

Scene Context, Object Reference, and Image Memorability: Insights into Visuo-Linguistic Processing in Humans and Models

Ece Takmaz



Utrecht
University

19.11.2025

- **Multimodal NLP**
- **Visual & linguistic processes** in deep neural networks
- Inspired by **cognitive science and psycholinguistics**
- Also using AI models to gain insight into human processes

- **Modelling Human Gaze in Language Use**

- Looking at images
- Looking at text
- Producing and understanding language

Generating Image Descriptions Using Human Gaze

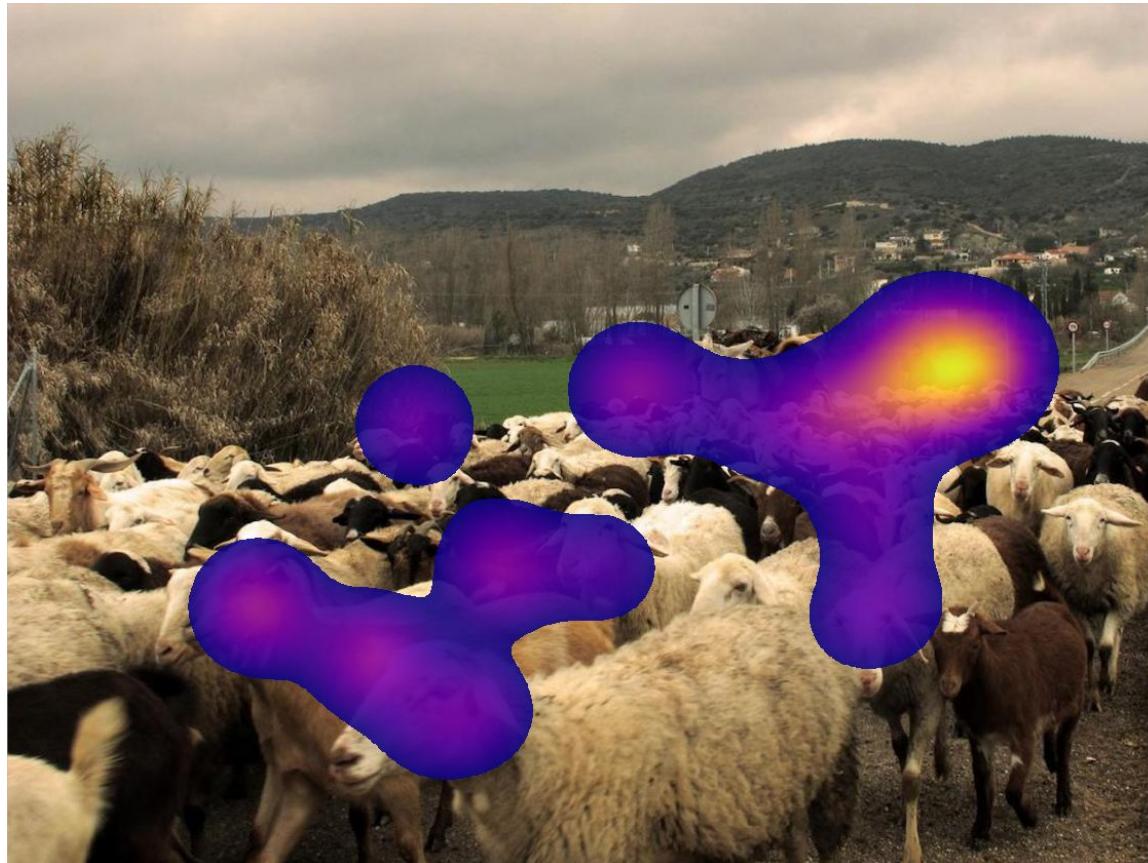


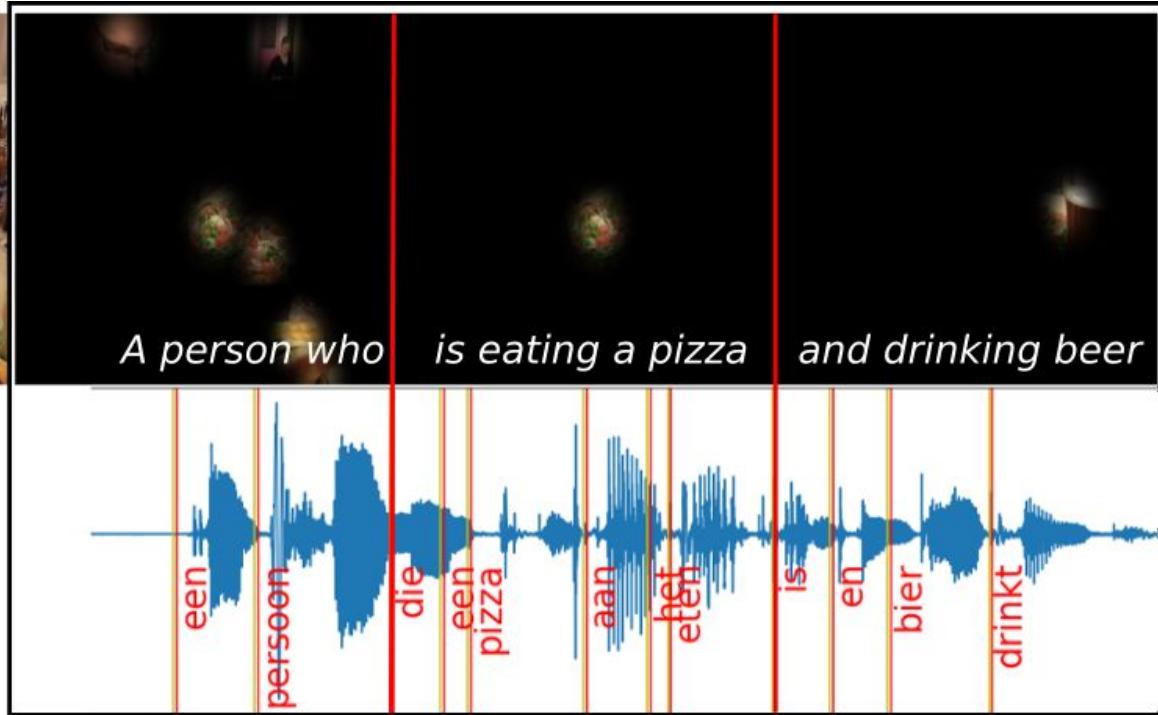
Een treinstation waarbij mensen op het perron aan het wachten zijn en waarbij net een goederentrein langsrijdt
(A train station where people are waiting on the platform and where a freight train is just passing by)

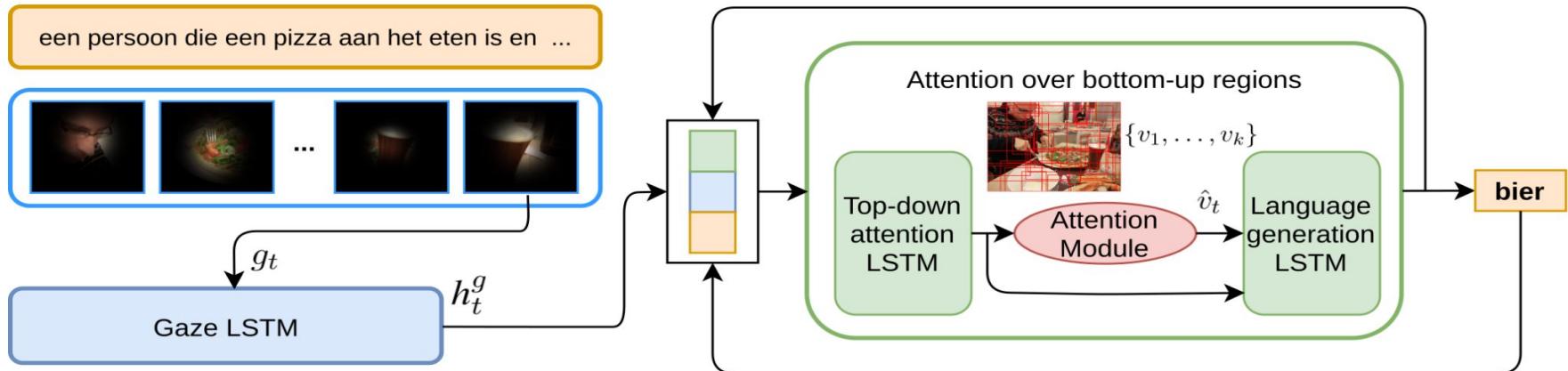
Een station waar een goederentrein voorbij komt
(A station where a freight train passes by)

Een vrachttrein op een Brits station
(A freight train at a British station)









Generated:

uh uh uh uh met een aantal vogels ...

Humans:

uh allemaal duiven

uh allemaal duiven die opvliegen of net
landen uh in een stadscentrum

uh een straat met heel veel duiven die
rond vliegen en heel veel
elektriciteitskabels in de lucht



Multi- and Cross-Lingual Prediction of Human Reading Behavior

Eye movements during reading

- **Communication strategies in dialogue that involves vision and language**
 - Referential tasks
 - Multimodal dialogue
 - Images in the context of
 - Other images
 - Dialogue
 - **PhotoBook Dataset (Haber et al., 2019)**

(Grice, 1975; Clark and Wilkes-Gibbs, 1986; Clark and Brennan, 1991; Clark, 1996, Garrod and Anderson, 1987; Brennan and Clark, 1996, Pickering and Garrod, 2004)

Player A



B: do you have plate of food on 2 pink bowls? Rice in one, veggies and ham other

A: I have that one



B: white square bowl with white rice and broccoli and yellow shreds?

A: I have a fry tray with hot dogs next to it

A: I have that



Player B



B: i do not have that

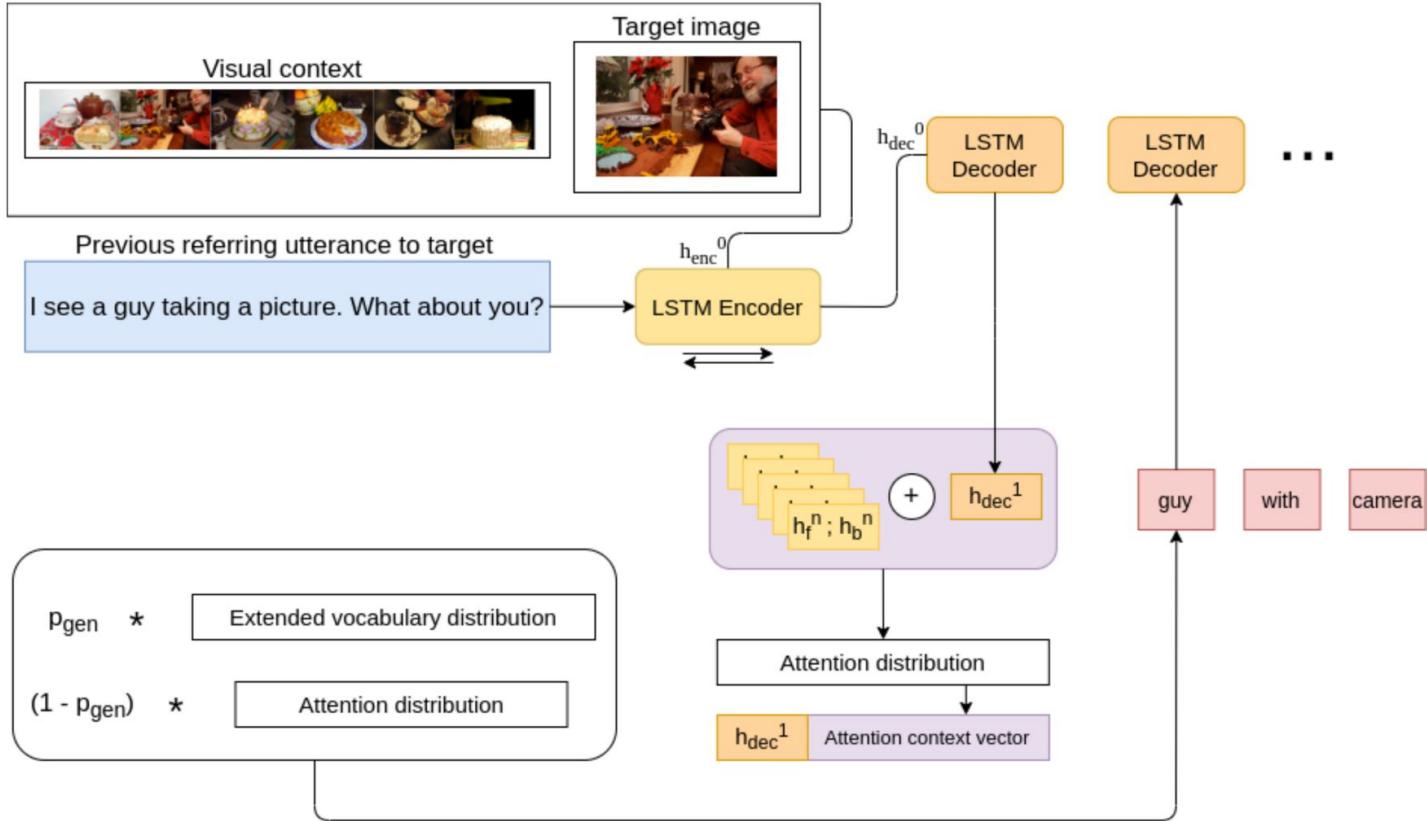
B: you have two bowls and an oval tray? one has rice. one has kimchi.

B: he oval tray has 3 slots of toppings



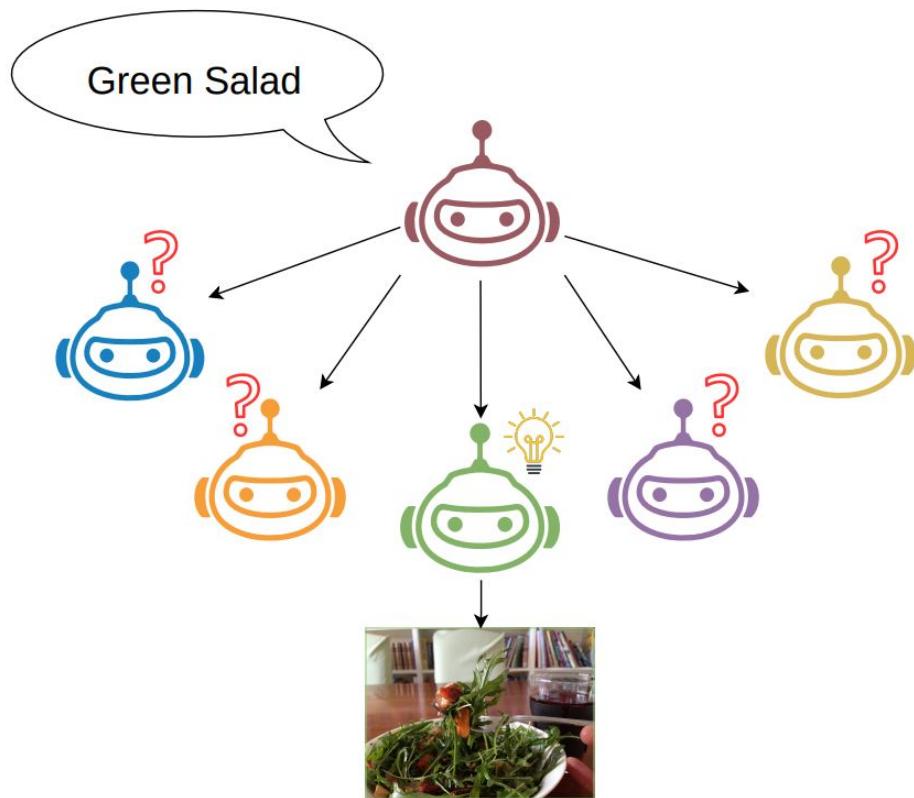
1. *girl on end of bed with computer, she has pigtails*
2. *Girl with pigtails?*
3. *Pigtail girl?*
4. *Pigtails? lol*

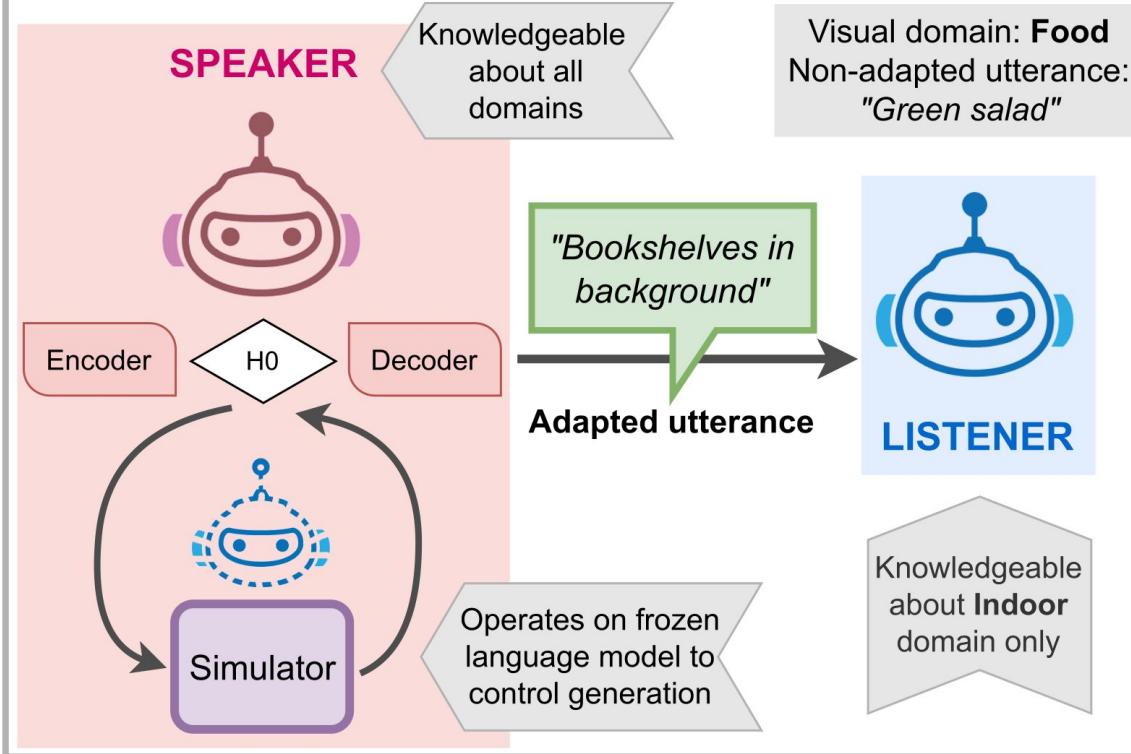
- Dialogue history and distilling the most important information
- **Less descriptive, yet discriminative**, as quantified by the CLIP model



Speaker Adaptation in Visually Grounded Dialogue

- Audience-aware adaptation of pretrained speaker models
- Adapting to domain-specific listeners with Theory of Mind









Describing Images *Fast and Slow*



Min: 1.69 sec



Max: 7.07 sec



Mean onset: 3.46 seconds
Variation in starting points: 11
Most common starting point: pier
Image specificity BLEU-2: 0.39
Variation in gaze: 38.47

een **pier** waar het heel erg druk is uh rechts is een vis aquarium waar je vissen kan aanraken
(a pier where it is very busy uh on the right is a fish aquarium where you can touch fish)

een drukke **straat** met een aantal restaurants pier 39
(a busy **street** with a number of restaurants pier 39)

pier waar veel mensen lopen
(**pier** where many people walk)

een drukbezette **pier**
(a busy **pier**)

een toeristische **plaats** waar veel verschillende entertainment dingen te doen zijn
(a touristic **place** where there are many different entertainment things to do)

de **ingang** van een aquarium met veel mensen op een plein
(the **entrance** to an aquarium with many people in a square)

Is variation in one signal correlated with variation in another?

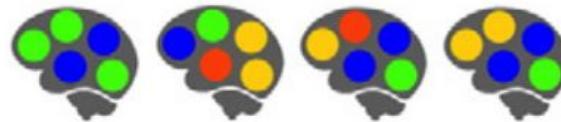
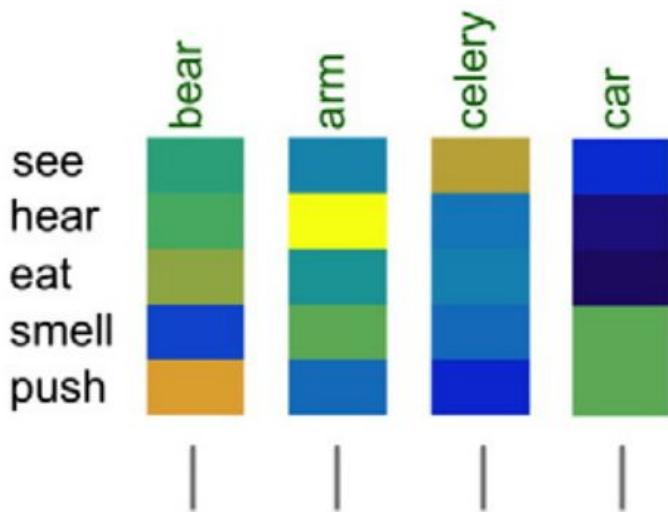
Can we predict variation using image representations obtained from pretrained encoders (CLIP and ViT)?

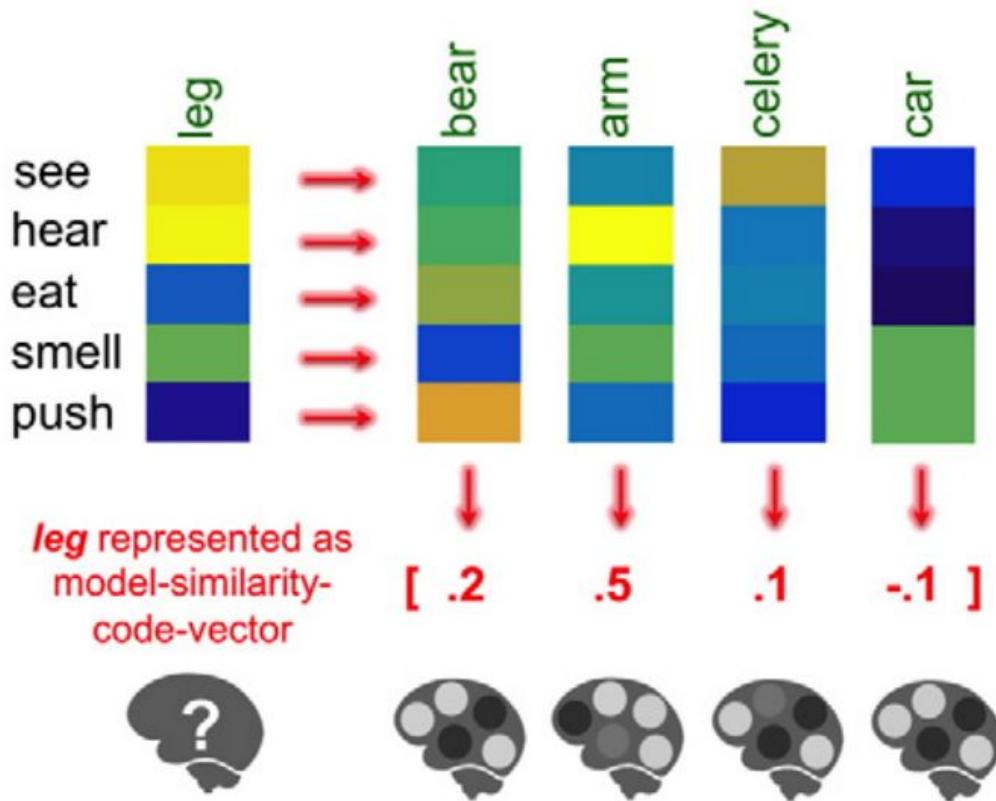


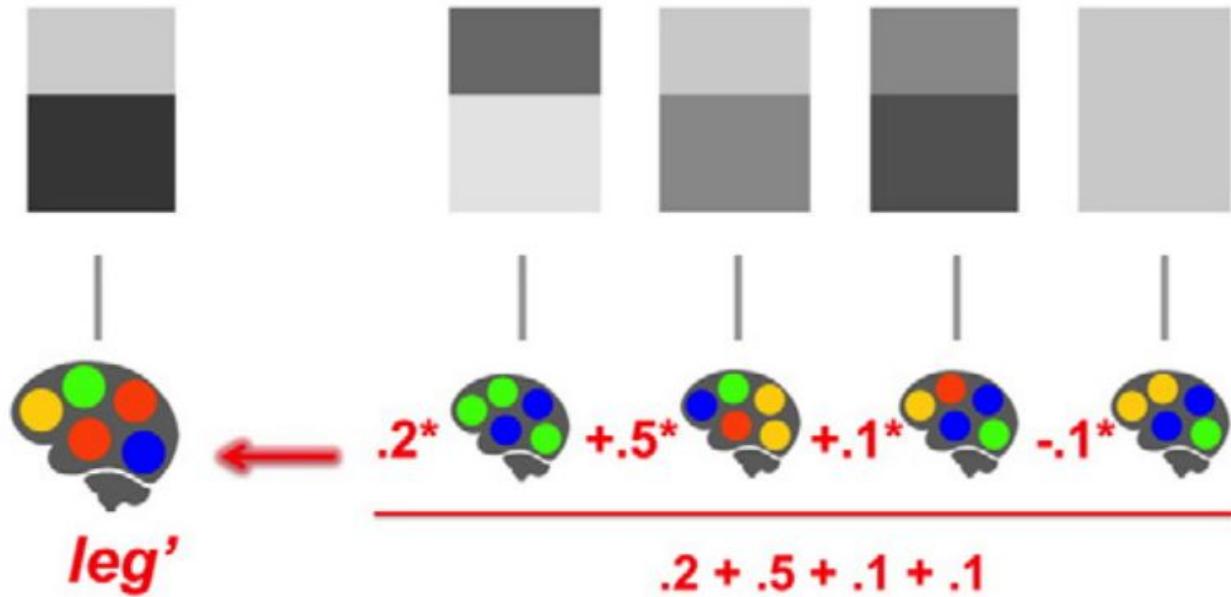
Min: 11.22



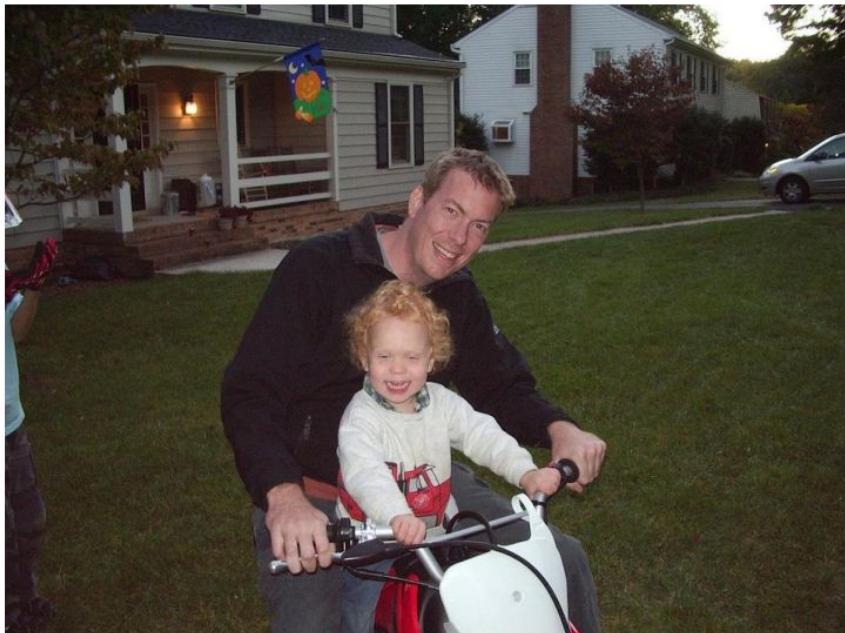
Max: 38.79







Predicted Gaze Variation



Min: 23.666



Max: 24.308

Leonidas_Smith1866 2013

Will a multimodal model demonstrate accuracy above chance level in predicting emotions for abstract art images and textual justifications?

Do the model's color-emotion associations align with those of humans found in existing literature?

Decoding Emotions in Abstract Art

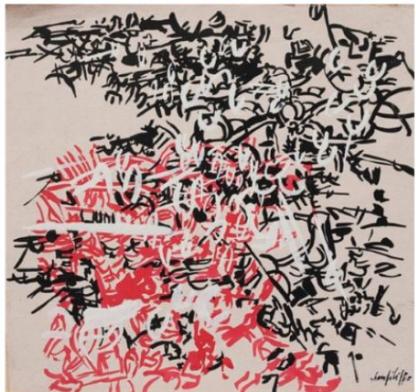
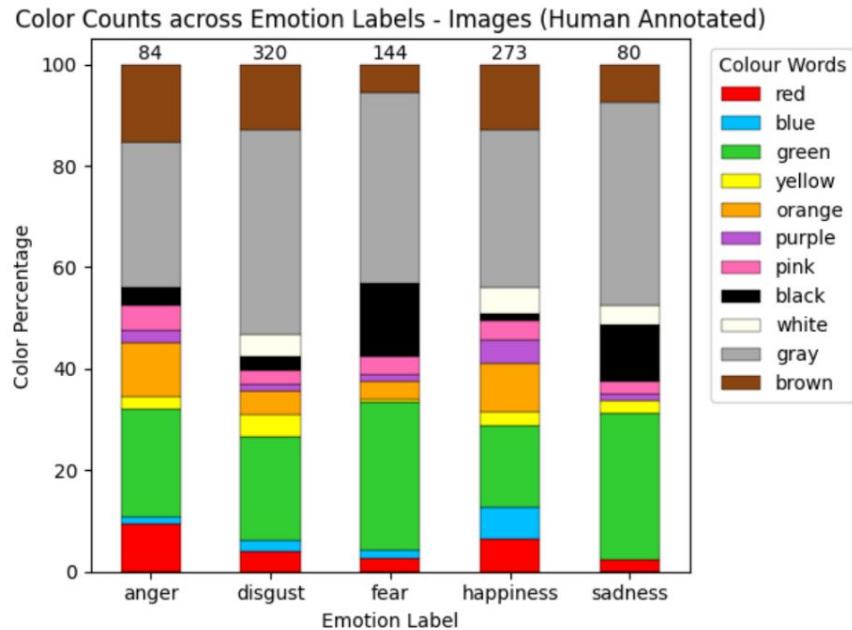
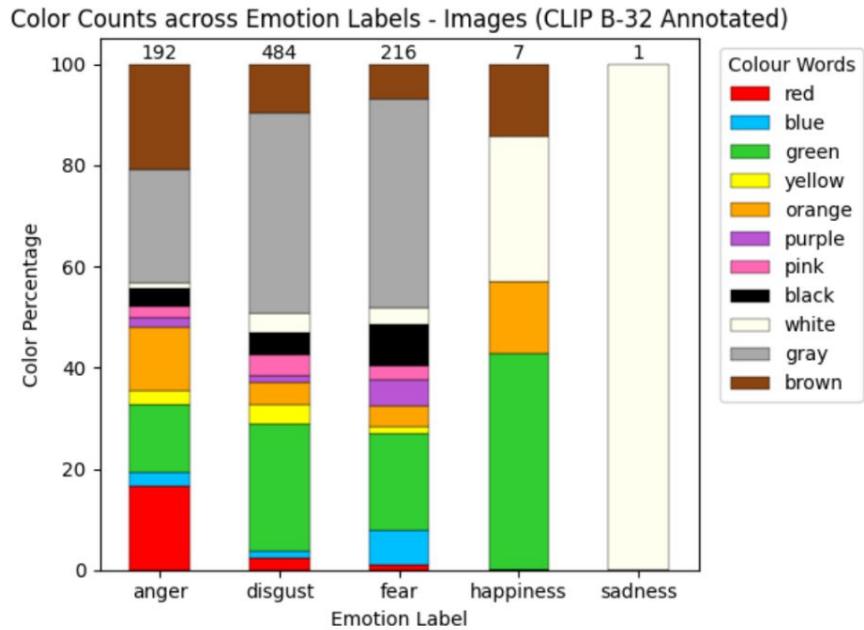


Image ID: image_328, Antonio-sanfilippo_composizione-1955

Annotator ID	Emotion	Rationale
28	anger	I see anger in the red and black masses arrayed against each other.
28	fear	Looks like the black shapes are running away in fear from the angry red shapes.
284	fear	I SAW SOME FIGURES AND BLACK RED COMBINE
285	anger	Looks like the black is swallowing the red paint.
260	anger	This looks like a bunch of idiots congregating during a pandemic and spreading the coronavirus. The red is the virus. It makes me angry that people are so selfish and stupid.
28	anger	The red is fighting the black in a vicious battle.
284	anger	BLACK LINES HAVE MORE FIGURES
284	fear	BLACK AND RED SHAPES

Number of annotators: 4

Dominant color of an image (CLIP- vs. human-annotated images)





**Utrecht
University**

- Postdoc - NGF call AiNed XS Europe
 - Albert Gatt (March-November 2025)

- Postdoc - ERC Consolidator project
 - Jakub Dotlačil (March 2025 - October 2027)



CON

SEM

SYN



SEM



SYN



SEMSYN



EXSYN



EXSEMSYN



	TRF _{tgt}	red van (A)
noise 0.0	TRF _{vis}	red truck (A)
	TRF _{sym}	red truck (A)
	TRF _{tgt}	left elephant (F)
noise 1.0	TRF _{vis}	white truck (A)
	TRF _{sym}	car on left (A)

Semantic Violation in Scenes: Investigating Multimodal Context in Referential Communication

COOCO - Common Objects Out-of-Context

<https://huggingface.co/datasets/fmerlo/COOCO>

Original



Clean



Generated



High



Medium

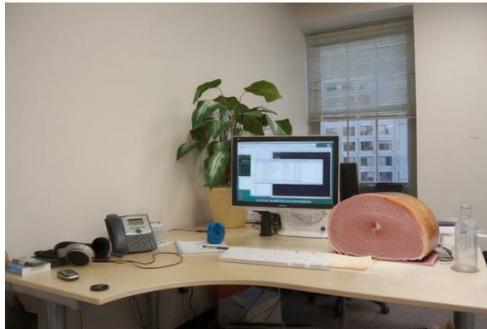


Low





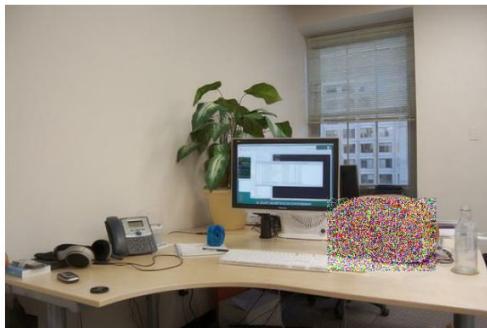
Original



Modified

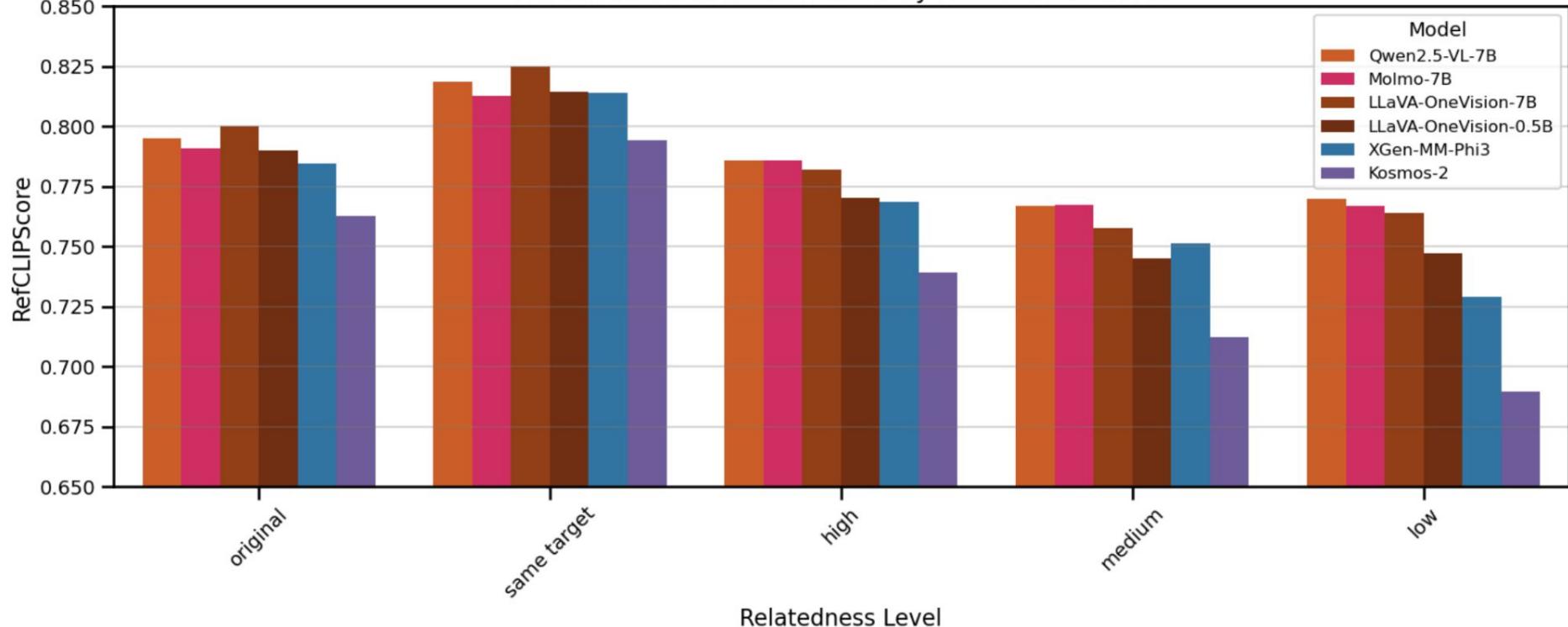


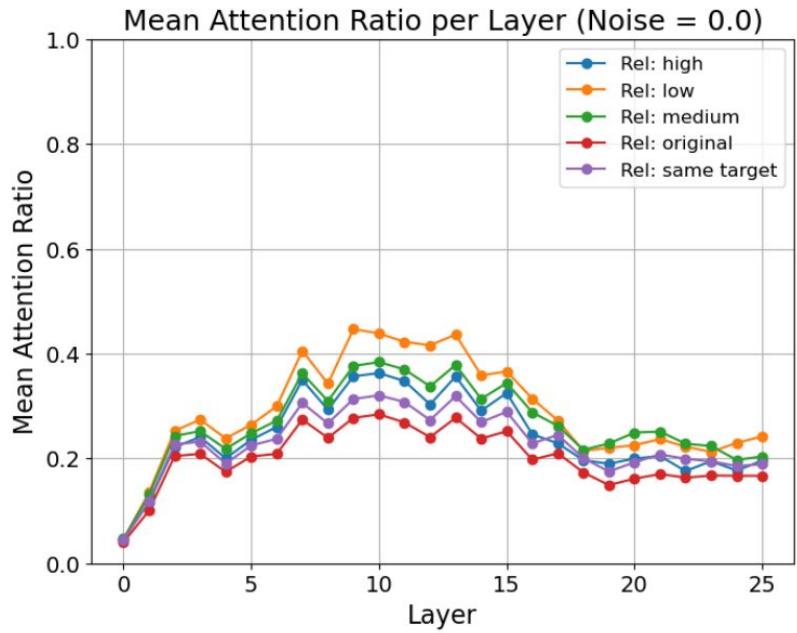
Context Noise



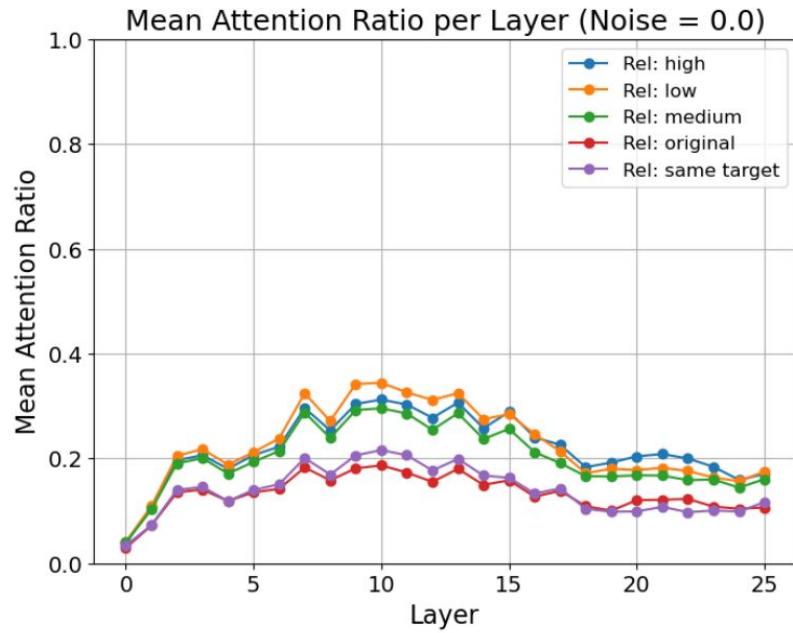
Target Noise

RefCLIPScore at Noise Level 0 by Relatedness Level

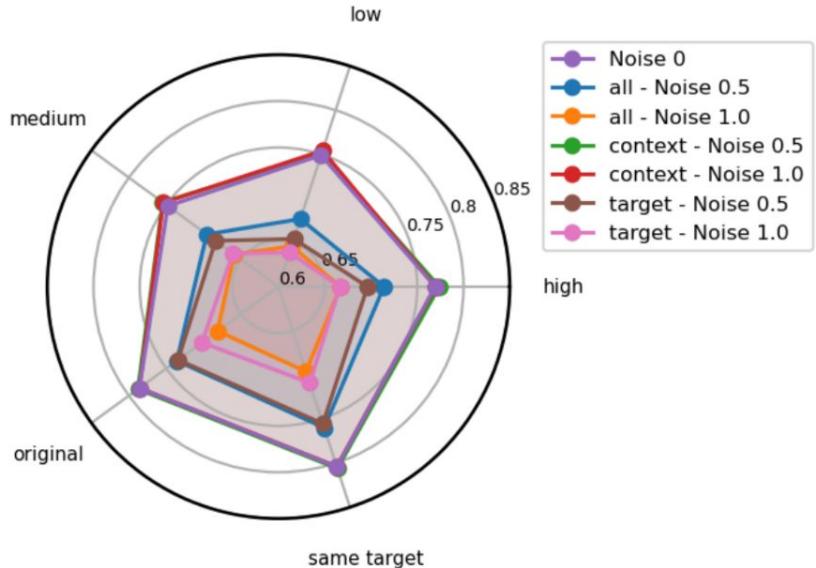




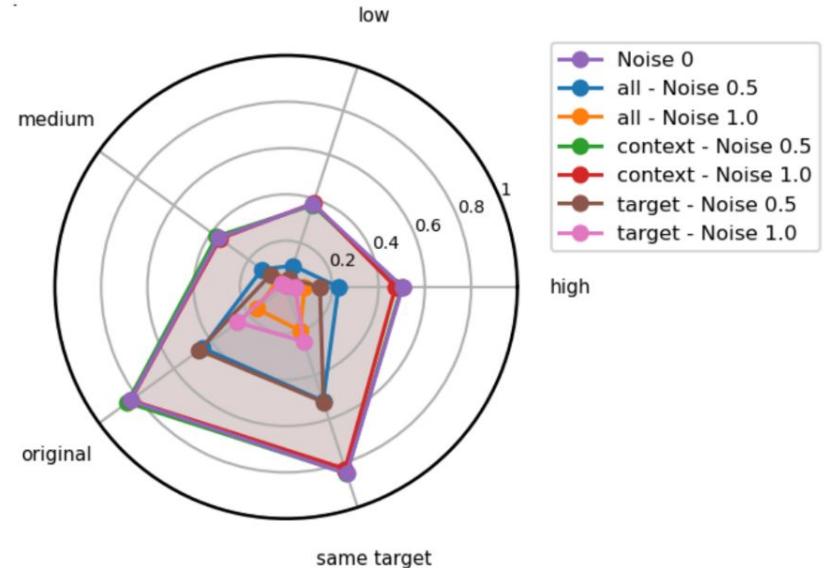
Correct Responses



Incorrect Responses



(a) RefCLIPScores by relatedness, noise area, and noise level.



(b) Accuracy across noise levels, noise areas, and relatedness conditions.

- - Scene context acts as a distractor for targets that violate scene semantics
 - Scene context acts as a facilitator when congruent targets are obscured
- Targets that violate scene semantics attract more attention
- Targets attract attention even under moderate noise conditions



human: sink (5); bowl (2); crab (2); bathtub shape (1)

LLaVA 7b: bathtub (6); rectangle (2); bathroom (2)

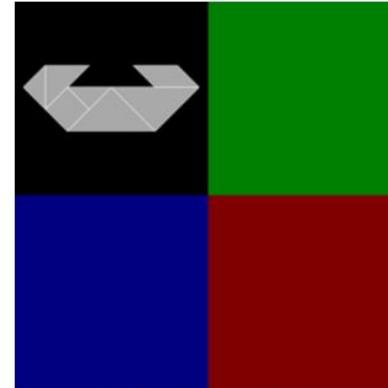
LLaVA 72b: house (8); boat (1); bathtub (1)



human: crab (7); bathtub (1); bowl (1); bull (1)

LLaVA 7b: sun (3); bird (2); diamond (2); boat (1); wave (1); house (1)

LLaVA 72b: sailboat (4); house (3); boat (3)



human: crab (4); bowl (2); dog (1); seal (1); letter c (1); space ship (1)

LLaVA 7b: house (3); square (2); diamond (2); triangle (1); parallelogram (1); box (1)

LLaVA 72b: house (7); boat (3)

Traces of Image Memorability in Vision Encoders

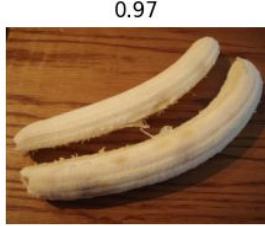


Image Memorability

Complex phenomenon, intrinsic property of images, consistent across individuals with some influence from extrinsic factors

More memorable - humans, faces, food, animals

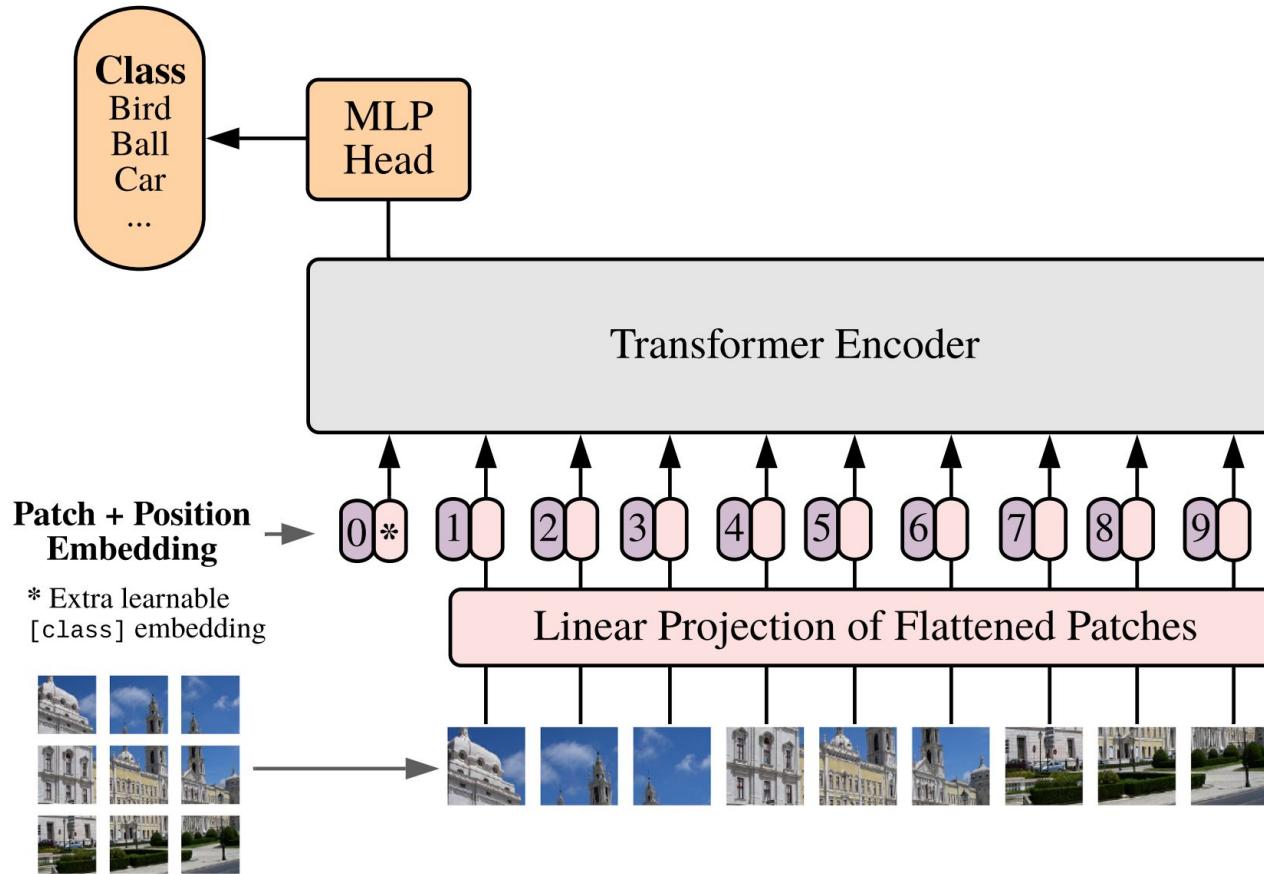
Less memorable - nature images, large uniform regions in images

Brain activations, deeper levels of processing during encoding, region uniformity, distribution of visual attention

Predicting memorability: mainly using CNNs

- Do internal proxies from transformer-based vision encoders capture information related to image memorability?
- Does autoencoder image reconstruction loss correlate with memorability?

Vision Transformer (ViT)



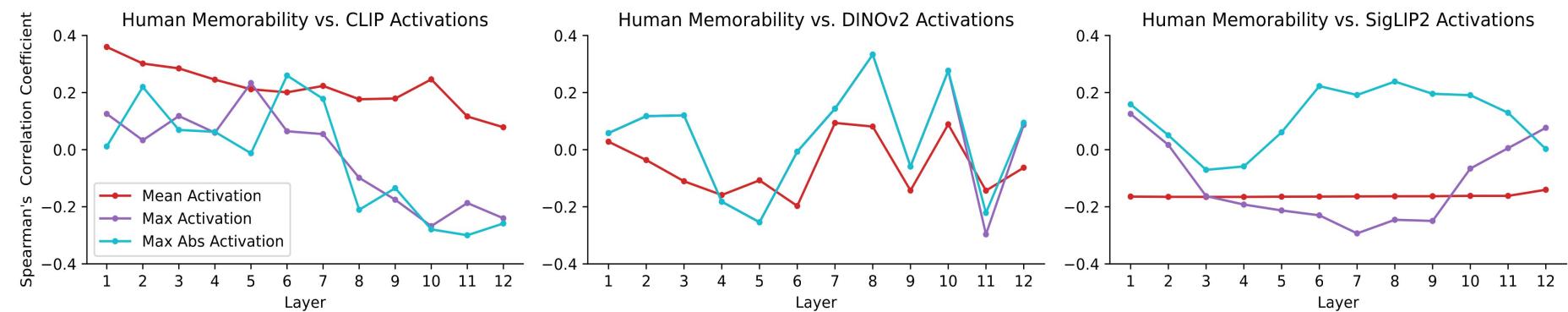
Model-Internal Features

[CLS] activations from CLIP, DINOv2, first image token from SigLIP2 over the layers

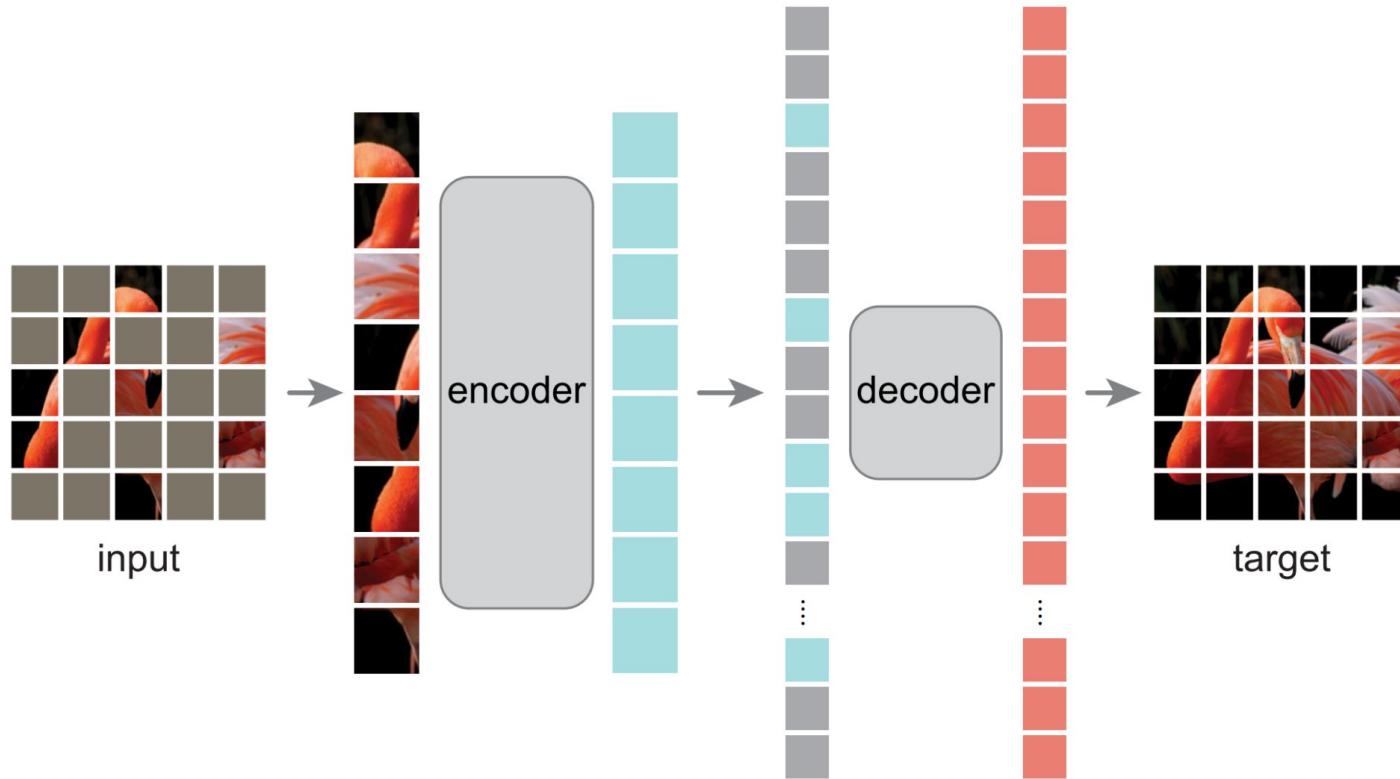
[CLS] delta: changes over the layers (cosine similarity)

Attention entropy of the attention applied by [CLS]

Patch uniformity: Variation in image token representations



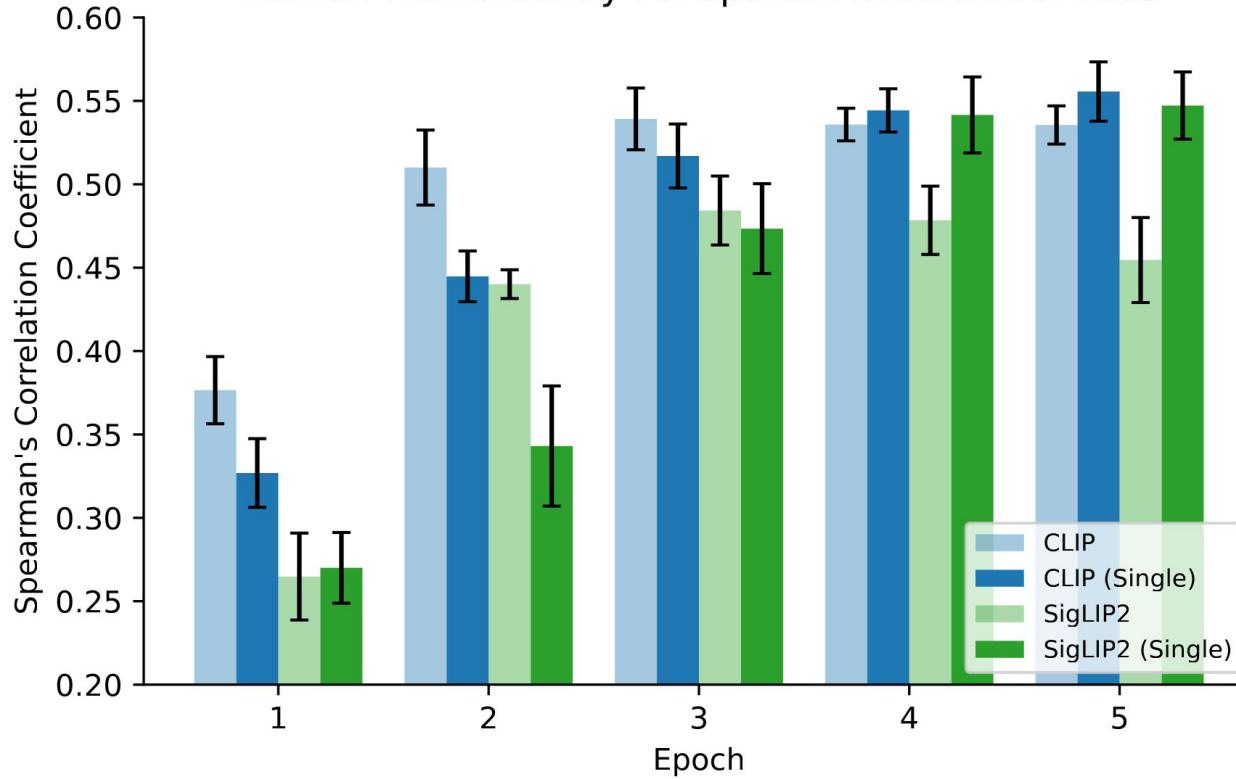
	Coef	Std Err	Z
Intercept	0.6932	0.001	577.806***
Activation max	0.0144	0.001	23.556***
Activation mean	0.0264	0.001	20.171***
Activation max abs	0.0144	0.001	23.556***
Patch uniformity	-0.0170	0.002	-10.574***
Attention entropy	0.0457	0.002	27.745***

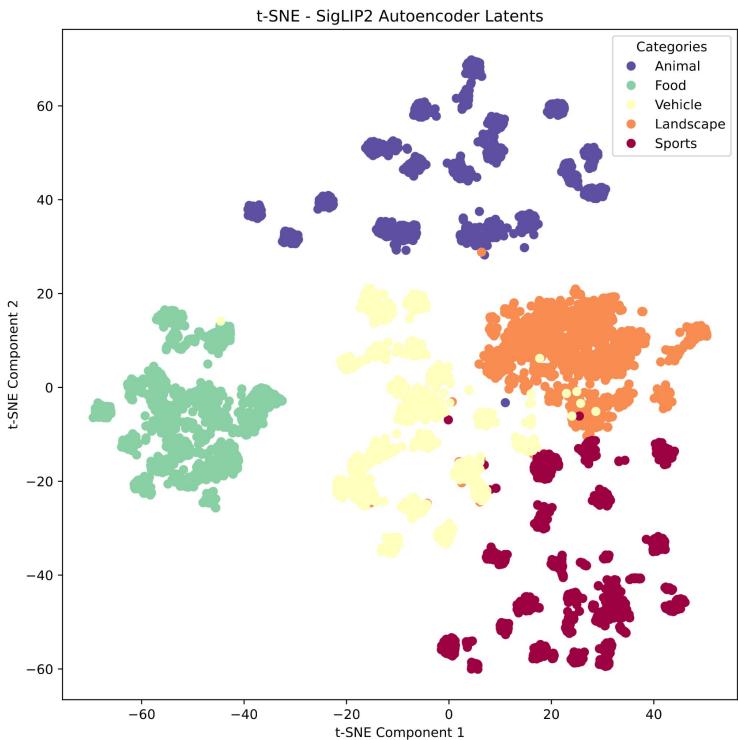
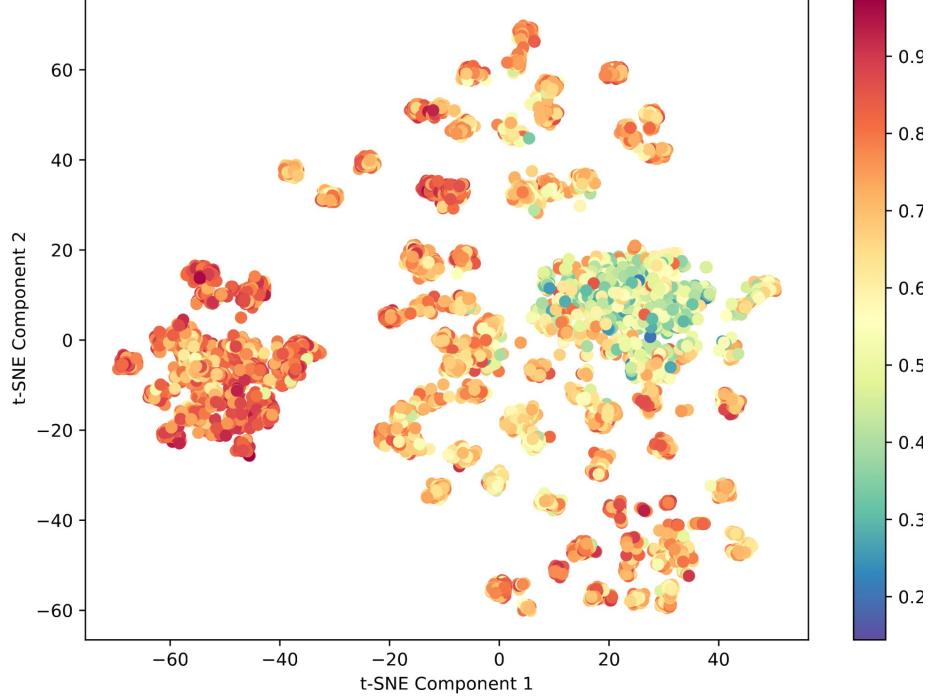


	ViTMAE - base	ViTMAE - large
Full MemCat	0.073***	0.056***
No ImageNet	0.118***	0.097***

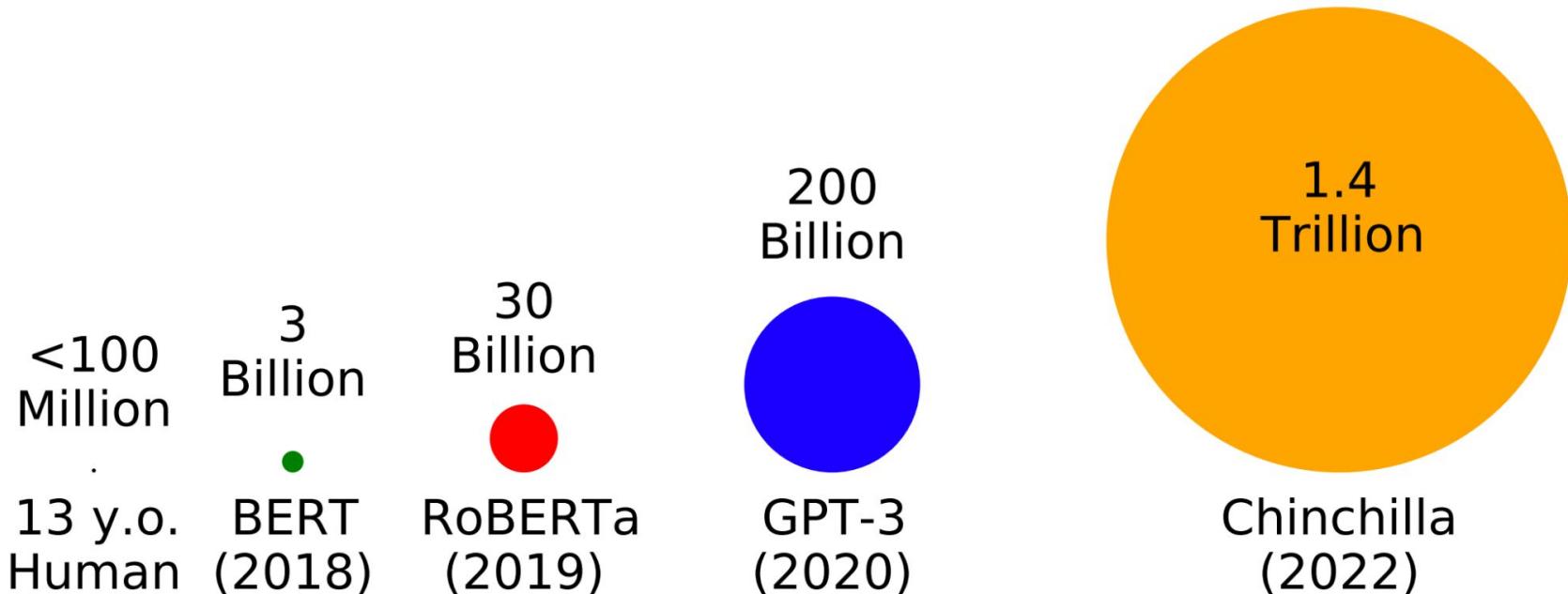
$$\mathcal{L} = \underbrace{\sum_{i=1}^N (\hat{y}_i - y_i)^2}_{\text{MSE loss (sum reduction)}} + \lambda \underbrace{\sum_{j=1}^M |z_j|}_{\text{sparsity loss}}$$

Human Memorability vs. Sparse Autoencoder Loss





BabyLM Challenge

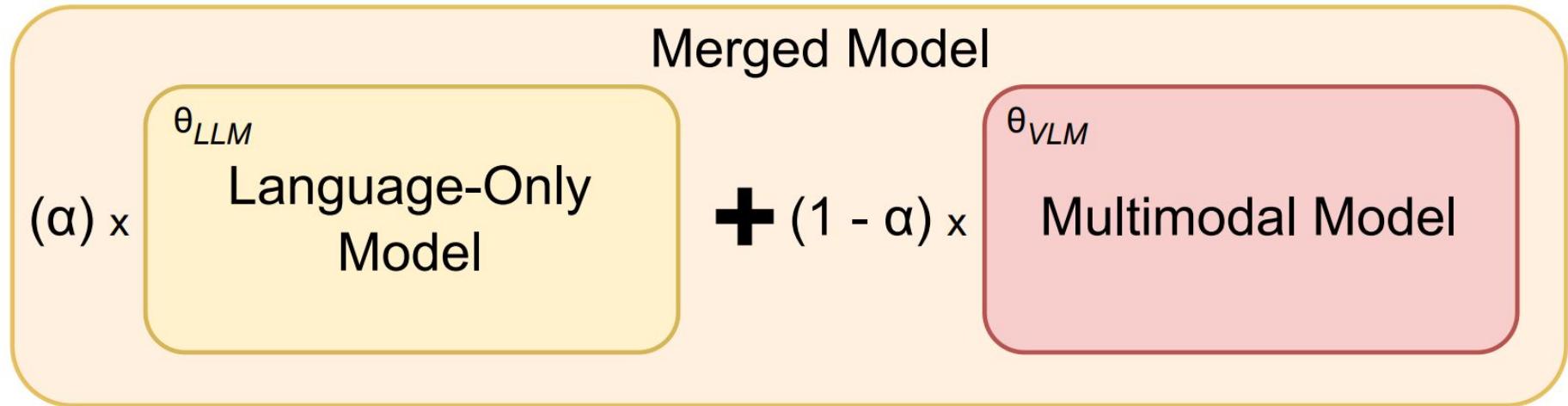


BabyLM Setup

- Data and training constraints
 - 100M words
 - 10 epochs
- Benchmarks
 - Language-only - BLiMP, EWOK, Wug, Entity Tracking ...
 - Multimodal - Winoground, DevBench

BabyLM - Shortcomings of Multimodal Models

BabyLM - Shortcomings of Multimodal Models



Conclusion

- Importance of exploring various aspects of **visuo-linguistic processes in humans** when modelling them with **deep neural networks**
- Current trends in AI (multimodal multilingual **large language models**)
- Further work in exploring computational approaches leveraging **human signals**
- Simultaneously benefit the development of **better AI models** and provide **insights into human cognition** itself

<https://ecekt.github.io>