

Robust evaluation of language–brain encoding experiments

Lisa Beinborn, Samira Abnar, Rochelle Choenni
Institute for Logic, Language and Computation
Universiteit van Amsterdam



Cognitively inspired
Language Processing

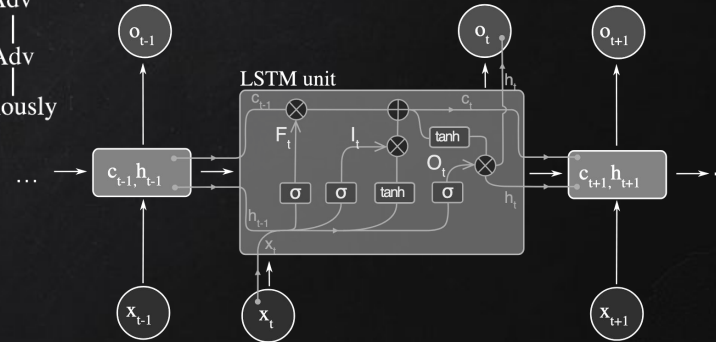
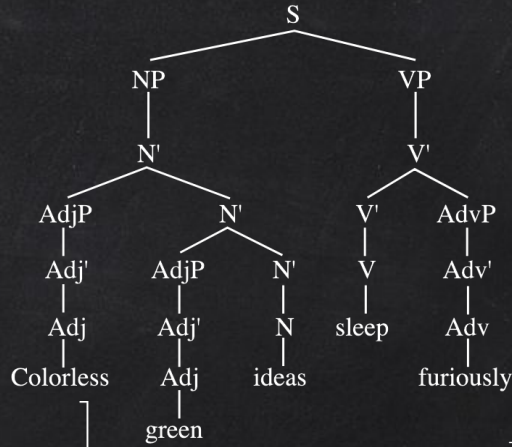
HOW TO MODEL LANGUAGE?

ۛ ۛ ۛ
 ۛ ۛ ۛ
 ۛ ۛ ۛ

Wiktionary
The free dictionary

CAT	<i>v</i>						
VFORM	<i>finite</i>						
TNS	<i>pres</i>						
SUBJ	<table border="1"> <tr> <td>CAT</td> <td><i>np</i></td> </tr> <tr> <td>PER</td> <td><i>1</i></td> </tr> <tr> <td>NUM</td> <td><i>sg</i></td> </tr> </table>	CAT	<i>np</i>	PER	<i>1</i>	NUM	<i>sg</i>
CAT	<i>np</i>						
PER	<i>1</i>						
NUM	<i>sg</i>						
OBJ1	[CAT <i>np</i>]						

→ writes



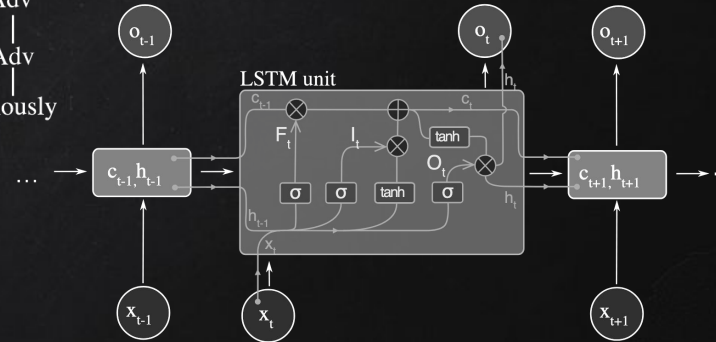
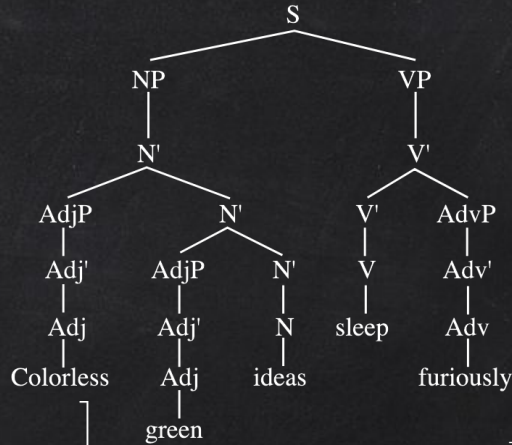
HOW TO MODEL LANGUAGE?

无 维 王
 入 W 山
 之 尸 昌

Wiktionary
The free dictionary

CAT	<i>v</i>						
VFORM	<i>finite</i>						
TNS	<i>pres</i>						
SUBJ	<table border="1"> <tr> <td>CAT</td> <td><i>np</i></td> </tr> <tr> <td>PER</td> <td><i>1</i></td> </tr> <tr> <td>NUM</td> <td><i>sg</i></td> </tr> </table>	CAT	<i>np</i>	PER	<i>1</i>	NUM	<i>sg</i>
CAT	<i>np</i>						
PER	<i>1</i>						
NUM	<i>sg</i>						
OBJ1	[CAT <i>np</i>]						

→ writes



[4 1 5 3 1 6 8 0 3 ...]

WHICH ONE?

HIERARCHICAL

BIDIRECTIONAL

CONTEXTUALIZED

LONG SHORT-TERM NETWORK

MULTILINGUAL

STACKED

ATTENTION

CHARACTER

UNIVERSAL

CONVOLUTION

TRANSFORMER



TASKS

Model	Avg	Single Sentence		Similarity and Paraphrase			Natural Language Inference			
		CoLA	SST-2	MRPC	QQP	STS-B	MNLI	QNLI	RTE	WNLI
Single-Task Training										
BiLSTM	62.0	15.7	85.9	69.3/79.4	81.7/61.4	66.0/62.8	70.3/70.8	60.8	52.8	62.3
+ELMo	<u>66.2</u>	35.0	<u>90.2</u>	69.0/80.8	85.7/65.6	64.0/60.2	72.9/73.4	<u>69.4</u>	50.1	65.1
+CoVe	62.4	14.5	88.5	<u>73.4/81.4</u>	83.3/59.4	<u>67.2/64.1</u>	64.5/64.8	64.8	<u>53.5</u>	61.6
+Attn	60.0	15.7	85.9	68.5/80.3	83.5/62.9	59.3/55.8	74.2/73.8	51.9	51.9	55.5
+Attn, ELMo	64.8	35.0	<u>90.2</u>	68.8/80.2	86.5/66.1	55.5/52.5	76.9/76.7	61.1	50.4	65.1
+Attn, CoVe	60.8	14.5	88.5	68.6/79.7	84.1/60.1	57.2/53.6	71.6/71.5	53.8	52.7	64.4
Multi-Task Training										
BiLSTM	63.5	24.0	85.8	71.9/82.1	80.2/59.1	68.8/67.0	65.8/66.0	71.1	46.8	63.7
+ELMo	64.8	<u>27.5</u>	89.6	76.2/83.5	78.5/57.8	67.0/65.9	67.1/68.0	66.7	55.7	62.3
+CoVe	62.2	16.2	84.3	71.8/80.0	82.0/59.1	68.0/67.1	65.3/65.9	70.4	44.2	65.1
+Attn	65.7	0.0	85.0	75.1/ 83.7	84.3/63.6	<u>73.9/71.8</u>	<u>72.2/72.1</u>	82.1	61.7	63.7
+Attn, ELMo	69.0	18.9	91.6	77.3/83.5	85.3/63.3	72.8/71.1	<u>75.6/75.9</u>	81.7	61.2	65.1
+Attn, CoVe	64.3	19.4	83.6	75.2/83.0	84.9/61.1	72.3/71.1	69.9/68.7	78.9	38.3	65.1
Pre-Trained Sentence Representation Models										
CBoW	58.9	0.0	80.0	73.4/81.5	79.1/51.4	61.2/58.7	56.0/56.4	75.1	54.1	62.3
Skip-Thought	61.5	0.0	81.8	71.7/80.8	82.2/56.4	71.8/69.7	62.9/62.8	74.7	53.1	65.1
InferSent	64.7	4.5	<u>85.1</u>	74.1/81.2	81.7/59.1	75.9/75.3	66.1/65.7	79.8	58.0	65.1
DisSent	62.1	4.9	83.7	74.1/81.7	82.6/59.5	66.1/64.8	58.7/59.1	75.2	56.4	65.1
GenSen	<u>66.6</u>	<u>7.7</u>	83.1	<u>76.6/83.0</u>	<u>82.9/59.8</u>	79.3/79.2	<u>71.4/71.3</u>	82.3	<u>59.2</u>	65.1

Table 3: Baseline performance on the GLUE tasks. For MNLI, we report accuracy on the matched and mismatched test sets. For MRPC and Quora, we report accuracy and F1. For STS-B, we report Pearson and Spearman correlation. For CoLA, we report Matthews correlation. For all other tasks we report accuracy. All values are scaled by 100. A similar table is presented on the online platform.

SAME SAME BUT DIFFERENT?

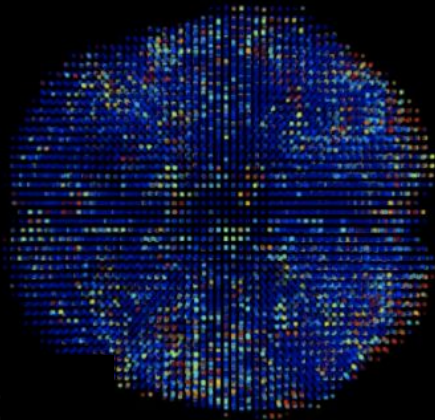
- ◆ Why is one model better than another?
- ◆ How do the representations differ?
- ◆ Which linguistic properties are encoded?
- ◆ Which phenomena cannot be modeled?
- ◆ How are they different from human language processing?

COMPARE COMPUTATIONAL MODELS
WITH THE SIGNAL THAT WE MEASURE
WHEN HUMANS PROCESS LANGUAGE.

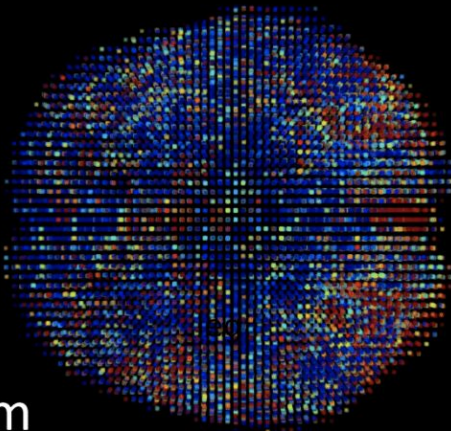


BUT HOW?

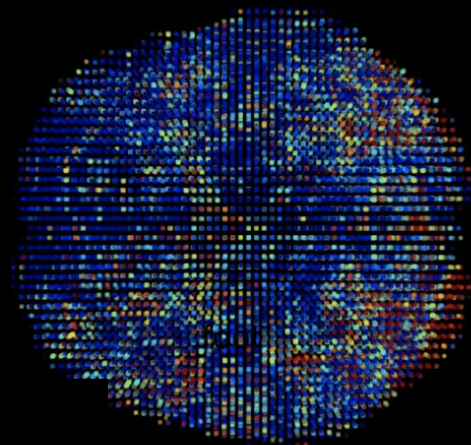
ant



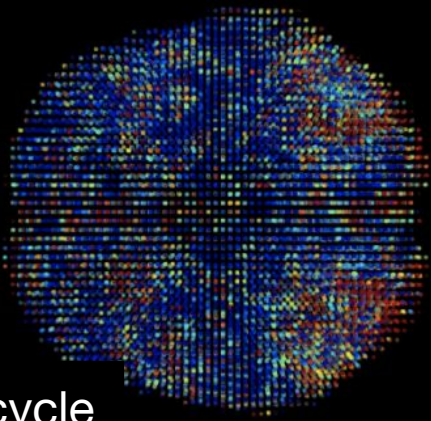
arm



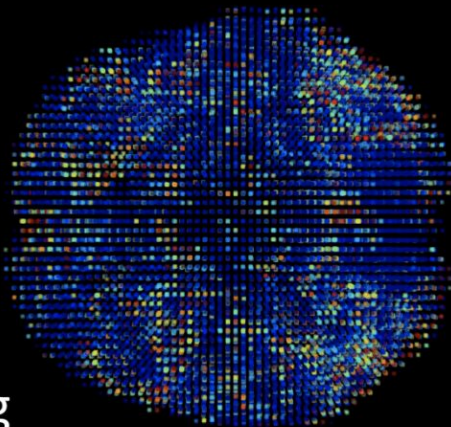
?



bicycle

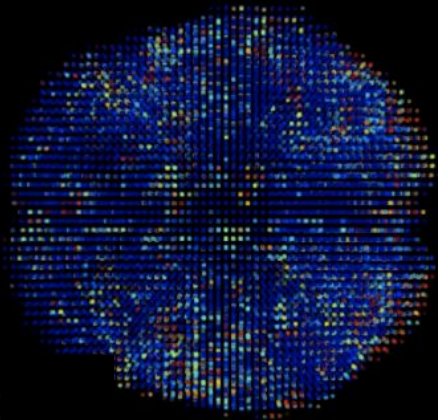


leg

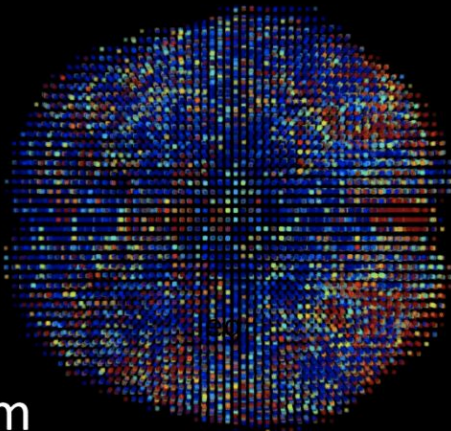


Data: Mitchell et al. 2008
Visualization: Samira Abnar

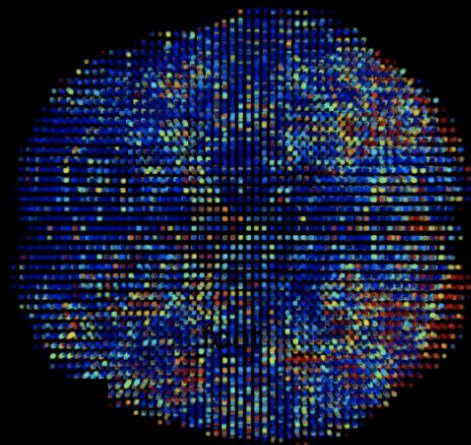
ant



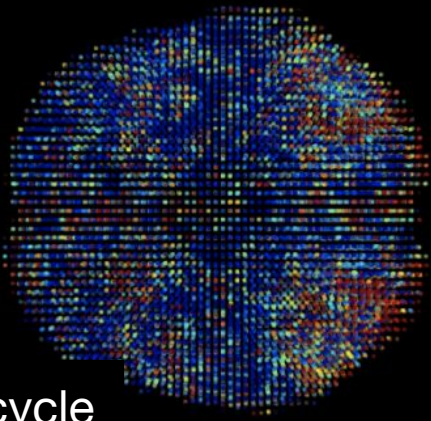
arm



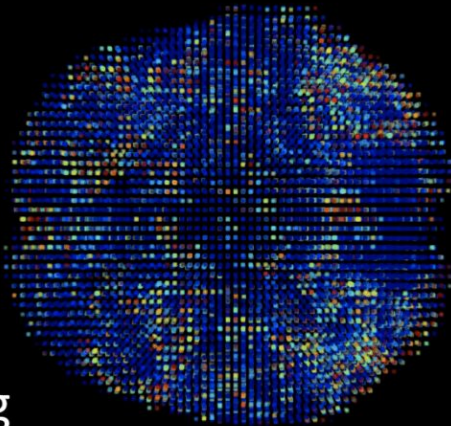
?



bicycle



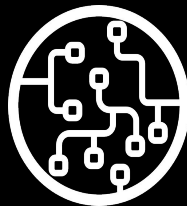
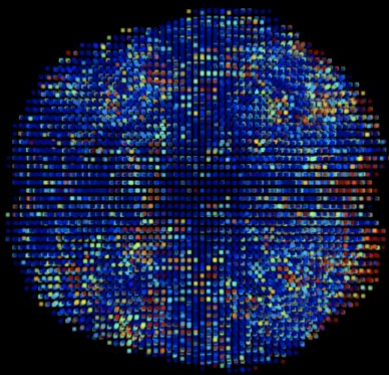
leg



Hand or foot?

Data: Mitchell et al. 2008
Visualization: Samira Abnar

LEARN MAPPING MODEL



cat: [2 6 8 1 3 1 0 0 ...]

WHAT'S NEXT?

WORDS IN
CONTEXT!



Harry had never believed he would meet a boy he hated more than Dudley.

[Wehbe et al. 2014]



A few weeks ago, a man I hardly know wrote me a really sweet love letter.

[Dehghani et al. 2017]



Alice was beginning to get very tired of sitting by her sister on the bank.

[Brennan et al. 2016]



DOES THAT WORK?

↪ WE DON'T REALLY NOW...

WHAT'S THE PROBLEM?

Researchers use

- different datasets

WHAT'S THE PROBLEM?

Researchers use

- different datasets
- different encoding models

WHAT'S THE PROBLEM?

Researchers use

- different datasets
- different encoding models
- different experimental parameters

WHAT'S THE PROBLEM?

Researchers use

- different datasets
- different encoding models
- different experimental parameters
- different evaluation metrics

WHAT'S THE PROBLEM?

Researchers use

- different datasets
- different encoding models
- different experimental parameters
- different evaluation metrics
- no comparison to baseline

WHAT'S THE PROBLEM?

Researchers use

- different datasets
- different encoding models
- different experimental parameters
- different evaluation metrics
- no comparison to baseline



AND THEY ARE OFTEN NOT
VERY TRANSPARENT
ABOUT THE DIFFERENCES:
NO DATA, NO CODE

OUR APPROACH

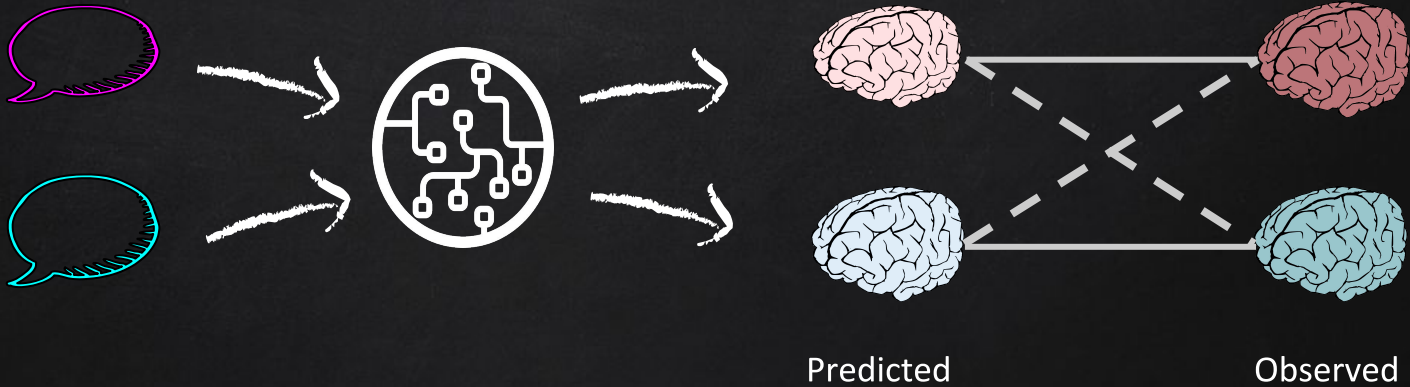
We use

- 4 datasets
- Constant encoding model
- Constant experimental parameters (if possible)
- Multiple evaluation metrics
- Comparison to a baseline with a “random language model”

We tried to standardize the procedure as much as possible and publish the experimental framework on github.


EVALUATION METHODS

1. Pairwise evaluation



EVALUATION METHODS

1. Pairwise evaluation
2. Voxelwise evaluation

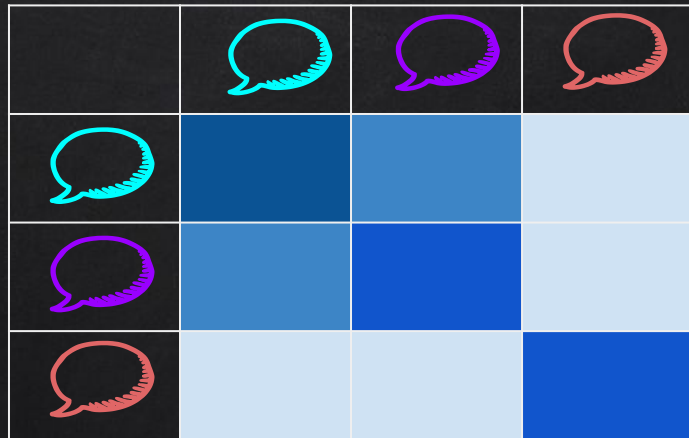
 **IDEA:** Not all voxels are related to language processing.
Evaluate the prediction for every voxel individually.

EVALUATION METHODS

1. Pairwise evaluation
2. Voxelwise evaluation
3. Representational similarity analysis
 - IDEA: directly compare relations between stimuli
 - No more prediction!

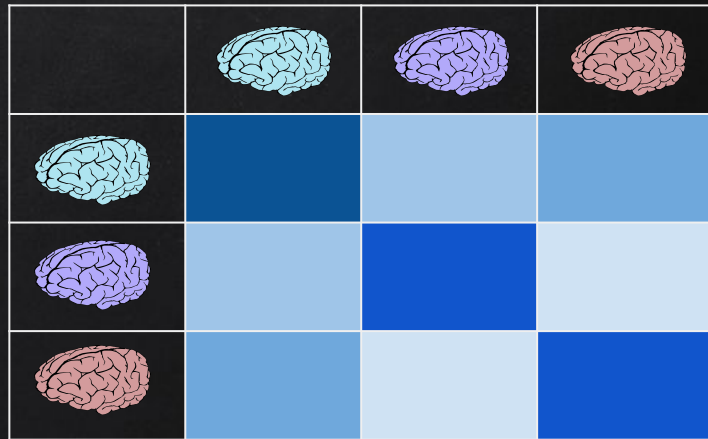
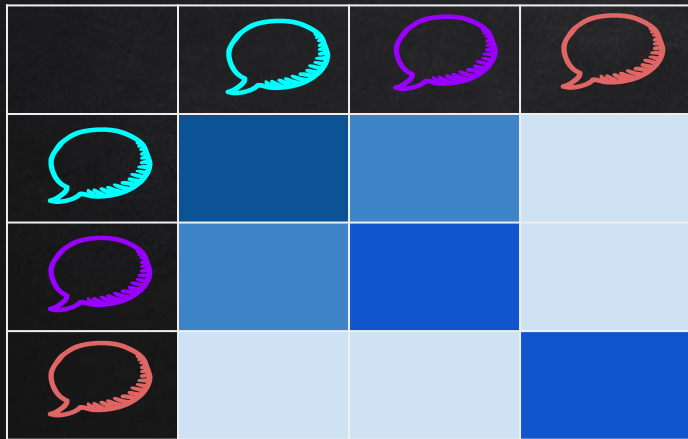
EVALUATION METHODS

1. Pairwise evaluation
2. Voxelwise evaluation
3. Representational similarity analysis



EVALUATION METHODS

1. Pairwise evaluation
2. Voxelwise evaluation
3. Representational similarity analysis



EVALUATION METHODS

1. Pairwise evaluation
2. Voxelwise evaluation
3. Representational similarity analysis

EACH METHOD CAN BE REALIZED WITH
DIFFERENT PARAMETERS.



WE COMPARE THEIR EFFECTS.

EXAMPLE PIPELINE

```
# Set the components
mitchell_reader = WordsReader(data_dir=mitchell_dir)
mapper = RegressionMapper()
stimuli_encoder = ElmoEncoder(save_dir)
random_encoder = RandomEncoder(save_dir)
```



TRY IT YOURSELF:

<https://github.com/beinborn/brain-lang>

```
# Try different language models
for encoder in [stimuli_encoder, random_encoder]:
```

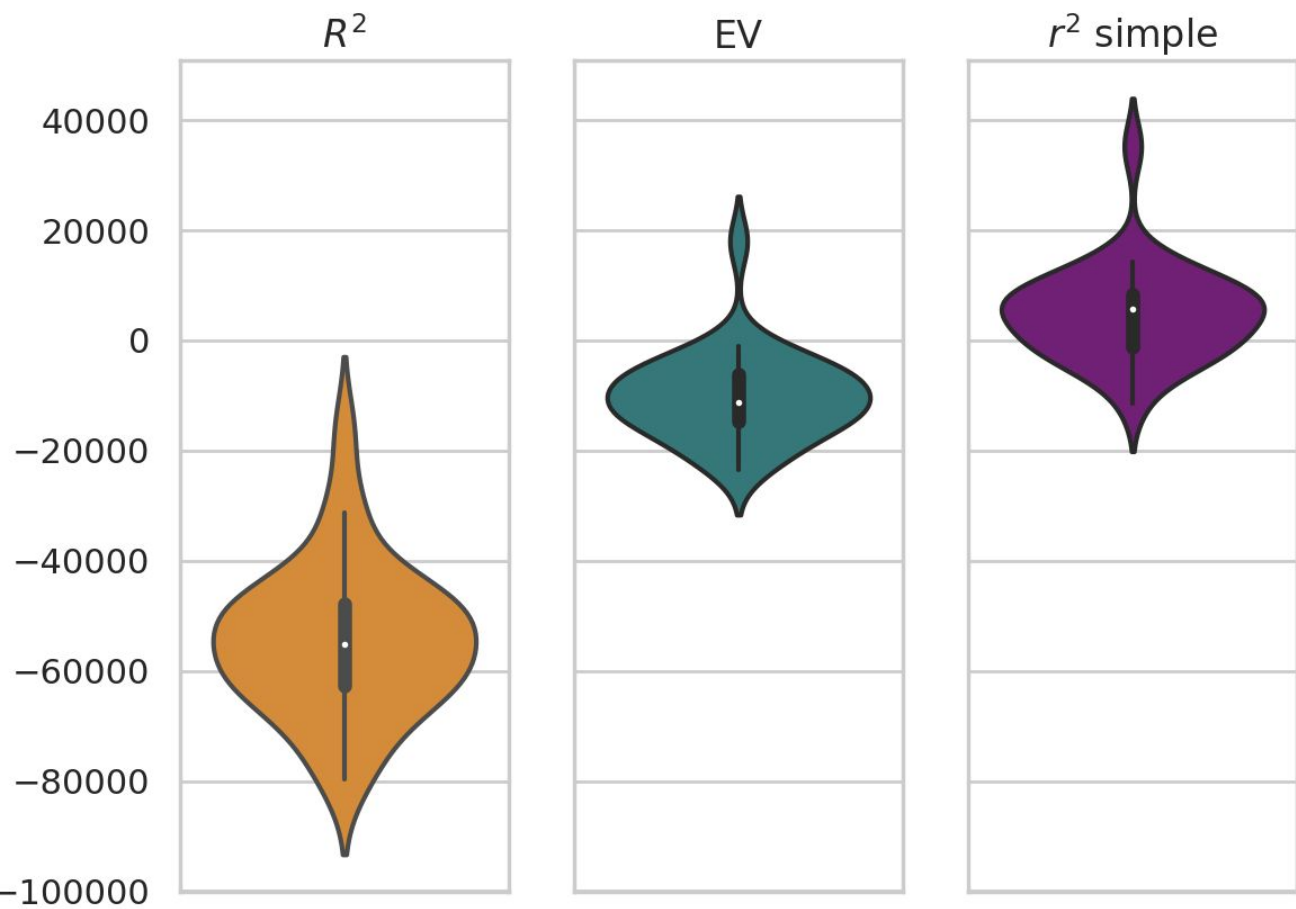
```
    # Set up the pipelines
    mitchell_pipeline_name = "Words" + encoder.__class__.__name__
    mitchell_pipeline = SingleInstancePipeline(mitchell_reader, encoder, mapper,
                                              mitchell_pipeline_name, save_dir=save_dir)
    mitchell_pipeline.voxel_selection = "none"

    mitchell_pipeline.pairwise_procedure("Mitchell_pairwise_noVS")
    mitchell_pipeline.run_standard_crossvalidation("Mitchell_CV_noVS")
    mitchell_pipeline.runRSA("Mitchell_RSA")
```



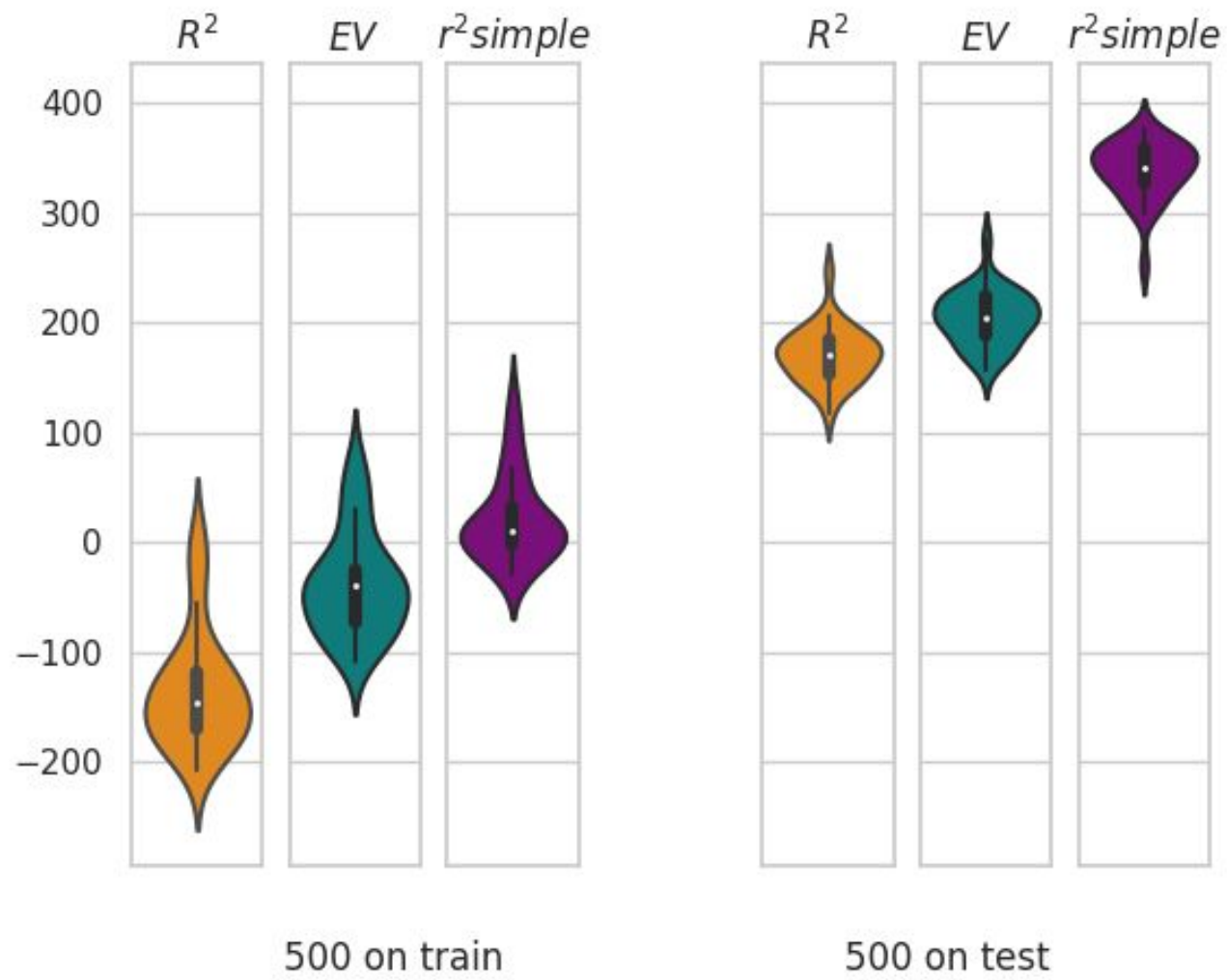
A PEAK
AT THE RESULTS?





Sum over all voxels

Voxel-wise results for
the data by Dehghani
et al. (2017)




Voxelwise results with
selected voxels for the
data by
Dehghani et al. (2017)

So?

Language-brain encoding is hard.

Many crucial design decisions:

preprocessing, language model, encoding parameters, ...

- 
1. Make these decisions transparent and reproducible.
 2. Analyze your hypothesis on several datasets with the same metric.
 3. Compare to reasonable baselines.
 4. Do not oversell your results! A tiny signal is already impressive.



QUESTIONS?



I.beinborn@uva.nl



BACK-UP SLIDES

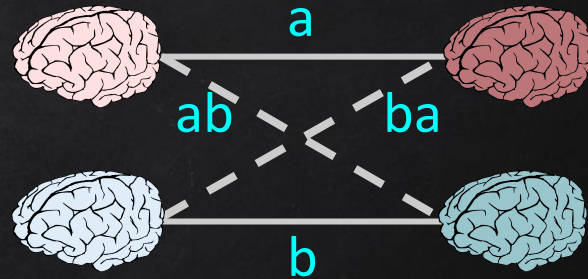
PAIRWISE EVALUATION

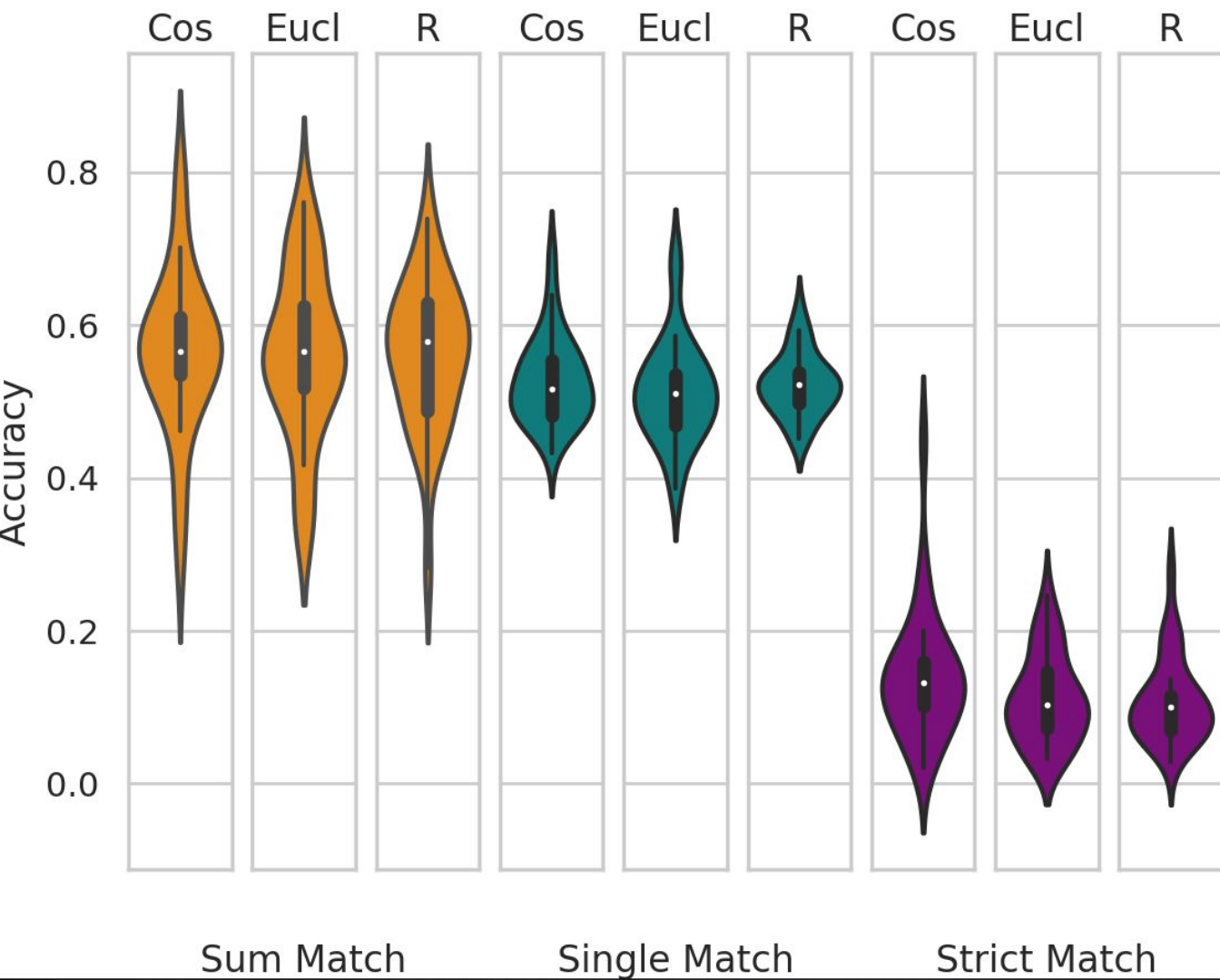
Match definition:

Sum: $(a + b) < (ab + ba)$

Single: $a < ab$

Strict: $(a < ab) \& (b < ba)$





Results for the data by Dehghani et al. (2017)

WHY?



ACROSS DATASETS?

		Encoding Model (Random LM)			
Match		WORDS	STORIES	ALICE	HARRY
	Sum	.67 (.54)	.57 (.53)	.54 (.53)	.50 (.49)
Cos	Single	.60 (.53)	.53 (.53)	.53 (.51)	.49 (.49)
	Strict	.26 (.13)	.14 (.02)	.28 (.27)	.25 (.24)



RSA RESULTS

	WORDS	STORIES	ALICE	HARRY
PEARSON	0.41	0.19	0.06	0.06
SPEARMAN	0.09	0.22	0.02	0.03



RSA RESULTS - RANDOM LM

	WORDS		STORIES		ALICE		HARRY	
PEARSON	0.41	0.44	0.19	0.21	0.06	0.02	0.06	0.03
SPEARMAN	0.09	0.05	0.08	0.09	0.03	0.01	0.00	0.01