# On visual grounding at the age of deep learning

Monthly online ILFC Seminar

Denis Paperno, 18.11.2024

#### Goals of the talk

Widespread assumptions:

- Models for language based on text inherently limited
- Grounding them e.g. in visual modality is qualitatively different
- In particular, representation of meaning qualitatively different
- How can we check this empirically?
- What changes when we switch to multimodality?
- How good/humanlike is multimodal models' grasp of meaning?

#### Transformer Language Models



## The Octopus thought experiment (Bender and Koller 2020)



alone and we should not expect machines to do so either"



# Grounding and representations: controlled model learning

Timothee Mickus, Elaine Zosa, and Denis Paperno. <u>Grounded and well-rounded</u>: A Methodological Approach to the Study of Cross-modal and Cross-lingual Grounding. In Findings of EMNLP 2023.

#### Vatex dataset

#### • Wang et al.2019



#### **10 English Descriptions:**

- A person wearing a bear costume is inside an inflatable play area as they lose their balance and fall over.
- A person in a bear costumer stands in a bounce house and falls down as people talk in the background.
- A person dressed in a cartoon bear costume attempts to walk in a bounce house.
- A person in a mascot uniform trying to maneuver a bouncy house.
- A person in a comic bear suit falls and rolls around in a moonbounce.
- A person dressed as a teddy bear stands in a bouncy house and then falls over.
  - Someone dressed in a bear costume falling over in a bouncy castle. (一) 有个穿着熊装的人在充气城堡摔倒了。
- A person dressed up as a bear is standing in a bouncy castle and falls down.
- A man in a bear costume is balancing in a bouncy castle before they tumble to the floor.
- A man in costume was trying to stand straight on a bouncy castle but fell.

#### **10 Chinese Descriptions:**

- 一个人穿着熊的布偶外套倒在了蹦床上。
- 一个人穿着一套小熊服装在充气蹦蹦床上摔 倒了。
- 一个穿着熊外衣的人在充气垫子上摔倒了。
- -个穿着深色衣服的人正在蹦蹦床上。
- 在一个充气大型玩具里,有一个人穿着熊的 衣服站了一下之后就摔倒了。
- ➡○ 一个打扮成泰迪熊的人站在充气房上, 然后 摔倒了。

  - 一个装扮成熊的人站在充气蹦床里, 然后摔 0 倒了。
- 一个穿着熊服装的人在一个有弹性的城堡里 平衡,然后他们就倒在了地板上。
- 一个穿着布偶熊的人试图站在一个充气城堡 上,但却摔倒了。



Multi-modal

#### Experimental setup

- For each type, train 25 model instances
- All models share the same architecture (64M parameters):
  - 6 layers
  - 8 heads per multihead sublayer
  - hidden dimensions of 512
  - latent feedforward dimensions of 2048
- Add varying degrees of noise to ensure comparable accuracy



## Results: attention patterns



#### Results: concreteness in word embeddings

	Silhouette		Pu	ırity	Inverse Purity			
Model	input	last hidden state	input	last hidden state	input	last hidden state		
Single-task								
Р	$0.021(\pm 0.004)$	$0.054(\pm 0.007)$	$0.746(\pm 0.021)$	$0.887(\pm 0.005)$	$0.062(\pm 0.015)$	$0.049(\pm 0.005)$		
C	<b>0.026</b> (±0.002)	$0.049(\pm 0.006)$	<b>0.760</b> (±0.013)	<b>0.890</b> (±0.005)	<b>0.080</b> (±0.025)	<b>0.068</b> (±0.014)		
Т	$0.023(\pm 0.004)$	<b>0.056</b> (±0.006)	$0.748(\pm 0.016)$	$0.884(\pm 0.005)$	$0.060(\pm 0.014)$	$0.047(\pm 0.006)$		
			Multitas	K		1		
Р	$0.023(\pm 0.004)$	0.051(±0.005)	$0.757(\pm 0.017)$	$0.888(\pm 0.005)$	$0.073(\pm 0.021)$	$0.050(\pm 0.008)$		
P∨C	$0.024(\pm 0.004)$	$0.056(\pm 0.009)$	$0.752(\pm 0.021)$	$0.891(\pm 0.007)$	$0.070(\pm 0.021)$	<b>0.061</b> (±0.011)		
<b>P</b> ∨C∨T	<b>0.027</b> (±0.004)	$0.055(\pm 0.007)$	0.765(±0.016)	$0.889(\pm 0.006)$	$0.073(\pm 0.016)$	0.052(±0.009)		
P∨T	$0.024(\pm 0.004)$	$0.053(\pm 0.008)$	$0.754(\pm 0.019)$	$0.889(\pm 0.006)$	$0.069(\pm 0.016)$	$0.050(\pm 0.007)$		
Multimodal								
Р	$0.025(\pm 0.004)$	$0.051(\pm 0.007)$	$0.759(\pm 0.022)$	0.890(±0.006)	$0.074(\pm 0.023)$	0.054(±0.010)		
P∧C	$0.026(\pm 0.003)$	0.052(±0.006)	$0.761(\pm 0.018)$	<b>0.891</b> (±0.007)	$0.073(\pm 0.019)$	0.052(±0.006)		
$P \land C \land T$	$0.024(\pm 0.004)$	<b>0.056</b> (±0.007)	$0.759(\pm 0.017)$	0.889(±0.0067)	$0.072(\pm 0.032)$	$0.048(\pm 0.005)$		
$P \wedge T$	$0.026(\pm 0.004)$	$0.051(\pm 0.008)$	$0.765(\pm 0.017)$	$0.886(\pm 0.007)$	$0.065(\pm 0.014)$	$0.046(\pm 0.005)$		

#### Interim conclusions

- Grounding leads to quantitatively and qualitatively different models
- Multimodal vs multilingual grounding: qualitatively different
- Multimodality: same input, qualitatively different representations

## Grounding and representations: pretrained models

Aleksey Tikhonov, Lisa Bylinina and Denis Paperno. <u>Leverage Points In Modality Shifts: Comparing</u> <u>Language-Only and Multimodal Word Representations</u>. Proceedings of the 12th Joint Conference on Lexical and Computational Semantics: 11–17.

#### Goals

- There are pretrained models for both text and multimodal input.
- Models used in this study:
  - CLIP, OpenCLIP, Multilingual CLIP
    - For each: 2 ways of extracting word embeddings iso (word in isolation), avg (over 10 sentences)
    - fastText, multilingual BERT, XLM-RoBERTa For BERT and RoBERTa: 3 ways of extracting word embeddings iso, avg-last, avg-bottom
- Unimodal vs. multimodal representation spaces: different structure?
- What semantic factors account for differences in structure?

#### Methodology

- Use a sample of word pairs
- For each model M, calculate cosine distance of pairs  $dist^{M}(w_{i}, w_{j})$
- Control for differences in embedding space, e.g. anisotropy:
  - Rank pairs by shift of distance between two models  $\frac{dis^{M_1}(w_i, w_j)}{dist^{M_2}(w_i, w_j)}$
  - Regression analysis: predict rank of shifts from linguistic properties
- Hypothesis: some properties explain more of a difference when models belong to different classes (multimodal vs textual)

#### Factors considered

- Frequency
- Concreteness (Ghent norms)
- Ontological category (WordNet supersenses)
- Affective meaning (VAD: Valence, Arousal, Dominance)
- Relation within word pair
  - WordNet relations (e.g. hypernymy)
  - ConceptNet relations (e.g. *is used for*)

#### Illustration: method of model comparison

• Variance of distance ratio rank explained:

CLIP-iso vs.	XLMR-iso	mBERT-iso	<b>BERT-avg-last</b>	fastText
Baselines				
concreteness	9.5	11.68	2.27	8.71
frequency	5.43	7.81	1.91	0.45
concreteness+frequency	16.73	17.16	3.65	9.54
+taxonomic	21 (+4.27)	20.35 (+3.19)	5.43 (+1.78)	19.50 (+9.96)
+VAD	17.36 (+0.63)	17.49 (+0.33)	4.62 (+0.97)	10.78 (+1.24)
+WordNet relations	18.47 (+1.74)	17.36 (+0.2)	10.05 (+6.4)	10.34 (+0.8)
+ConceptNet relations	19.8 (+3.07)	17.47 (+0.31)	8.84 (+5.19)	10.26 (+0.72)

#### Illustration: specific features



#### Results: global view



#### Results summary

- **Concreteness** plays a major role in explaining modality shifts, in line with results of previous studies.
- Combined WordNet supersenses have a significant effect for many, although different subsets of features prove significant in different pairs of models.
- WordNet and ConceptNet relations tend to be significant when aggregated, although no individual relation has a systematic effect across model pairs.
- VAD features produce varied effects, with valence showing the most consistent modality difference.

### Case study: derivational contrast

Claudia Tagliaferri, Sofia Axioti, Albert Gatt and Denis Paperno. <u>The Scenario Refiner: Grounding</u> <u>subjects in images at the morphological level</u>. LIMO workshop (Linguistic Insights from and for Multimodal Language Processing @KONVENS 2023).

#### Methodology

- Collect human ratings on image-text match on a Likert scale
- -er nouns vs -ing forms of verbs from four domains:

professional	sports	artistic	general
baker vs bake	runner vs run	painter vs paint	supporter vs support
teacher vs teach	surfer vs surf	singer vs sing	reader vs read
cleaner vs clean	skier vs ski	dancer vs dance	lover vs love
	•••		

• Contrast human judgments to VLM predictions of match

### The task:

the woman with pink gloves is a driver the woman with pink gloves is driving

**4.4** 4.3



the man and the woman are supporters6.2the man and the woman are supporting6.4



### Human vs human (correlation, two groups)

Domain	<b>Derived noun</b>	Verb	Morphological contrast
Professional domain	0.76	0.84	0.75
Sport domain	0.69	0.70	0.60
Artistic domain	0.79	0.31	0.51
General	0.92	0.88	0.94
All domains	0.80	0.81	0.78

### Statistics of preference

	<b>Derived noun</b>	Verb
Humans	8.3%	91.7%
CLIP ViT-L/14@336px	51.9%	48.1%
CLIP RN50x64	52.8%	47.2%
CLIP ViT-B/32	49.1%	50.9%
ViLT	47.2%	52.8%
LXMERT	51.9%	48.1%

#### Model vs human correlations overall

Model	<b>Derived noun</b>	Verb	Morphological contrast
CLIP ViT-L/14@336px	0.13	0.08	0.15
CLIP RN50x64	0.09	0.08	-0.01
CLIP ViT-B/32	0.09	0.18	0.08
ViLT	0.07	0.26	0.32
LXMERT	0.16	0.03	0.21

#### Model vs. human correlations per domain

<del>))</del>	Sport domain			Professional domain		
	Deriv. noun	Verb	Morph. contrast	Deriv. noun	Verb	Morph. contrast
CLIP ViT-L/14@336p	0.02	-0.25	-0.06	-0.04	0.19	0.33
CLIP RN50x64	0.02	-0.31	-0.11	-0.22	0.33	0.23
CLIP ViT-B/32	-0.04	-0.26	-0.11	-0.16	0.23	0.40
ViLT	-0.01	-0.04	0.45	0.30	0.45	0.68
LXMERT	0.08	-0.32	0.17	-0.10	0.18	0.40
-	A	Artistic domain		General		
	Derived noun	Verb	Contrast	Derived noun	Verb	Contrast
CLIP ViT-L/14@336p	-0.03	0.01	0.22	0.28	0.18	-0.06
CLIP RN50x64	0.06	0.21	0.27	0.08	0.23	0.10
CLIP ViT-B/32	-0.09	-0.04	-0.007	0.23	0.25	-0.06
ViLT	0.44	0.15	0.39	-0.06	0.05	0.25
LXMERT	0.29	0.26	0.42	-0.01	-0.25	0.12

#### Conclusions from the three studies

- 1. Methodology for testing qualitative differences between modalities
- 2. Identify which dimensions of meaning affected by modality
- 3. Meaning in multimodal models tested in a novel way, not humanlike

#### Thank you!

• Collaborators:

Albert Gatt Aleksey Tikhonov Claudia Tagliaferri Elaine Zosa Lisa Bylinina Sofia Axioti Timothee Mickus