self-formed, you will not be assigned to a group
- You need to come up with a project proposal and then work on the project as a group (each person in the group gets the same grade for the project)
- Project needs to be **research** oriented (not simply implementing an existing paper); you can use code of existing paper as a starting point though

Project proposal + class presentation: 15%
Project + final presentation (during finals week): 25%

Correlated Representations vs. Joint Embeddings

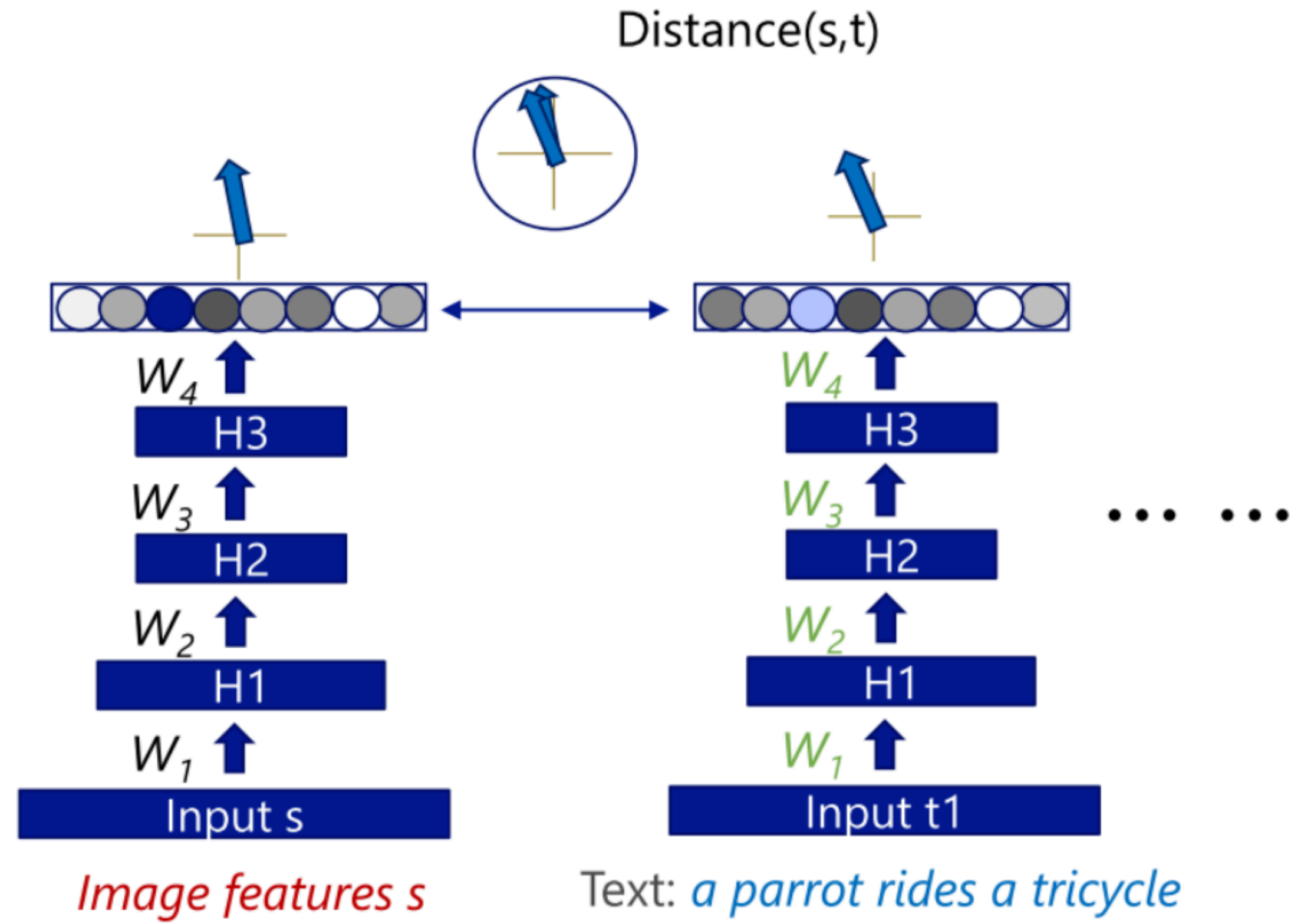
Correlated Representations: Find representations $f_1(\mathbf{x}_1)$, $f_2(\mathbf{x}_2)$ for each view that maximize correlation:

$$\mathbf{corr}(f_1(\mathbf{x}_1), f_2(\mathbf{x}_2)) = \frac{\mathbf{cov}(f_1(\mathbf{x}_1), f_2(\mathbf{x}_2))}{\sqrt{\mathbf{var}(f_1(\mathbf{x}_1)) \cdot \mathbf{var}(f_2(\mathbf{x}_2))}}$$

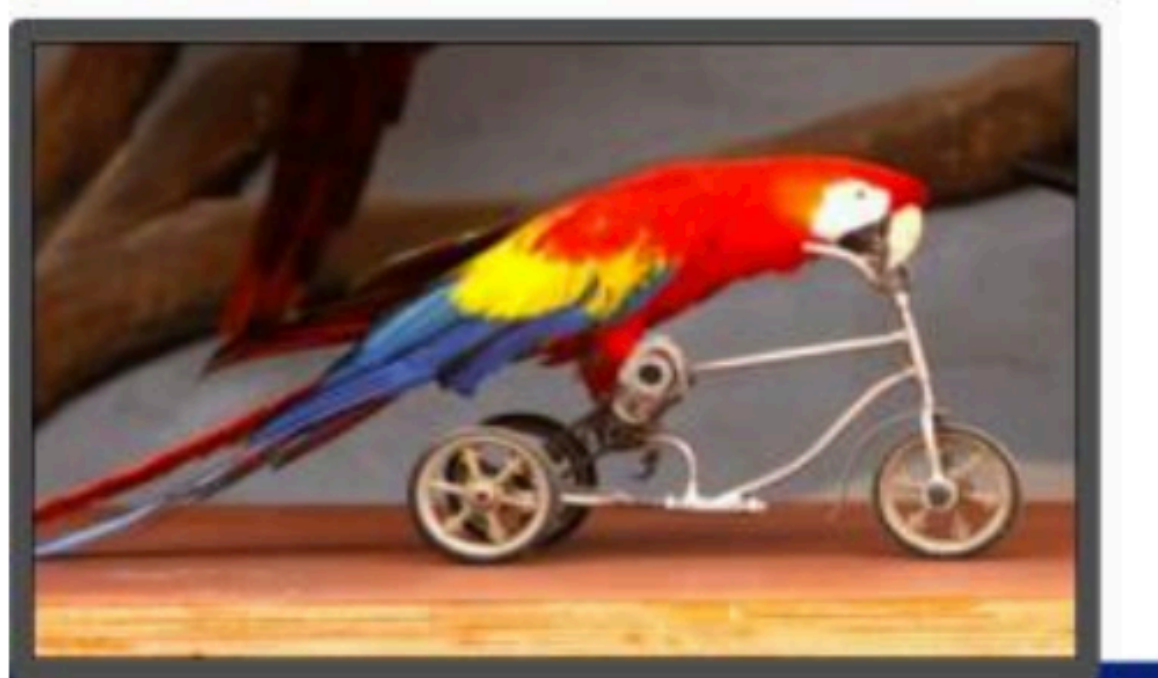
Joint Embeddings: Models that minimize distance between ground truth pairs of samples:

$$\min_{f_1, f_2} D \left(f_1(\mathbf{x}_1^{(i)}), f_2(\mathbf{x}_2^{(i)}) \right)$$

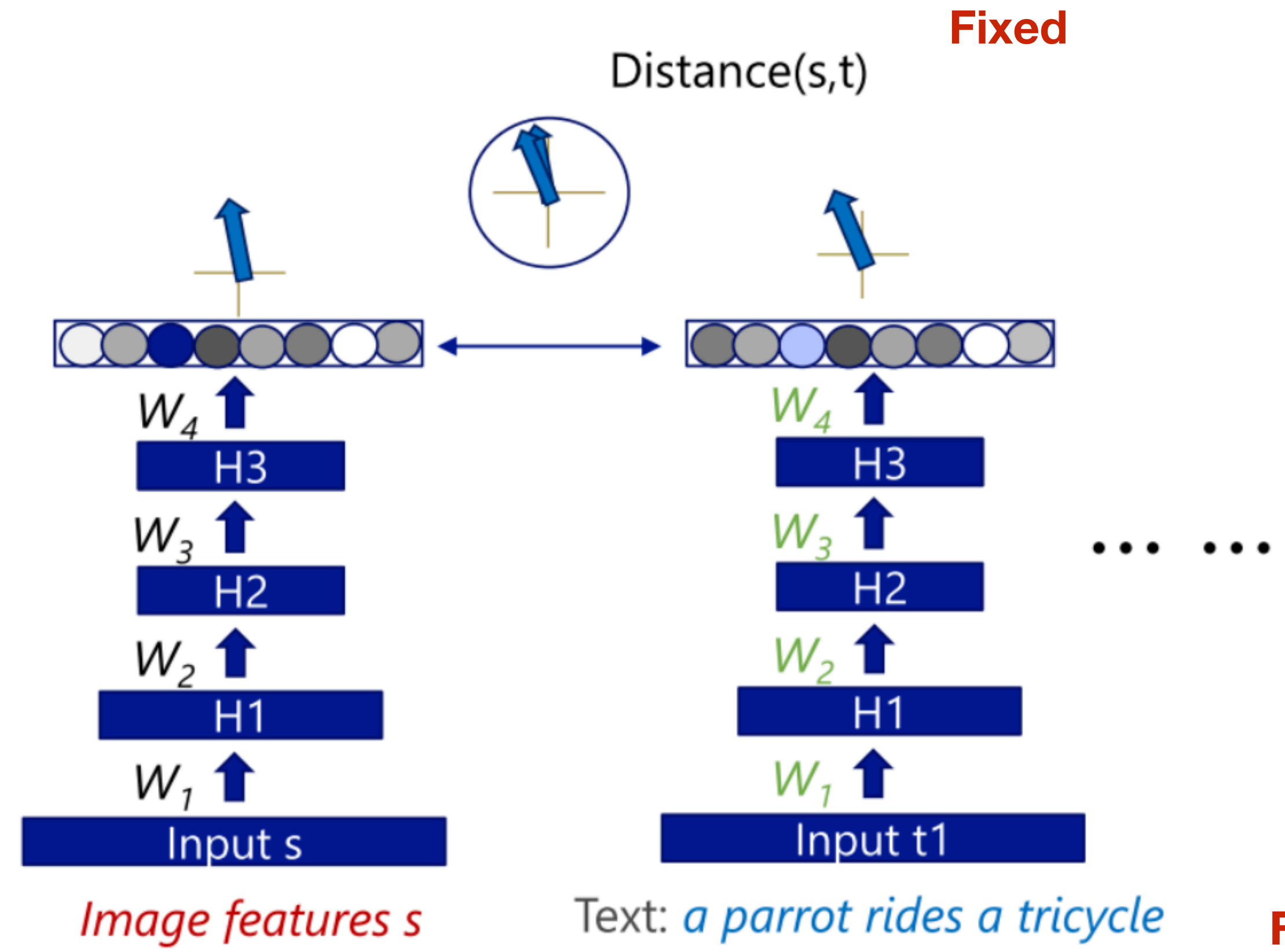
Joint Embeddings



Joint Embeddings











Fixed



Joint Embeddings

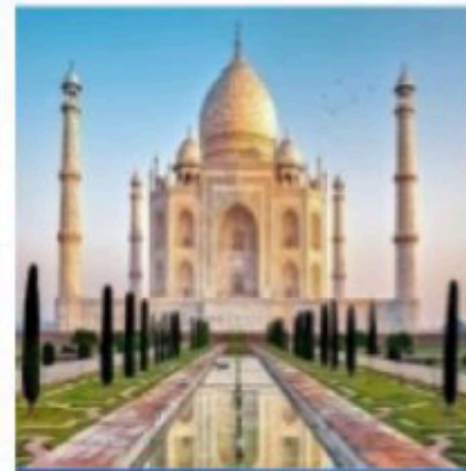
Nearest images

	- blue + red =	
	- blue + yellow =	
	- yellow + red =	
	- white + red =	

[Kiros et al., Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, 2014]

Joint Embeddings

Nearest images



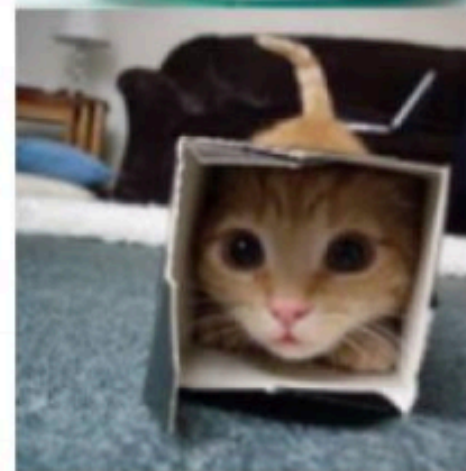
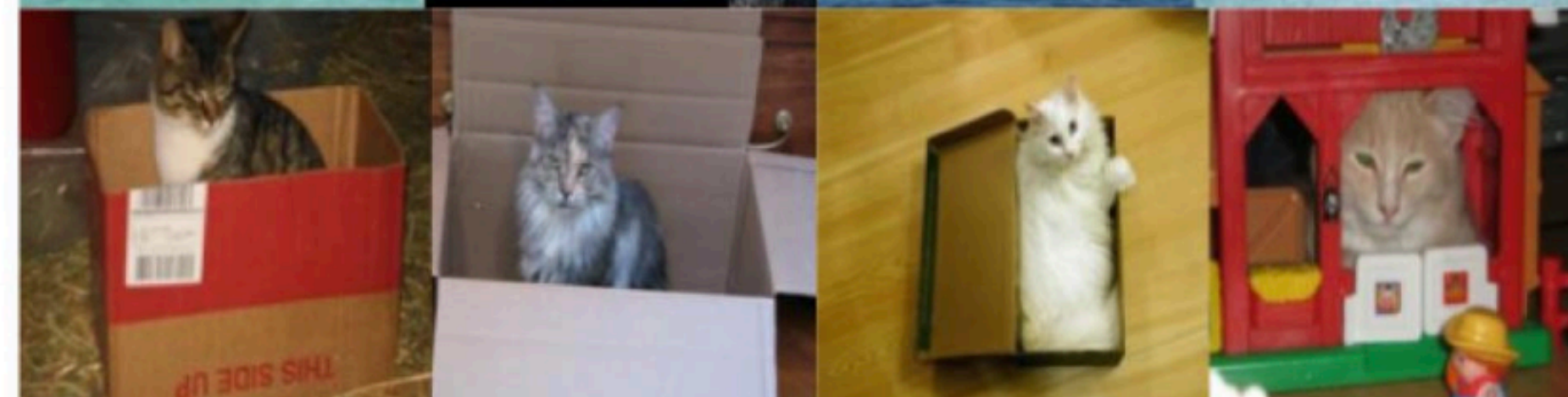
- day + night =



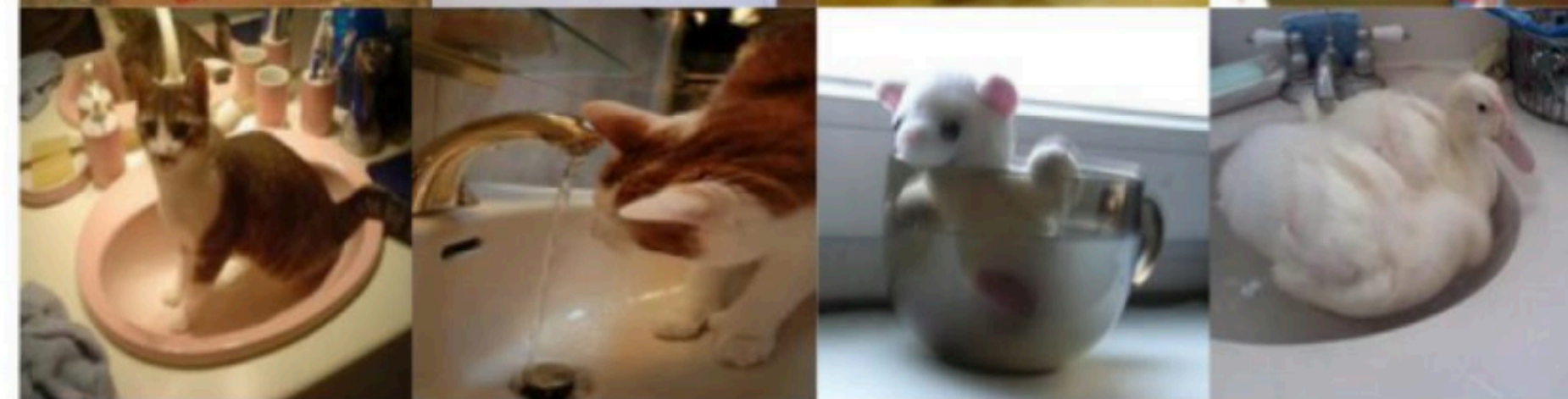
- flying + sailing =



- bowl + box =



- box + bowl =



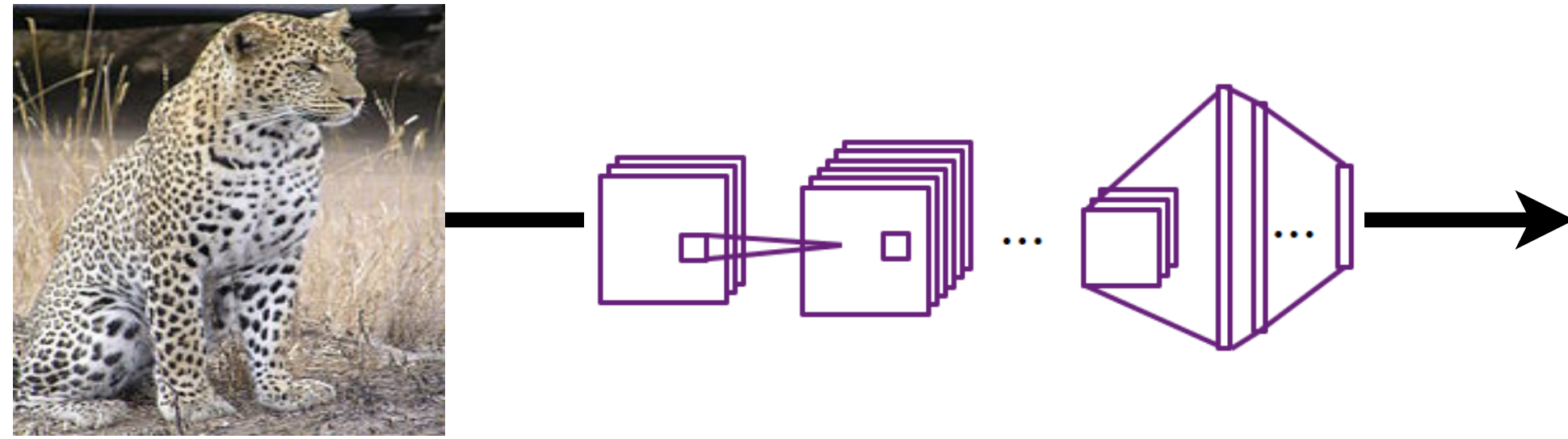
Object Classification



Category	Prediction
Dog	No
Cat	No
Couch	No
Flowers	No
Leopard	Yes
...	...

Problem: For each image predict which category it belongs to out of a fixed set

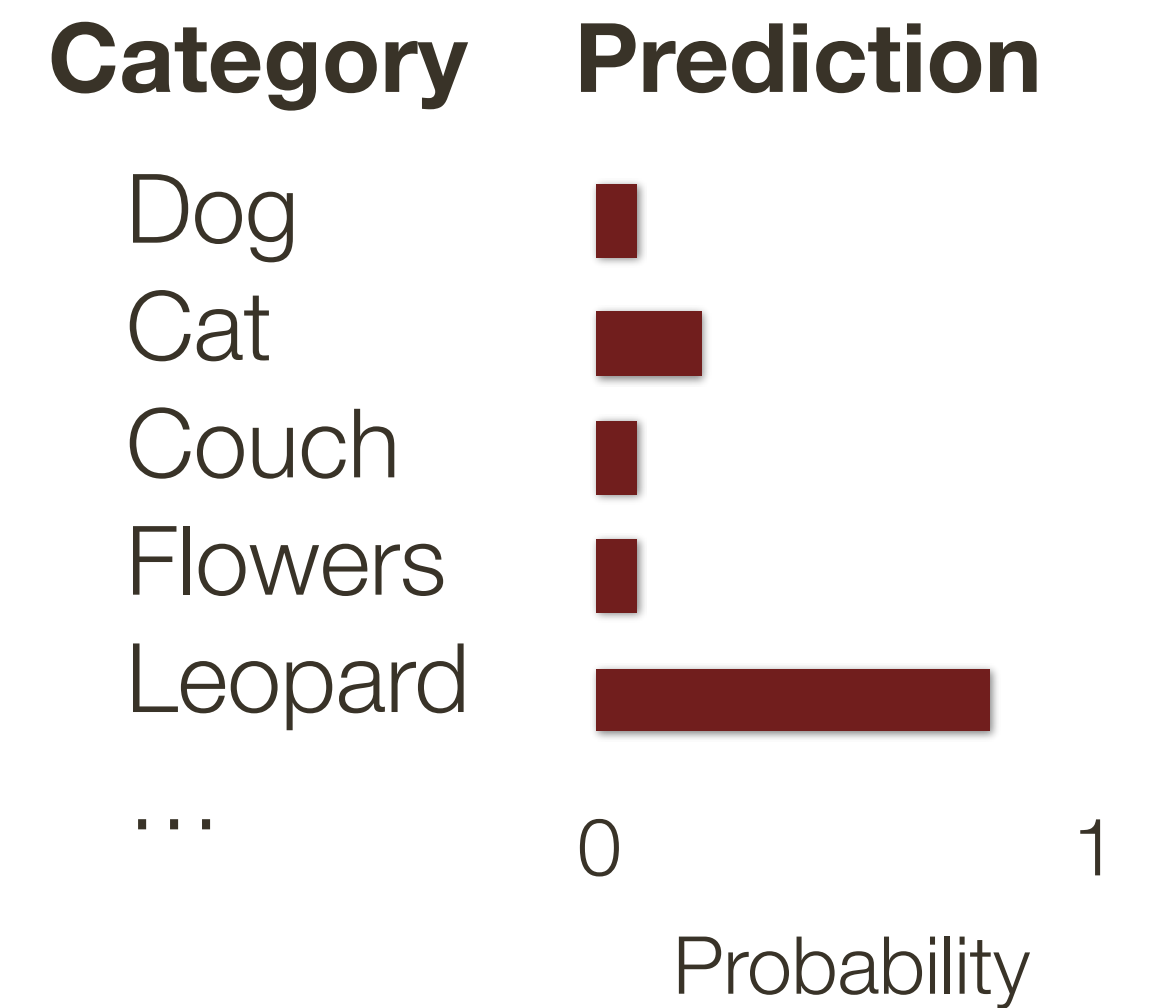
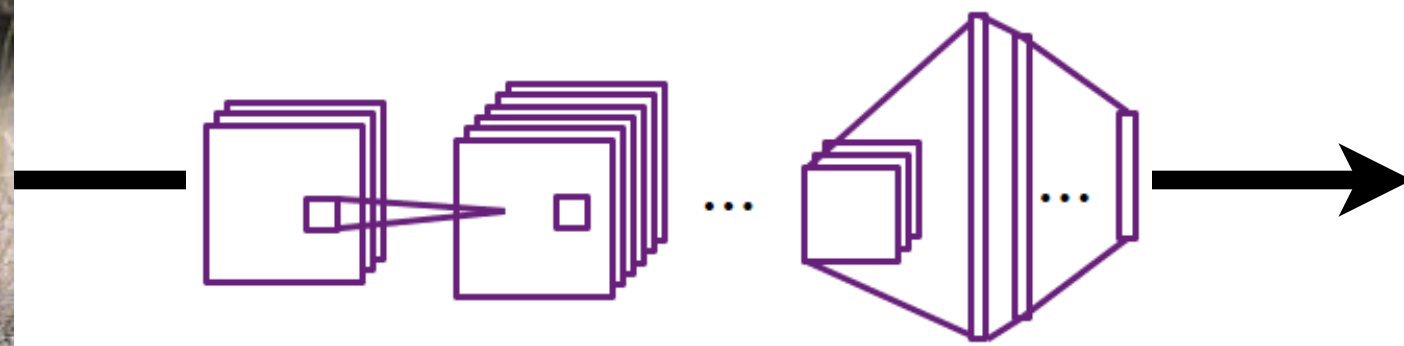
Object Classification



Category	Prediction
Dog	No
Cat	No
Couch	No
Flowers	No
Leopard	Yes
...	...

Problem: For each image predict which category it belongs to out of a fixed set

Object Classification



Problem: For each image predict which category it belongs to out of a fixed set

Discriminative Embeddings

Images and **class labels** are embedded into the same space

\mathbb{R}^d

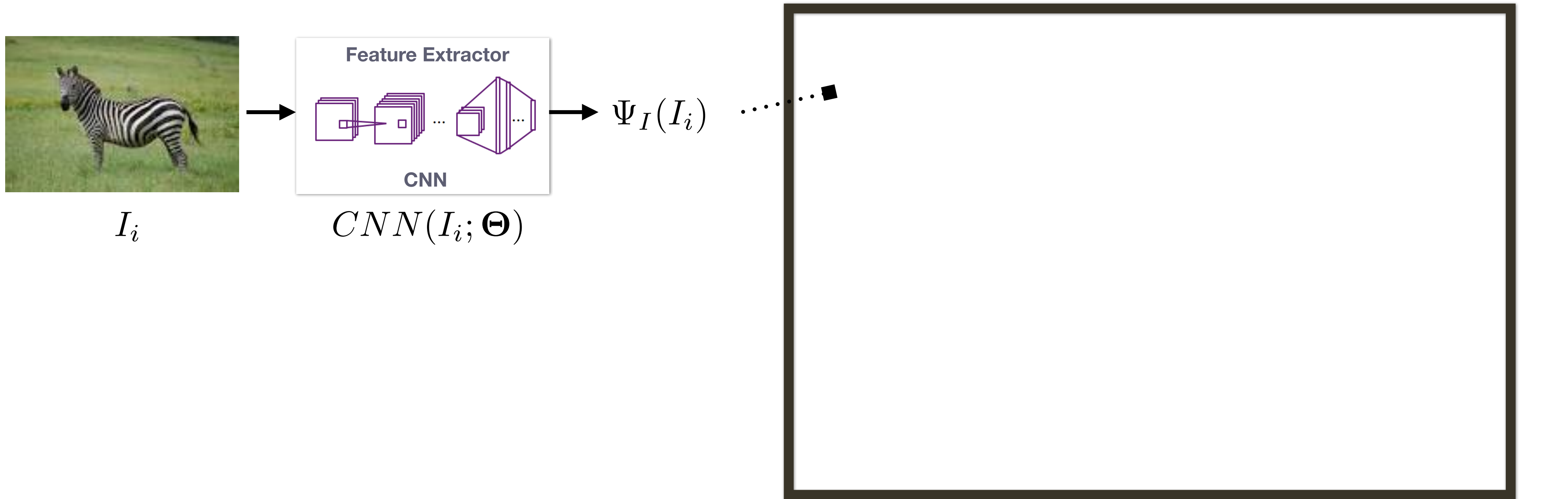


Discriminative Embeddings

Images and **class labels** are embedded into the same space

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$



Discriminative Embeddings

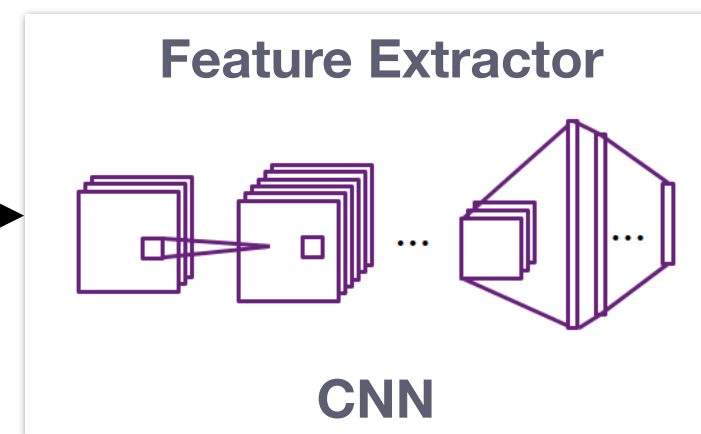
Images and **class labels** are embedded into the same space

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$



I_i



$\text{CNN}(I_i; \Theta)$

$\Psi_I(I_i)$



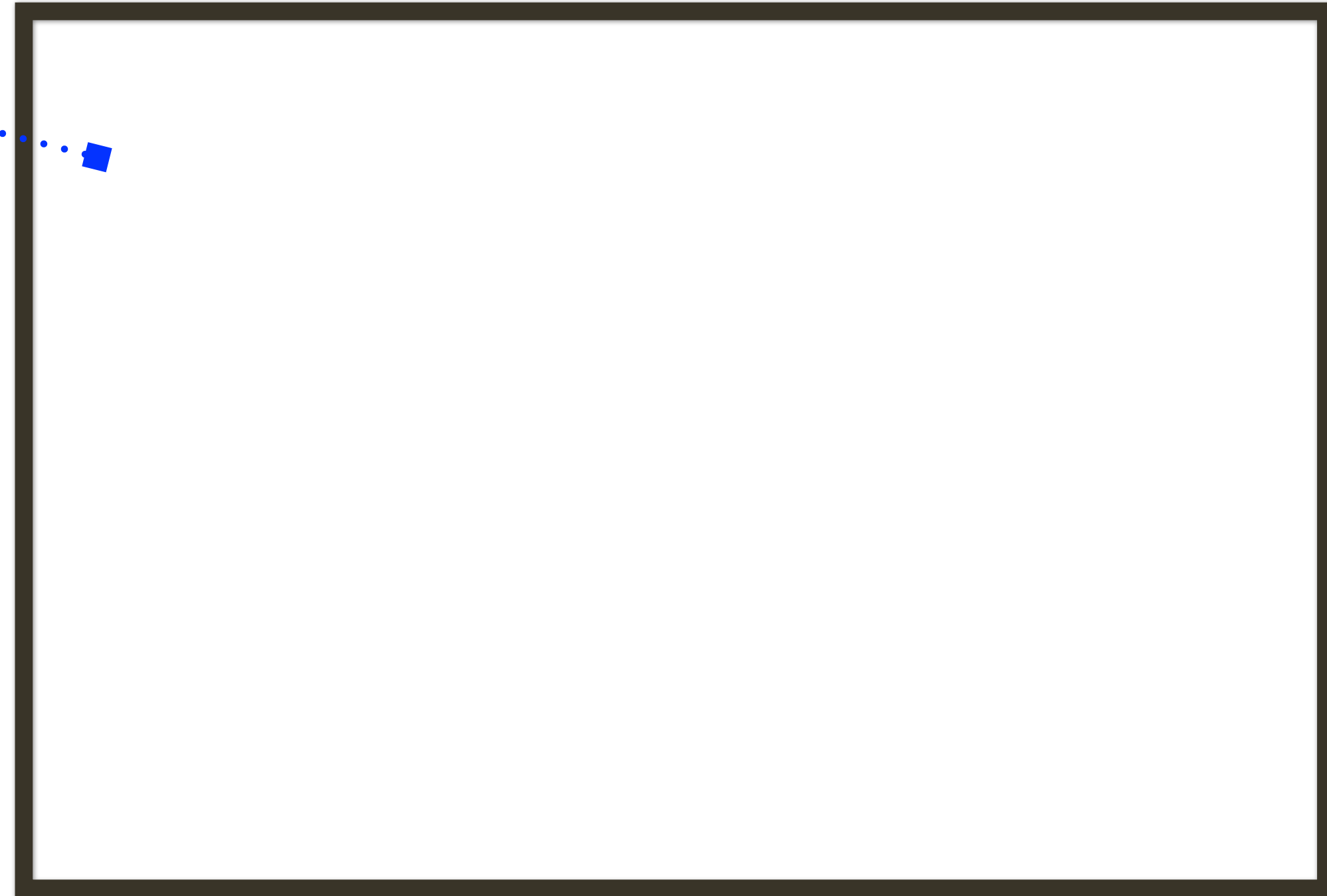
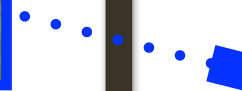
\mathbb{R}^d

Discriminative Embeddings

Images and **class labels** are embedded into the same space

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$



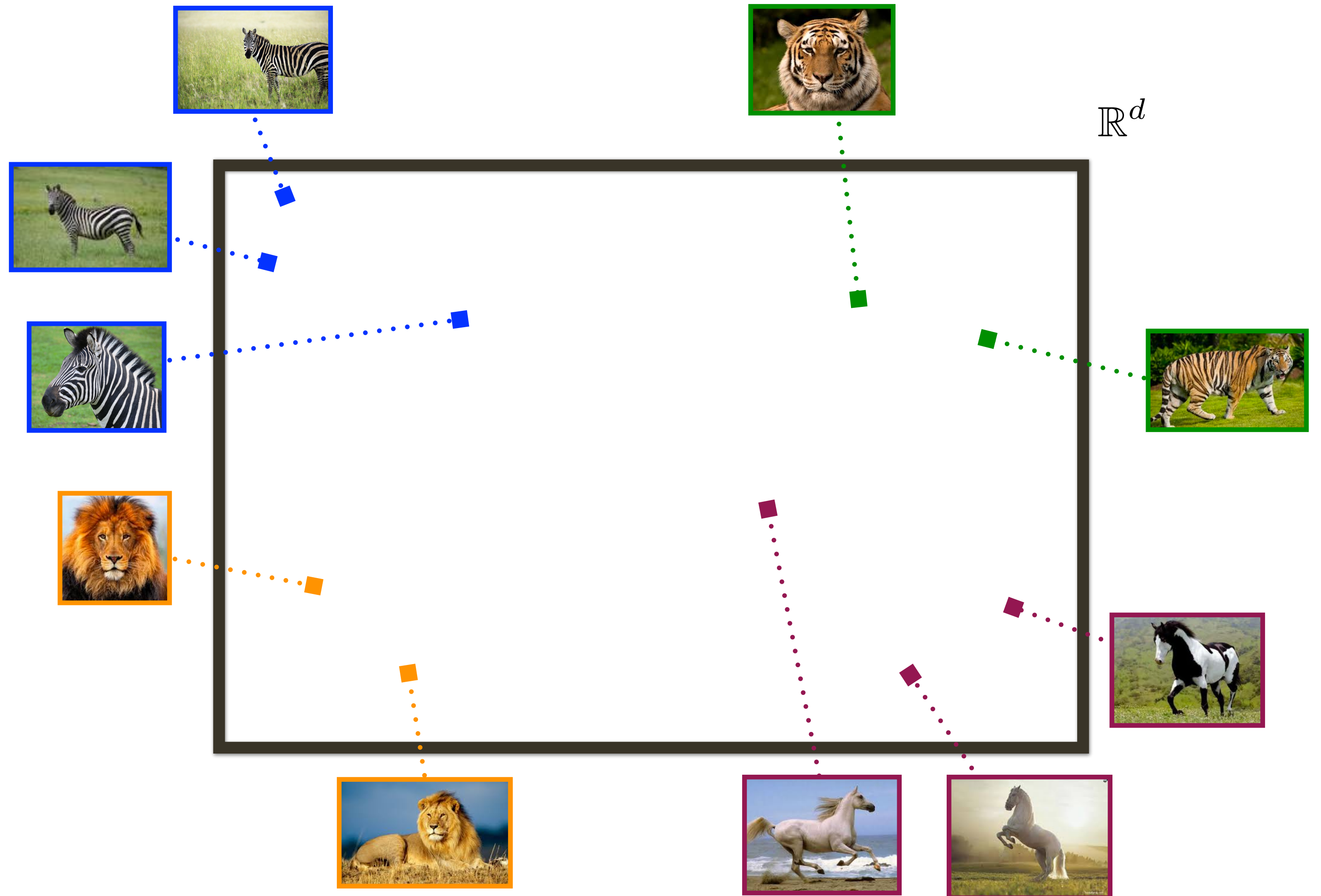
\mathbb{R}^d

Discriminative Embeddings

Images and **class labels** are embedded into the same space

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$



Discriminative Embeddings

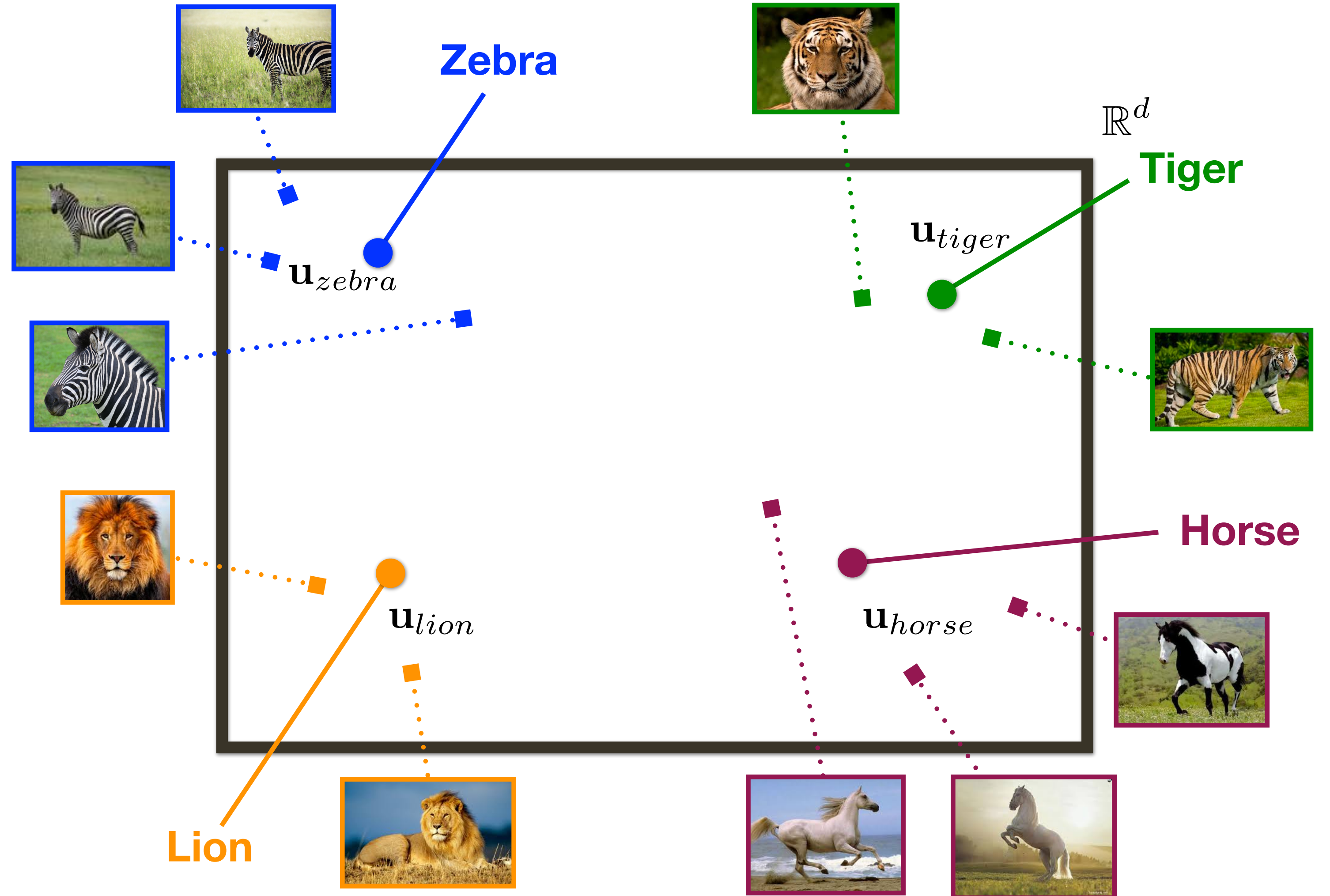
Images and **class labels** are embedded into the same space

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$



Discriminative Embeddings

Images and **class labels** are embedded into the same space

Image Embedding 

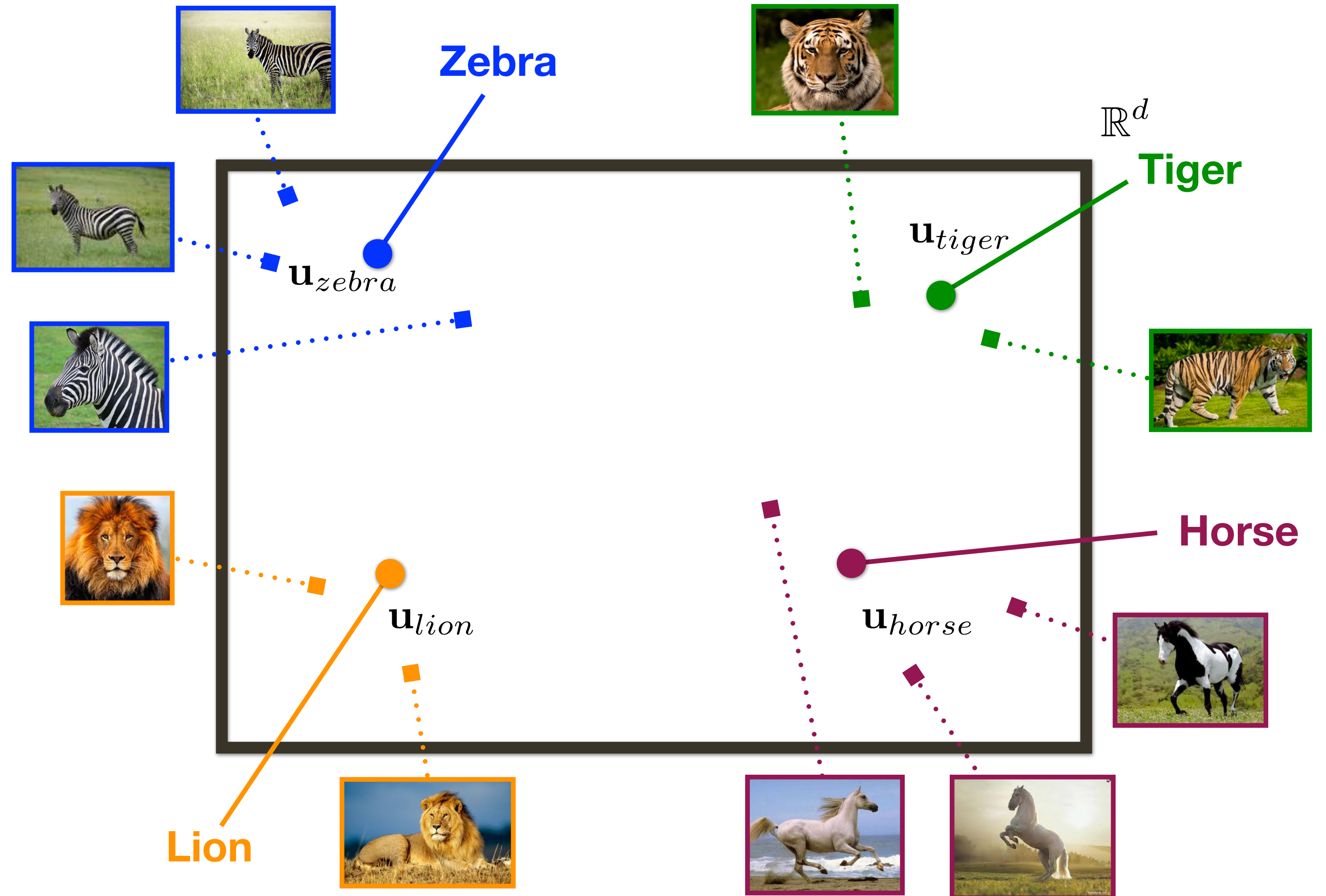
$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$



Discriminative Embeddings

Images and **class labels** are embedded into the same space

Image Embedding 

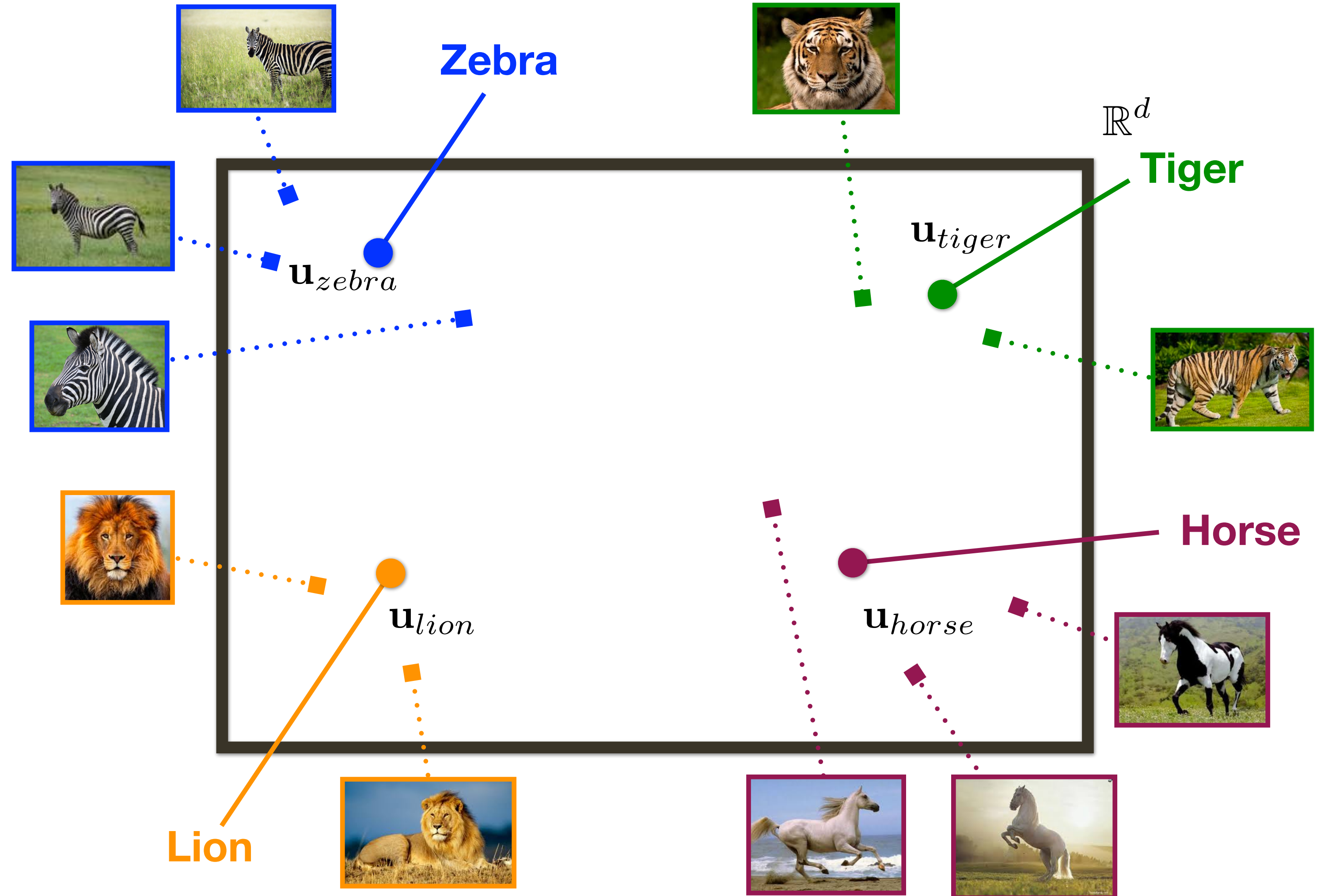
$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \frac{\mathbf{u}}{\|\mathbf{u}\|} \cdot \frac{\mathbf{u}'}{\|\mathbf{u}'\|}$$



Discriminative Embeddings

Image Categorization / Annotation

which object category does image belong to?

Image Embedding



$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

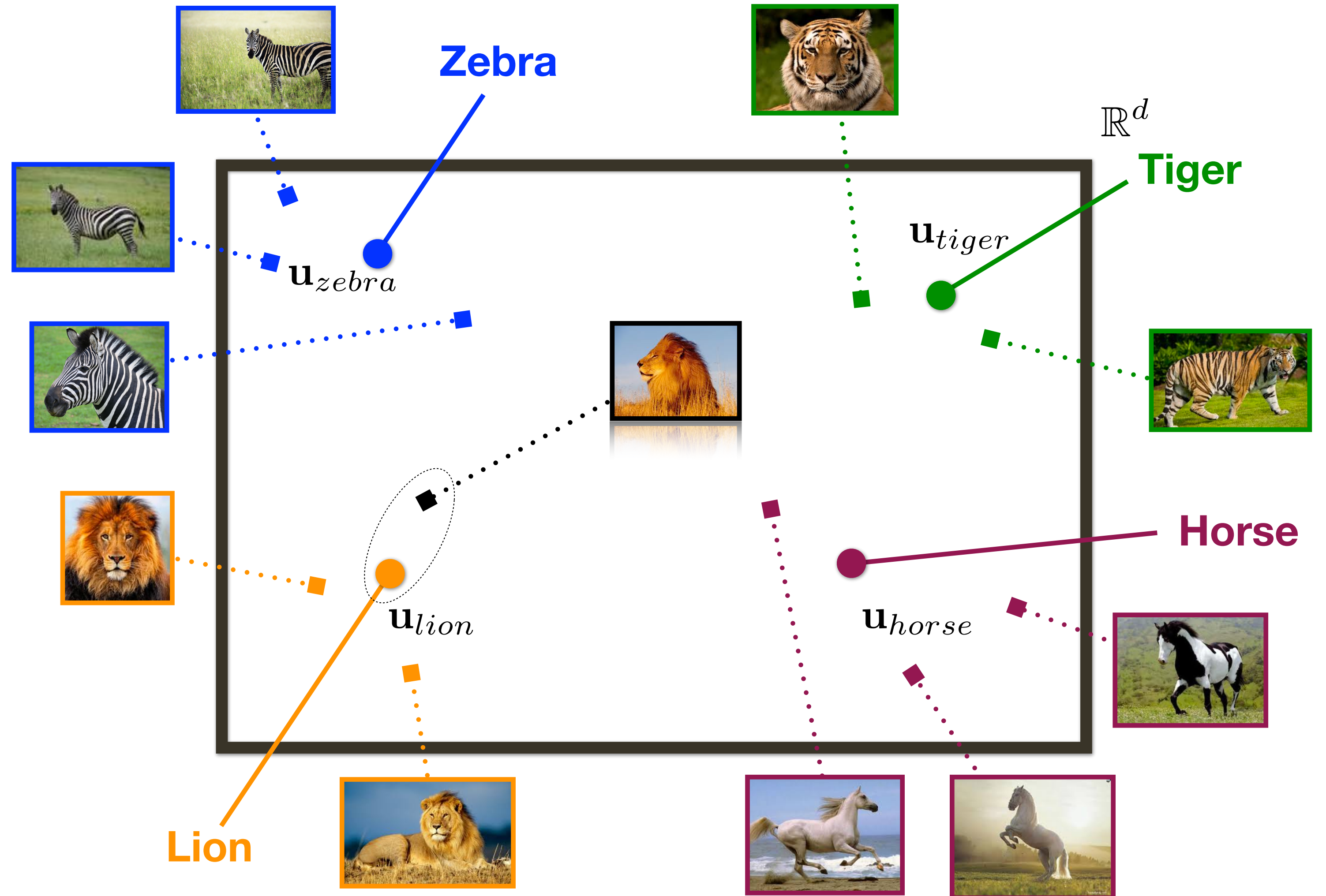
Label Embedding



$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$



Discriminative Embeddings

Image Categorization / Annotation

which object category does image belong to?

Image Embedding



$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding

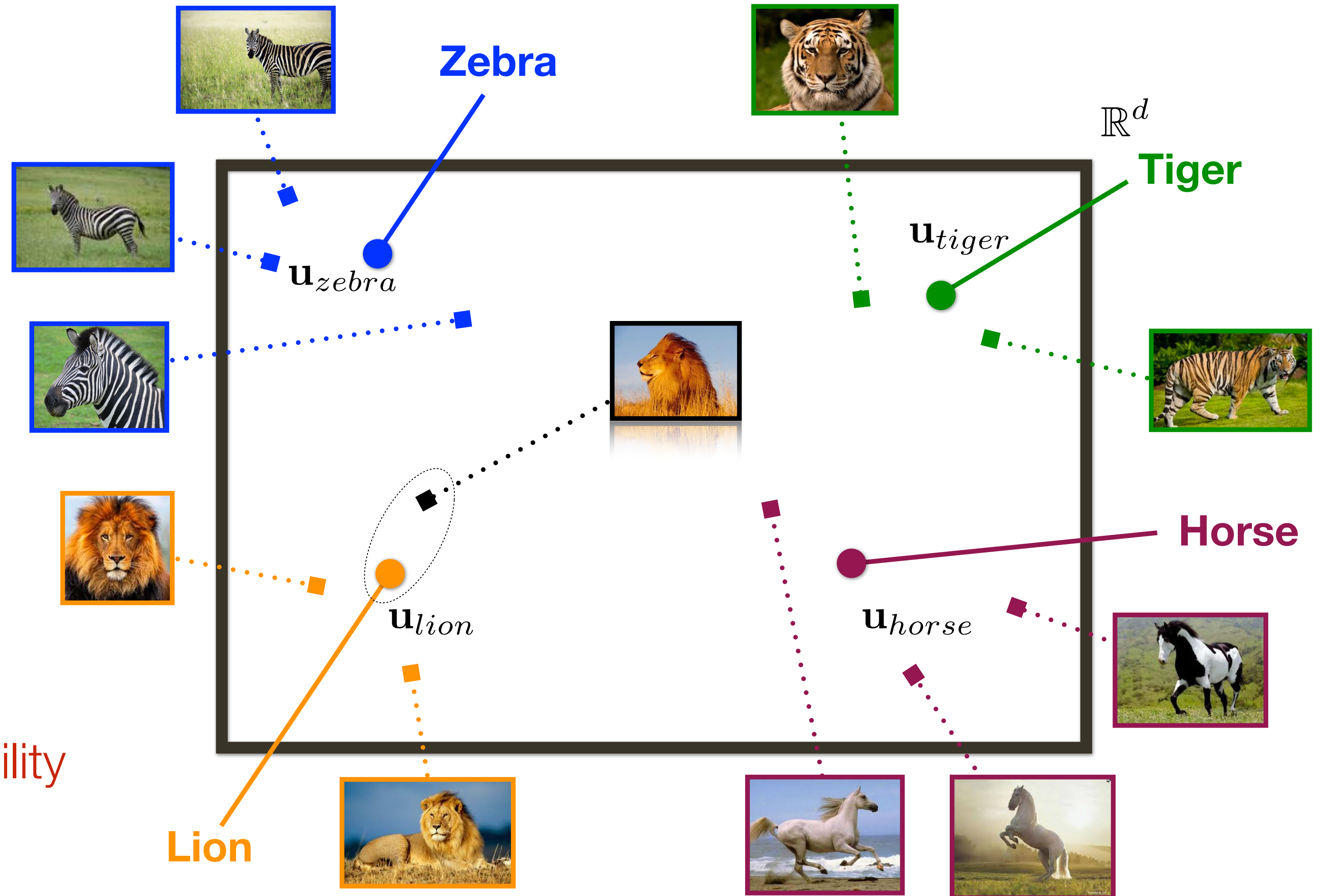


$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Distance can be interpreted as probability



Discriminative Embeddings

Image Categorization / Annotation

which object category does image belong to?

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

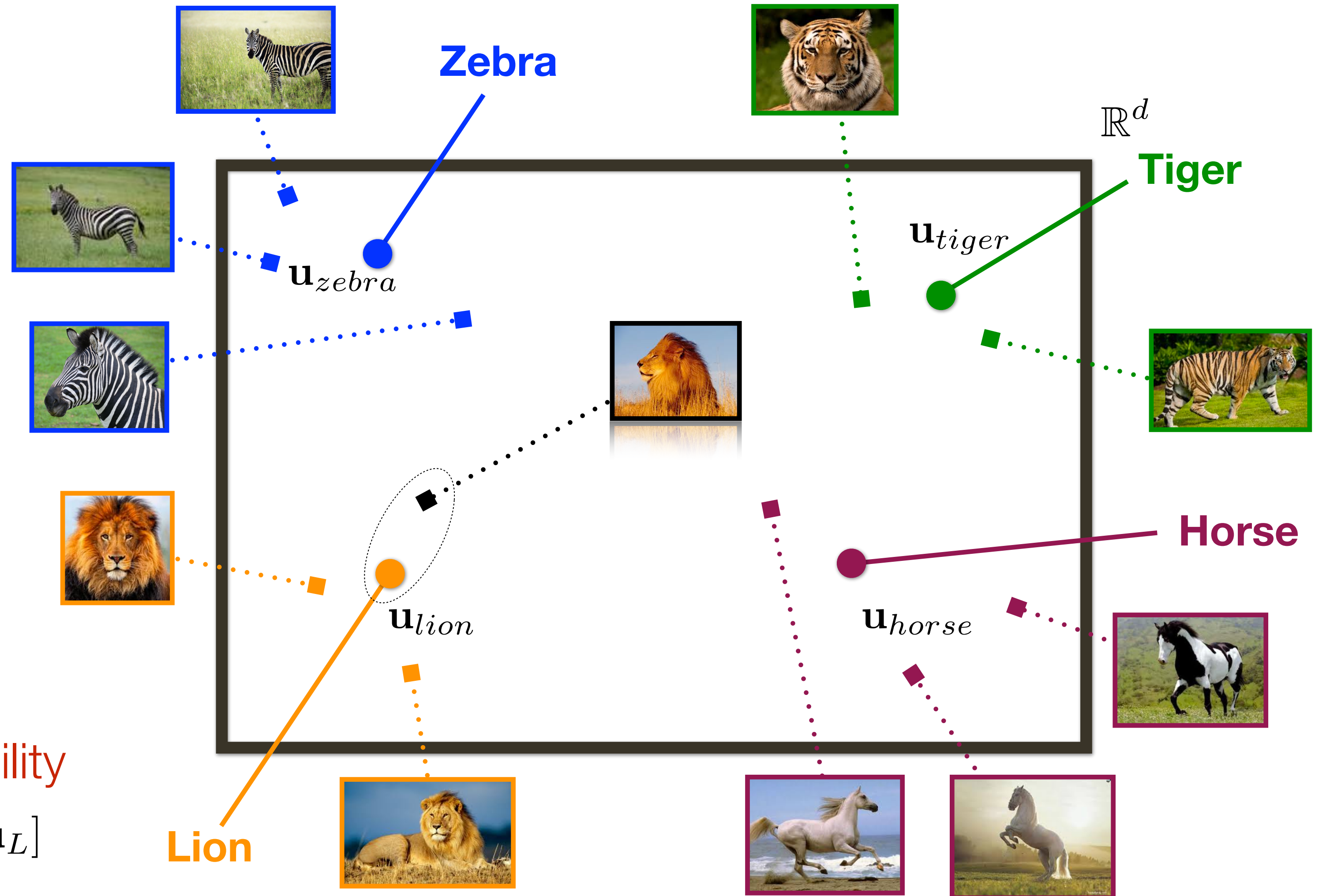
$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

Similarity in Embedding Space

$$D(\mathbf{u}_i, \mathbf{u}') = \mathbf{u}_i \cdot \mathbf{u}'$$

Distance can be interpreted as probability

$$\text{Softmax}(\mathbf{U}\mathbf{u}'), \text{ where } \mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_L]$$



Discriminative Embeddings

Search by Image
most similar image to a query?

Image Embedding 

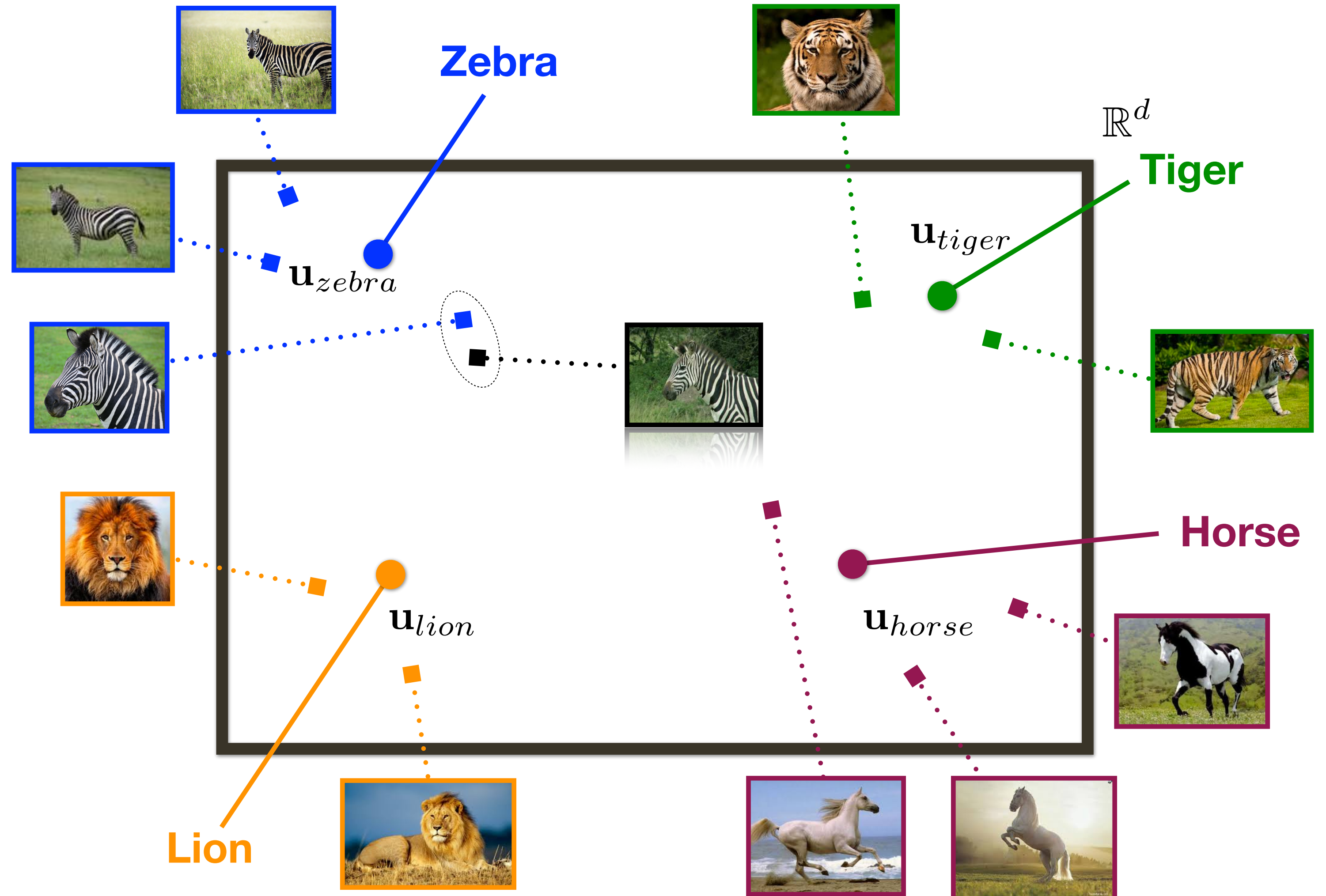
$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$



Discriminative Embeddings

Search by Label

most representative image for a label?

Image Embedding



$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

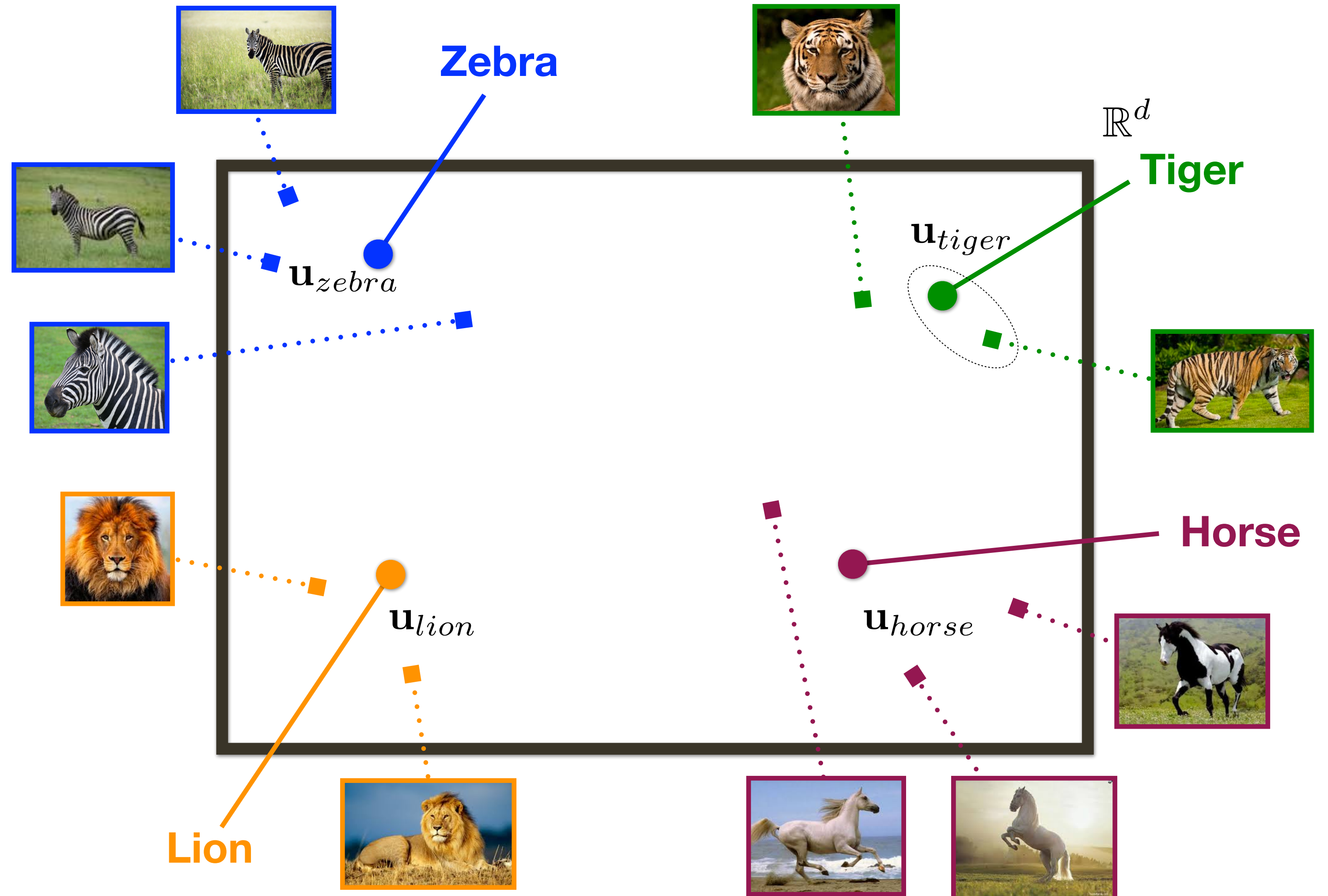
Label Embedding



$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$



Discriminative Embeddings

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

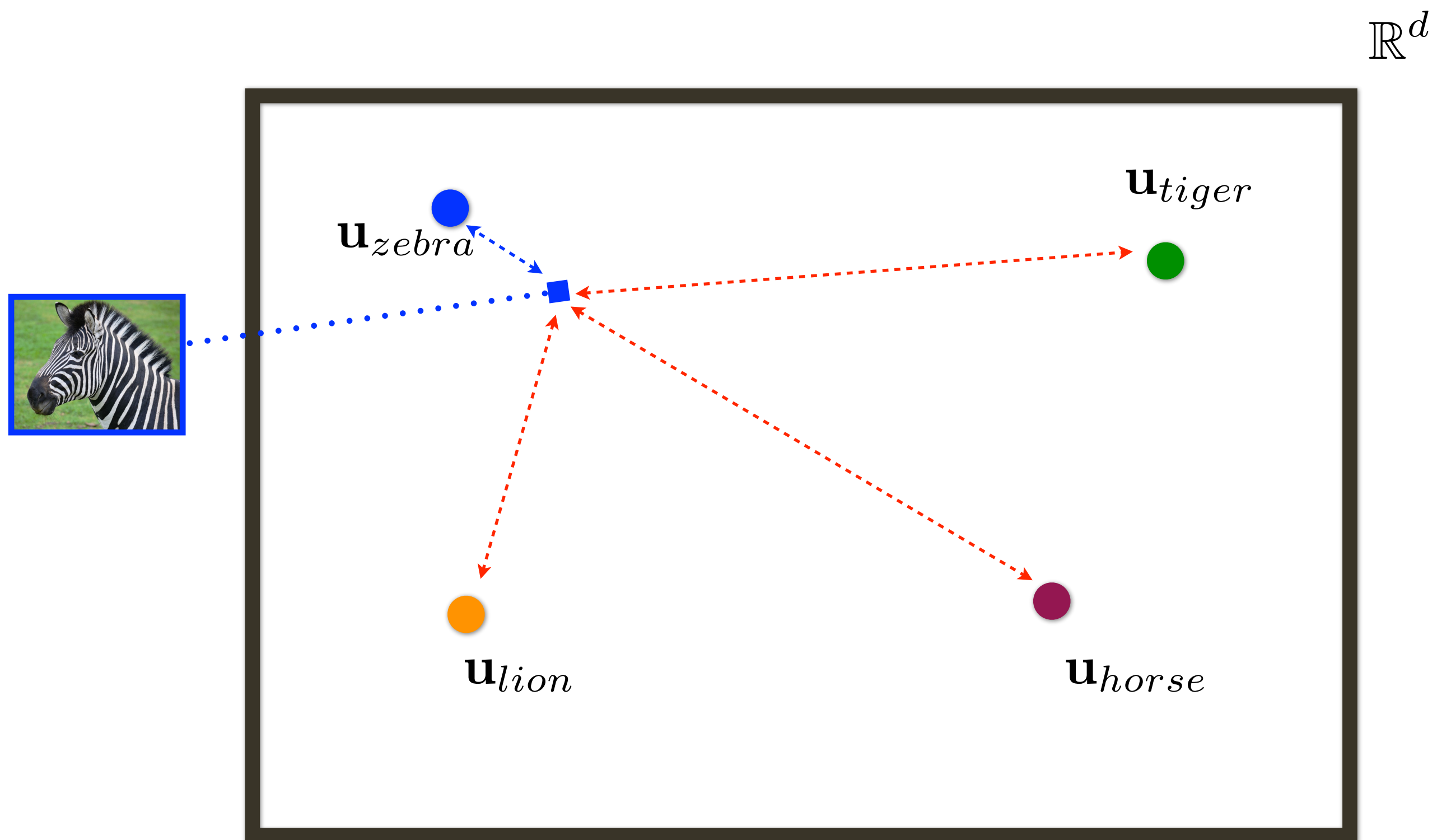
Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:

$$\min_{\mathbf{W}, \mathbf{U}} \sum_i^N \mathcal{L}_C(\mathbf{W}, \mathbf{U}, I_i, y_i) + \lambda_1 \|\mathbf{W}\|_F^2 + \lambda_2 \|\mathbf{U}\|_F^2$$

$$\mathcal{L}_C(\mathbf{W}, \mathbf{U}, I_i, y_i) = \sum [1 + \underbrace{D(\Psi(I_i), \mathbf{u}_{y_i})}_{\text{blue}} - \underbrace{D(\Psi(I_i), \mathbf{u}_{y_c})}_{\text{red}}]$$



Discriminative Embeddings

Why not minimize distance directly?

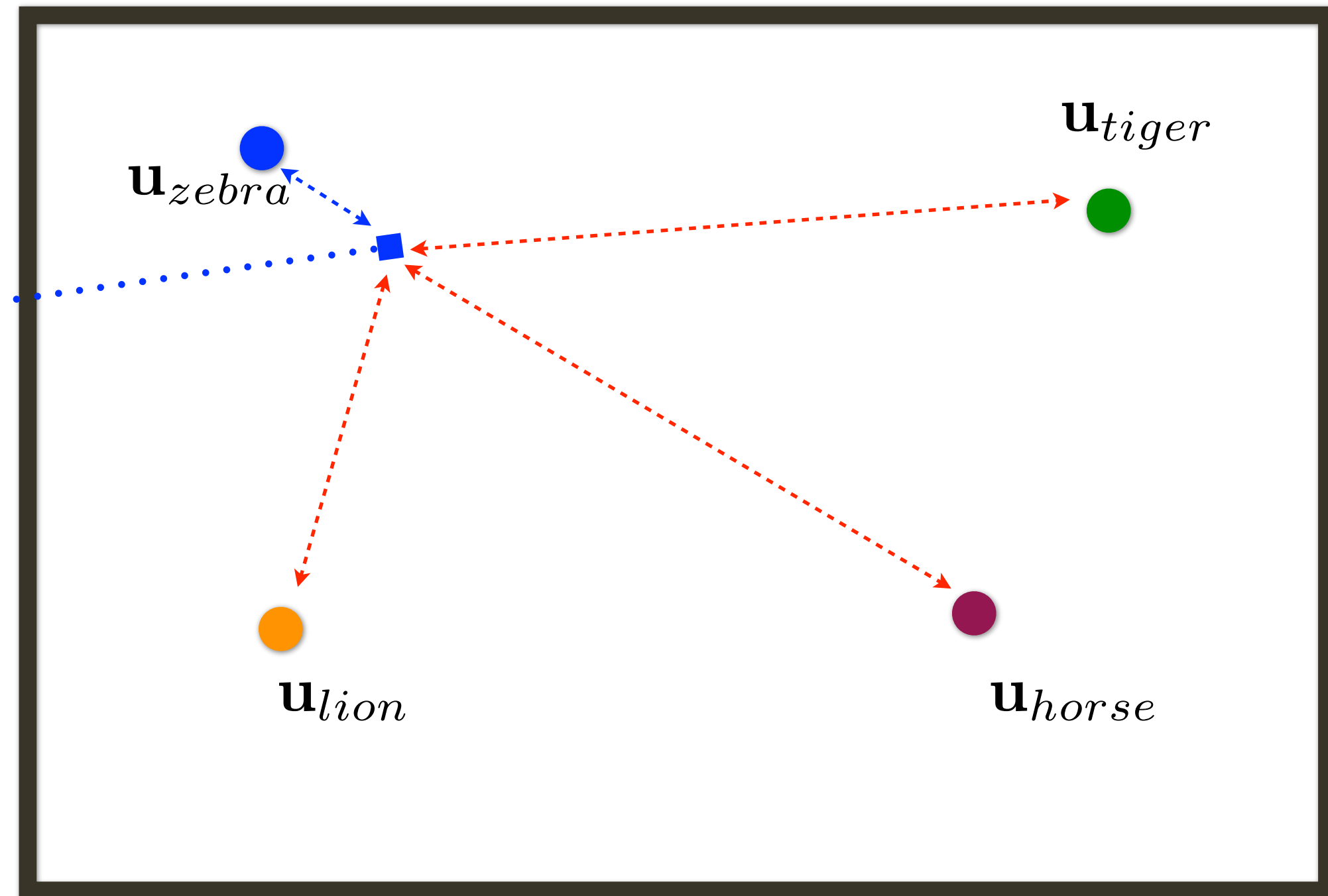
$$\mathcal{L}_C(\mathbf{W}, \mathbf{U}, I_i, y_i) = \sum [1 + \underbrace{D(\Psi(I_i), \mathbf{u}_{y_i})}_{\text{blue}} - \underbrace{D(\Psi(I_i), \mathbf{u}_{y_c})}_{\text{red}}]$$

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$



Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:

$$\min_{\mathbf{W}, \mathbf{U}} \sum_i^N \mathcal{L}_C(\mathbf{W}, \mathbf{U}, I_i, y_i) + \lambda_1 \|\mathbf{W}\|_F^2 + \lambda_2 \|\mathbf{U}\|_F^2$$

[Bengio et al., NIPS'10]

[Weinberger, Chapelle, NIPS'09]

Discriminative Embeddings

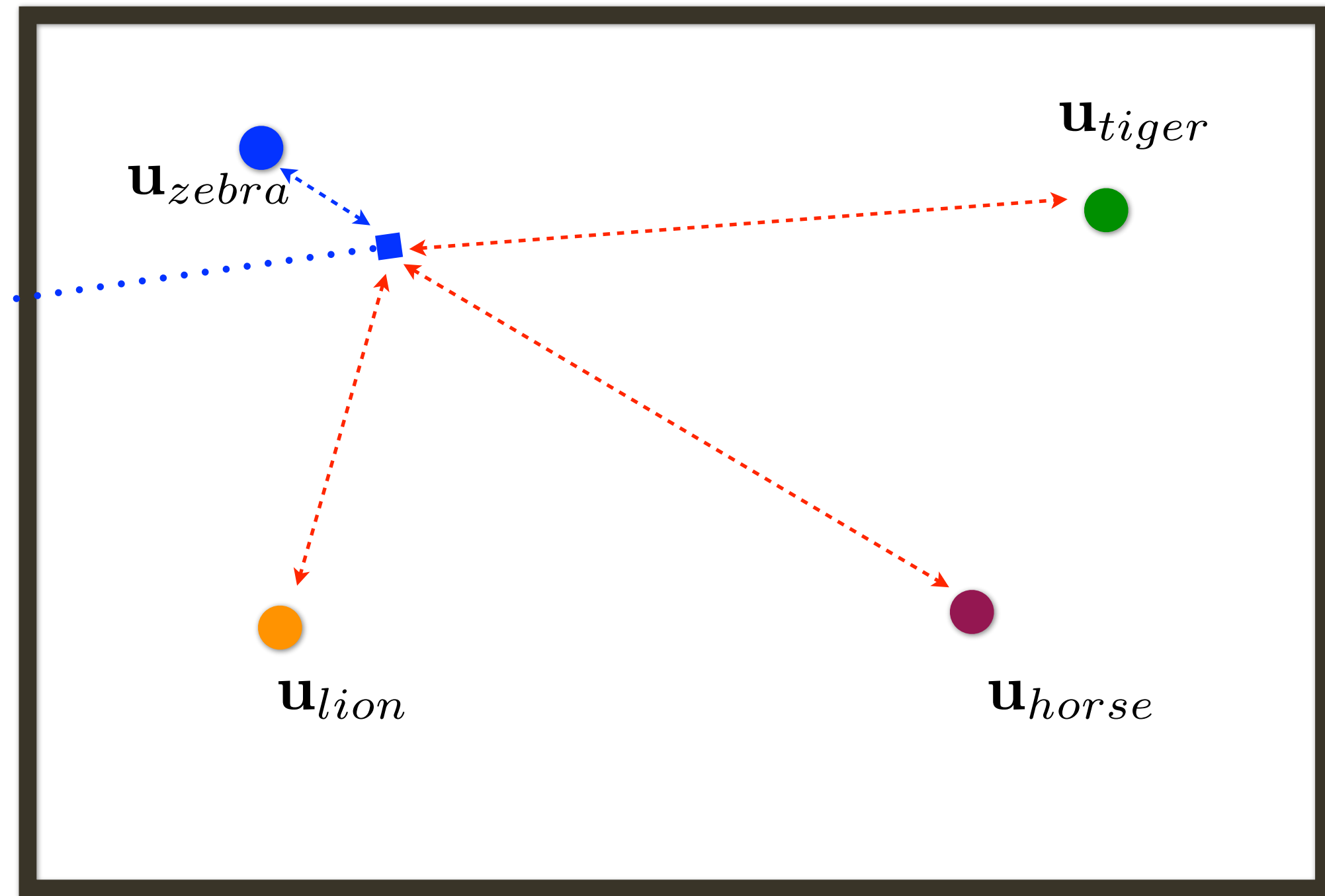
Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

$$\mathcal{L}_C(\mathbf{W}, \mathbf{U}, I_i, y_i) = \sum \max\{0, \alpha - \underbrace{D(\Psi(I_i), \mathbf{u}_{y_i})}_{\text{blue}} + \underbrace{D(\Psi(I_i), \mathbf{u}_{y_c})}_{\text{red}}\}$$



Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \frac{\mathbf{u}}{\|\mathbf{u}\|} \cdot \frac{\mathbf{u}'}{\|\mathbf{u}'\|}$$

Objective Function:

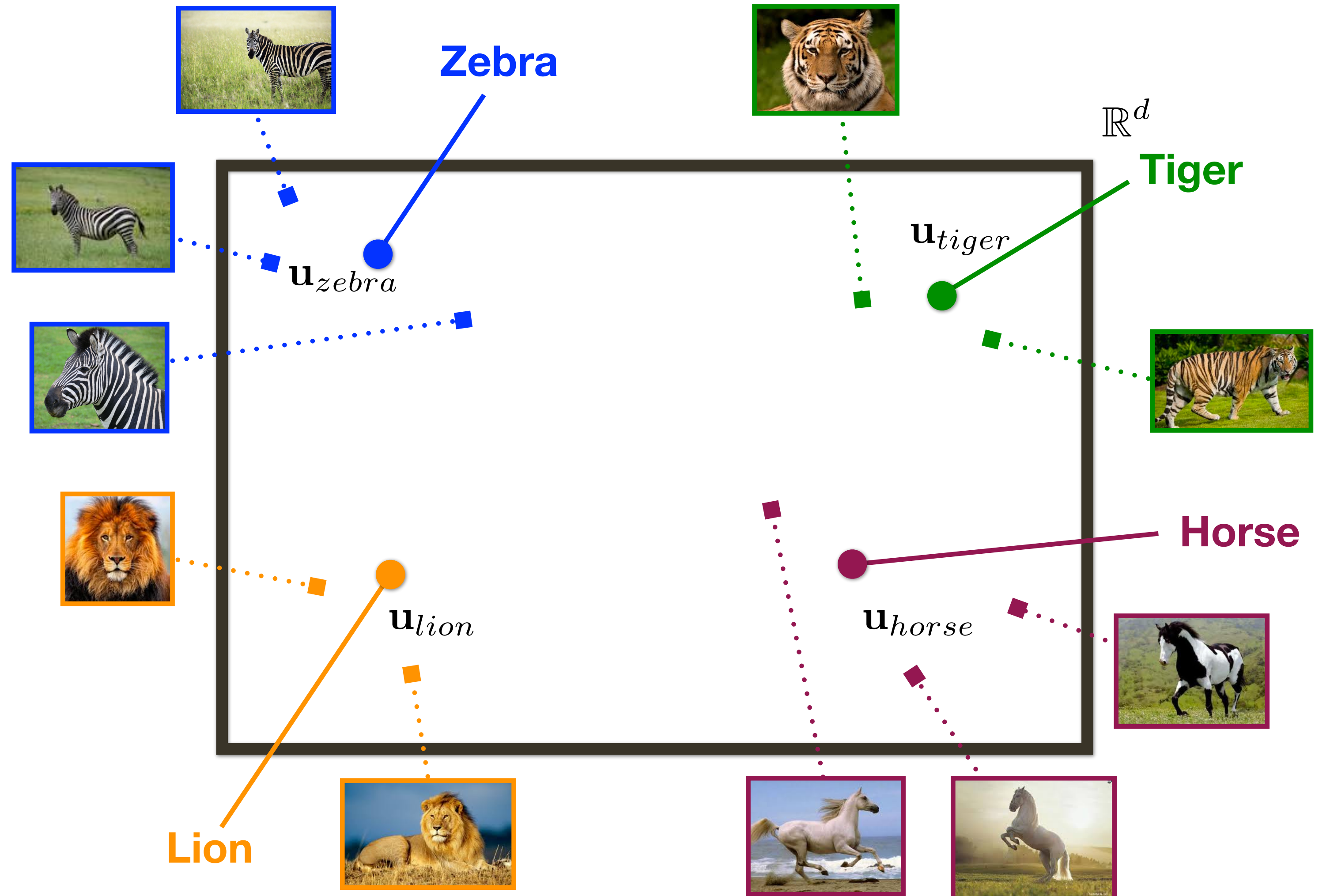
$$\min_{\mathbf{W}, \mathbf{U}} \sum_i^N \mathcal{L}_C(\mathbf{W}, \mathbf{U}, I_i, y_i) + \lambda_1 \|\mathbf{W}\|_F^2 + \lambda_2 \|\mathbf{U}\|_F^2$$

[Bengio et al., NIPS'10]

[Weinberger, Chapelle, NIPS'09]

Discriminative Embeddings

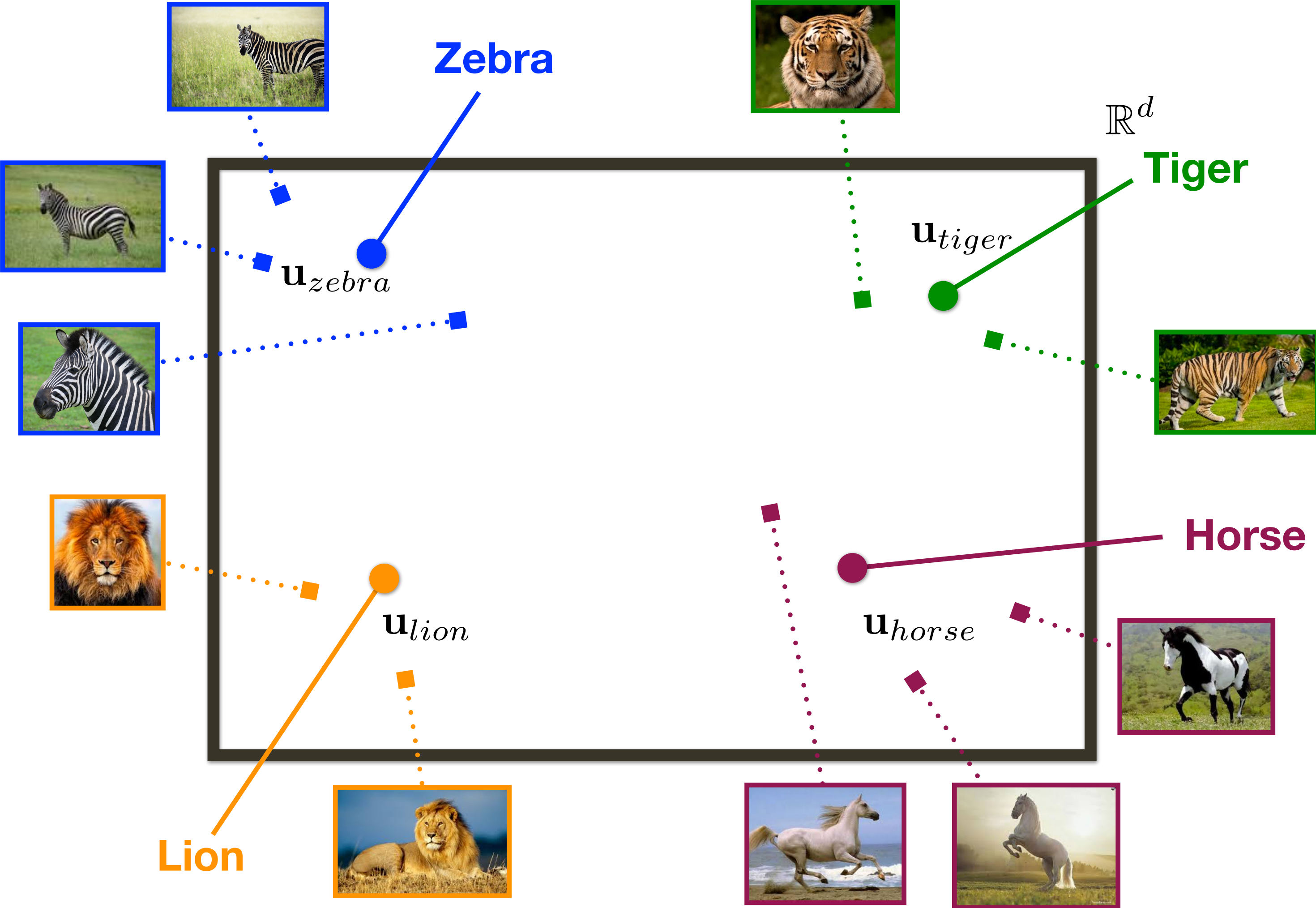
This is a very **convenient model**



Discriminative Embeddings

This is a very **convenient model**

Inducing semantics on the embedding space

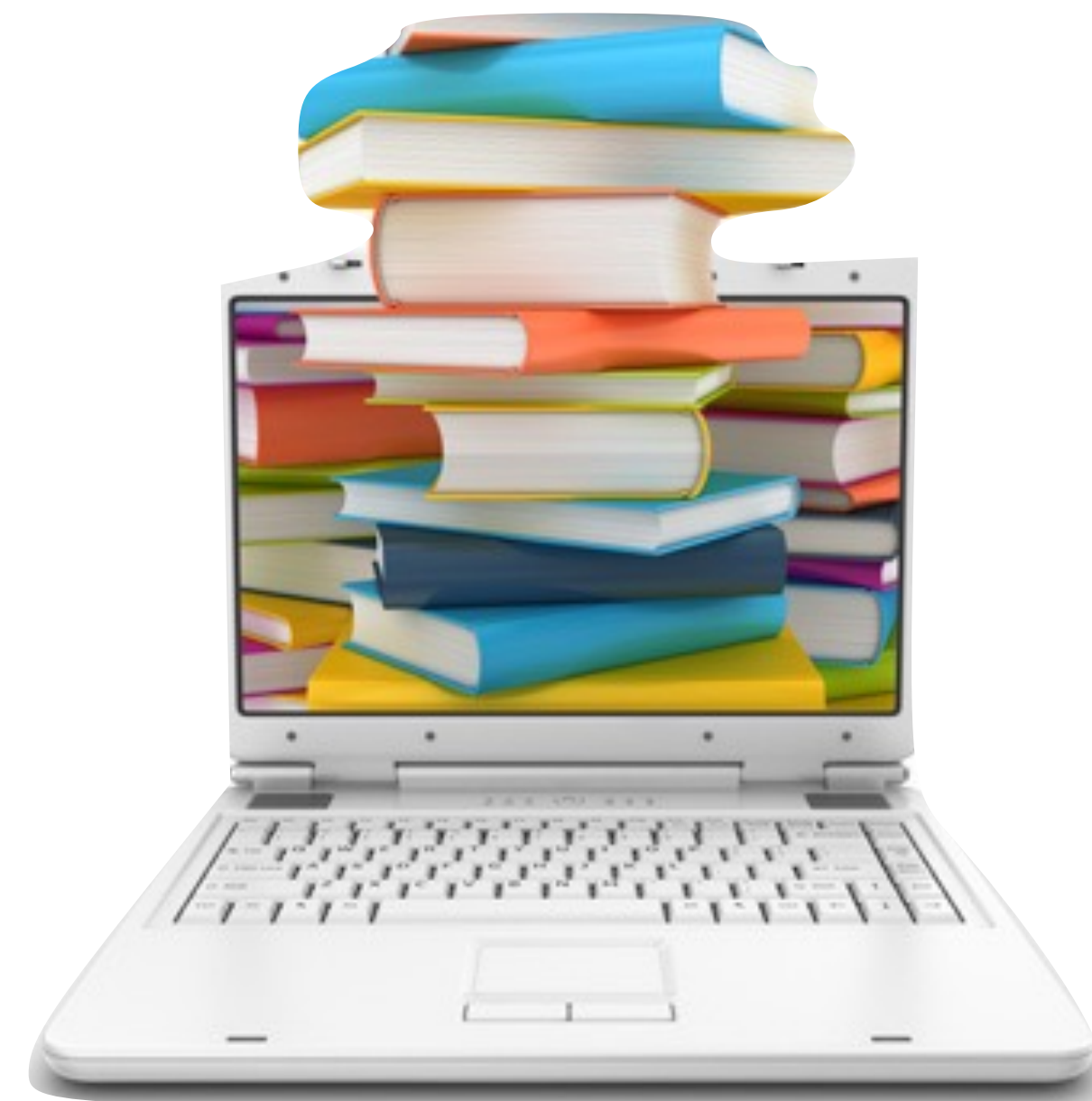


word2vec: Unsupervised Word Embedding

Distributional Semantics Hypothesis: words that are used and occur in the same context tend to have similar meaning

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i : \{1, \dots, L\} \rightarrow \mathbb{R}^d$$



word2vec: Unsupervised Word Embedding

[Fu et al., 2016]

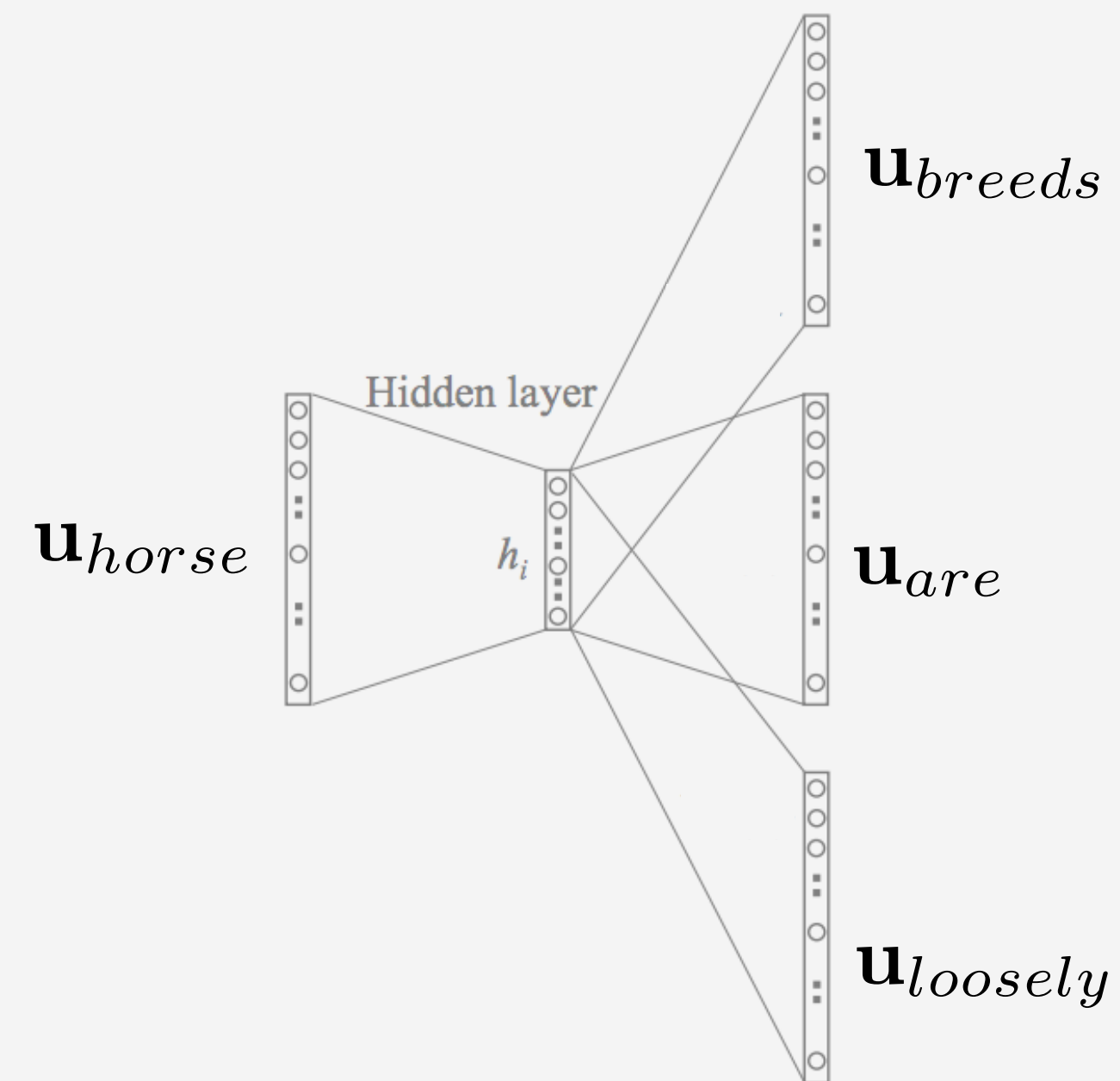
Distributional Semantics Hypothesis: words that are used and occur in the same context tend to have similar meaning

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i : \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

$$L = 310,000$$

e.g., Horse breeds are loosely divided into three categories



Skip-gram Model: unsupervised semantic representation for words
(trained from 7 billion word linguistic corpus)

Semi-supervised Vocabulary Informed Learning

[Fu et al., 2016]

Image Embedding



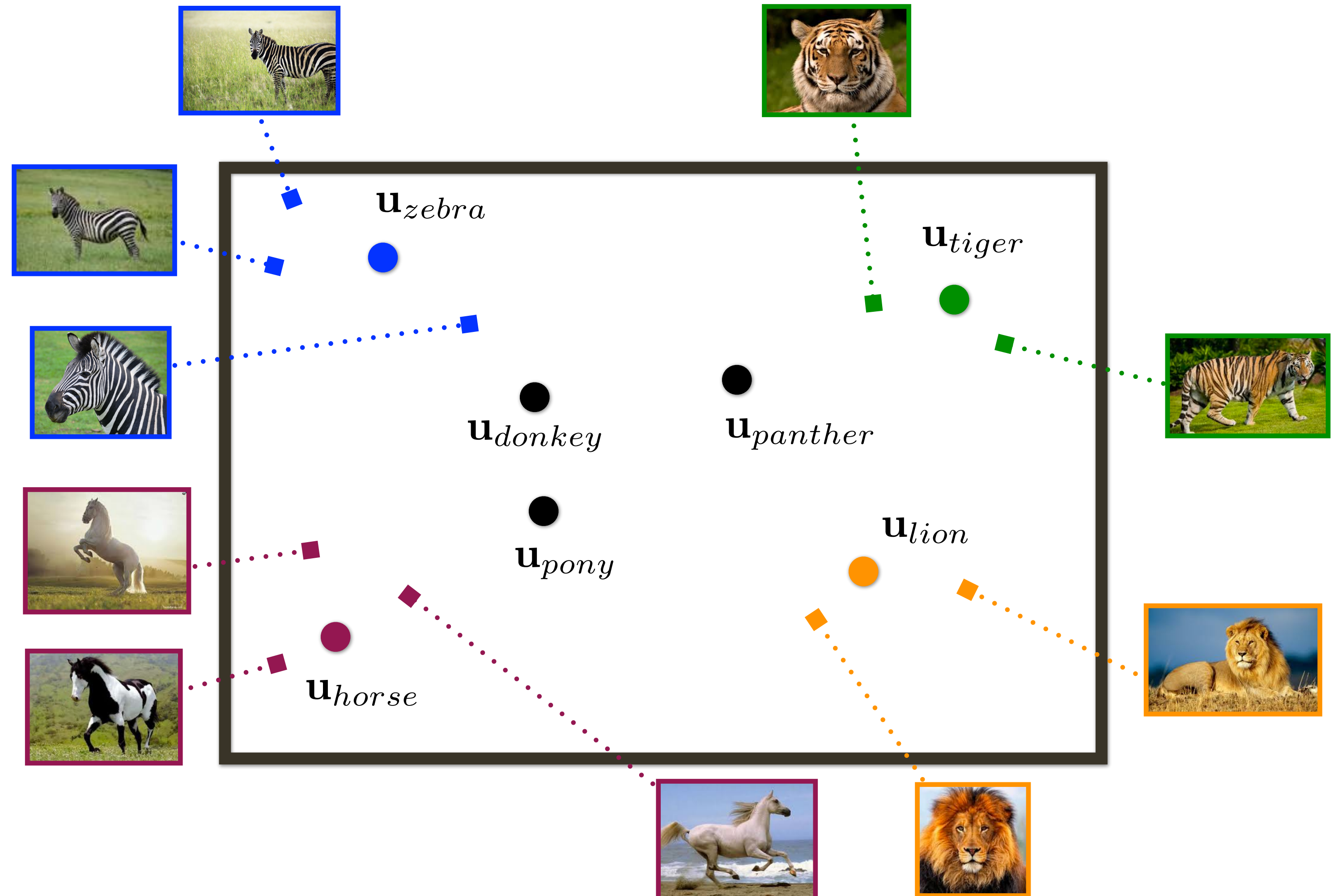
$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding



$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

$L = 310,000$



Semi-supervised Vocabulary Informed Learning

[Fu et al., 2016]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

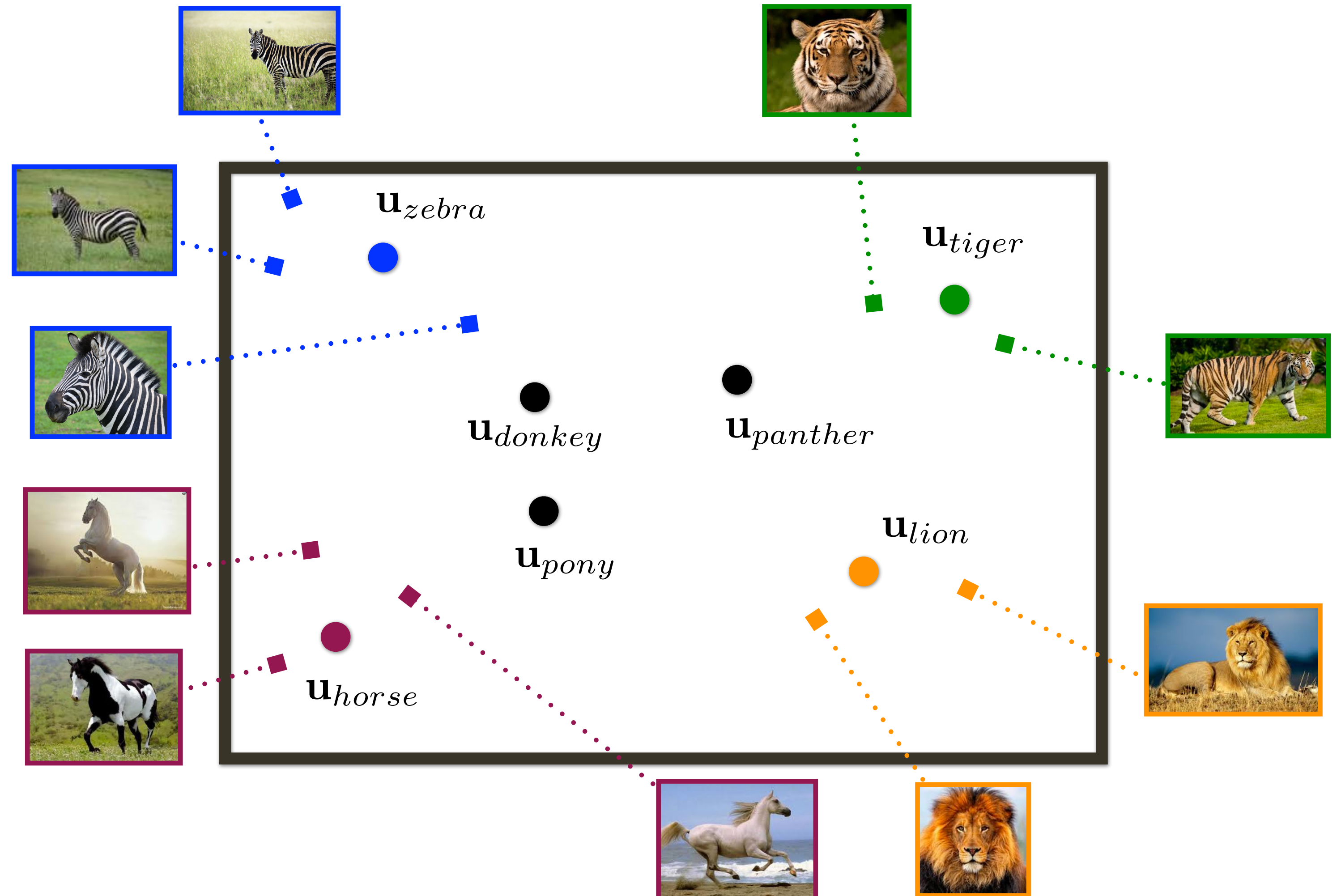
Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

$L = 310,000$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$



Semi-supervised Vocabulary Informed Learning

[Fu et al., 2016]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

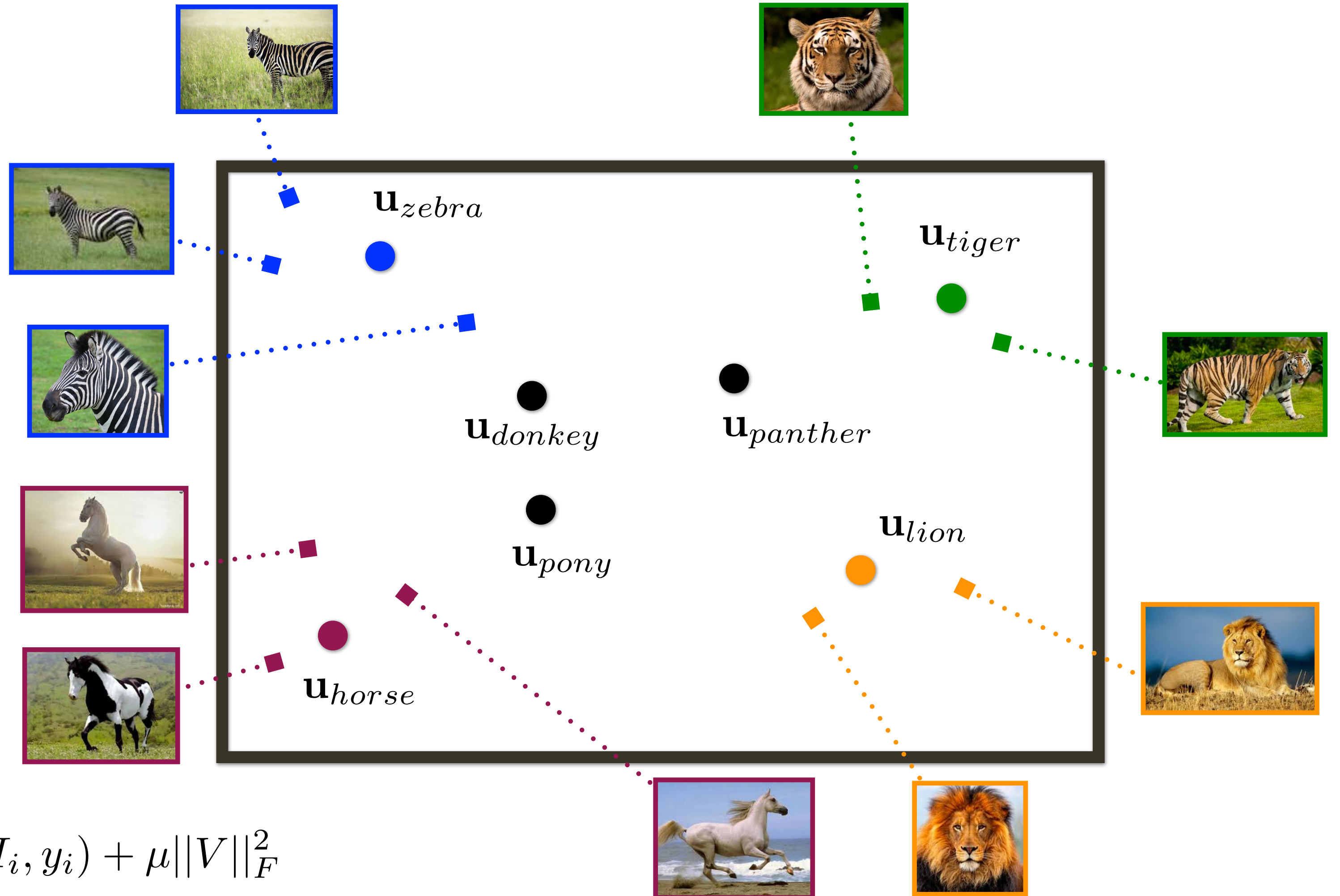
$L = 310,000$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:

$$\min_{\mathbf{W}} \sum_i^N \mathcal{L}_C(\mathbf{W}, \mathbf{V}, I_i, y_i) + \mathcal{L}_R(\mathbf{W}, \mathbf{V}, I_i, y_i) + \mu \|\mathbf{V}\|_F^2$$



Semi-supervised Vocabulary Informed Learning

[Fu et al., 2016]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

$L = 310,000$



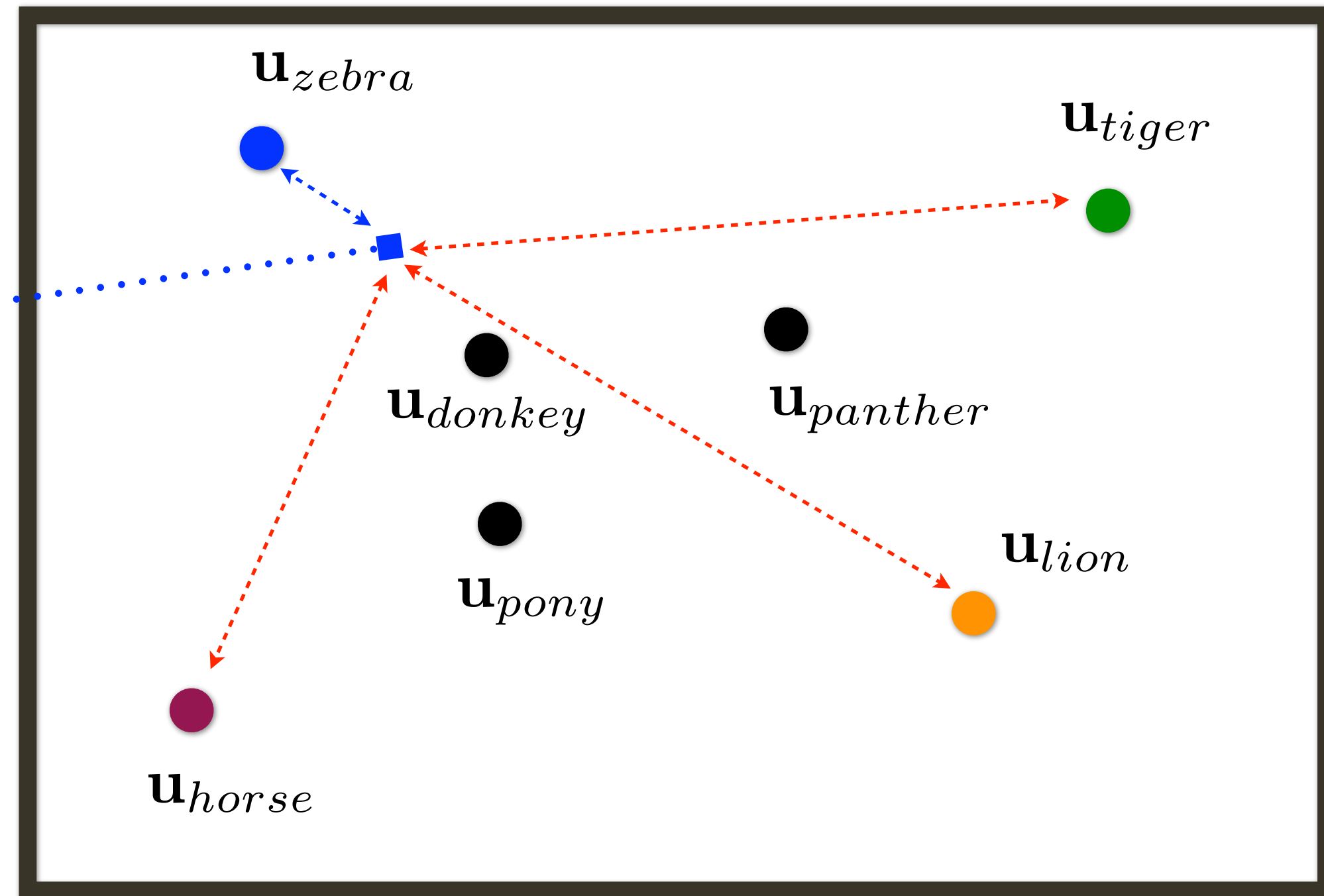
$$\mathcal{L}_C(\mathbf{W}, \mathbf{U}, \mathbf{x}_i, y_i) = \sum [1 + \underbrace{D(\mathbf{W}\mathbf{x}_i, \mathbf{u}_{y_i})}_{\text{blue}} - \underbrace{D(\mathbf{W}\mathbf{x}_i, \mathbf{u}_c)}_{\text{red}}]$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:

$$\min_{\mathbf{W}} \sum_i^N \mathcal{L}_C(\mathbf{W}, \mathbf{V}, I_i, y_i) + \mathcal{L}_R(\mathbf{W}, \mathbf{V}, I_i, y_i) + \mu \|\mathbf{V}\|_F^2$$

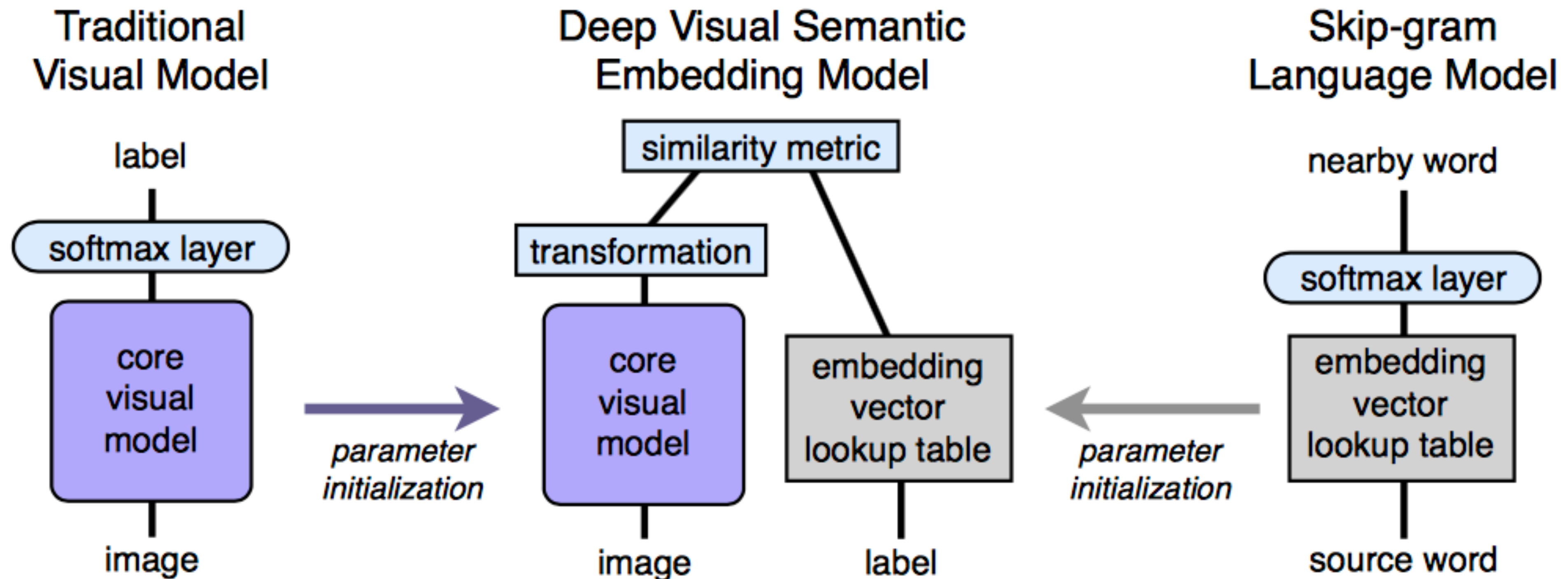


Intuition



DeViSE: A Deep Visual-Semantic Embedding Model

[Frome et al., 2013]



$$loss(image, label) = \sum_{j \neq label} \max[0, margin - \vec{t}_{label} M \vec{v}(image) + \vec{t}_j M \vec{v}(image)]$$

DeViSE: A Deep Visual-Semantic Embedding Model

[Frome et al., 2013]

Supervised Results

Model type	dim	Flat hit@ k (%)				Hierarchical precision@ k			
		1	2	5	10	2	5	10	20
Softmax baseline	N/A	55.6	67.4	78.5	85.0	0.452	0.342	0.313	0.319
DeViSE	500	53.2	65.2	76.7	83.3	0.447	0.352	0.331	0.341
	1000	54.9	66.9	78.4	85.0	0.454	0.351	0.325	0.331
Random embeddings	500	52.4	63.9	74.8	80.6	0.428	0.315	0.271	0.248
	1000	50.5	62.2	74.2	81.5	0.418	0.318	0.290	0.292
Chance	N/A	0.1	0.2	0.5	1.0	0.007	0.013	0.022	0.042

Zero-shot Results

Model	200 labels	1000 labels
DeViSE	31.8%	9.0%
Mensink et al. 2012 [12]	35.7%	1.9%
Rohrbach et al. 2011 [17]	34.8%	-

Semi-supervised Vocabulary Informed Learning

[Fu et al., 2016]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

$L = 310,000$



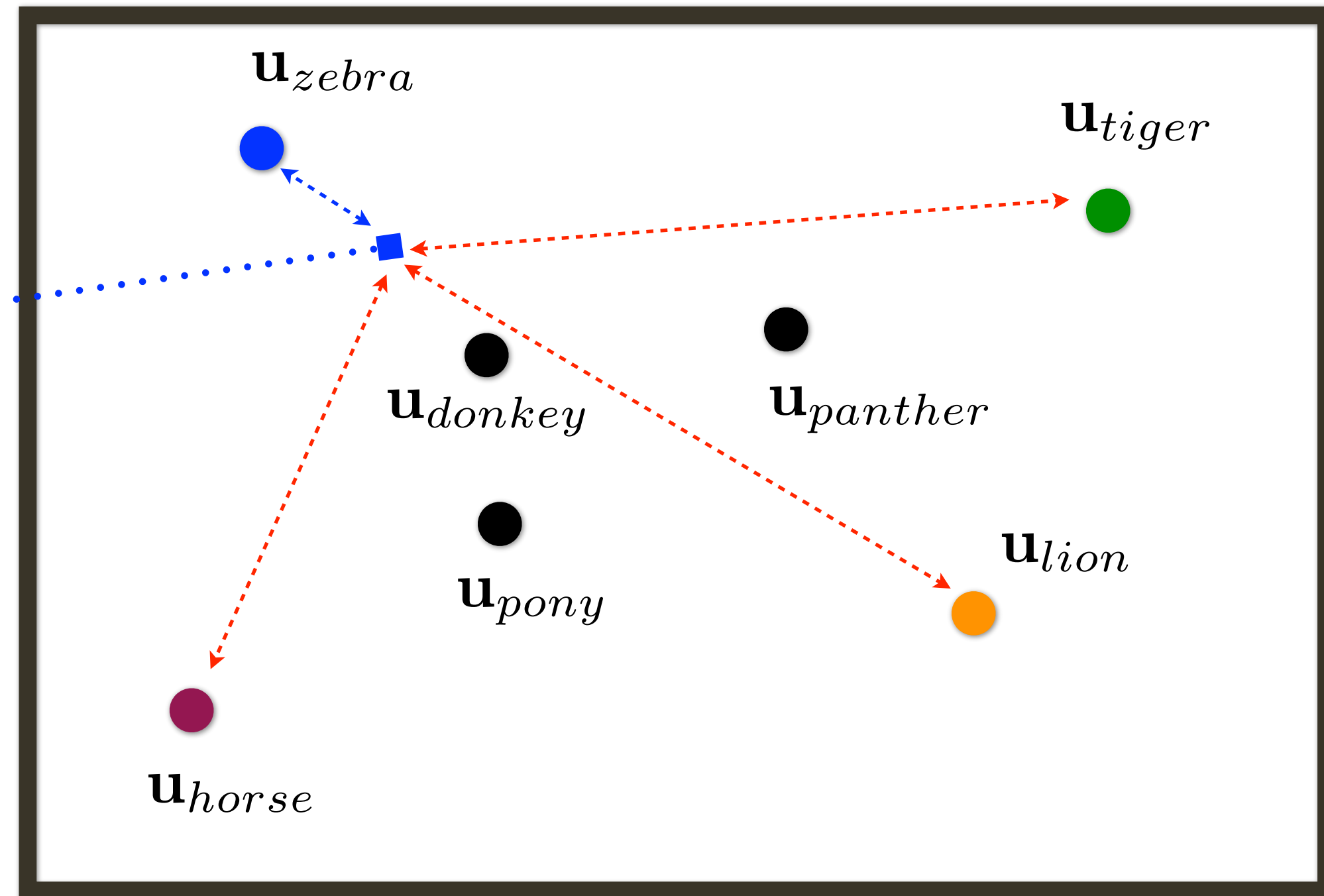
$$\mathcal{L}_C(\mathbf{W}, \mathbf{U}, \mathbf{x}_i, y_i) = \sum [1 + \underbrace{D(\mathbf{W}\mathbf{x}_i, \mathbf{u}_{y_i})}_{\text{blue}} - \underbrace{D(\mathbf{W}\mathbf{x}_i, \mathbf{u}_c)}_{\text{red}}]$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:

$$\min_{\mathbf{W}} \sum_i^N \mathcal{L}_C(\mathbf{W}, \mathbf{V}, I_i, y_i) + \mathcal{L}_R(\mathbf{W}, \mathbf{V}, I_i, y_i) + \mu \|\mathbf{V}\|_F^2$$



Semi-supervised Vocabulary Informed Learning

[Fu et al., 2016]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta): \mathbb{R}^D \rightarrow \mathbb{R}^d$$

Label Embedding 

$$\Psi_L(\text{word}_i) = \mathbf{u}_i: \{1, \dots, L\} \rightarrow \mathbb{R}^d$$

$L = 310,000$



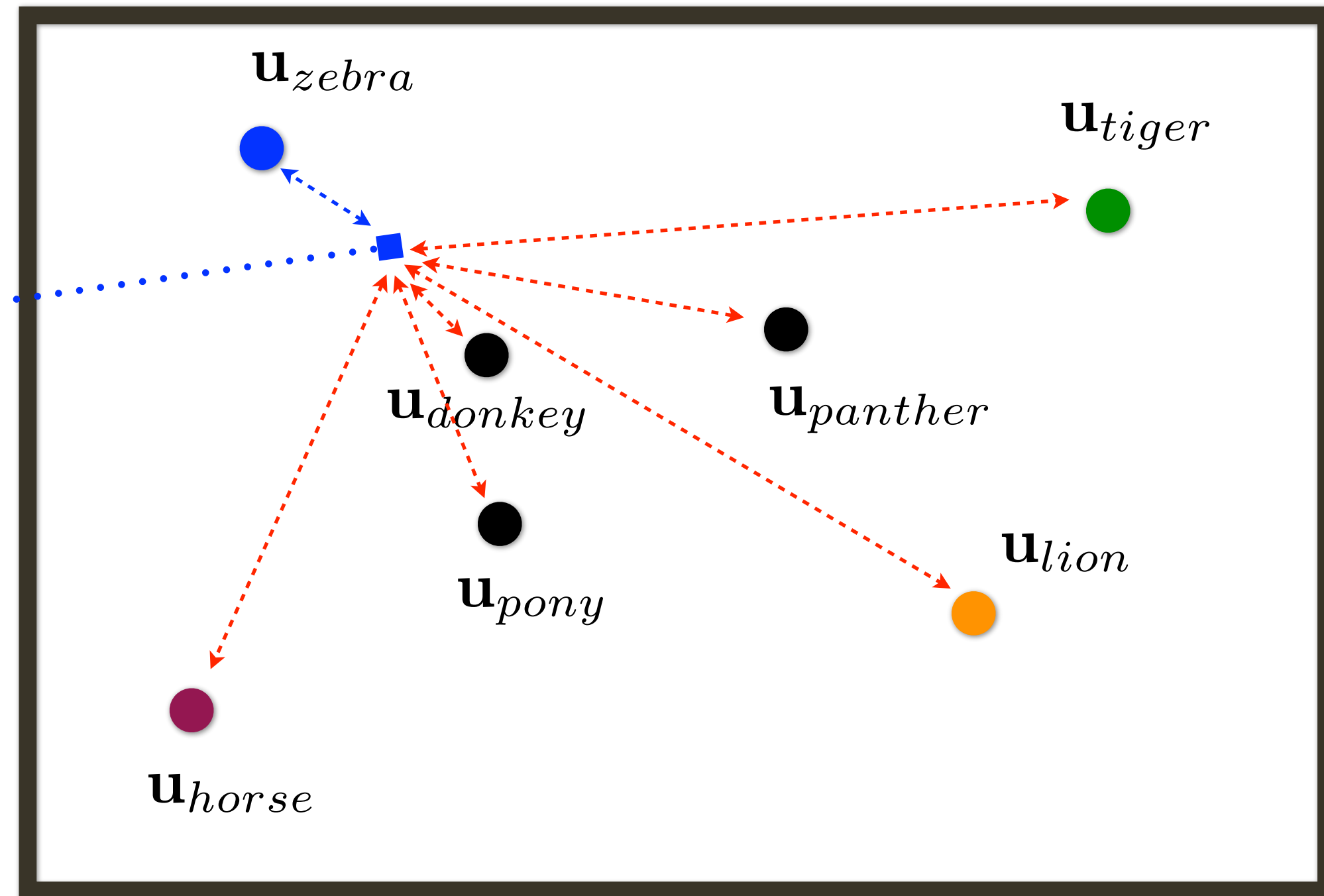
$$\mathcal{L}_C(\mathbf{W}, \mathbf{U}, \mathbf{x}_i, y_i) = \sum [1 + \underbrace{D(\mathbf{W}\mathbf{x}_i, \mathbf{u}_{y_i})}_{\text{blue}} - \underbrace{D(\mathbf{W}\mathbf{x}_i, \mathbf{u}_c)}_{\text{red}}]$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

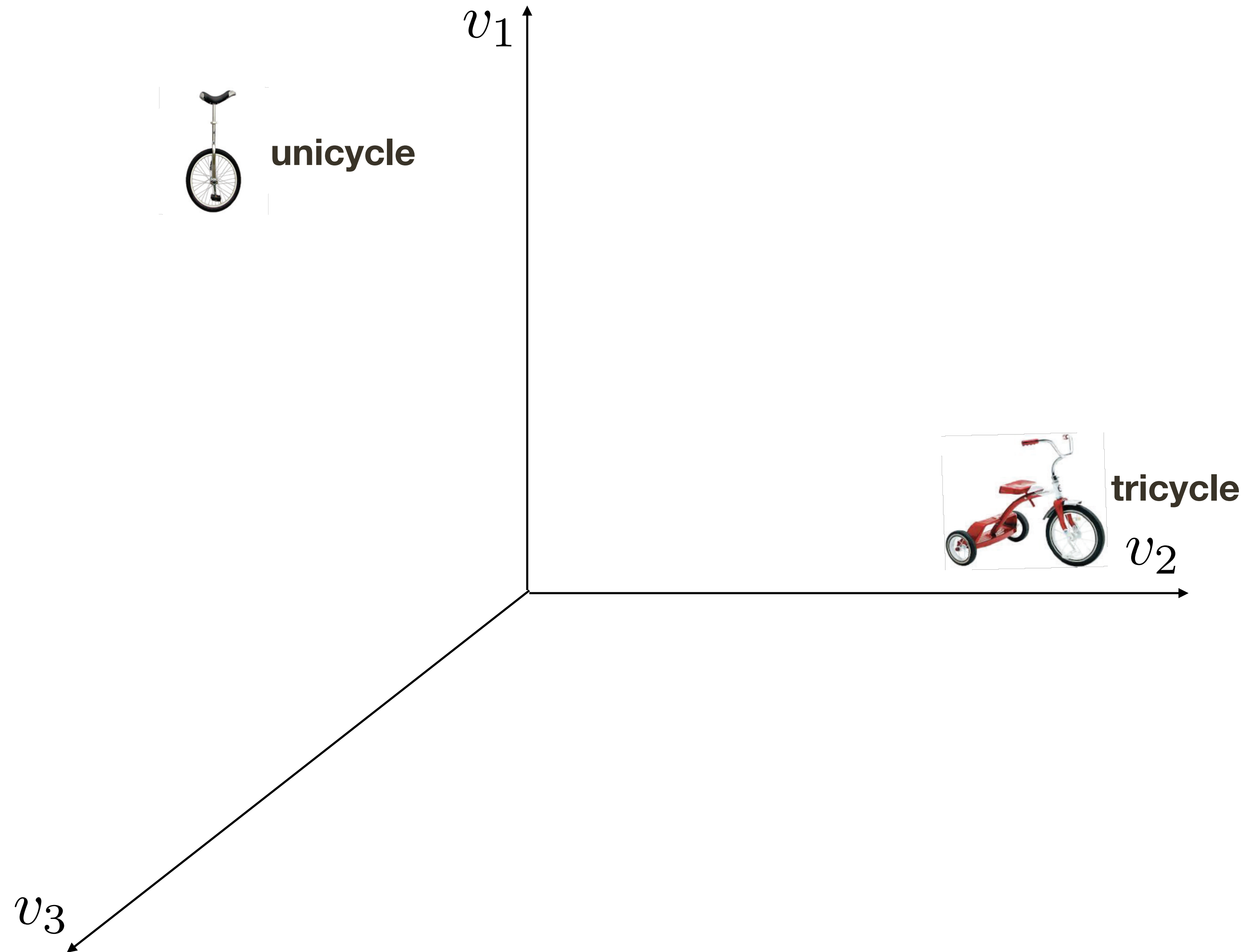
Objective Function:

$$\min_{\mathbf{W}} \sum_i^N \mathcal{L}_C(\mathbf{W}, \mathbf{V}, I_i, y_i) + \mathcal{L}_R(\mathbf{W}, \mathbf{V}, I_i, y_i) + \mu \|\mathbf{V}\|_F^2$$



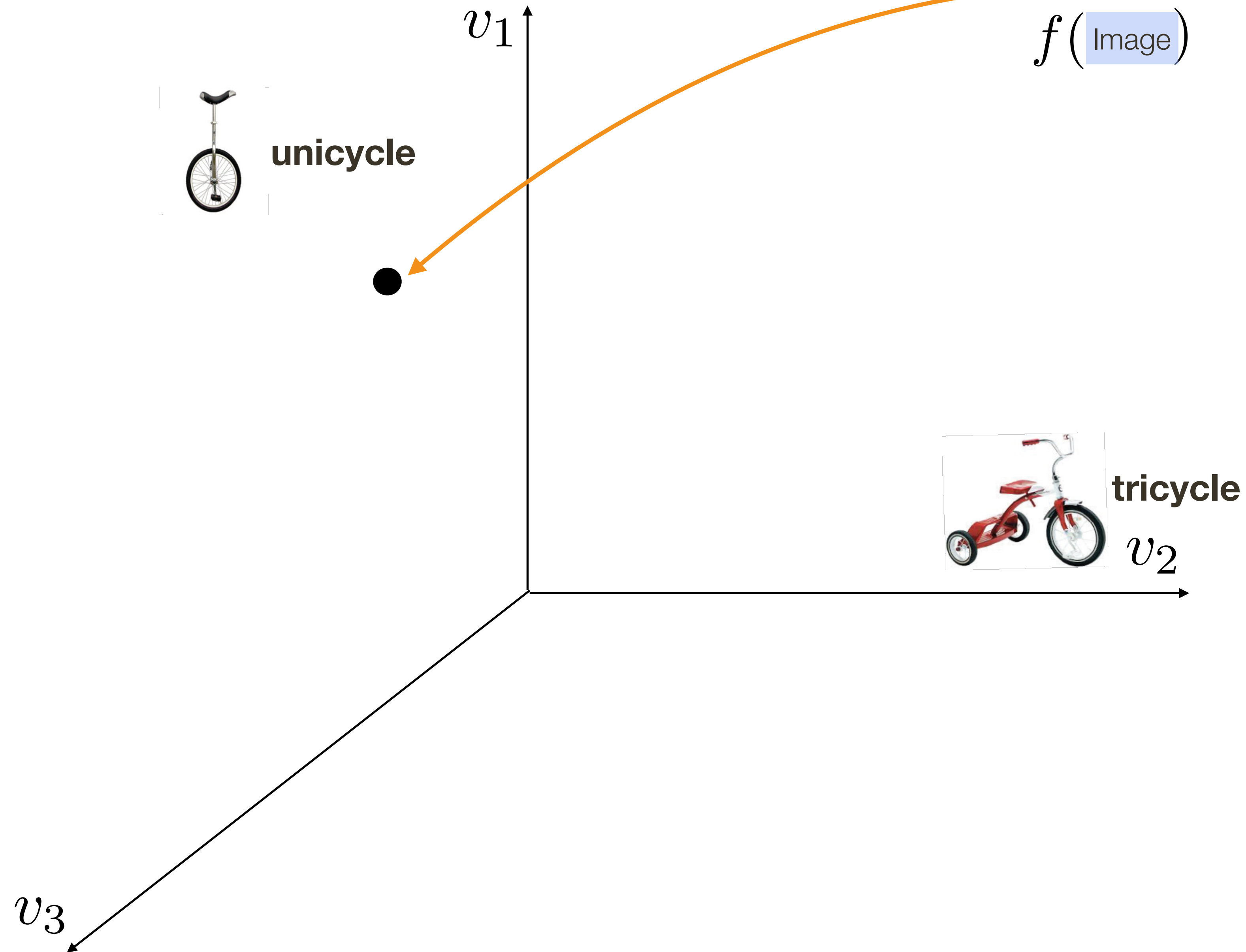
Vocabulary Informed Recognition

[Fu et al., 2016]



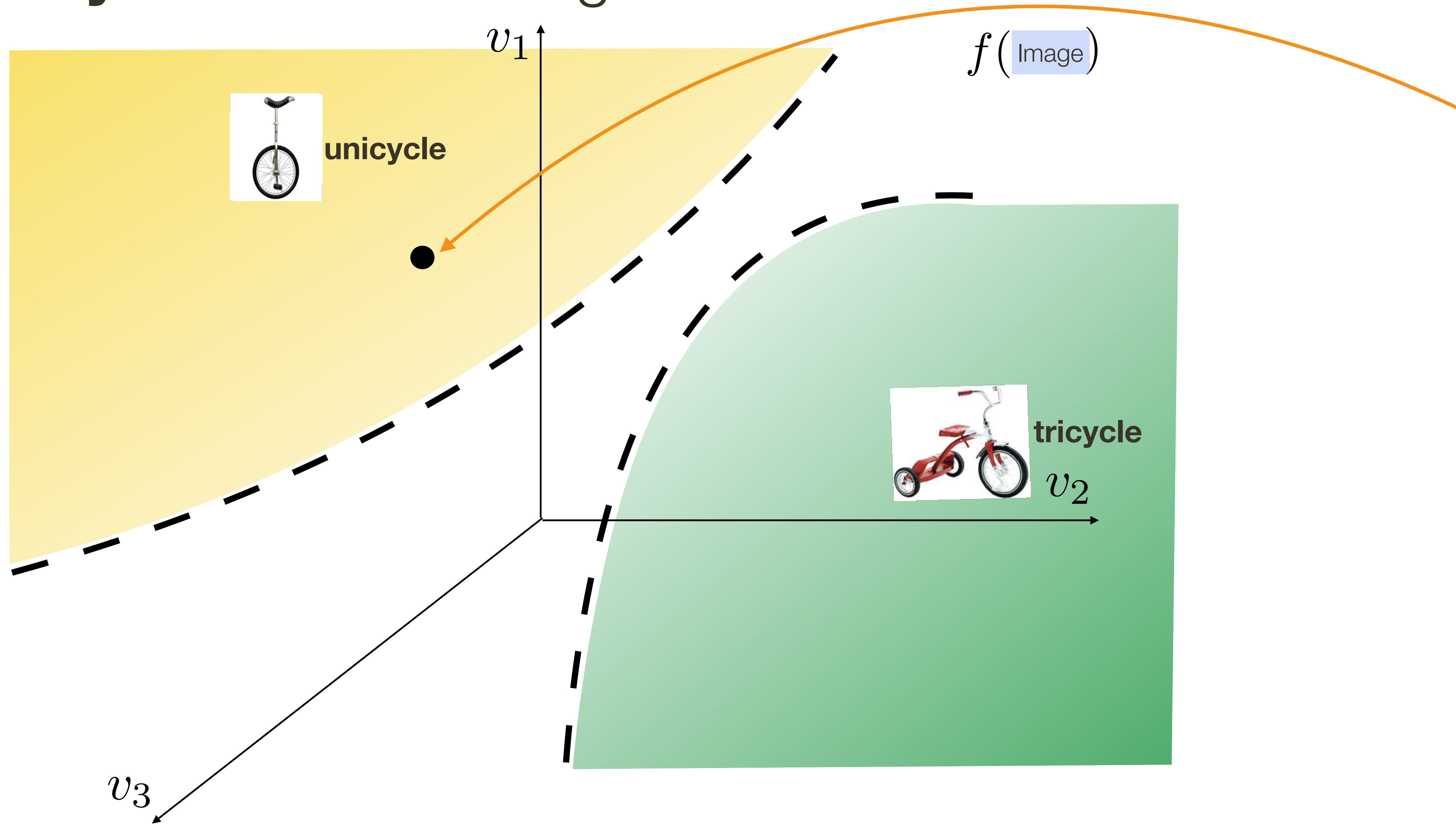
Vocabulary Informed Recognition

[Fu et al., 2016]



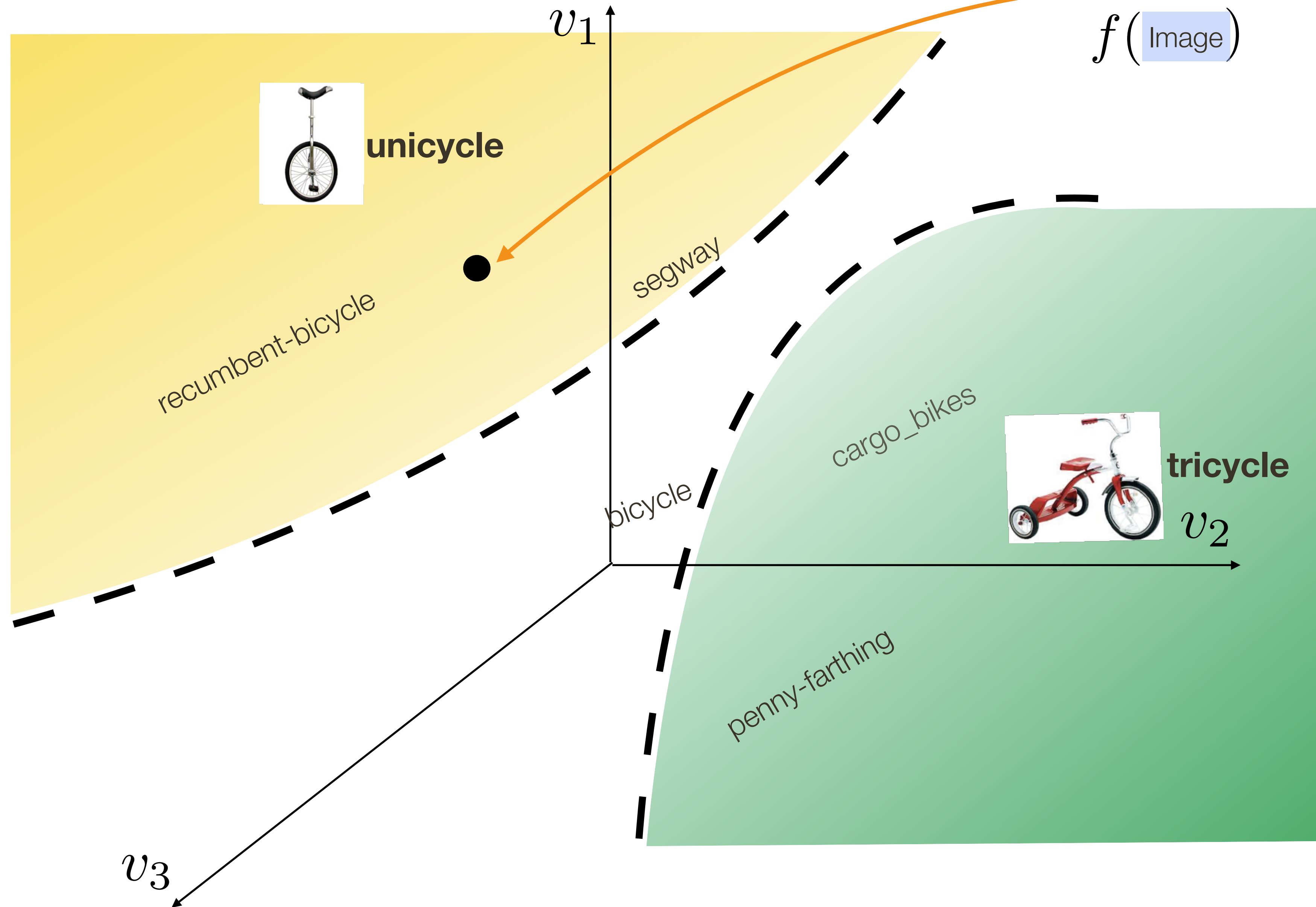
Vocabulary Informed Recognition

[Fu et al., 2016]



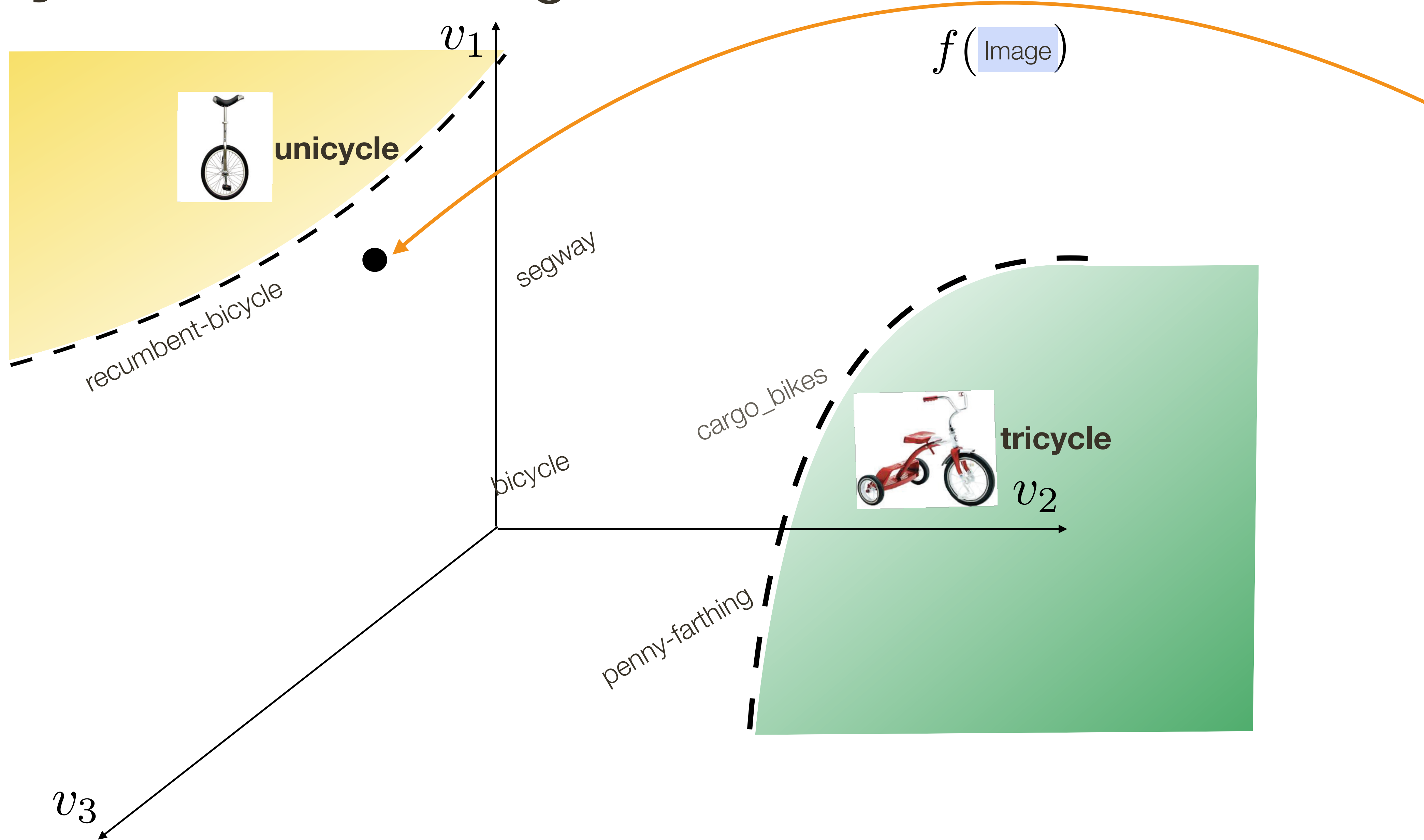
Vocabulary Informed Recognition

[Fu et al., 2016]



Vocabulary Informed Recognition

[Fu et al., 2016]



Zero-shot Results

[Fu et al., 2016]

Results with AWA

Method	Features	Accuracy	
SS-Voc: full instances	CNN _{OverFeat}	78.3	+4.4%
Akata <i>et al.</i> CVPR 2015	CNN _{GoogLeNet}	73.9	
TMV-BLP (Fu <i>et al.</i> ECCV 2014)	CNN _{OverFeat}	69.9	
AMP (SR+SE) (Fu <i>et al.</i> CVPR 2015)	CNN _{OverFeat}	66.0	
DAP (Lampert <i>et al.</i> TPAMI 2013)	CNN _{VGG19}	57.5	
PST (Rohrbach <i>et al.</i> NIPS 2013)	CNN _{OverFeat}	53.2	
DS (Rohrbach <i>et al.</i> CVPR 2010)	CNN _{OverFeat}	52.7	
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN _{OverFeat}	44.5	
HEX (Deng <i>et al.</i> ECCV 2014)	CNN _{DECAF}	44.2	

Zero-shot Results

[Fu et al., 2016]

Results with AWA

**3.3% of
training data**

Method	Features	Accuracy
SS-Voc: full instances	CNN _{OverFeat}	78.3
800 instances (20 inst*40 class);	CNN _{OverFeat}	74.4
		+0.5%
Akata et al. CVPR 2015	CNN _{GoogLeNet}	73.9
TMV-BLP (Fu et al. ECCV 2014)	CNN _{OverFeat}	69.9
AMP (SR+SE) (Fu et al. CVPR 2015)	CNN _{OverFeat}	66.0
DAP (Lampert et al. TPAMI 2013)	CNN _{VGG19}	57.5
PST (Rohrbach et al. NIPS 2013)	CNN _{OverFeat}	53.2
DS (Rohrbach et al. CVPR 2010)	CNN _{OverFeat}	52.7
IAP (Lampert et al. TPAMI 2013)	CNN _{OverFeat}	44.5
HEX (Deng et al. ECCV 2014)	CNN _{DECAF}	44.2

Zero-shot Results

[Fu et al., 2016]

Results with AWA

**0.82% of
training data**

Method	Features	Accuracy
SS-Voc: full instances	CNN _{OverFeat}	78.3
800 instances (20 inst*40 class);	CNN _{OverFeat}	74.4
200 instances (5 inst*40 class);	CNN _{OverFeat}	68.9
Akata et al. CVPR 2015	CNN _{GoogLeNet}	73.9
TMV-BLP (Fu et al. ECCV 2014)	CNN _{OverFeat}	69.9
AMP (SR+SE) (Fu et al. CVPR 2015)	CNN _{OverFeat}	66.0
DAP (Lampert et al. TPAMI 2013)	CNN _{VGG19}	57.5
PST (Rohrbach et al. NIPS 2013)	CNN _{OverFeat}	53.2
DS (Rohrbach et al. CVPR 2010)	CNN _{OverFeat}	52.7
IAP (Lampert et al. TPAMI 2013)	CNN _{OverFeat}	44.5
HEX (Deng et al. ECCV 2014)	CNN _{DECAF}	44.2

Weakly-supervised **Visual Grounding** of Phrases [Xiao et al., 2017]

Given **image-sentence pairs** learn how to **localize** arbitrary language phrase or sentence in new images



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

Weakly-supervised **Visual Grounding** of Phrases [Xiao et al., 2017]

Given **image-sentence pairs** learn how to **localize** arbitrary language phrase or sentence in new images



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

a man



Weakly-supervised **Visual Grounding** of Phrases [Xiao et al., 2017]

Given **image-sentence pairs** learn how to **localize** arbitrary language phrase or sentence in new images



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

a man



Weakly-supervised **Visual Grounding** of Phrases [Xiao et al., 2017]

Given **image-sentence pairs** learn how to **localize** arbitrary language phrase or sentence in new images



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

a table



Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

Label Embedding 

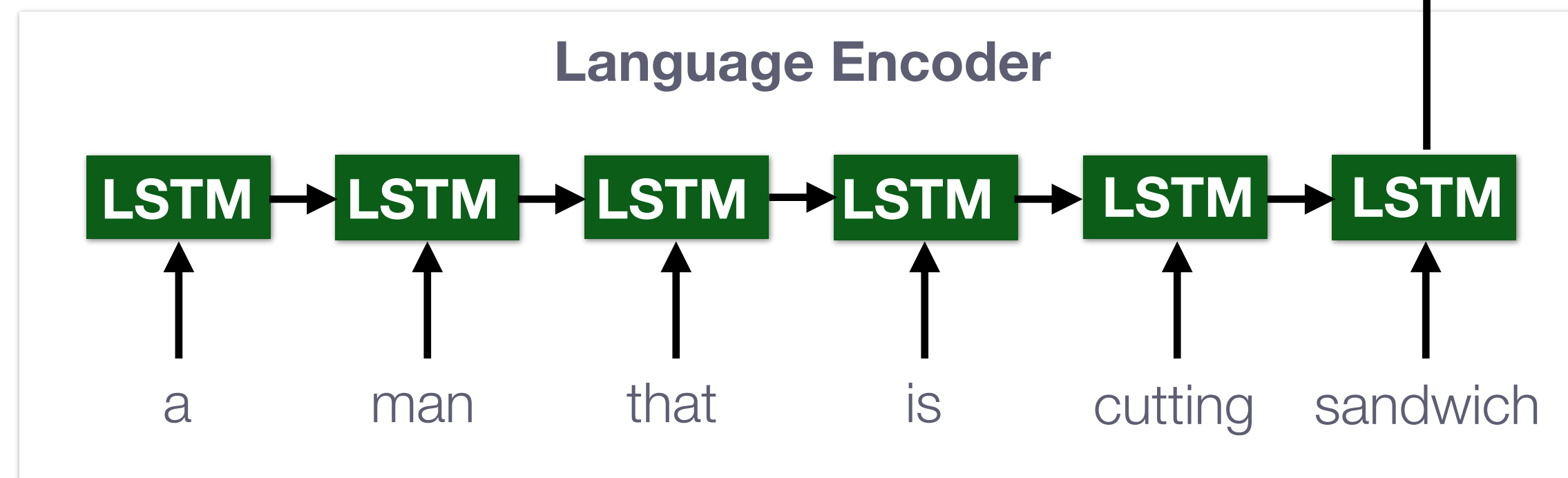
$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

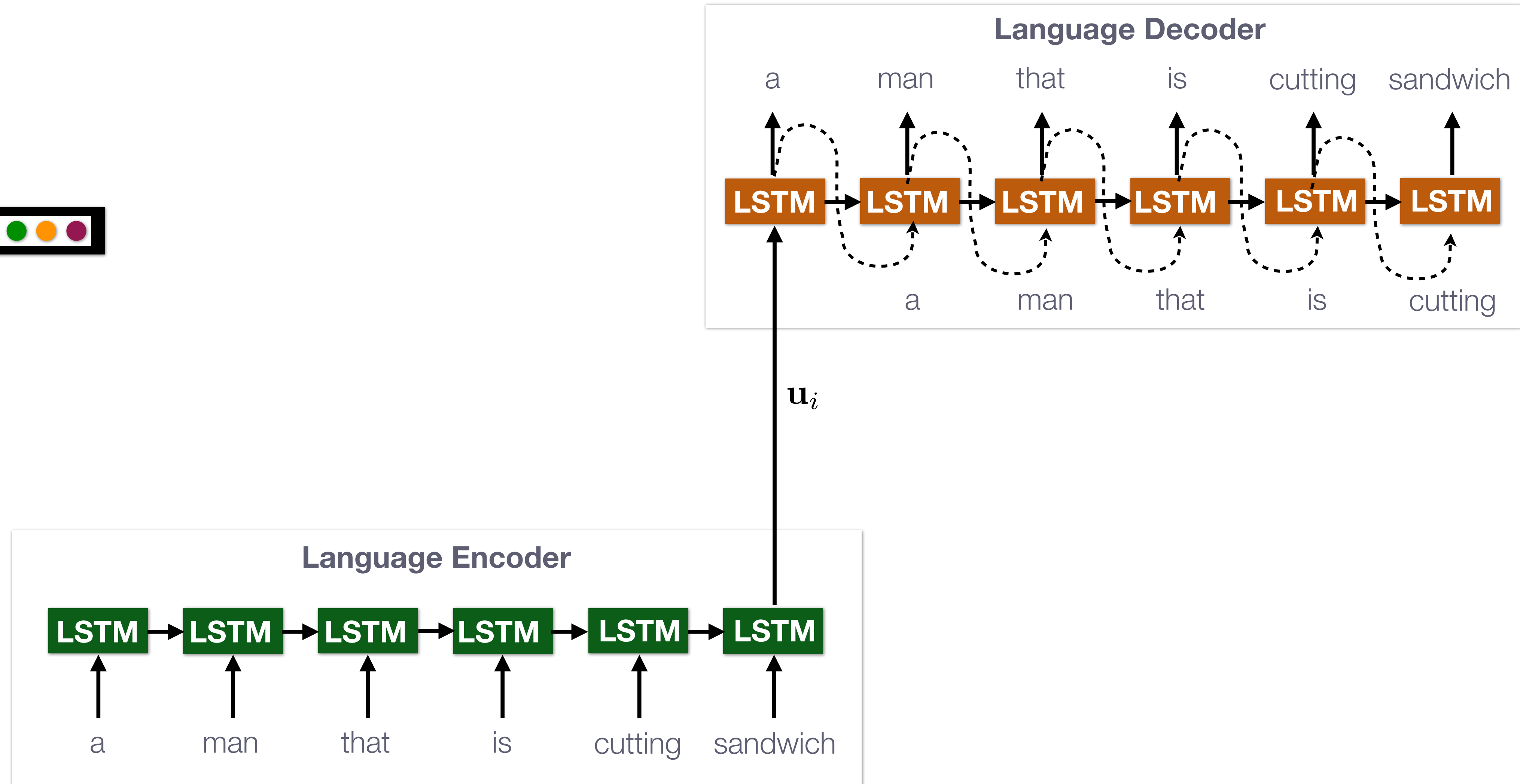


Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

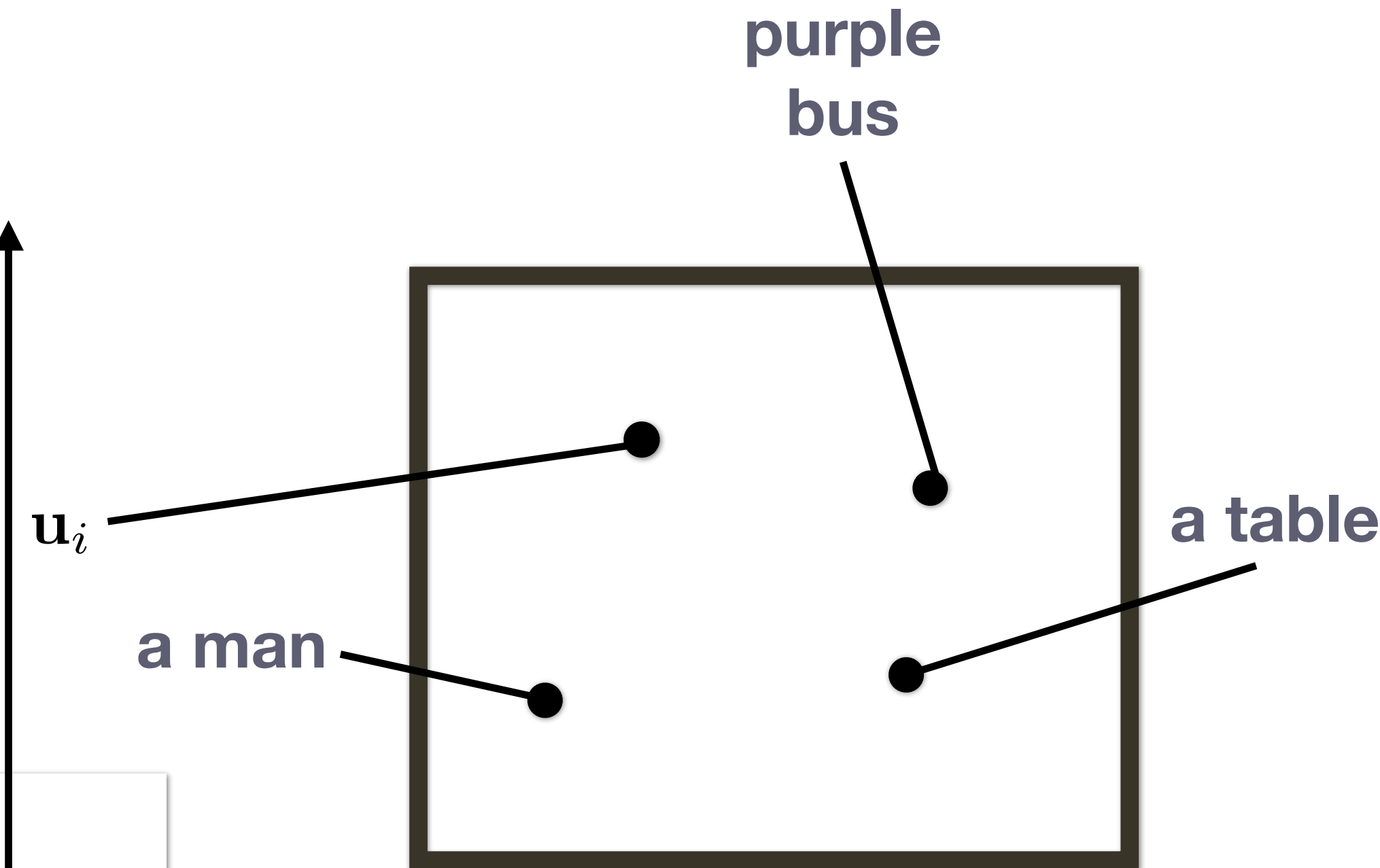
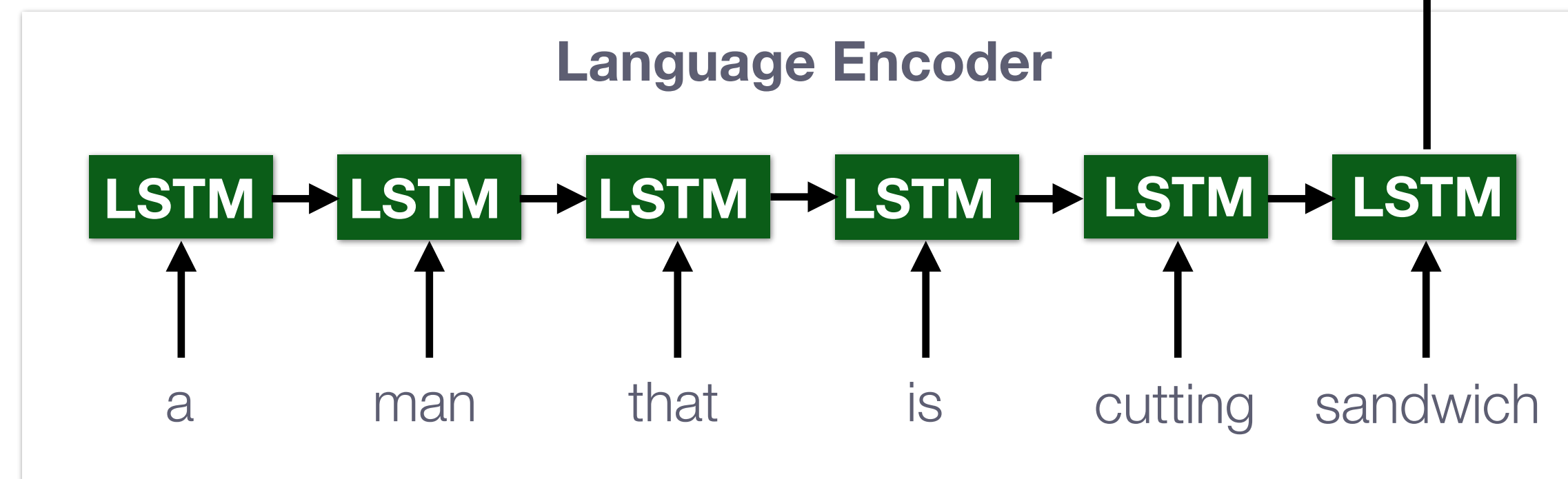


Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

Label Embedding ●●●●

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

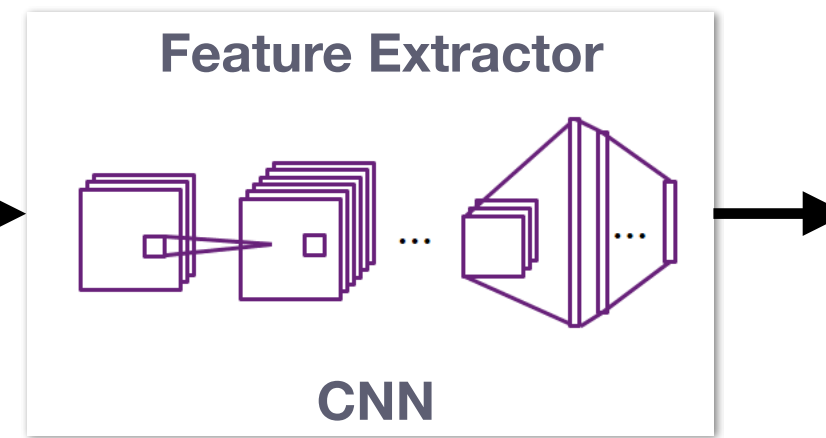


Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

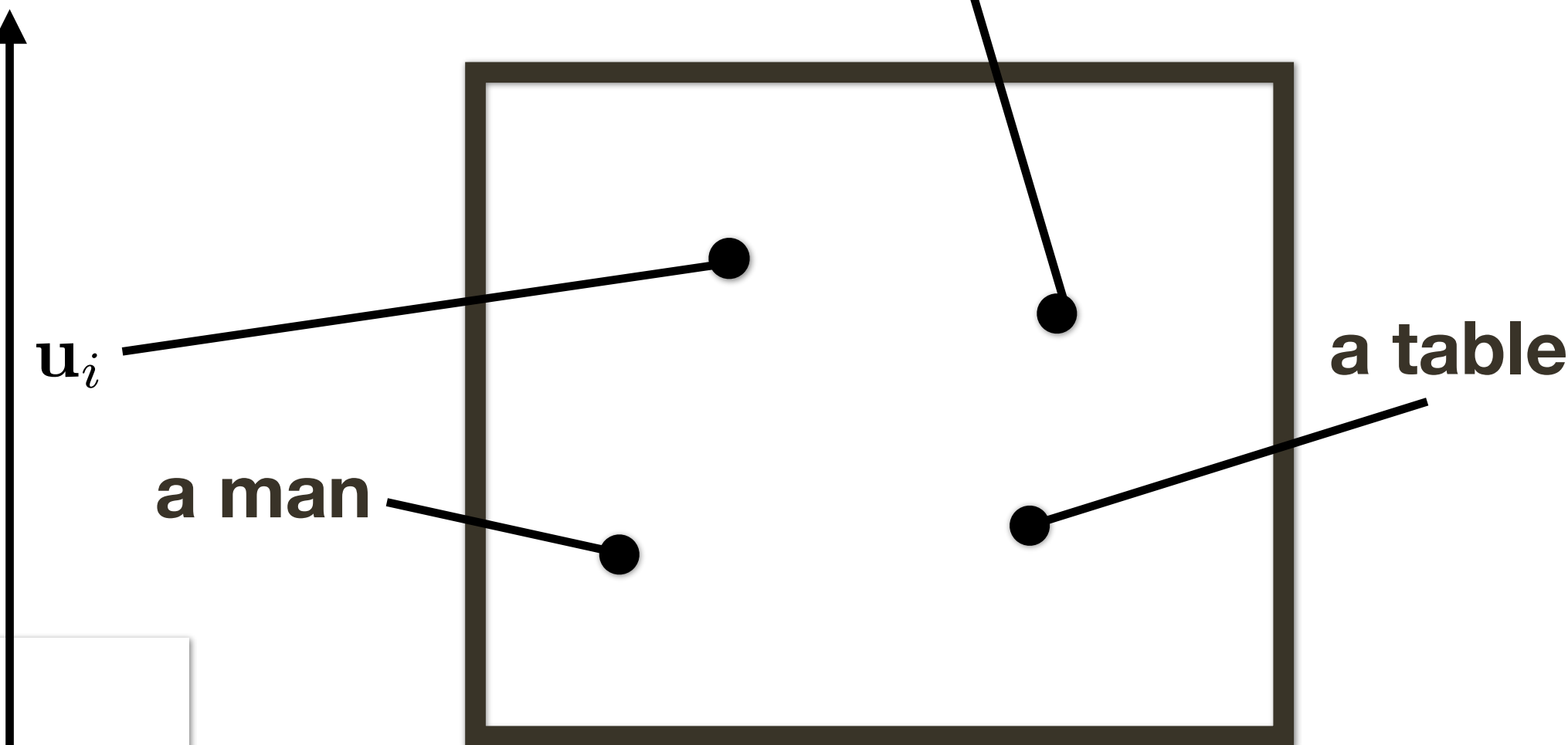
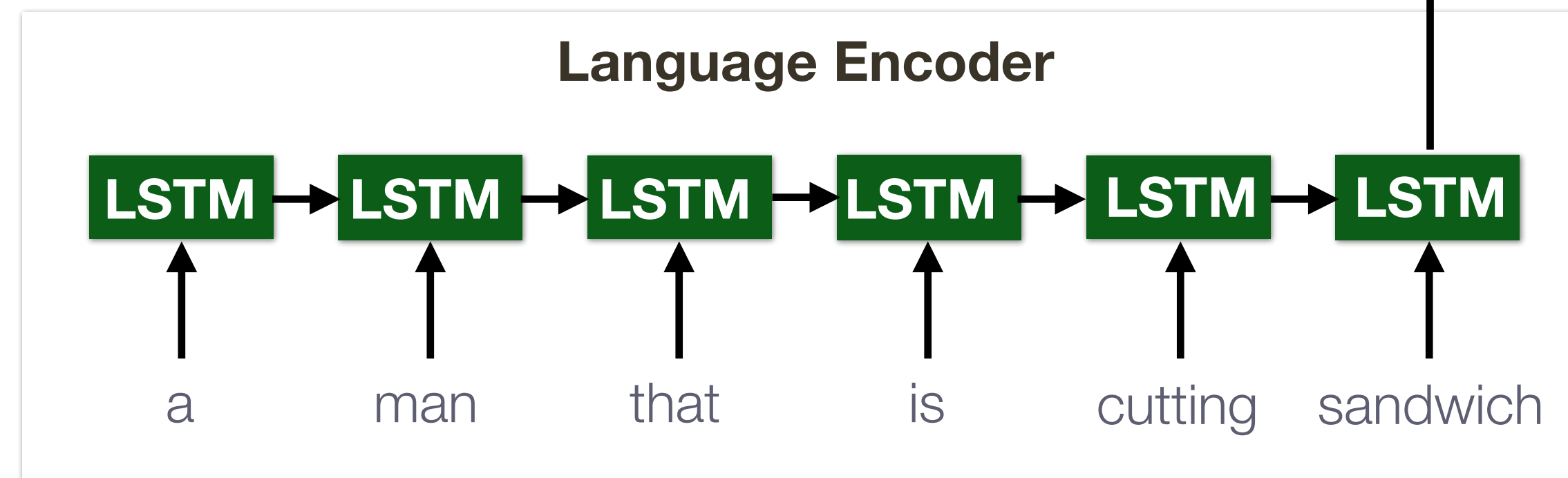
Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$



Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$



Weakly-supervised **Visual Grounding** of Phrases

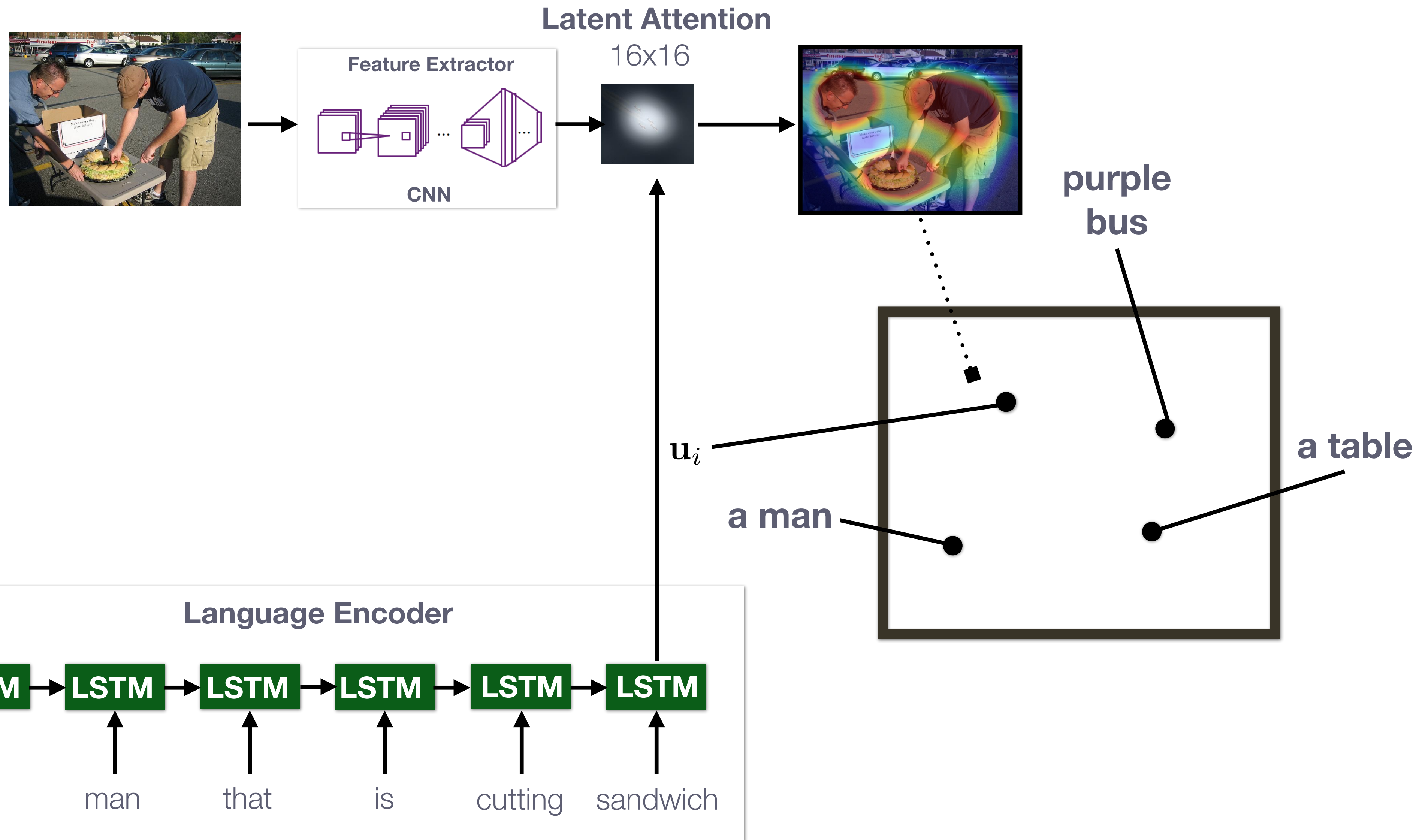
[Xiao et al., 2017]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$

Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$



Weakly-supervised **Visual Grounding** of Phrases

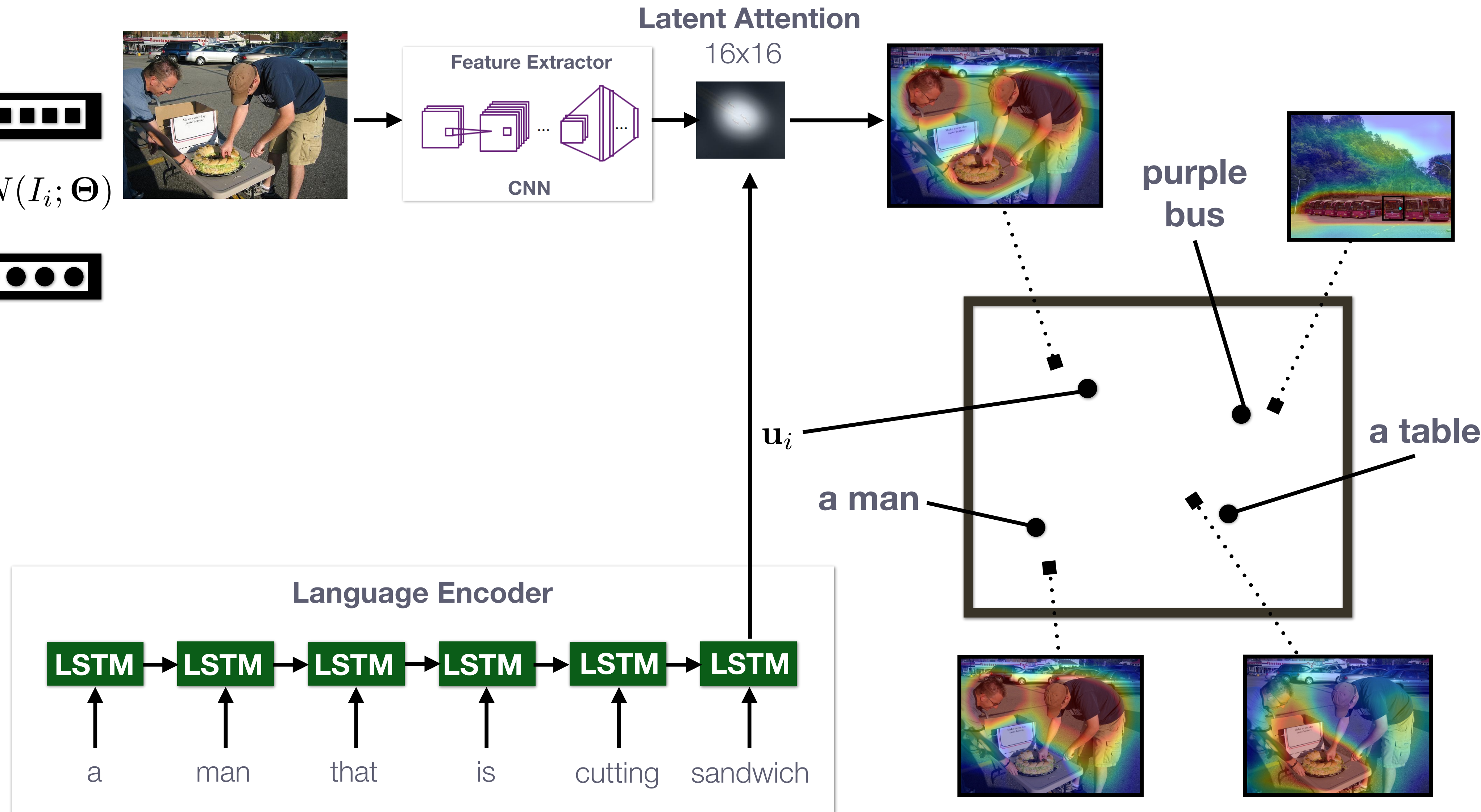
[Xiao et al., 2017]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$

Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$



Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$

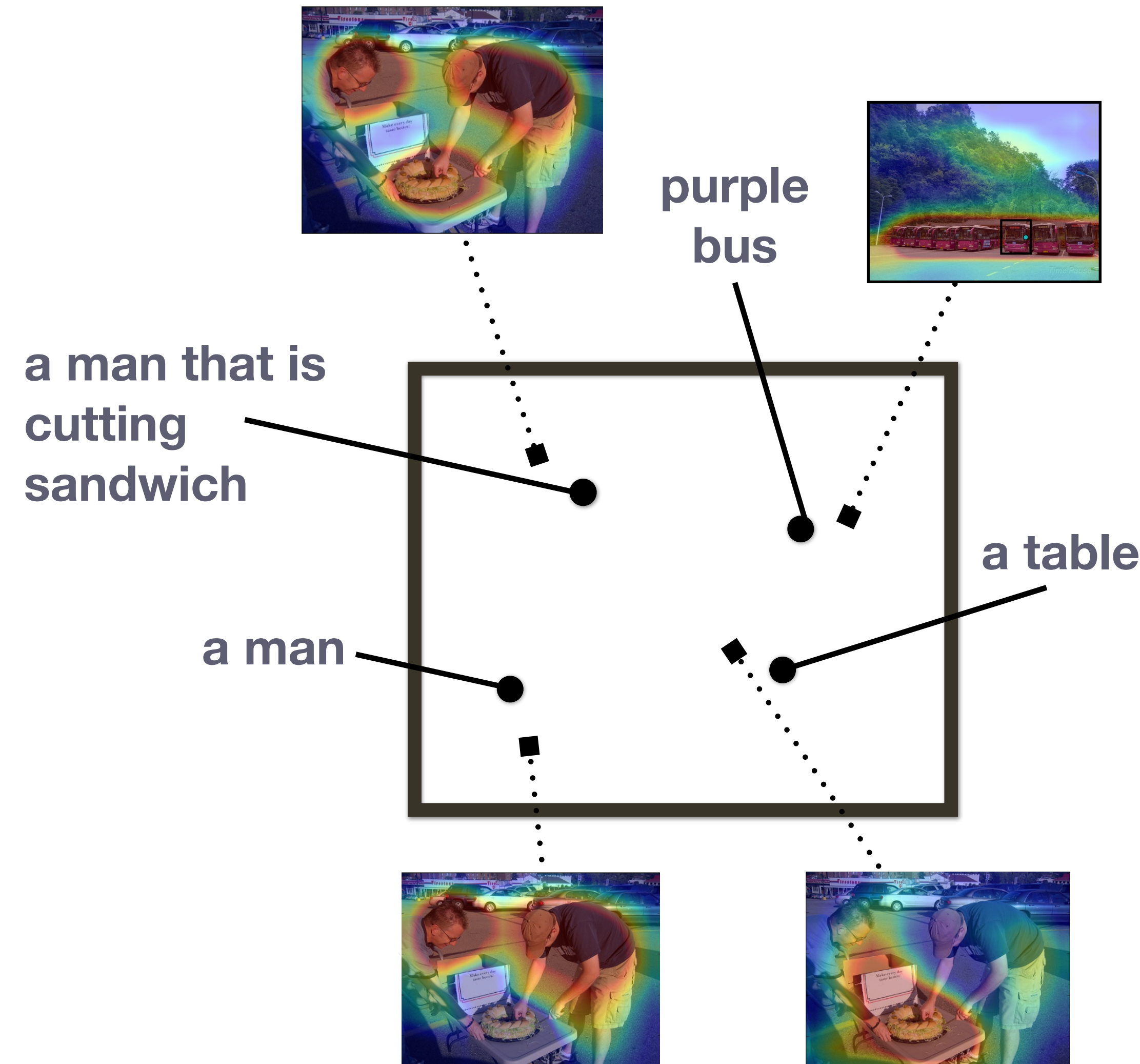
Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:



Combination of previous discriminative similarity and **linguistic regularization**

Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$

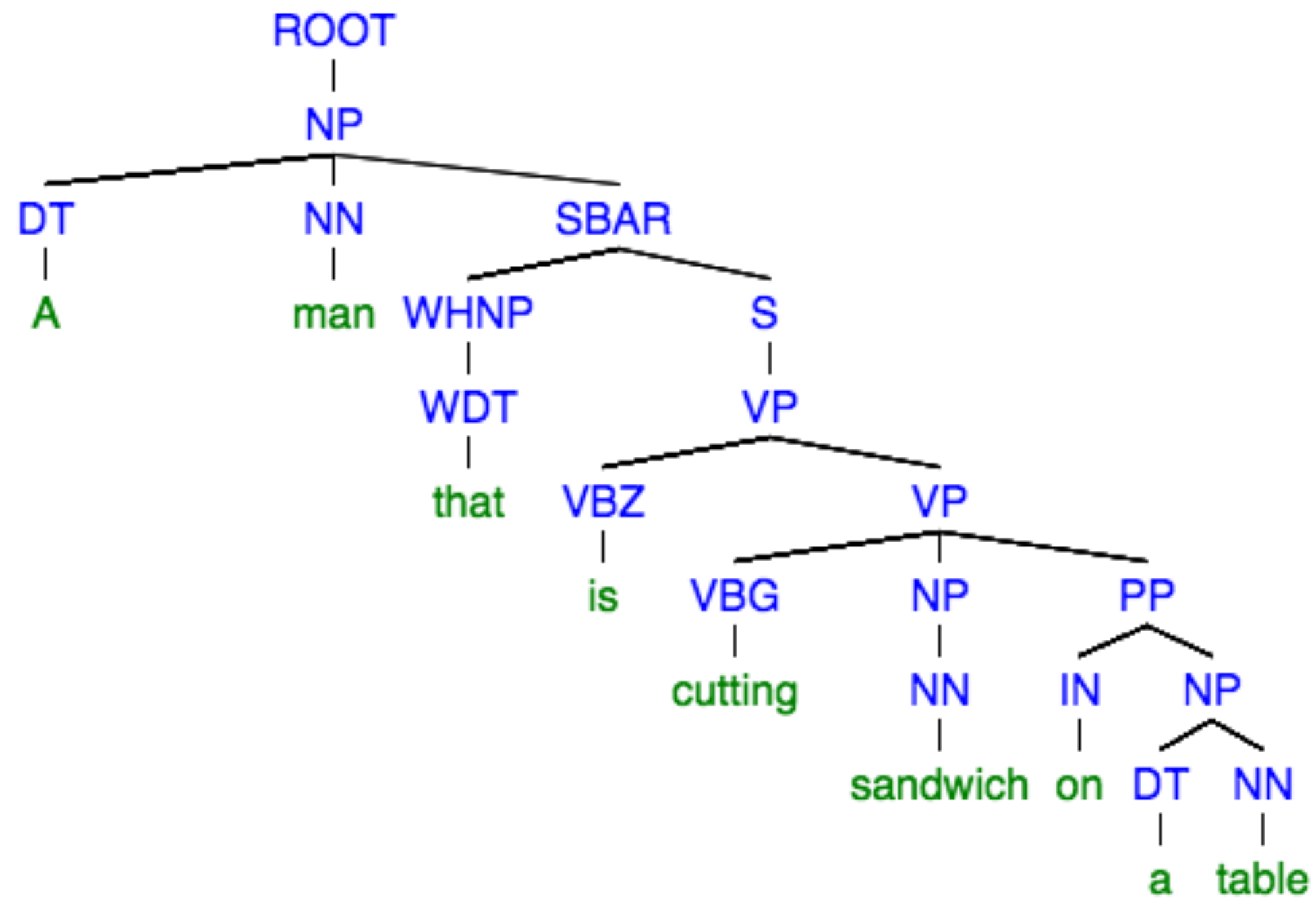
Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:



Combination of previous discriminative similarity and **linguistic regularization**

Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

For **noun phrases**:

- **siblings** should have **disjoint** masks

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$

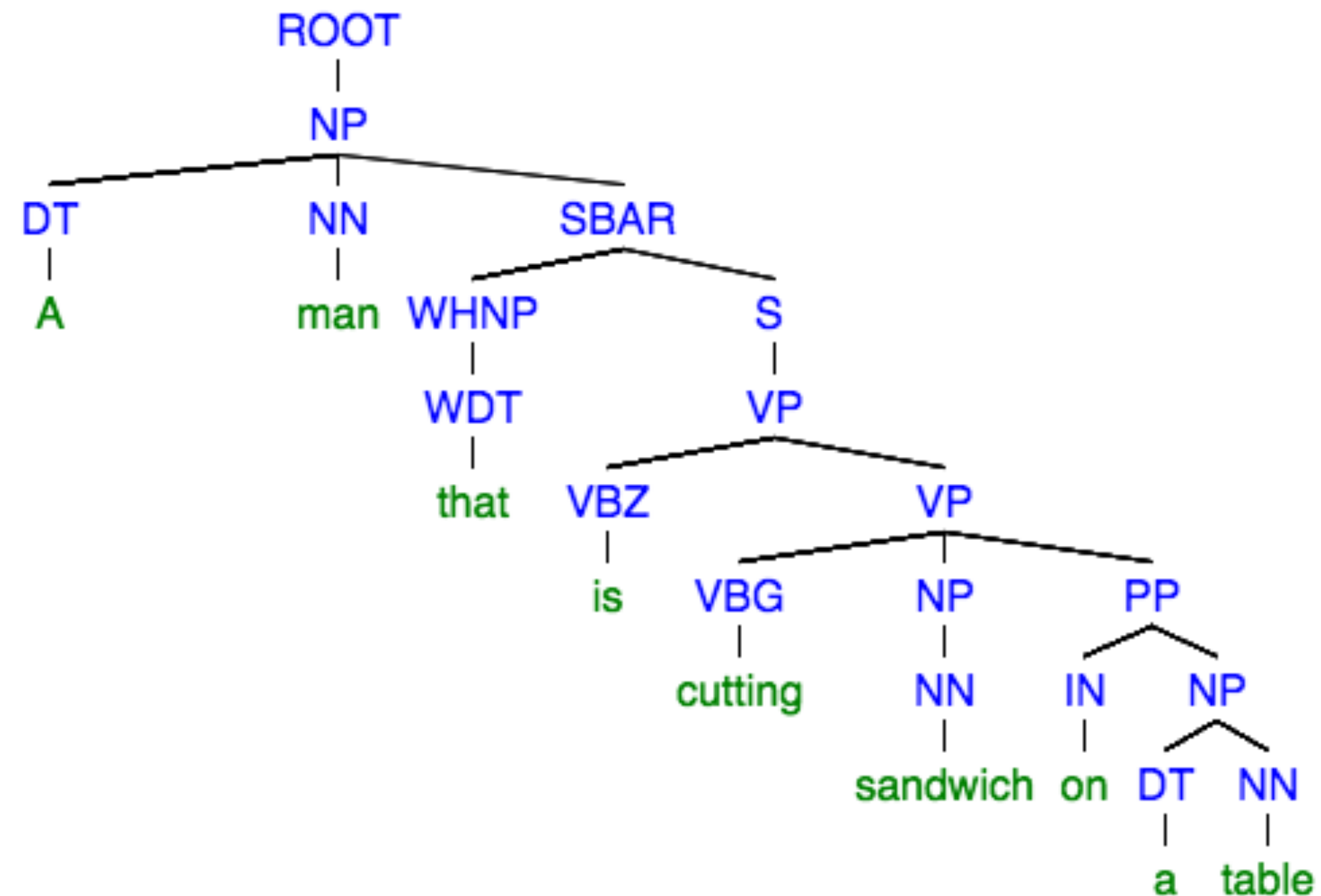
Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:



Combination of previous discriminative similarity and **linguistic regularization**

Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

For **noun phrases**:

- **siblings** should have **disjoint** masks

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$

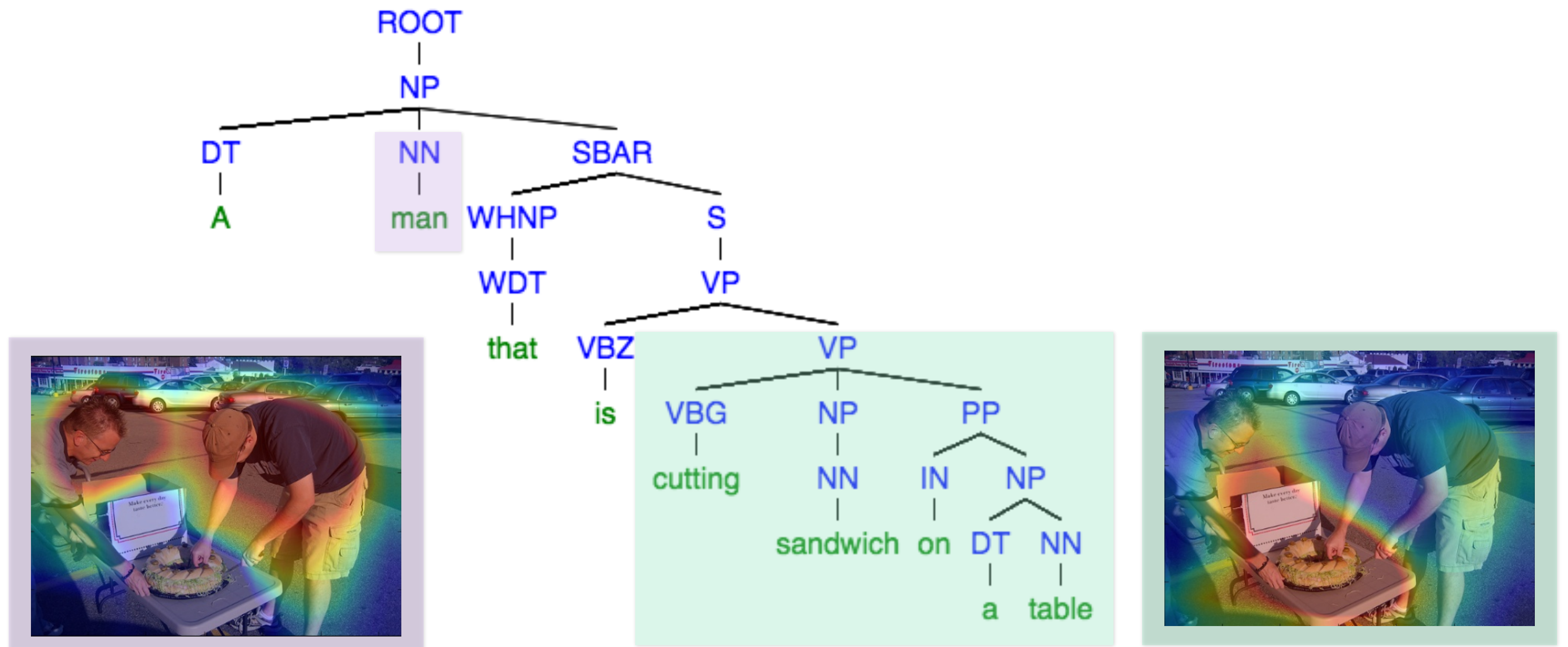
Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:



Combination of previous discriminative similarity and **linguistic regularization**

Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

For **noun phrases**:

- **siblings** should have **disjoint** masks
- **parents** should be **union of children** masks

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$

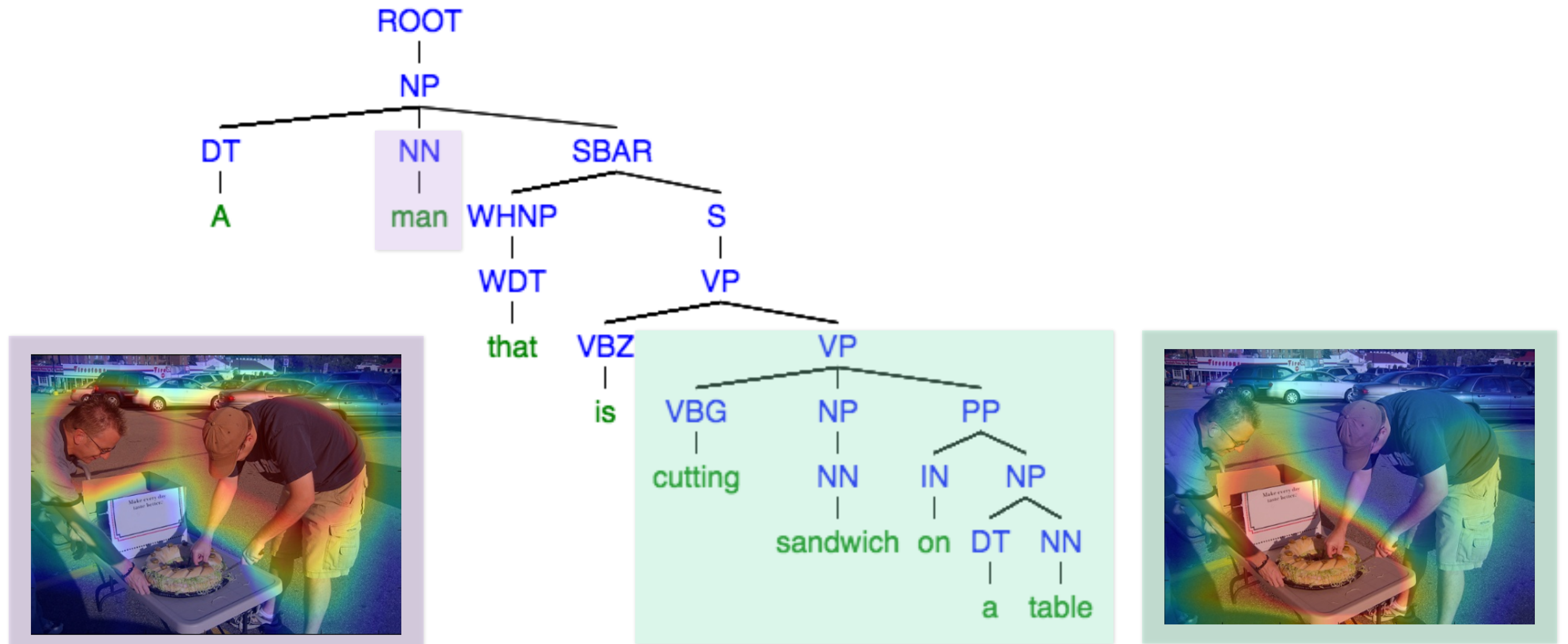
Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:



Combination of previous discriminative similarity and **linguistic regularization**

Weakly-supervised **Visual Grounding** of Phrases

[Xiao et al., 2017]

For **noun phrases**:

- **siblings** should have **disjoint**
- **parents** should be **union of**

Image Embedding 

$$\Psi(I_i) = \mathbf{W} \cdot \text{CNN}(I_i; \Theta)$$

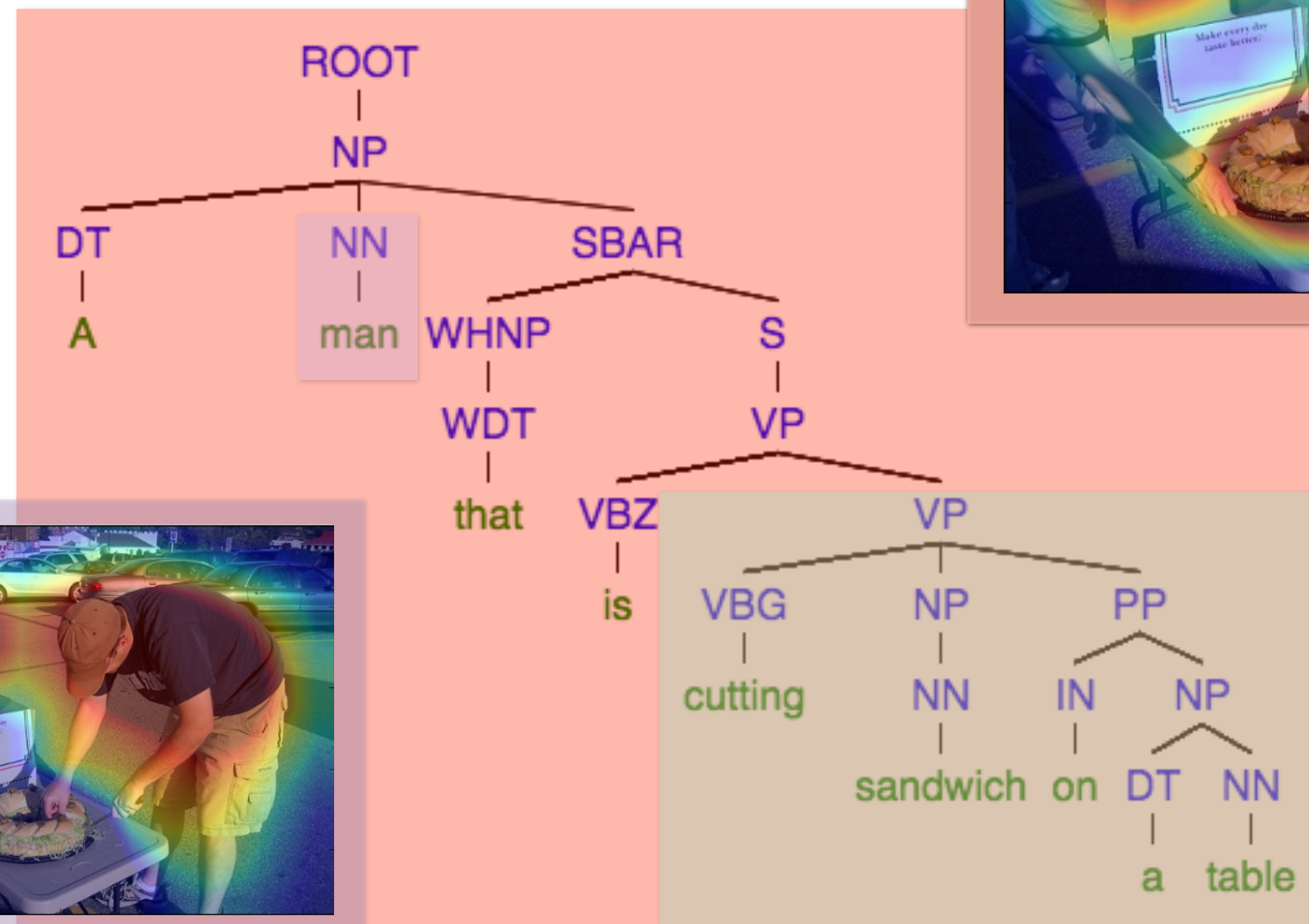
Label Embedding 

$$\Psi_L(\text{phrase}_i) = \mathbf{u}_i$$

Similarity in Embedding Space

$$D(\mathbf{u}, \mathbf{u}') = \|\mathbf{u} - \mathbf{u}'\|_2^2$$

Objective Function:



Combination of previous discriminative similarity and **linguistic regularization**

Qualitative Results

[Xiao et al., 2017]

Input:



guy in green t-shirt holding
skateboard

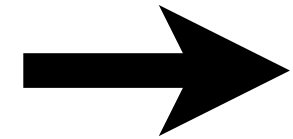
Qualitative Results

[Xiao et al., 2017]

Input:



NO linguistic constraints



guy in green t-shirt holding
skateboard

Qualitative Results

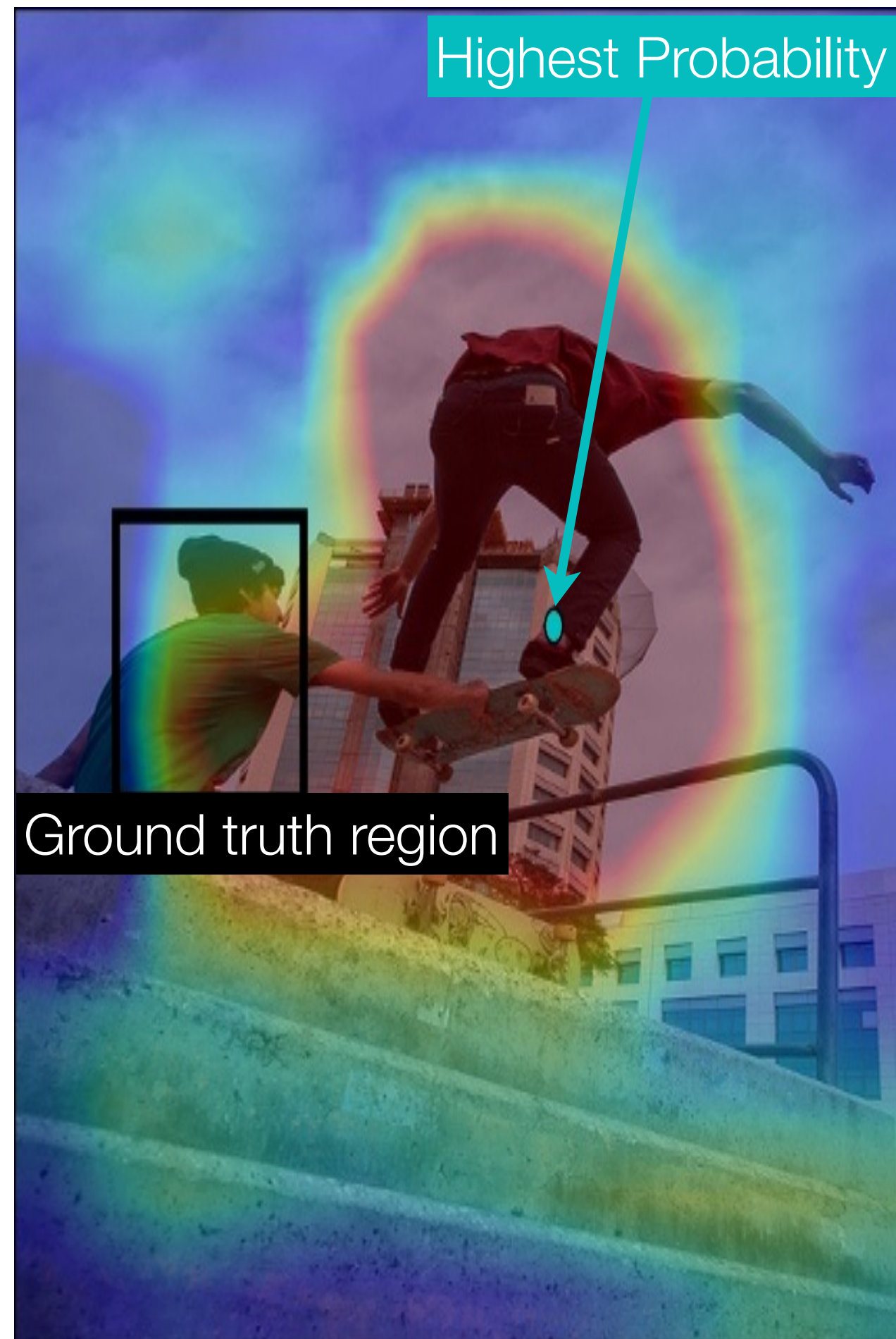
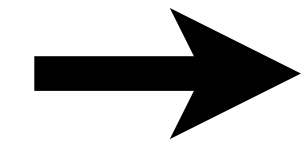
[Xiao et al., 2017]

Input:



guy in green t-shirt holding
skateboard

NO linguistic constraints



Ground truth region

Qualitative Results

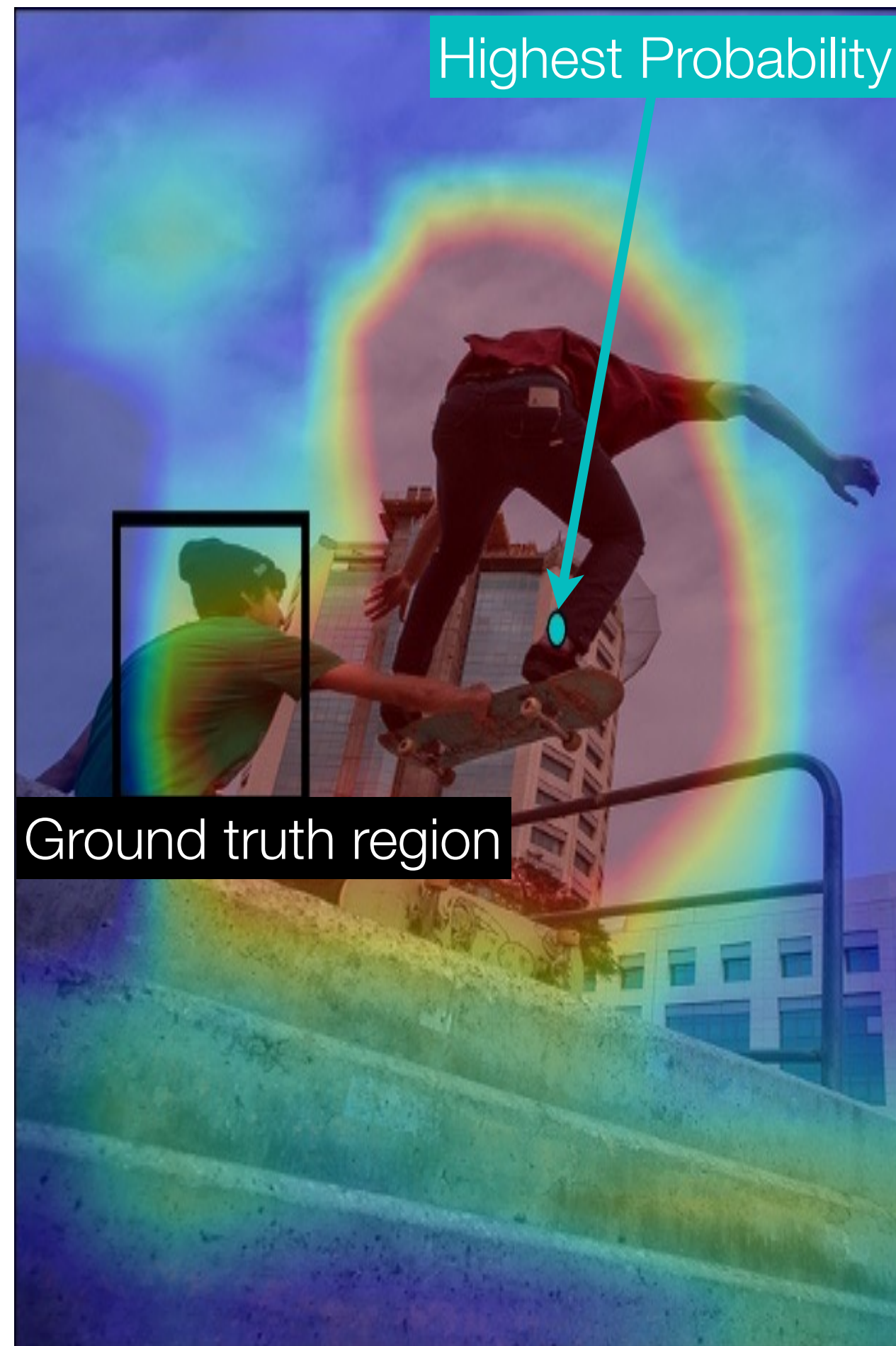
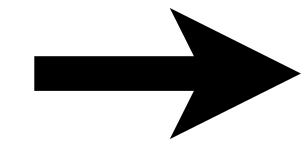
[Xiao et al., 2017]

Input:

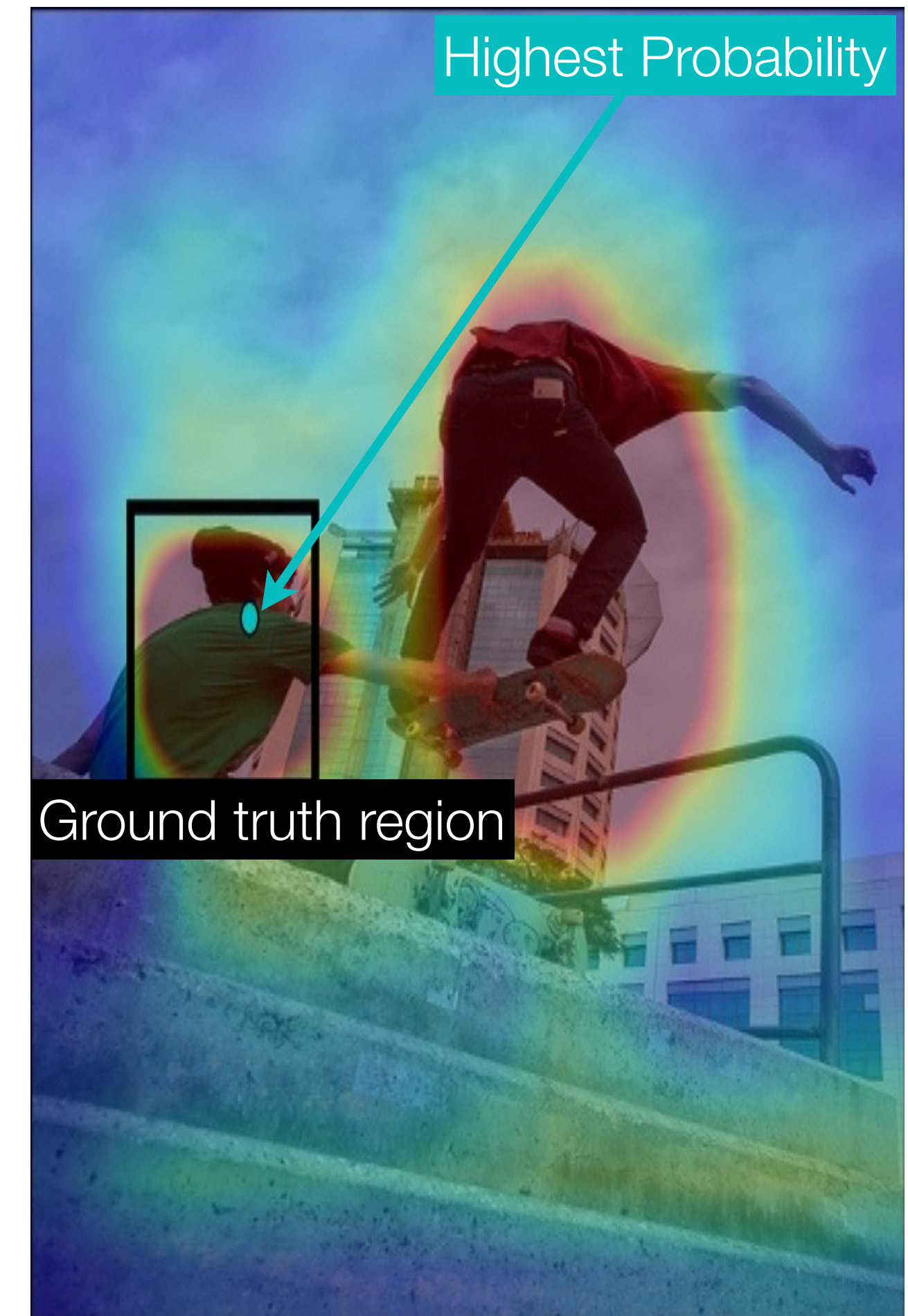


guy in green t-shirt holding
skateboard

NO linguistic constraints



Our Model



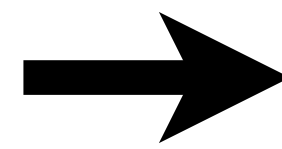
Qualitative Results

[Xiao et al., 2017]

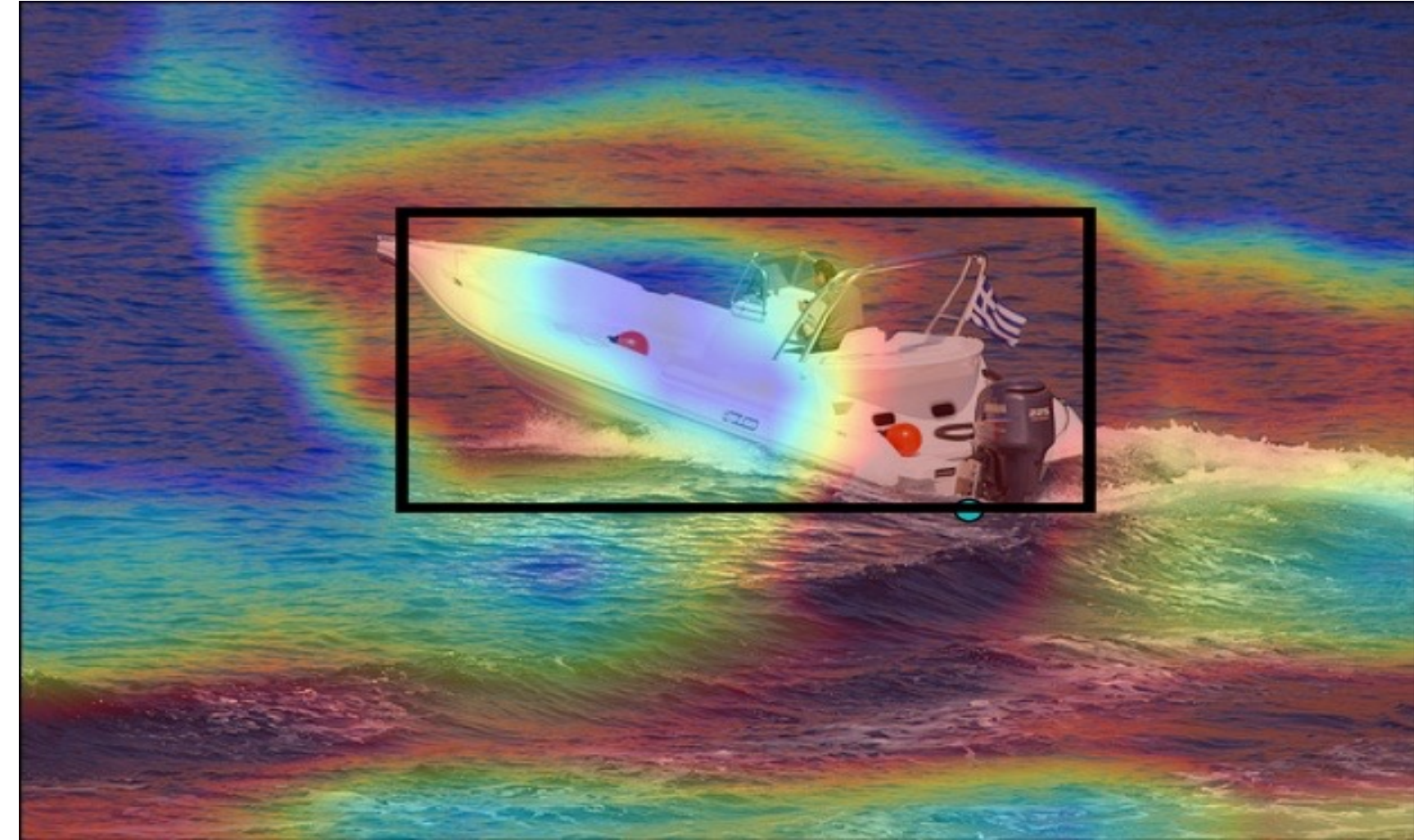
Input:



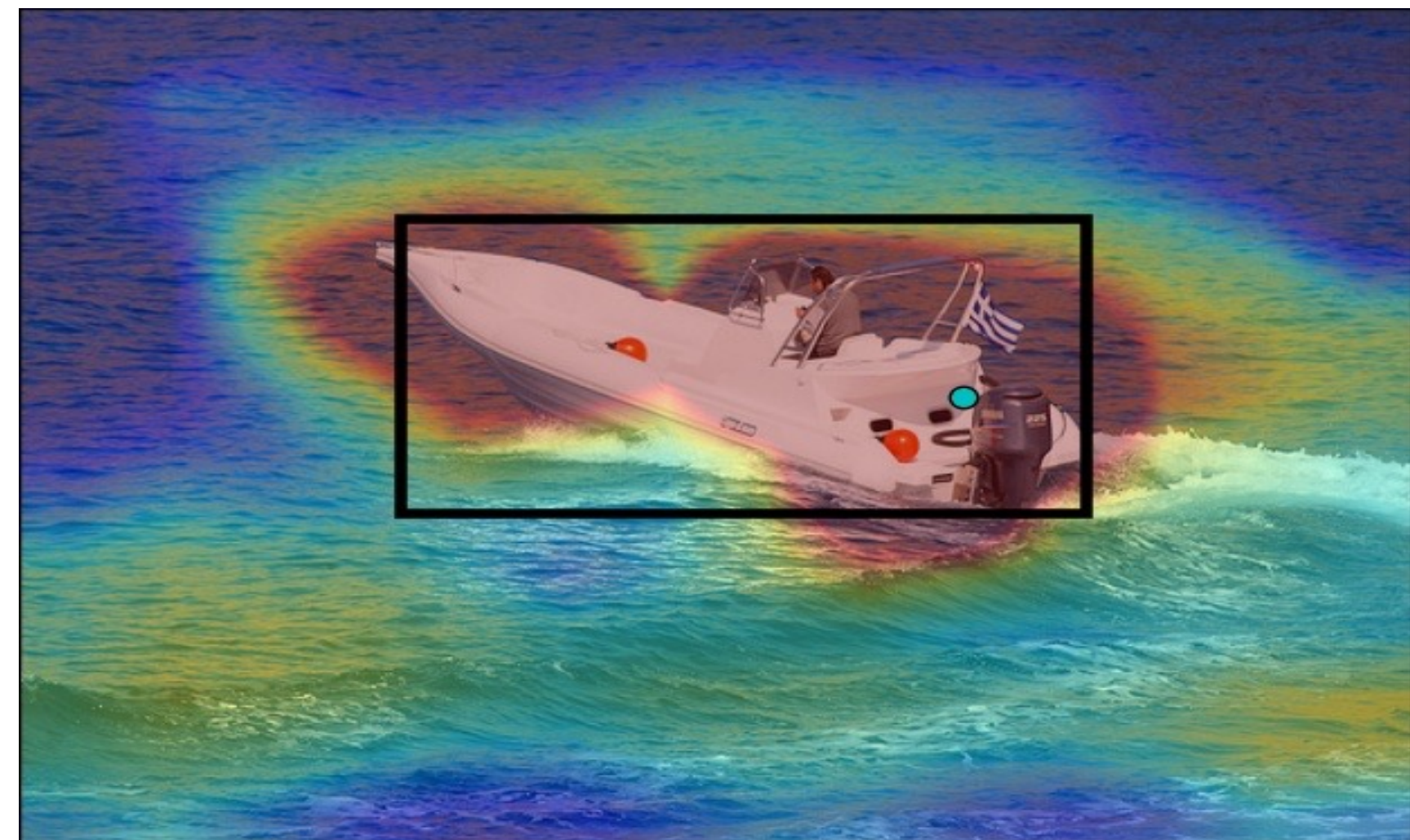
a person driving a boat



NO linguistic constraints



Our Model



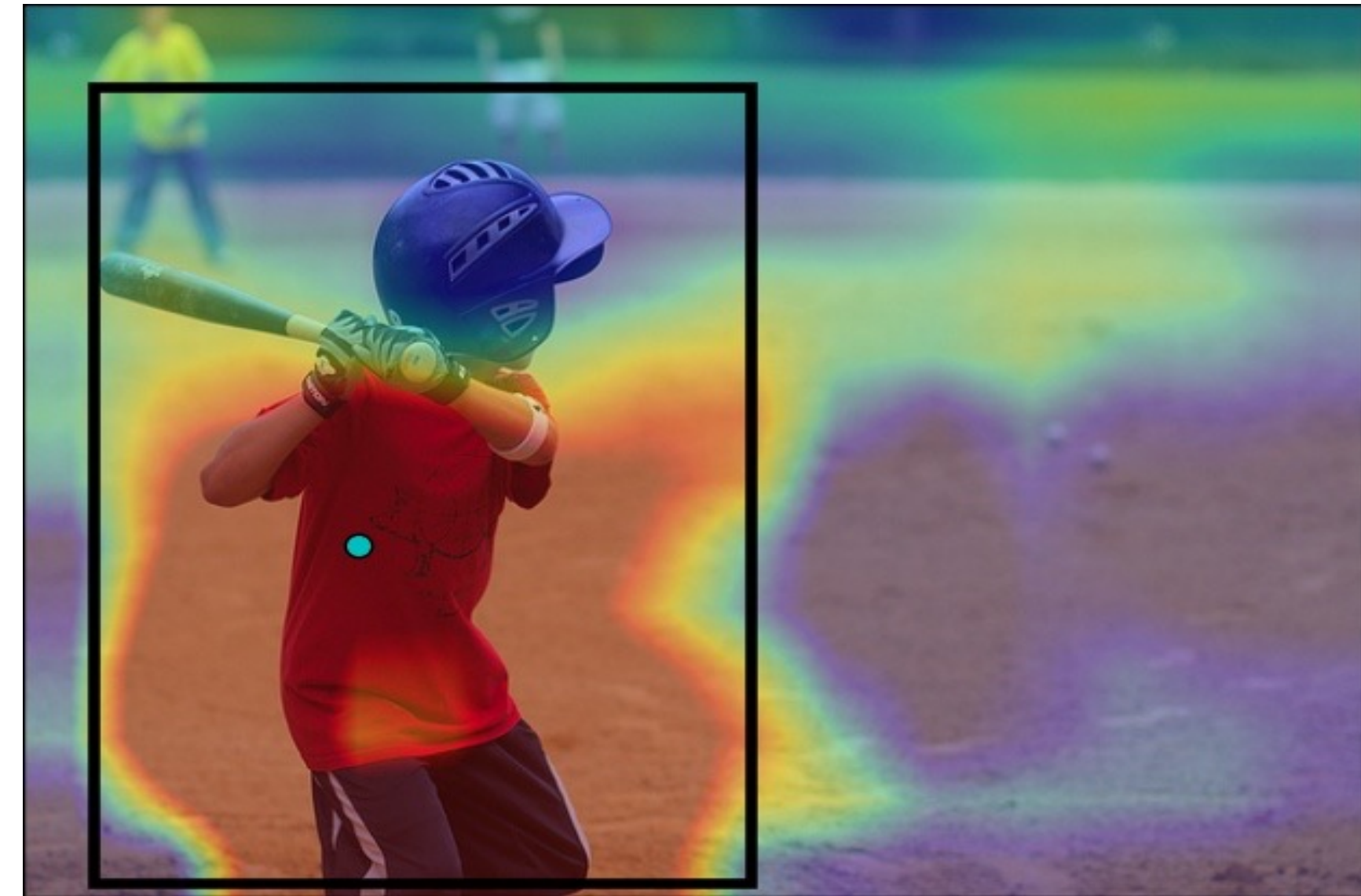
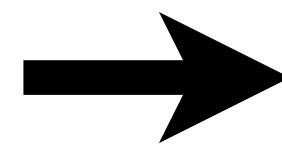
Qualitative Results

NO linguistic constraints [Xiao et al., 2017]

Input:



a child wearing black protective helmet



Our Model



Quantitative Results

[Xiao et al., 2017]

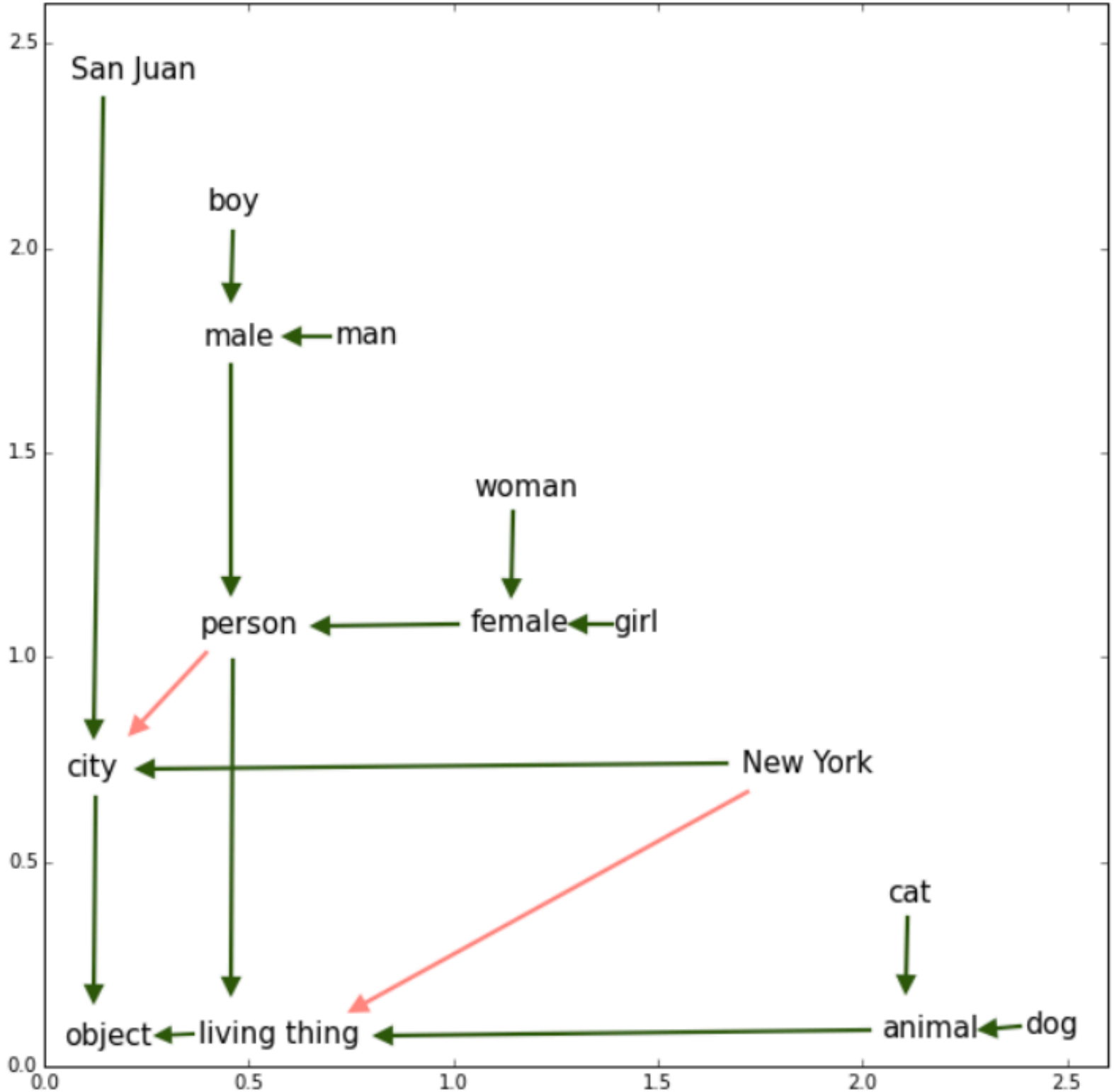
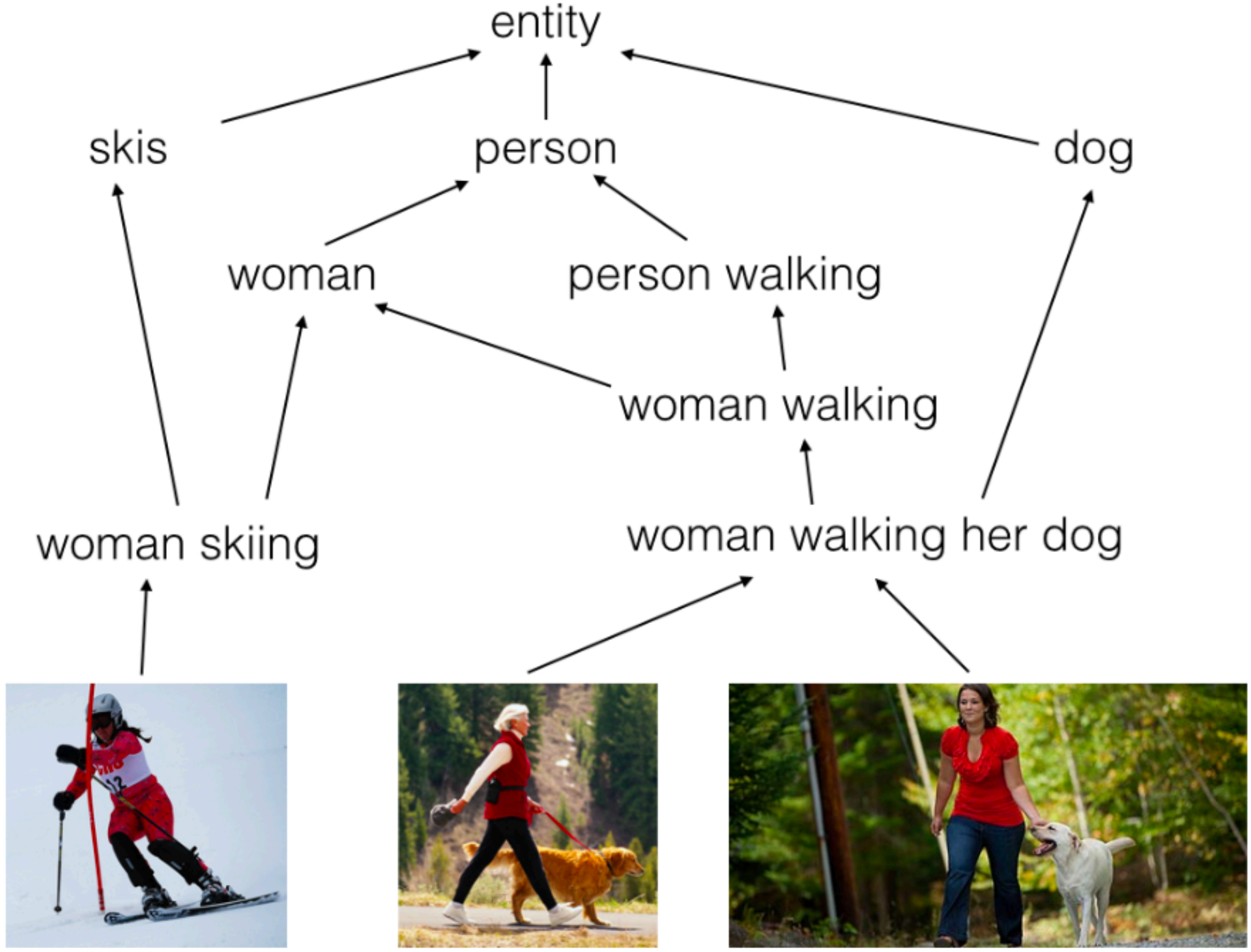
Segmentation performance on COCO dataset

[Lin, Maire, Belongie, Hays, Perona, Ramanan, Dollar, Zitnick, ECCV'14]

	IoU@0.3	IoU@0.4	IoU@0.5	Avg mAP
Non-structured	0.302	0.199	0.110	0.203
Parent-Child	0.327	0.213	0.118	0.219
Sibling	0.316	0.203	0.114	0.211
Ours	0.347	0.246	0.159	0.251

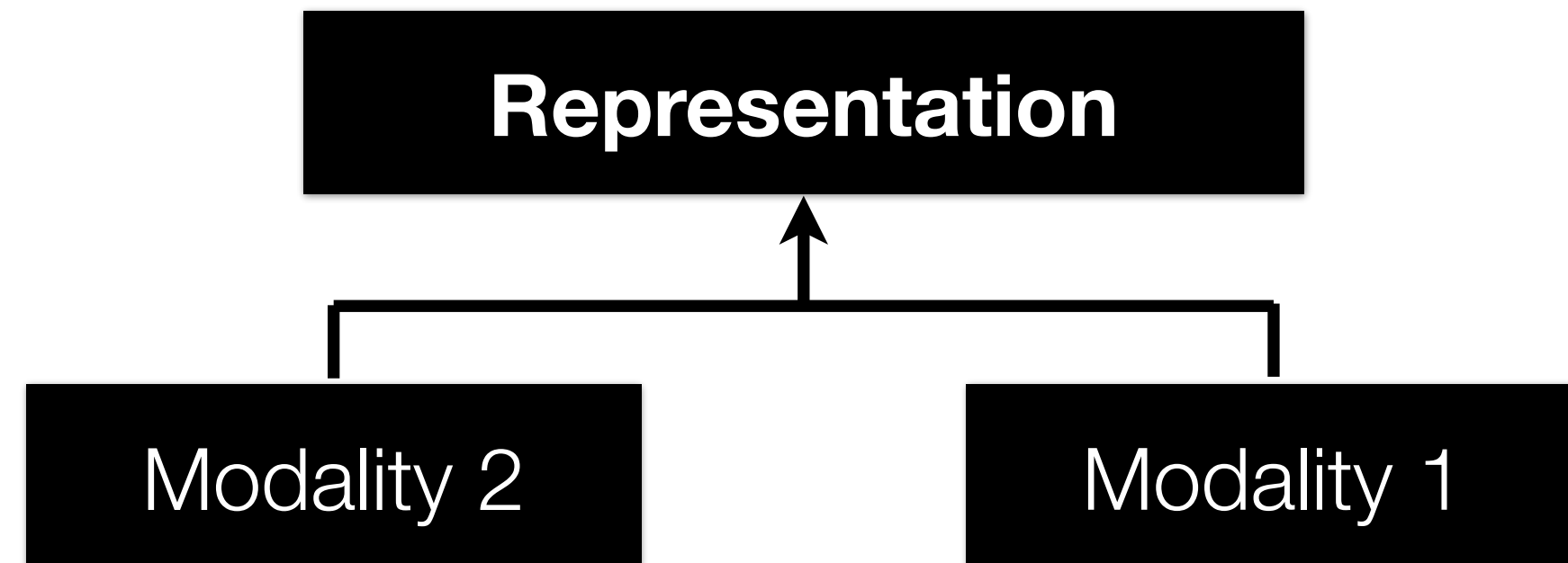
Order Embeddings

[Vendrov et al., 2016]



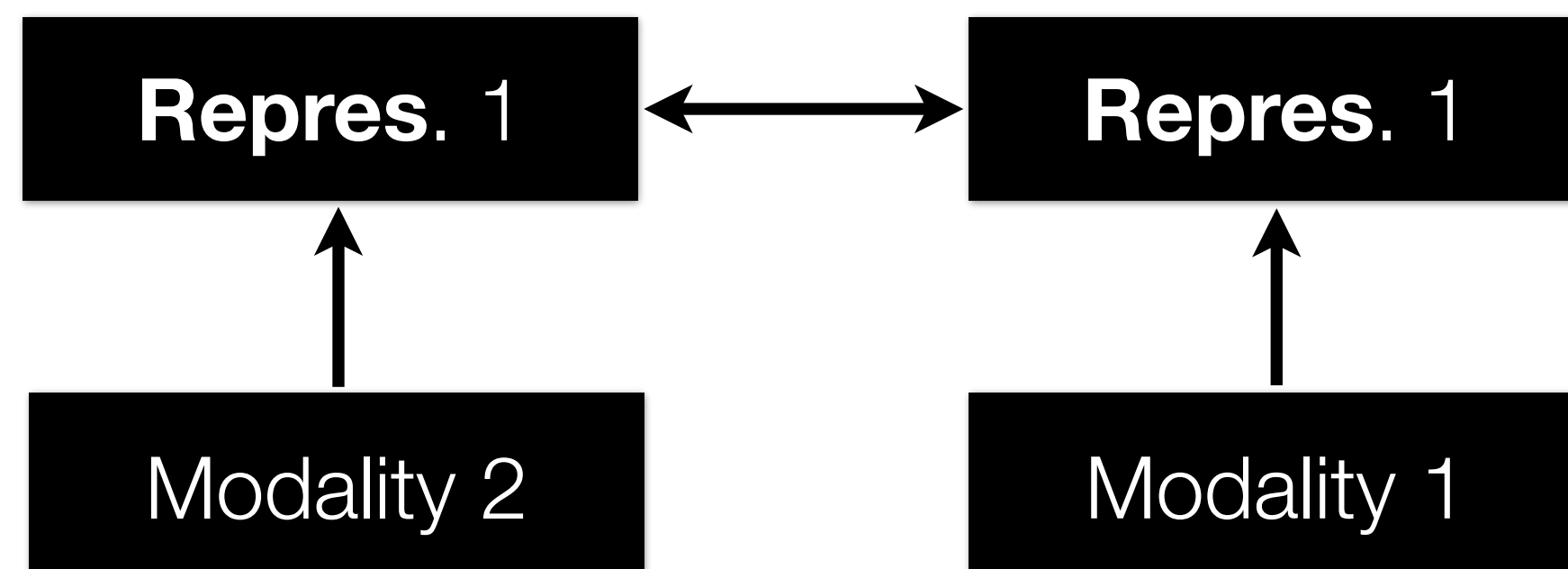
Multimodal Representation Types

Joint representations:



- Simplest version: modality concatenation (early fusion)
- Can be learned supervised or unsupervised

Coordinated representations:



- Similarity-based methods (e.g., cosine distance)
- Structure constraints (e.g., orthogonality, sparseness)
- CCA (unsupervised), joint embeddings (supervised)

Final Words ...

Joint representations

- Project modalities to the same space
- Use when all the modalities are present during test time
- Suitable for multi-model fusion

Coordinated representations

- Project modalities to their own coordinated spaces
- Use when only one of the modalities is present during test-time
- Suitable for multimodal translation
- Good for multimodal retrieval