

Explainable Semantic Space by Grounding Language to Vision with Cross-Modal Contrastive Learning

Yizhen Zhang^{1,2}, Minkyu Choi¹, Kuan Han¹, Zhongming Liu^{1,3}

¹ Department of Electrical Engineering and Computer Science, University of Michigan

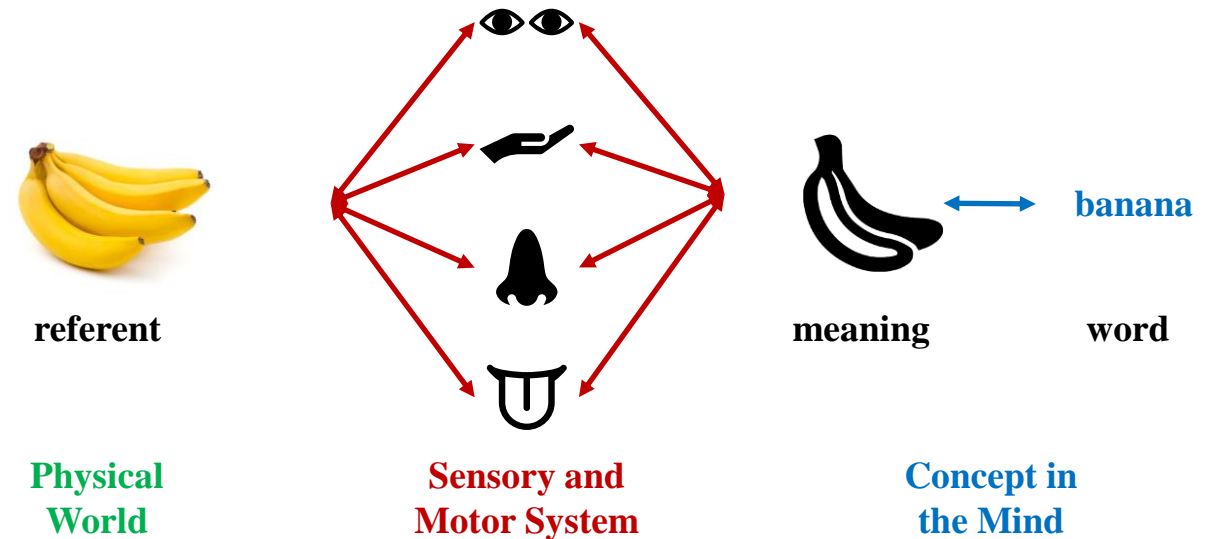
² Department of Neurological Surgery, University of California San Francisco

³ Department of Biomedical Engineering, University of Michigan

Motivation

How humans learn language?

- We learn the **meaning of a word** (e.g., banana) by associating it with the **sensory features** (e.g., shape, color, smell, taste etc.) of its **referent** (a real banana in the physical world).
- We learn **concepts** from **real world experiences**.



Motivation

Predominant language learning models learn word representations only from **textual** context instead of **multimodal** context.



The latest news from Google AI

[Bert: Pre-training of deep bidirectional transformers for language understanding](#)

J Devlin, [MW Chang](#), [K Lee](#), [K Toutanova](#) - arXiv preprint arXiv ..., 2018 - arxiv.org

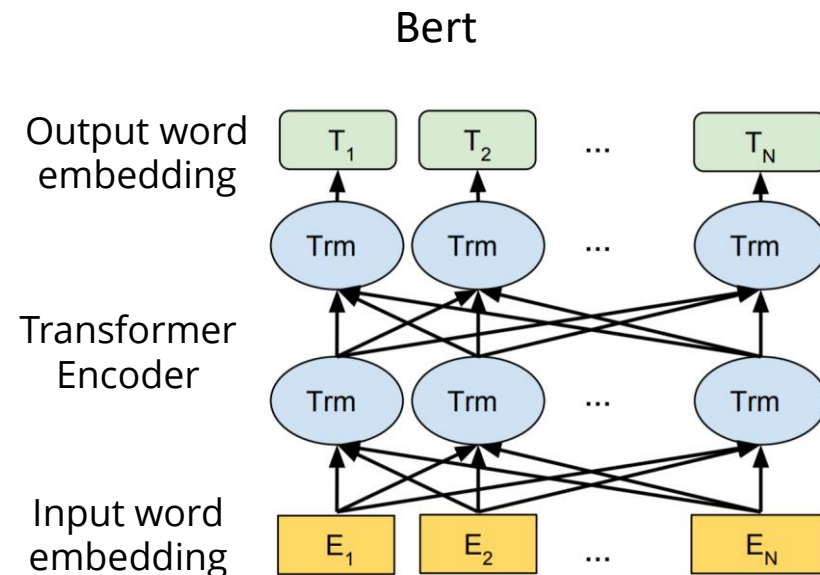
We introduce a new language representation model called **BERT**, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models, **BERT** is designed to pre-train deep bidirectional representations ...

☆ [Cited by 20810](#) [Related articles](#) [All 26 versions](#) [↔](#)

Open Sourcing BERT: State-of-the-Art Pre-training for Natural Language Processing

Friday, November 2, 2018

Posted by Jacob Devlin and Ming-Wei Chang, Research Scientists, Google AI Language

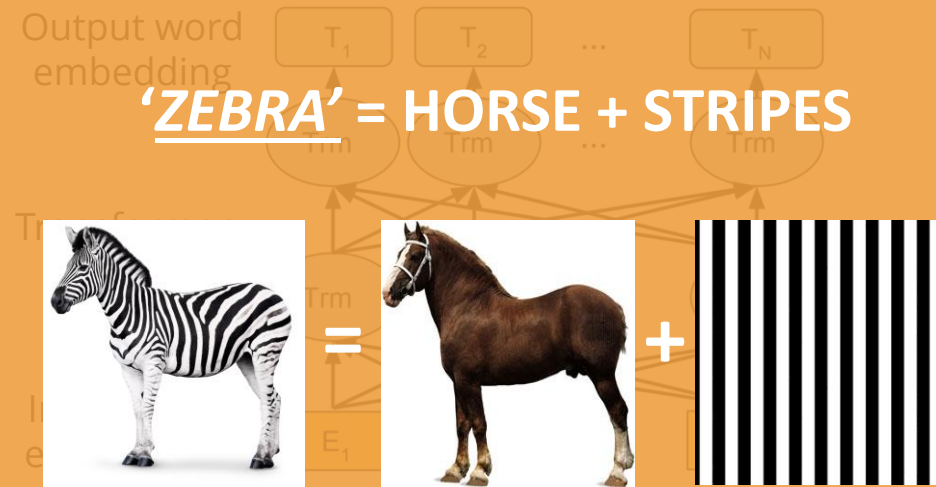


Devlin et al. 2018

Motivation

Predominant language learning models learn word representations only from **textual** context instead of **multimodal** context.

The machine is doing a task like deciphering an ancient language – since ‘word symbols’ in these language models are not grounded in real world experience.

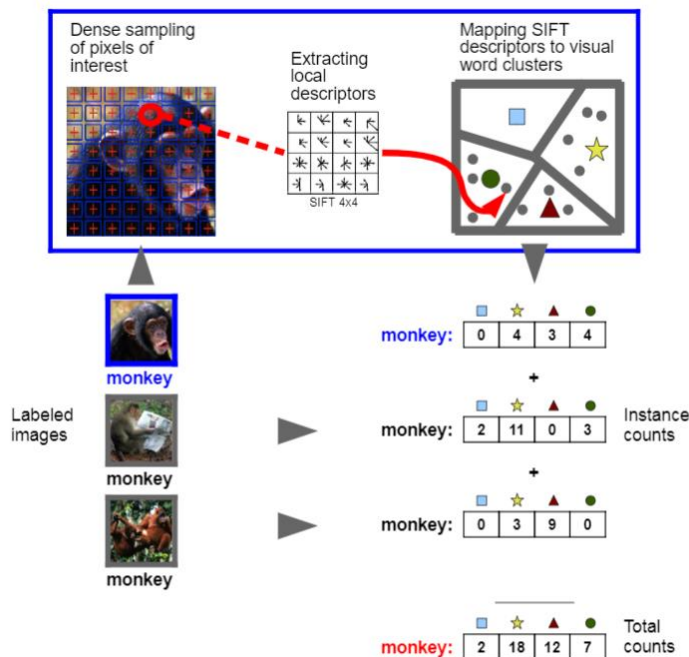


Devlin et al. 2018

Background

Early works for grounding language in vision:

bag-of-visual-word



Bruni et al. 2014

capturing word similarity

Model	all	adjs	nouns	verbs	conc-q1	conc-q2	conc-q3	conc-q4	hard
Glove	40.8	62.2	42.8	19.6	43.3	41.6	42.3	40.2	27.2
Picturebook	37.3	11.7	48.2	17.3	14.4	27.5	46.2	60.7	28.8
Glove + Picturebook	45.5	46.2	52.1	22.8	36.7	41.7	50.4	57.3	32.5
Picturebook (Visual)	31.3	11.1	38.8	<u>20.4</u>	13.9	26.1	38.7	47.7	23.9
Picturebook (Semantic)	<u>37.3</u>	<u>11.7</u>	<u>48.2</u>	17.3	<u>14.4</u>	<u>27.5</u>	<u>46.2</u>	<u>60.7</u>	<u>28.8</u>
Picturebook (1)	24.5	2.6	33.5	12.1	4.7	17.8	32.8	47.8	13.6
Picturebook (2)	28.4	6.5	38.9	9.0	5.0	21.3	34.3	55.1	15.7
Picturebook (3)	30.3	<u>11.9</u>	41.9	3.1	2.6	24.3	37.5	58.3	18.4
Picturebook (5)	34.4	6.8	44.5	<u>18.0</u>	9.0	<u>27.9</u>	42.8	58.3	25.9
Picturebook (10)	<u>37.3</u>	11.7	<u>48.2</u>	17.3	<u>14.4</u>	27.5	<u>46.2</u>	<u>60.7</u>	<u>28.8</u>

Table 3: SimLex-999 results (Spearman's ρ). Best results overall are bolded. Best results per section are underlined. Bracketed numbers signify the number of images used. Some rows are copied across sections for ease of reading.

Kiros and Chan et al. 2018

capturing word analogy

Analogy	Answer Candidates	GloVe	ViCo
car:land::aeroplane:?	ocean, sky, road, railway	ocean	sky
clock:circle::tv:?	triangle, square, octagon, round	triangle	square
park:bench::church:?	door, sofa, cabinet, pew	door	pew
sheep:fur::person:?	hair, horn, coat, tail	coat	hair
monkey:zoo::cat:?	park, house, church, forest	park	house
leg:trouser::wrist:?	watch, shoe, tie, bandana	bandana	watch
yellow:banana::red:?	strawberry, lemon, mango, orange	mango	strawberry
rice:white::spinach:?	blue, green, red, yellow	blue	green
train:railway::car:?	land, desert, ocean, sky	land	land
can:metallic::bottle:?	wood, glass, cloth, paper	glass	glass
man:king::woman:?	queen, girl, female, adult	queen	girl
can:metallic::bottle:?	wood, plastic, cloth, paper	plastic	wood
train:railway::car:?	road, desert, ocean, sky	road	ocean

Table 6. Answering Analogy Questions. Out of 30 analogy pairings tested, we found both GloVe and ViCo to be correct 19 times, only ViCo was correct 8 times, and only GloVe was correct 3 times. Correct answers are **highlighted**.

Gupta et al. 2019

better word clusters

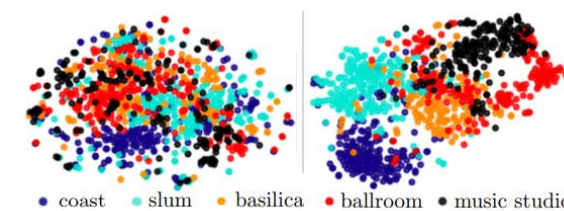
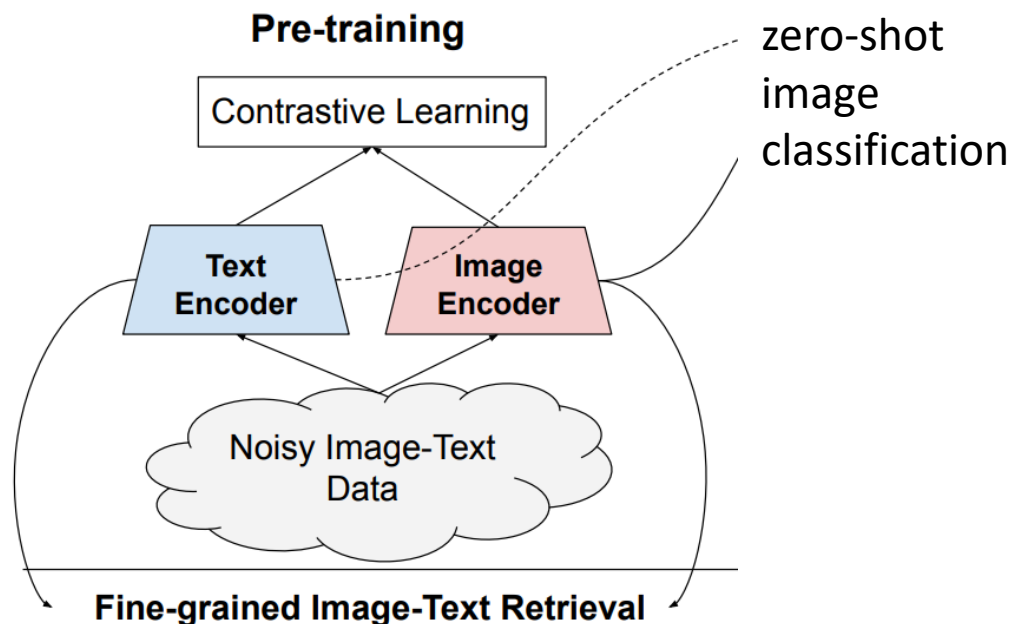


Figure 3: t-SNE visualization on CMPlaces sentences for a set of randomly sampled visual scenes. Left: textual model T . Right: grounded model $C_g + P_g$.

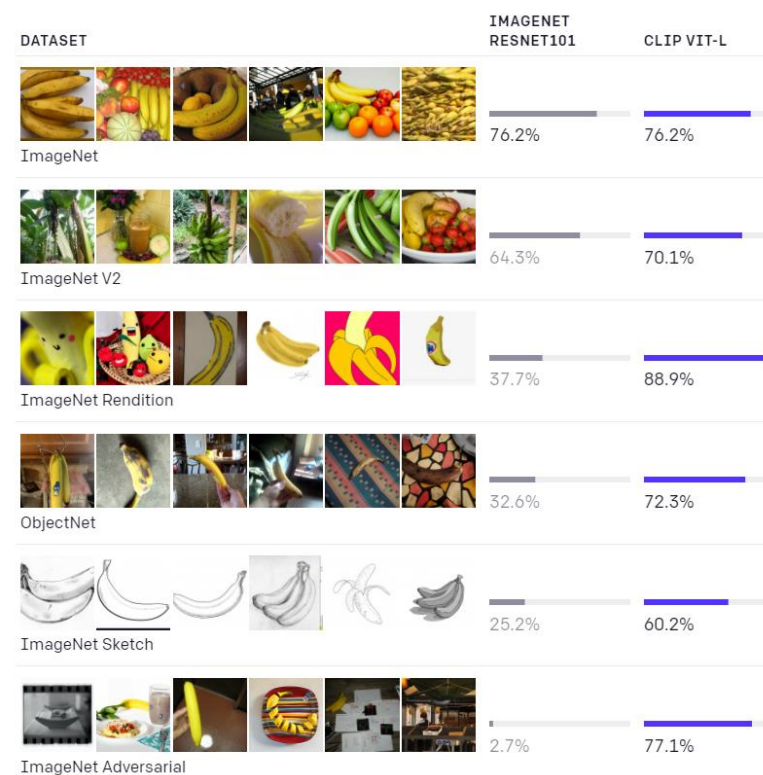
Bordes and Zablocki et al. 2020

Background

Visual-language cross-modal learning with contrastive loss

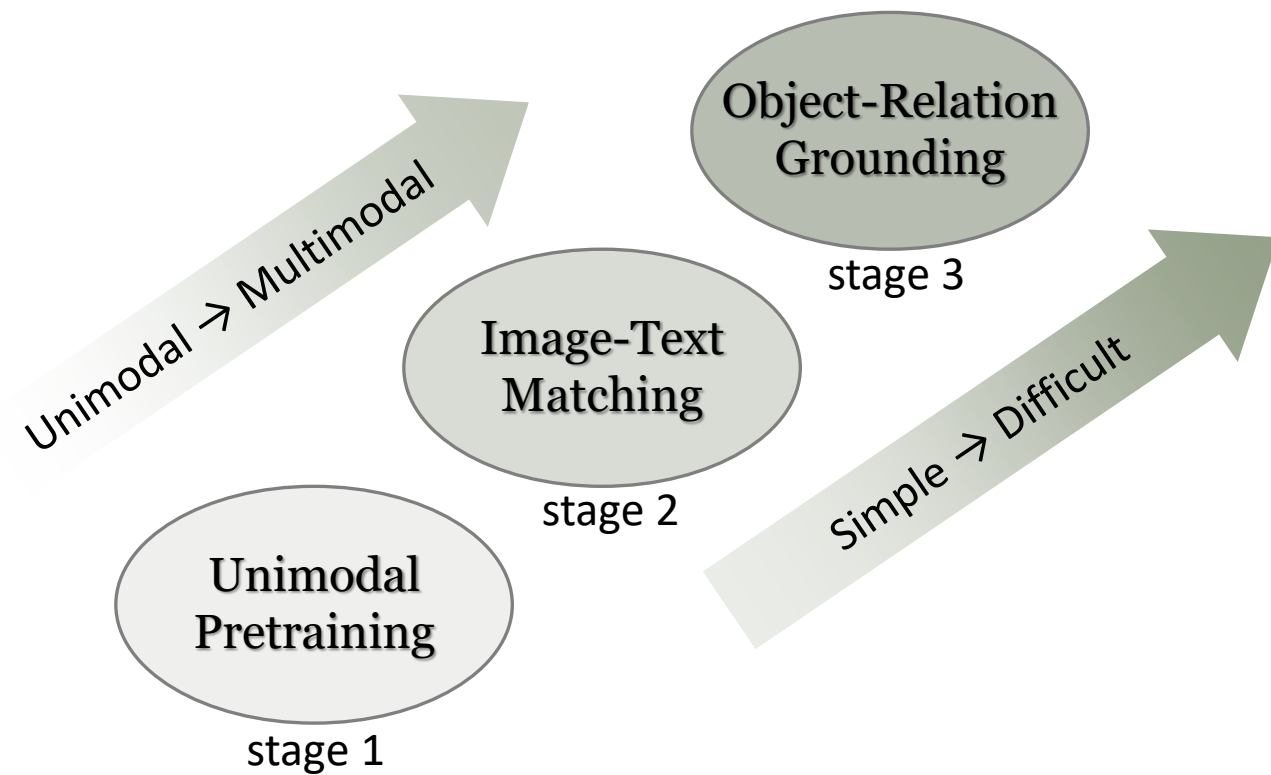
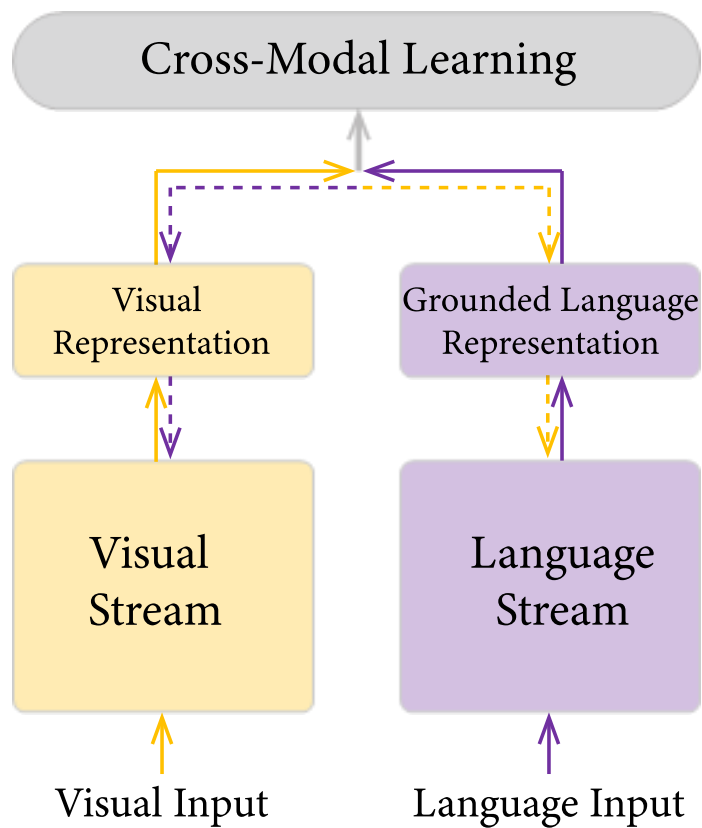


ALIGN: Jia et al. 2021



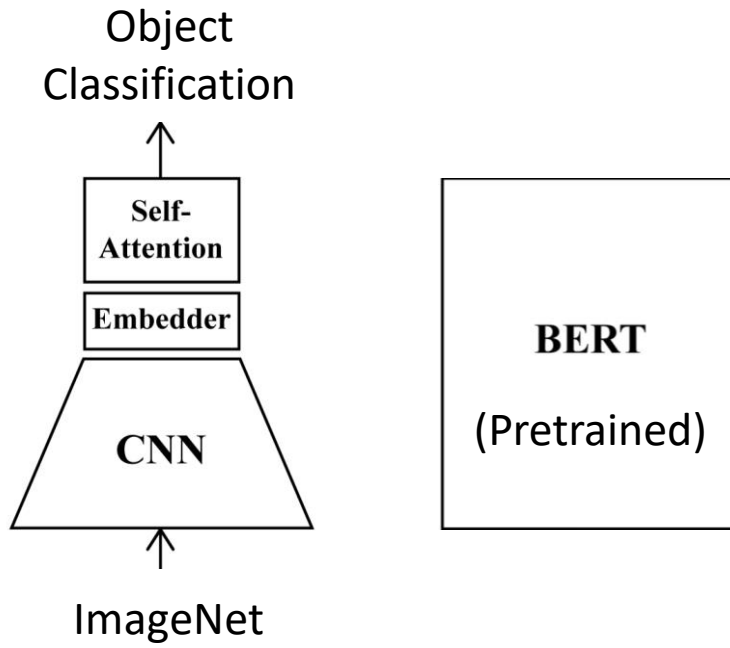
CLIP: Radford and Kim et al. 2021

Method



Model Design

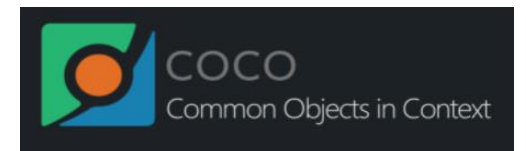
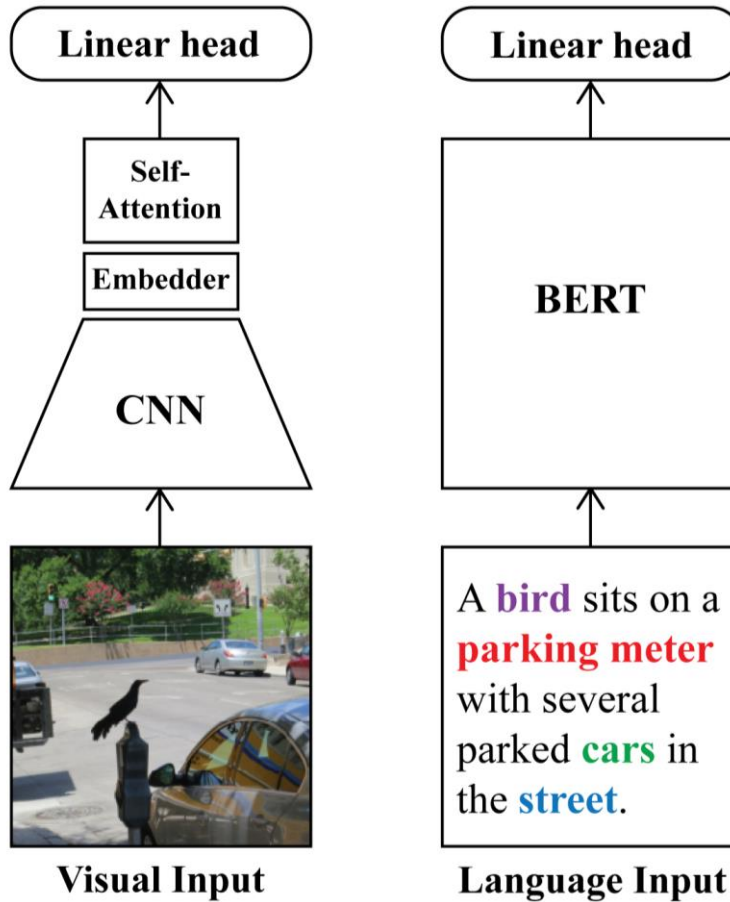
1. Unimodal Pretraining



Model Design

1. Unimodal Pretraining

2. Visual Grounding of Natural Language



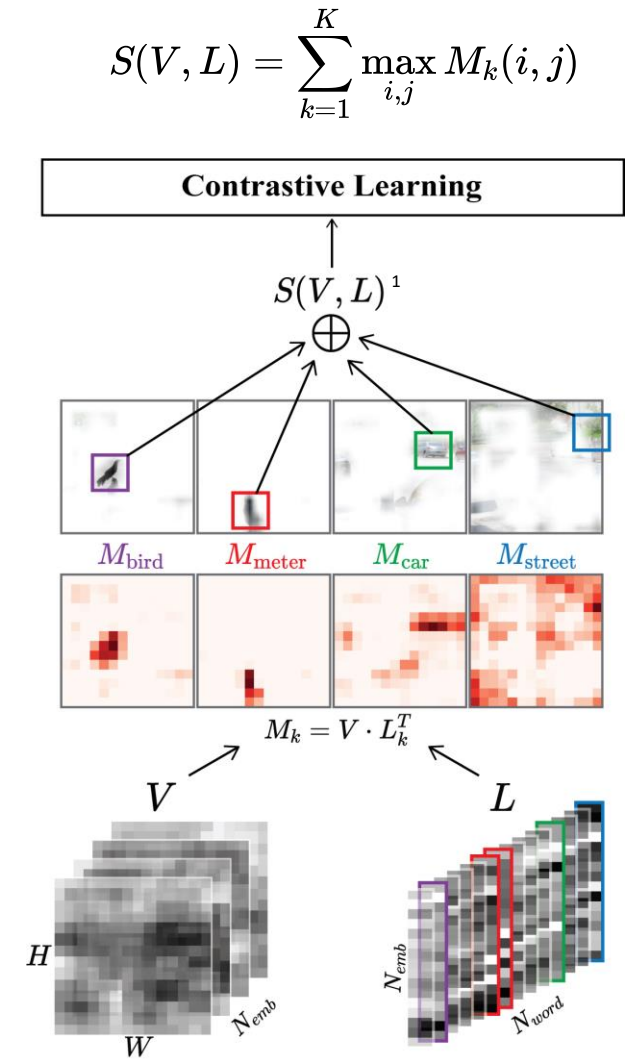
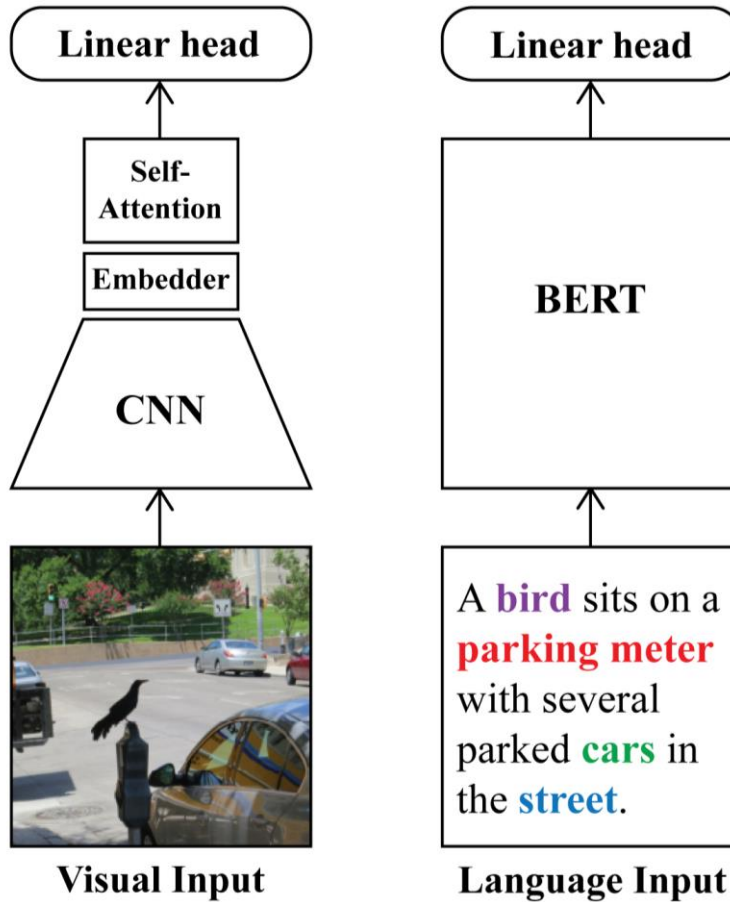
<https://cocodataset.org/>

Model Design

1. Unimodal Pretraining

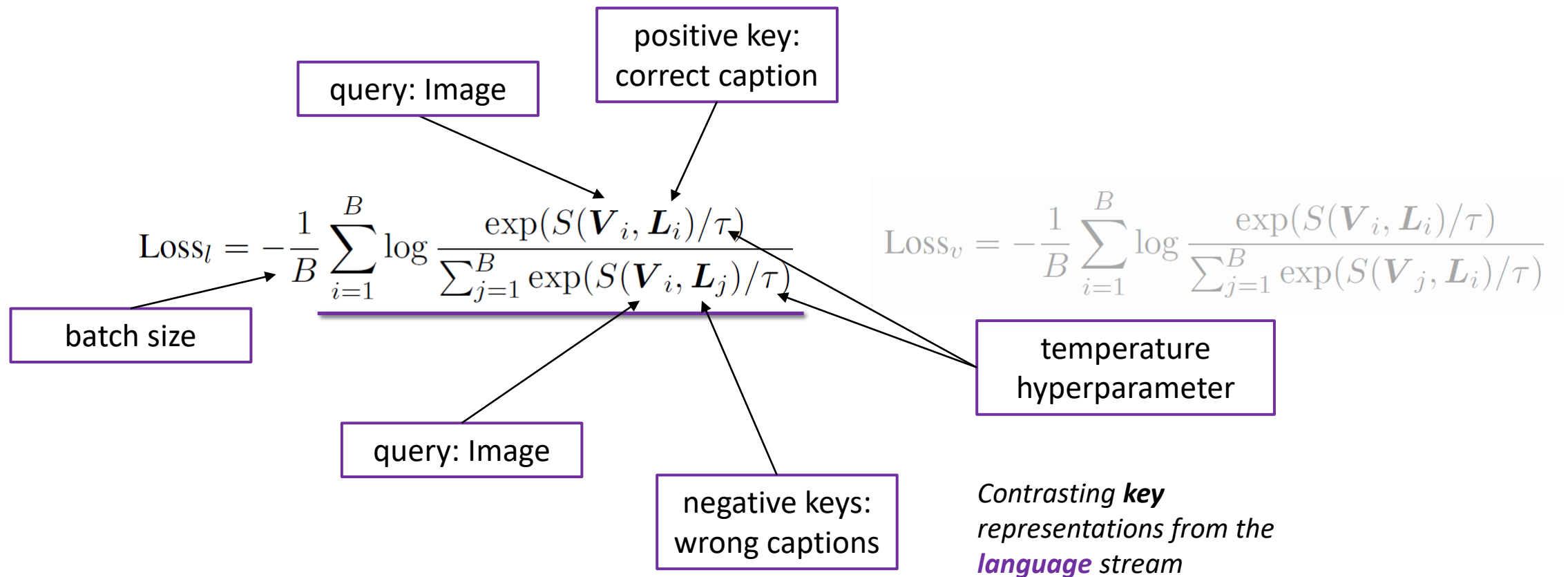
2. Visual Grounding of Natural Language

1. Harwath, D., Recasens, A., Surís, D., Chuang, G., Torralba, A., & Glass, J. (2018). Jointly discovering visual objects and spoken words from raw sensory input. In *Proceedings of the European conference on computer vision (ECCV)* (pp. 649-665).



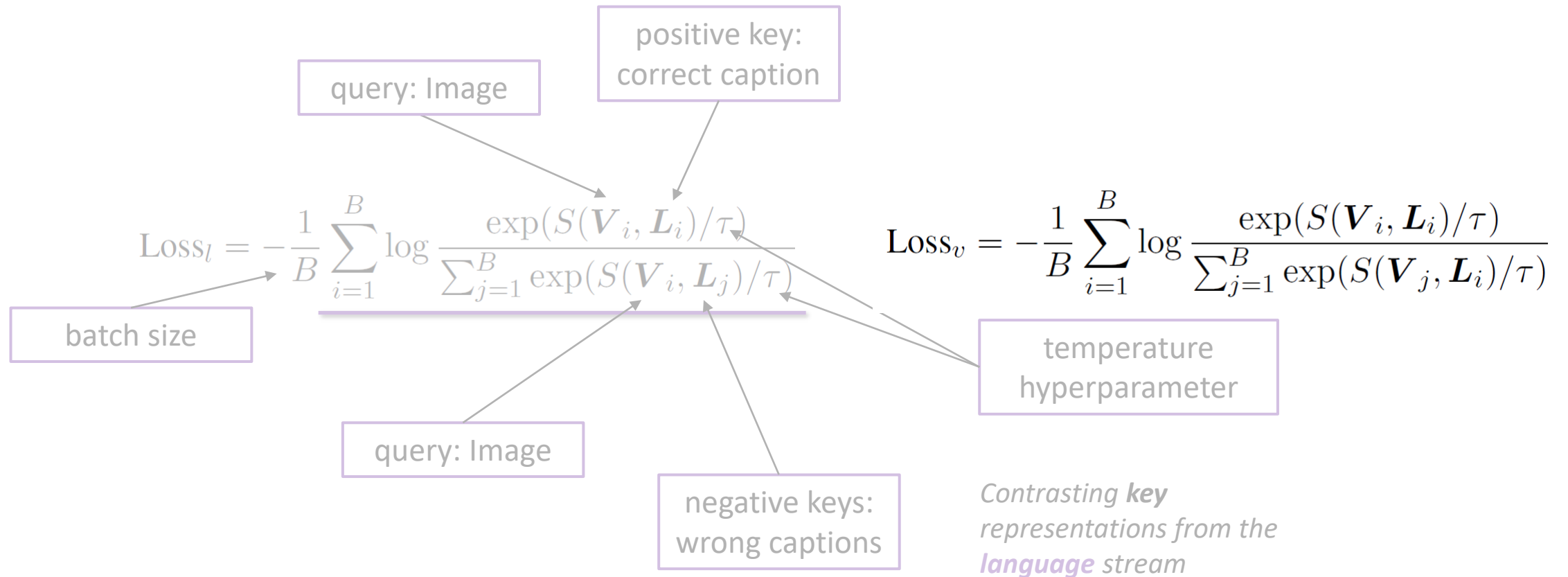
Training

NT-Xent¹ loss: Normalized Temperature-scaled Cross Entropy Loss



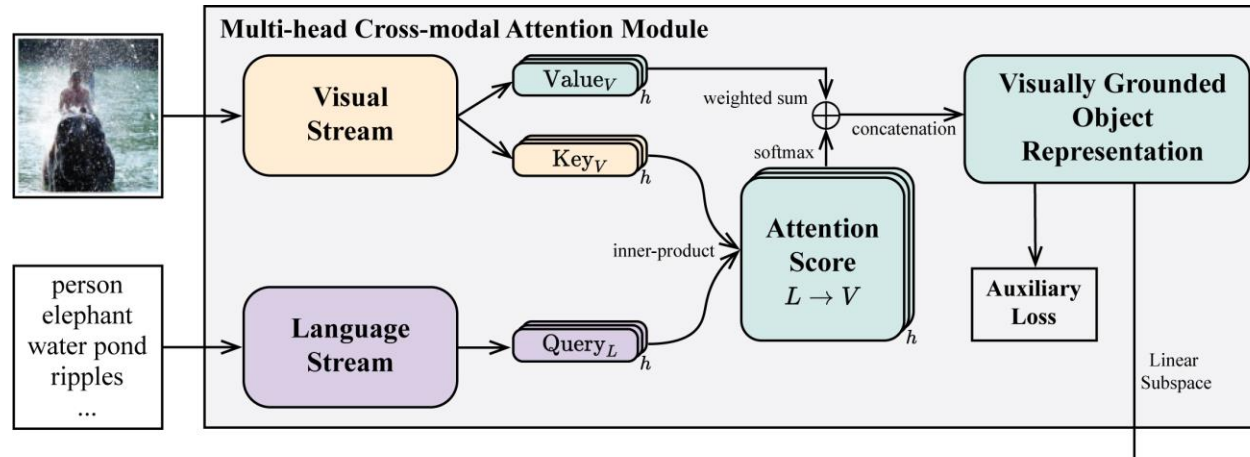
Training

NT-Xent¹ loss: Normalized Temperature-scaled Cross Entropy Loss



Model Design

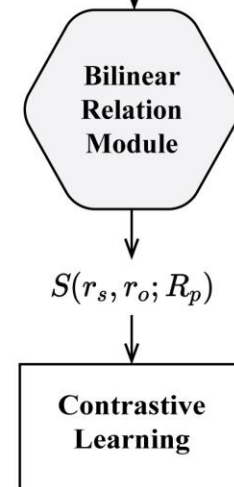
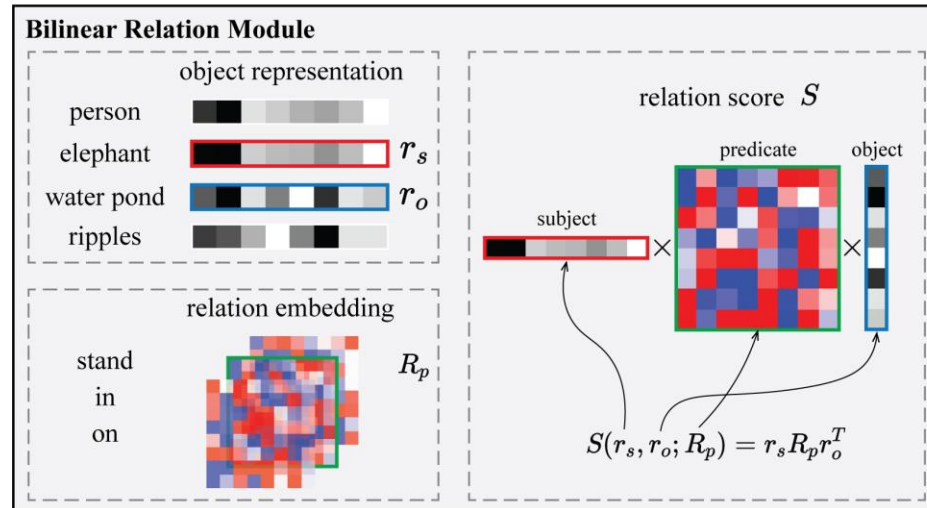
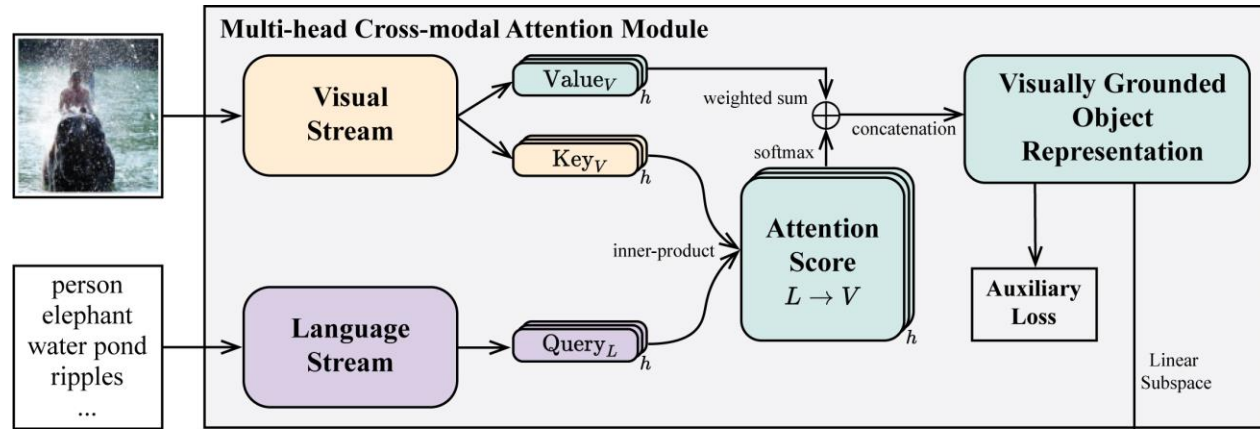
1. Unimodal Pretraining
2. Visual Grounding of Natural Language
3. Visual Grounding of object relations



$$\text{auxiliary loss}(\mathbf{x}, l) = -\log \frac{\exp(\mathbf{x}[l])}{\sum_{i=1}^N \exp(\mathbf{x}[i])}$$

Model Design

1. Unimodal Pretraining
2. Visual Grounding of Natural Language
3. Visual Grounding of object relations

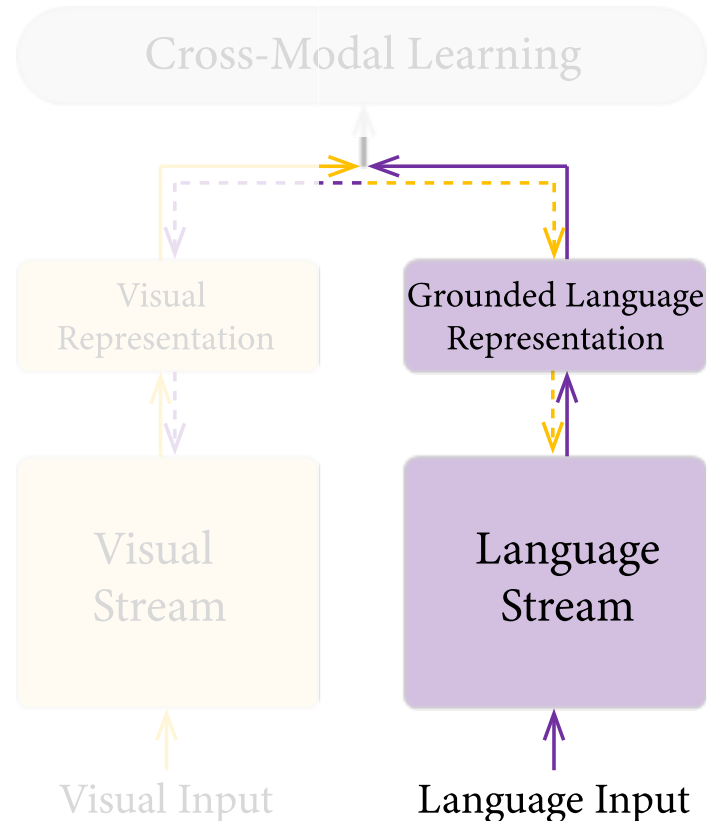


$$\text{Loss}_{\text{rel}} = -\frac{1}{|\mathcal{B}|} \sum_{(r_s, r_o; R_p) \in \mathcal{B}} \log \frac{\exp(S(r_s, r_o; R_p)/\tau)}{\sum_{k \in \mathcal{K}_{\text{rel}}} \exp(S(r_s, r_o; R_p^k)/\tau)}$$

$$\text{Loss}_{\text{obj}} = -\frac{1}{|\mathcal{B}|} \sum_{(r_s, r_o; R_p) \in \mathcal{B}} \log \frac{\exp(S(r_s, r_o; R_p)/\tau)}{\sum_{k \in \mathcal{K}_{\text{obj}}} \exp(S(r_s^k, r_o^k; R_p)/\tau)}$$

Interpreting the grounded semantic space

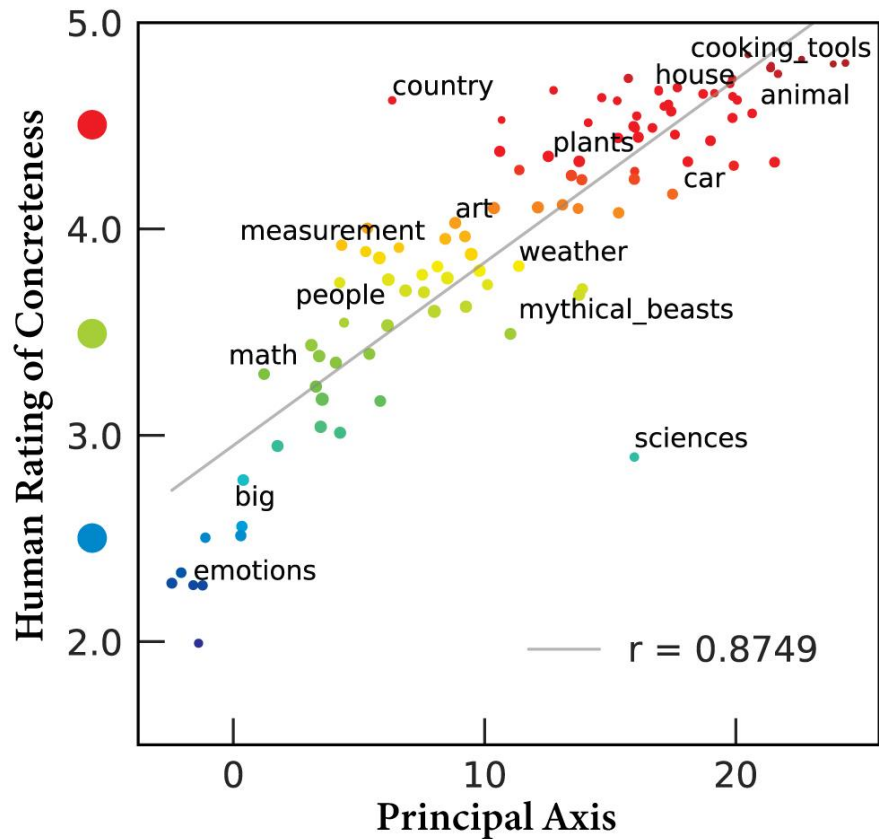
How does visual grounding reshape the semantic space in the language stream?



After visual grounding, we

- detach the language stream
- extract the grounded word embeddings
- apply intrinsic evaluations
 - Principal Component Analysis
 - Concreteness gradient
 - Clustering of Word Categories
 - Concept composition
 - Cross-modal image search
- interpret the semantic space by human intuition and neurobiological knowledge

Results - PCA



Examples:

word	category	principal axis	human rating
oven	cooking tool	30.86	4.97
zebra	animal	24.74	4.86
car	car	22.06	4.89
furniture	house	20.83	4.89
defense	country	4.65	4.19
chemistry	sciences	18.67	3.64
wood	plants	16.22	4.85
cartoon	art	11.84	4.33
humid	weather	11.58	3.48
angel	mythical beasts	7.22	3.82
lover	people	6.01	3.68
thousand	math	3.65	3.07
huge	big	0.17	3.54
cheerful	emotions	- 2.54	2.34

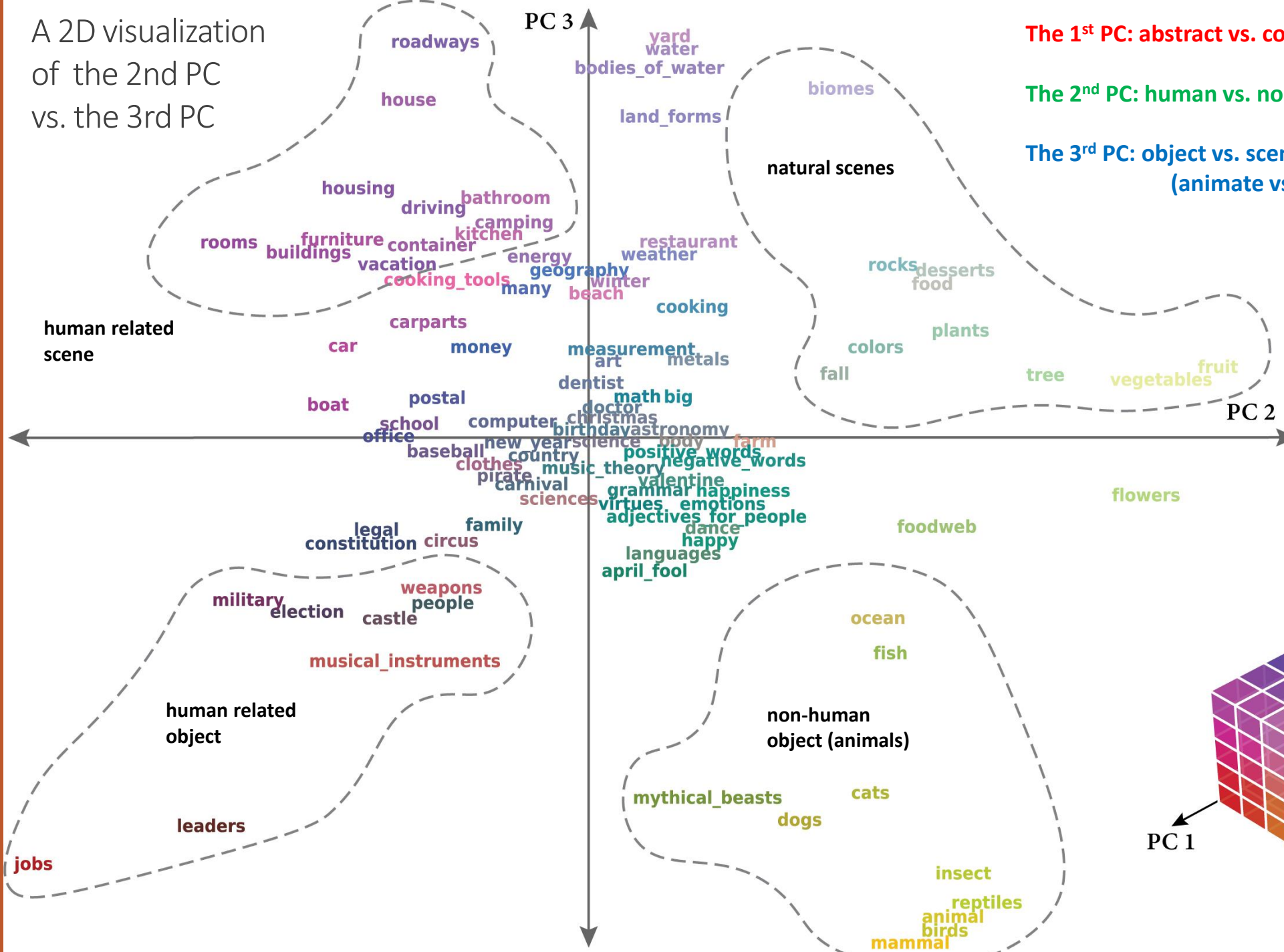
Human-rated Word Concreteness:

Brysbaert, M., Warriner, A. B., & Kuperman, V. (2014). Concreteness ratings for 40 thousand generally known English word lemmas. *Behavior research methods*, 46(3), 904-911.

Results - PCA

Group	Correlation (Pearson's r)		
	Bert	Grounded	Relational Grounded
word-level	0.1040	0.6615	0.6948
category-level	0.3538	0.8749	0.8001

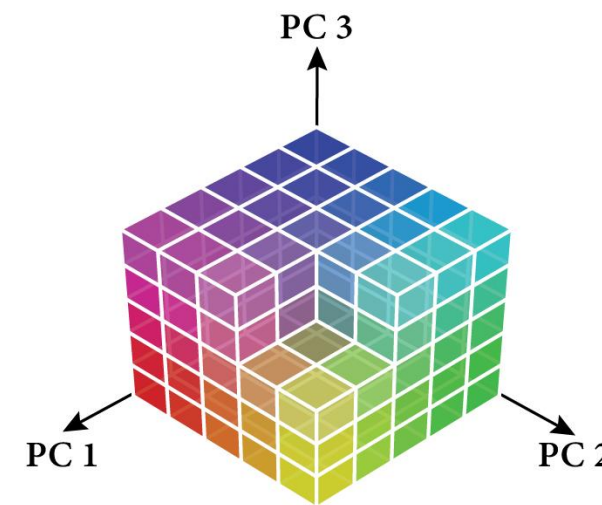
A 2D visualization of the 2nd PC vs. the 3rd PC



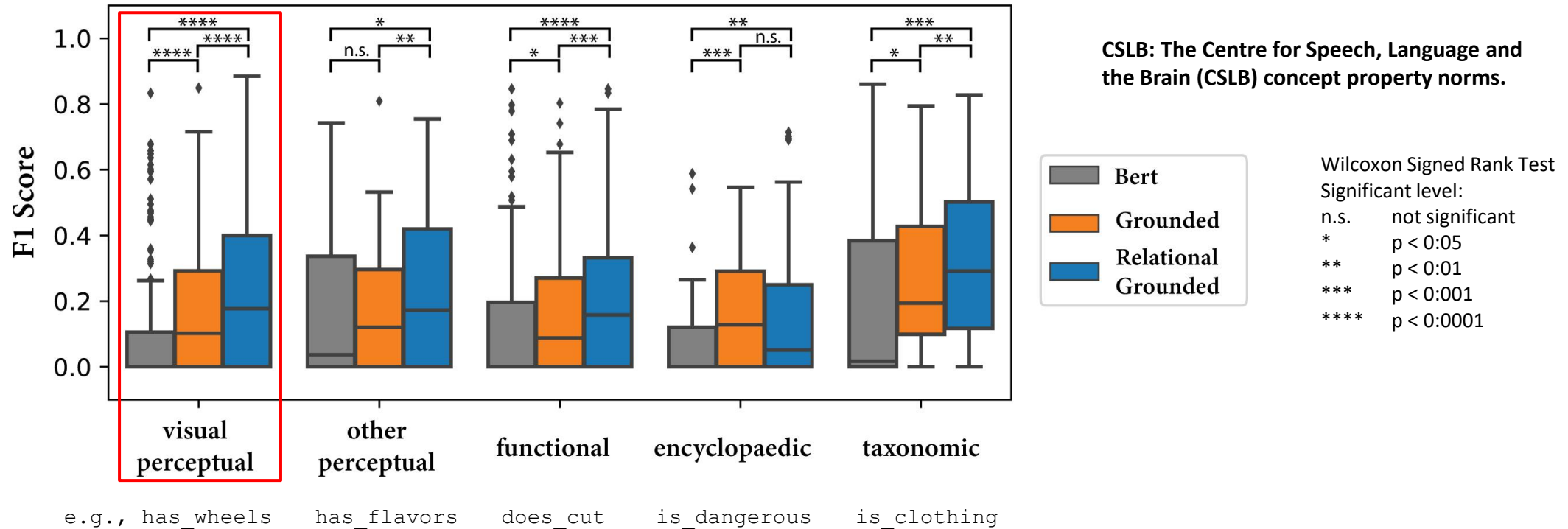
The 1st PC: abstract vs. concrete

The 2nd PC: human vs. non-human

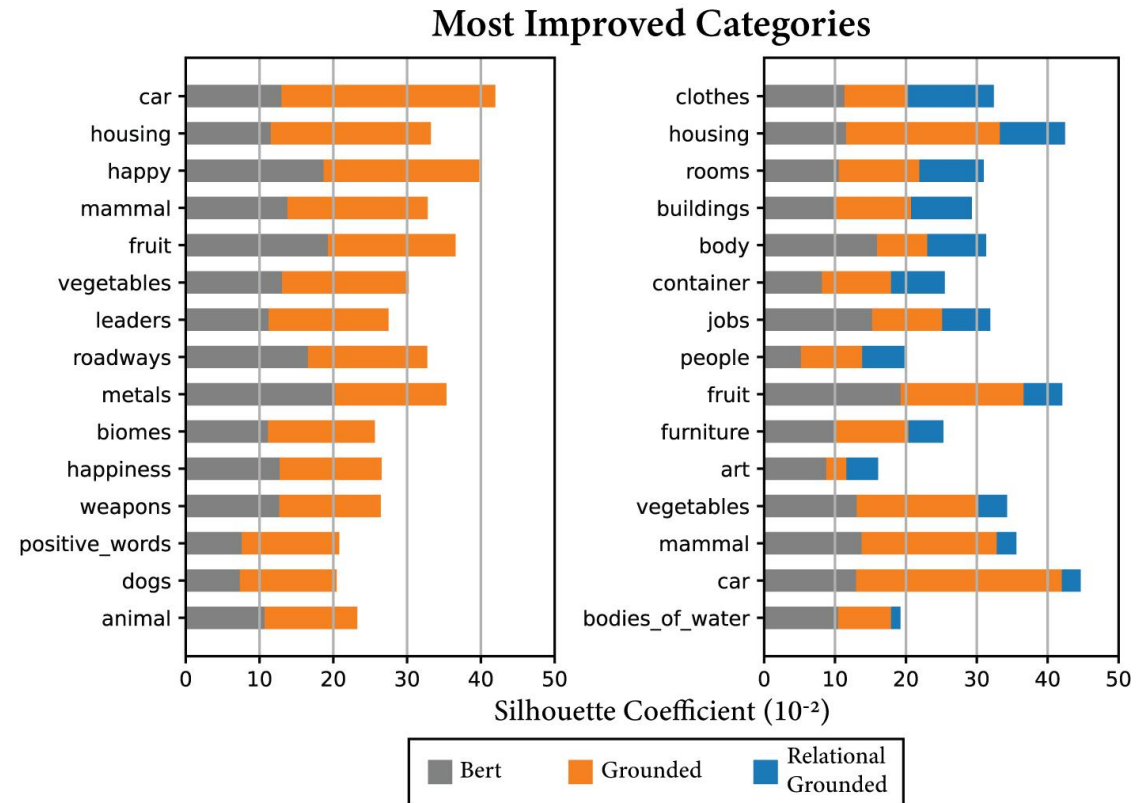
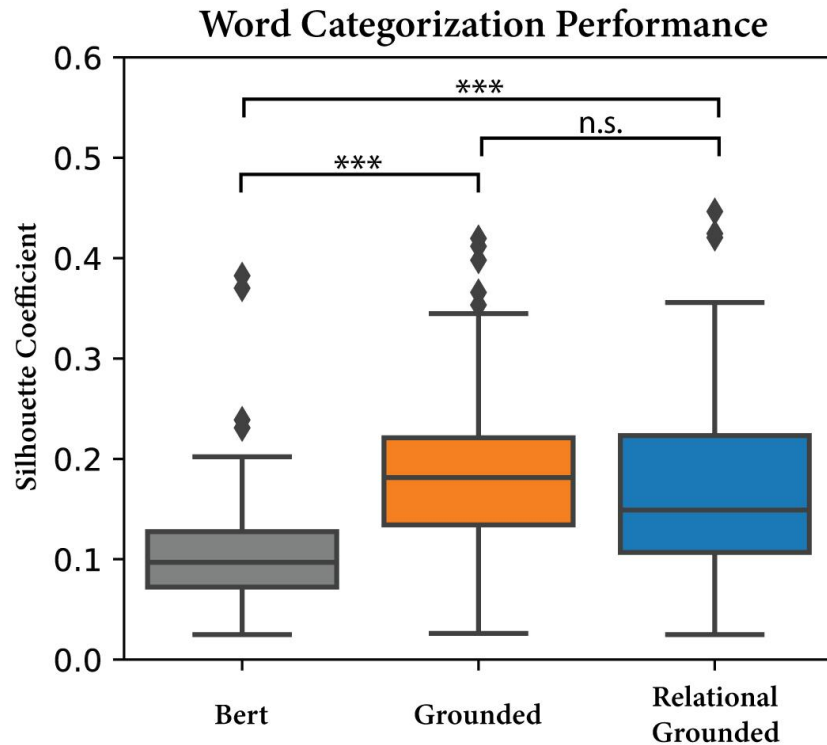
The 3rd PC: object vs. scene
(animate vs. inanimate)



Results – Semantic Norm Prediction



Results – word clustering



Results - Vision based compositional reasoning

Where is the phrase “striped horse” represented in the Semantic Space?

Is it a ZEBRA ?



What is a “Striped horse” in the grounded semantic space?

Representation Similarity of Visually Inferred Concepts

Striped Horse : Zebra



tomcat



seahorse



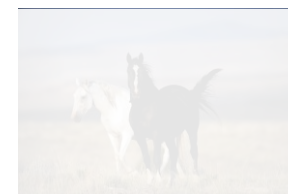
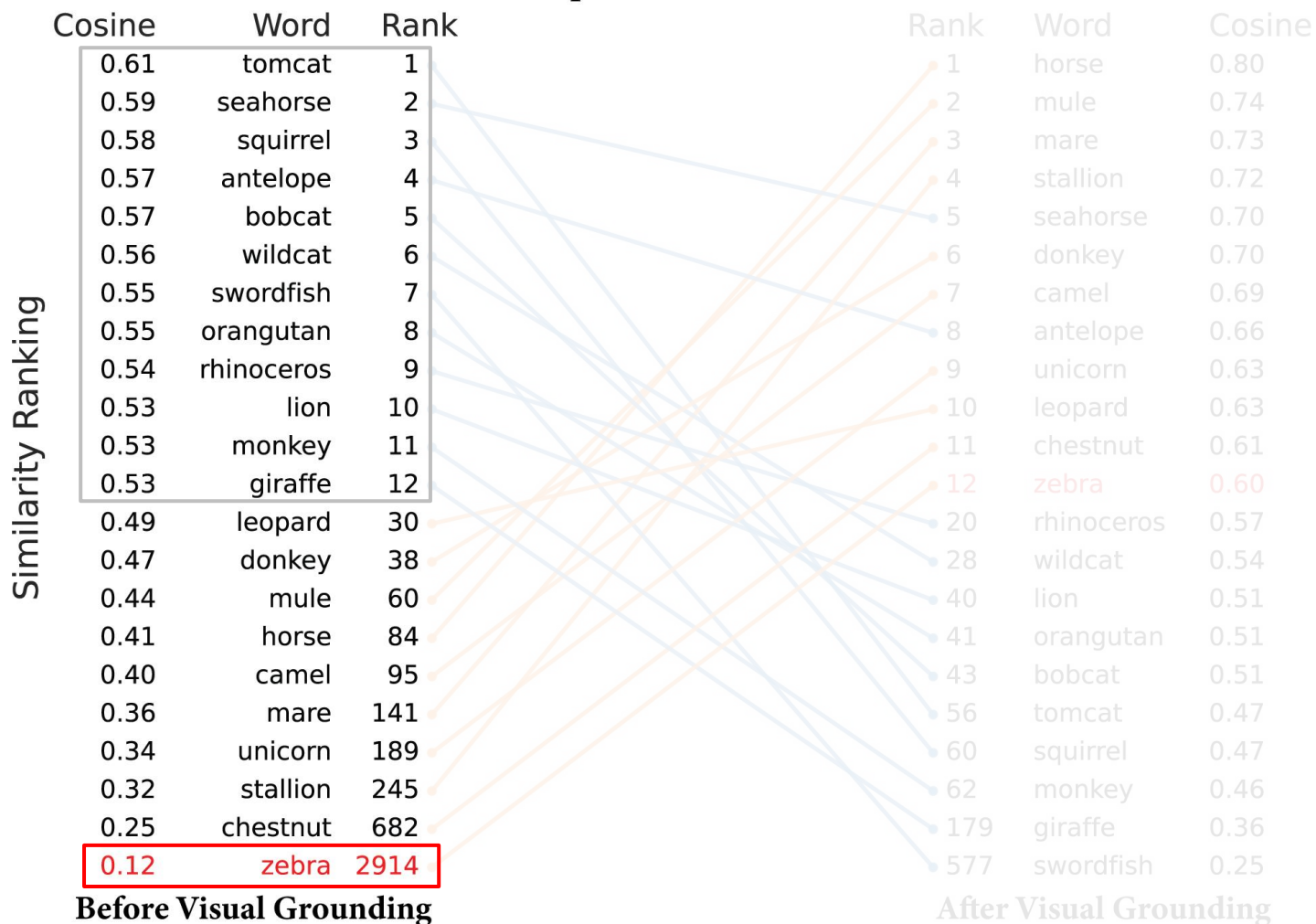
squirrel



antelope



swordfish



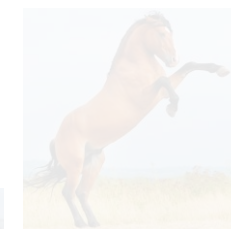
horse



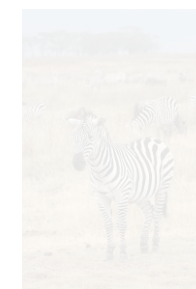
mule



mare



stallion



zebra

What is a “Striped horse” in the grounded semantic space?

Representation Similarity of Visually Inferred Concepts

Striped Horse : Zebra



tomcat



seahorse



squirrel

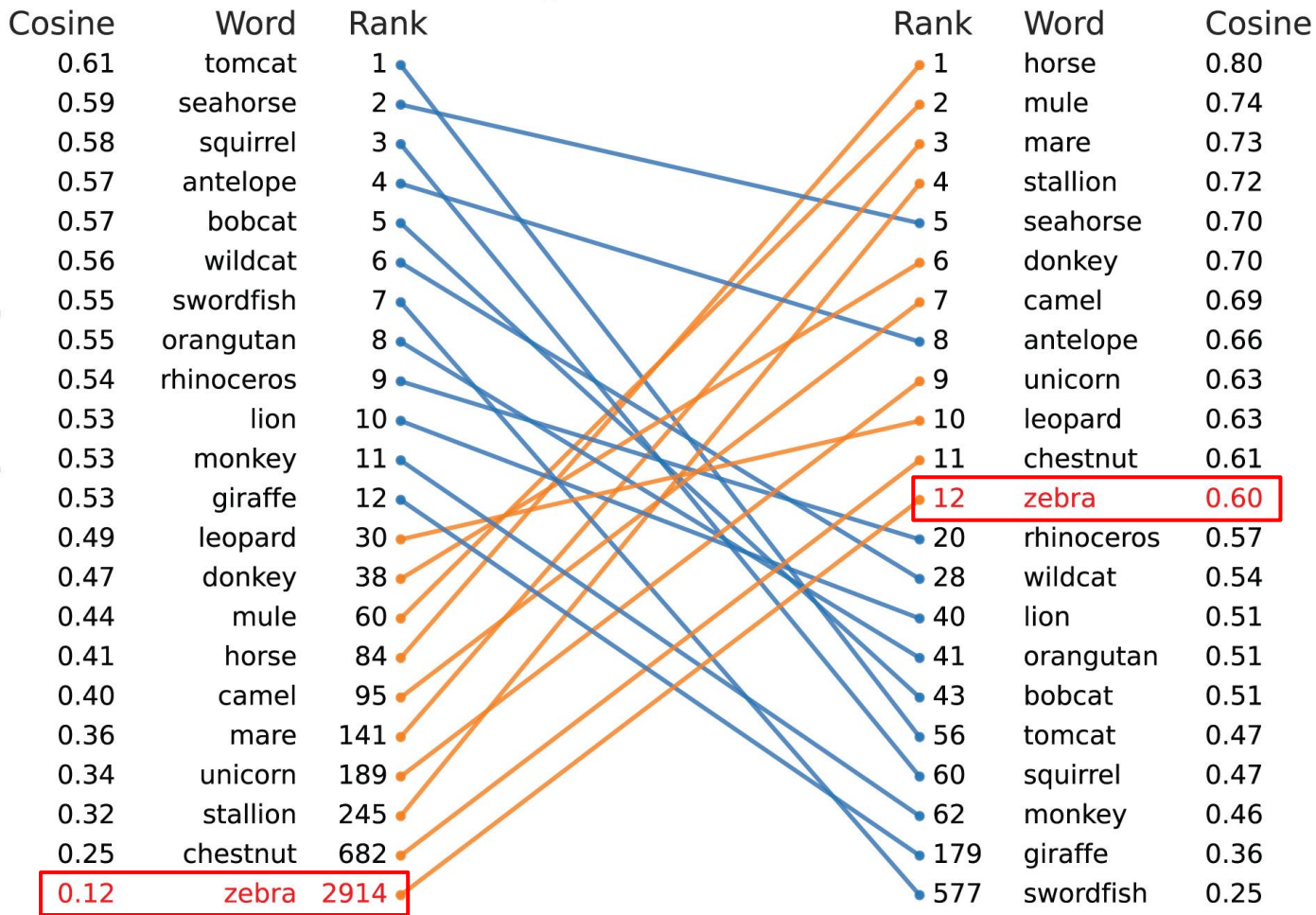


antelope



swordfish

Similarity Ranking



Before Visual Grounding

After Visual Grounding



horse



mare

mule



stallion



zebra

What is a “Red Fruit” in the grounded semantic space?



blueberry



watermelon



pomegranate

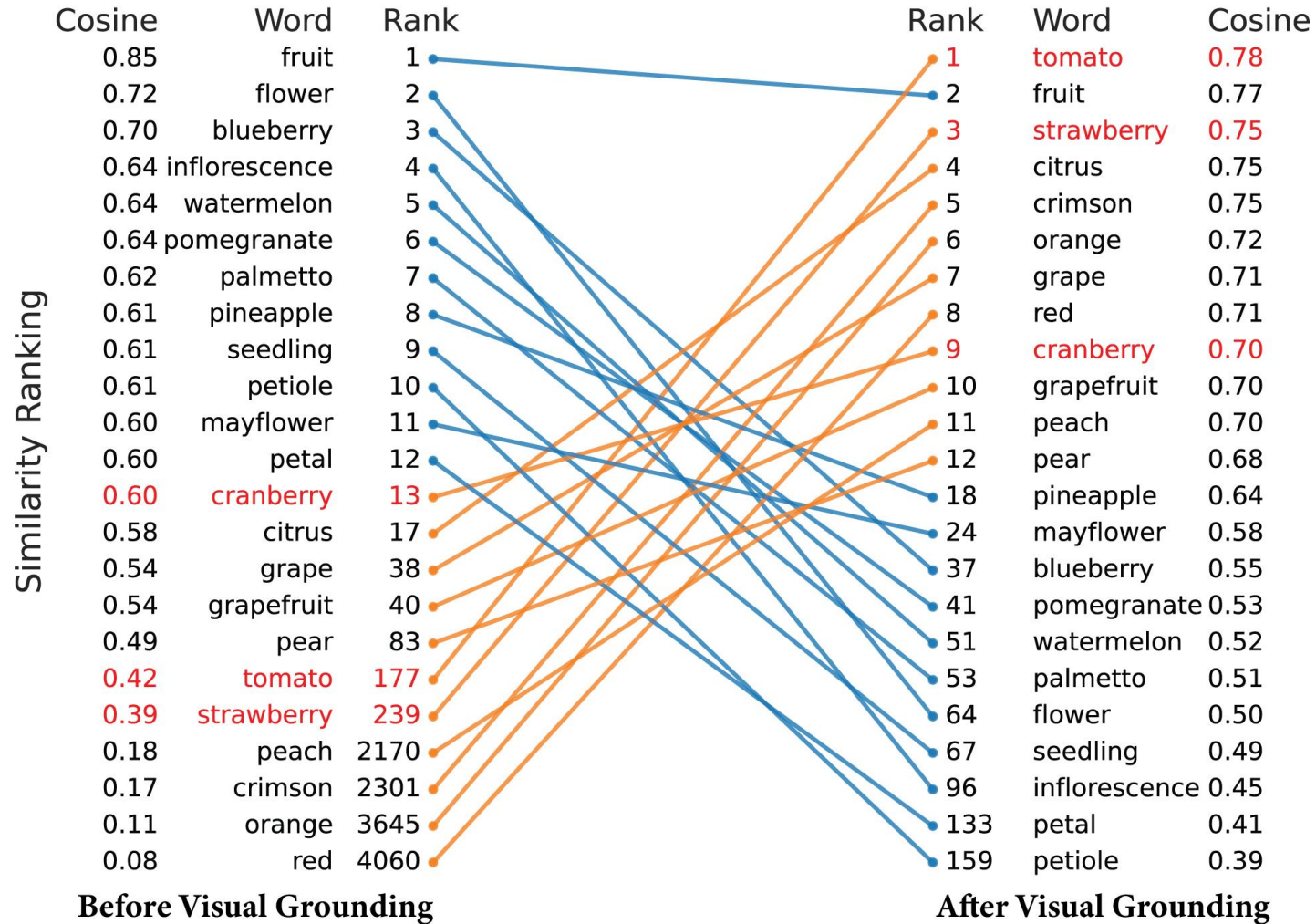


palmetto



pineapple

Representation Similarity of Visually Inferred Concepts Red Fruit



tomato



strawberry



citrus



grape



cranberry

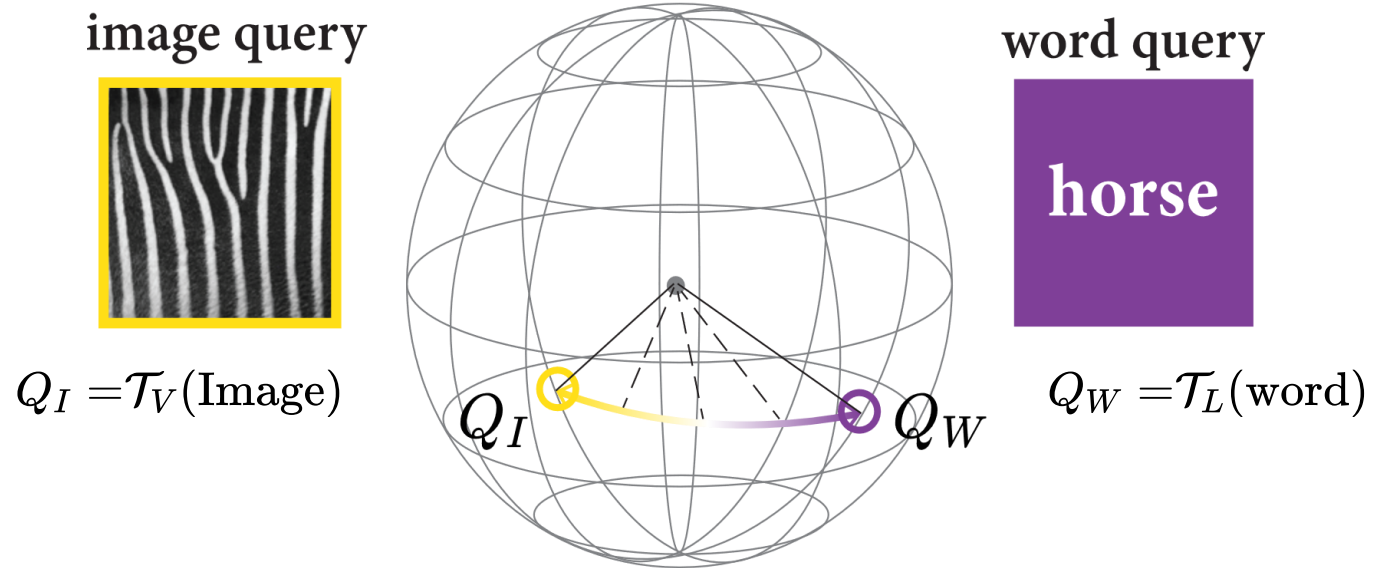


peach

Results - Vision based compositional reasoning

Query Phrase	Target Word	Similarity (cosine rank)					
			Bert		Grounded		Relational
striped horse	zebra	0.12	2914	0.60	12	0.63	8
black and white bear	panda	0.13	2478	0.69	2	0.81	2
flying car	plane	0.36	167	0.66	4	0.61	11
round container	bowl	0.25	489	0.56	8	0.67	2
red fruit	strawberry	0.39	239	0.75	3	0.85	3
young dog	puppy	0.40	94	0.92	2	0.93	2
iced mountain	glacier	0.44	20	0.86	1	0.73	5
clear sky	sunny	0.27	631	0.31	184	0.34	61
hot weather	summer	0.27	903	0.52	14	0.53	6

Results - Continuous Image Search from Text & Image



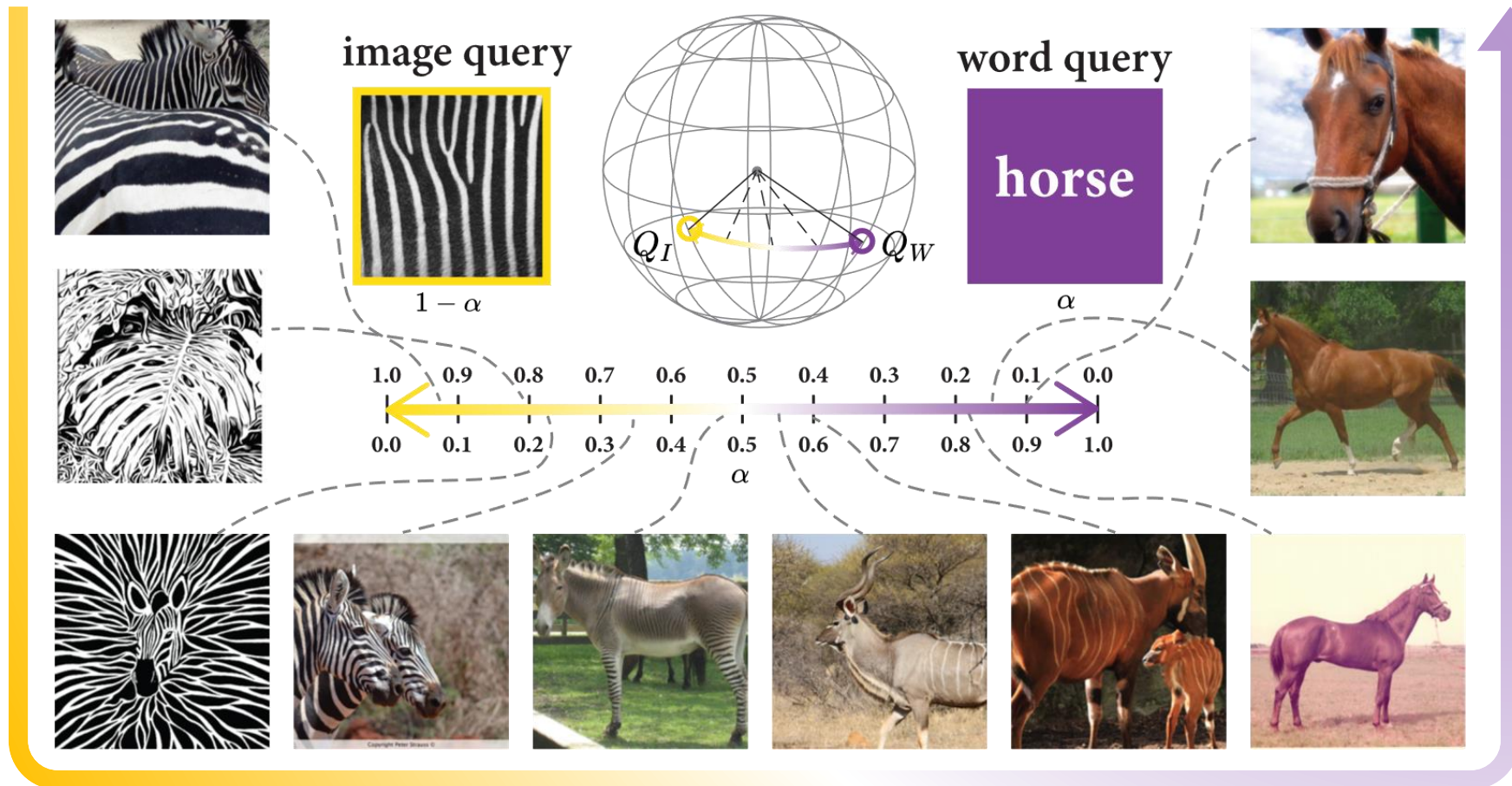
$\mathcal{T}_V, \mathcal{T}_L :$

The nonlinear transformation from the image / word space to the L2-normed multimodal representational space.


















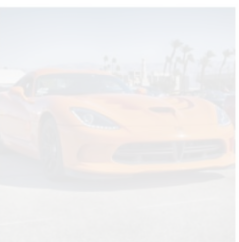



Search query: $Q_{\text{search}} = (1 - \alpha) \cdot Q_I + \alpha \cdot Q_W$

Searching: $\operatorname{argmax}_{K \in \text{Image Pool}} \cos(Q_{\text{search}}, \mathcal{T}(K))$

Results - Continuous Image Search from Text & Image



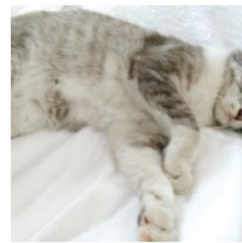
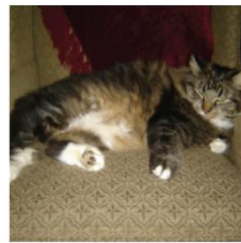
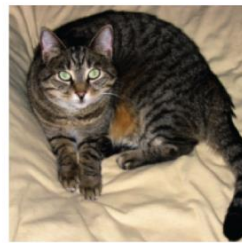
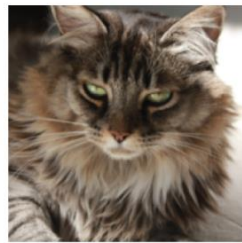
Results - Continuous Image Search from Text & Image

Image Query							Word Query
							sleep
							milk
							lego

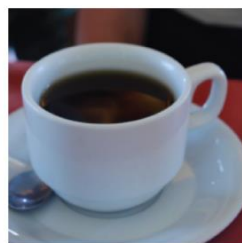
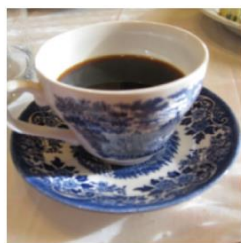
Results - Continuous Image Search from Text & Image

Image Query

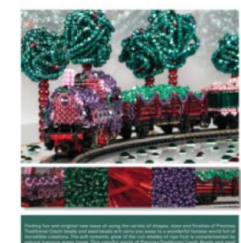
Word Query



sleep



milk



lego

Summary

We design a two-stream model for grounding language learning in vision:

- Progressive training
- cross-modal contrastive learning

After training, we analyze the language model as a stand-alone system. In this grounded word embedding space:

- The first principal axis = concrete vs. abstract gradient
- Principal axes are explainable by human intuition.
- Word representation captures human-defined semantic feature norms.
- Concepts are better clustered.

Besides, “zebra = striped horse” in both word embedding space and joint representational space.

Thank you!

QUESTIONS & COMMENTS?