

---

---

# Causal Attributions in Language Models

— Yanai Elazar —

---

---

*ETH Zürich, 23rd February, 2022*

# Hi There

Yanai Elazar, PhD student, Bar-Ilan University



With Yoav Goldberg



# Hi There



**NLP with Friends**

*By students, for students,  
where everyone is invited!*

[Home](#)

[Upcoming](#)

[Past](#)

[Calendar](#)

[Guidelines](#)

[FAQ](#)



## Welcome!

This is the home of **NLP with Friends**, an **online seminar series** made by students, for students, where everyone is invited!

## About the Seminar

We meet [Wednesdays](#) on a bi-weekly basis to talk about interesting work in NLP and related areas. The presenters are [students](#), who will talk about their **own work** (both ongoing and already published). Links are distributed through our [mailing list](#).

## About the Organizers



[Yanai Elazar](#) is a PhD candidate at Bar-Ilan University, where he works on neural representations, model analysis and missing elements. In his spare time he can be found nourishing flour-based organisms and converting them into bread.



[Abhilasha Ravichander](#) is a PhD candidate at Carnegie Mellon University, where she works on robust language understanding, including problems in interpretability, evaluation and computational reasoning. In her spare time she talks her plants into staying alive.



[Liz Salesky](#) is a PhD student at Johns Hopkins University, where she works on machine translation and computational linguistics. In her spare time she can be found biking to ice cream and bingeing Duolingo.



[Zeerak Waseem](#) is a PhD candidate at the University of Sheffield, where he works on abusive language detection and fairness in machine learning, and in his spare time he can be found napping.

# My Research

## Commonsense Reasoning

ACL19

### How Large Are Lions? Inducing Distributions over Quantitative Attributes

**Yanai Elazar\***  
Bar Ilan University  
yanaiela@gmail.com

**Abhijit Mahabal†**  
Pinterest  
amahabal@gmail.com

**Deepak Ramachandran**  
Google Research  
ramachandrand@google.com

**Tania Bedrax-Weiss**  
Google Research  
tbedrax@google.com

**Dan Roth**  
University of Pennsylvania  
danroth@seas.upenn.edu

### Back to Square One: Artifact Detection, Training and Commonsense Disentanglement in the Winograd Schema

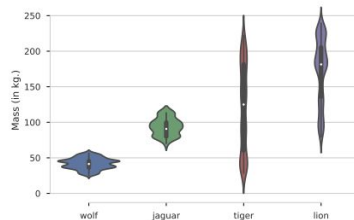
**Yanai Elazar<sup>1,2</sup> Hongming Zhang<sup>3,4</sup> Yoav Goldberg<sup>1,2</sup> Dan Roth<sup>4</sup>**

<sup>1</sup>Bar Ilan University, <sup>2</sup>Ai2, <sup>3</sup>HKUST, <sup>4</sup>UPenn

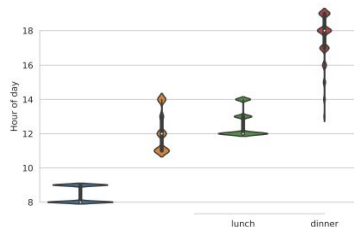
{yanaiela, yoav.goldberg}@gmail.com

hzhangal@cse.ust.hk, danroth@seas.upenn.edu

EMNLP21



(a) Mass distributions for multiple animals.



Setup	Example	Answer
<u>Original</u>		
twin-1	<i>The trophy doesn't fit into the brown suitcase because it is too <u>large</u>.</i>	🏆 trophy
twin-2	<i>The trophy doesn't fit into the brown suitcase because it is too <u>small</u>.</i>	👛 suitcase
<u>Baselines</u>		
<i>no-cands</i>	doesn't fit into because it is too <u>large</u> .	?
<i>part-sent</i>	because it is too <u>large</u> .	?
<u>Zero-shot</u>		
twin-1	<i>The trophy doesn't fit into the brown suitcase because <b>the trophy</b> is too [MASK].</i>	🟩 large
twin-2	<i>The trophy doesn't fit into the brown suitcase because <b>the brown suitcase</b> is too [MASK].</i>	🟦 small

# My Research

## Commonsense Reasoning V2: Missing Elements

Where  
for Nun

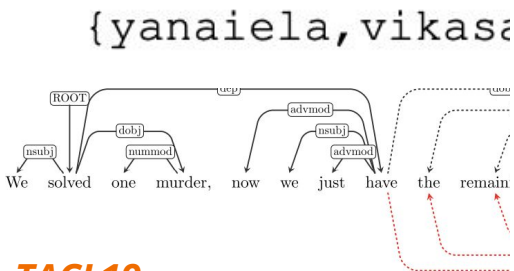
### Text-based NP Enrichment

Yanai Elazar\* Victoria Basmov\* Yoav Goldberg Reut Tsarfaty

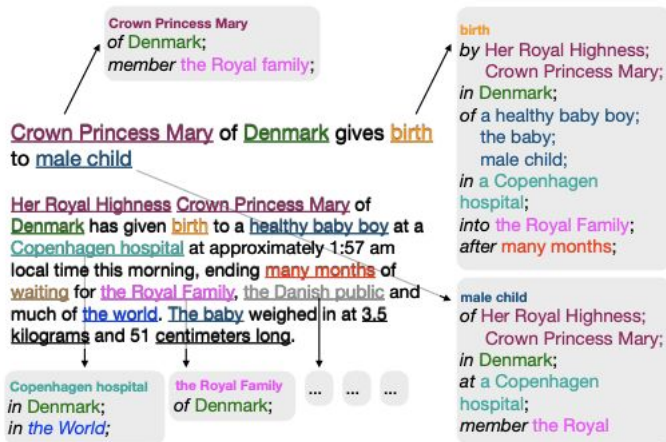
Computer Science Department, Bar Ilan University

Allen Institute for Artificial Intelligence

arfaty}@gmail.com



TAACL19



	Annotations
!t vow, she agrees	
egin kissing as the preacher	{officiating}, φ
ue — the ceremony.	
for max to finish	
ore asking him again. ~>	
for max to finish swallowing	ENT NEU CON
ore asking him again.	

# My Research

*and a bunch of BERTology...*

Do Language E

oLMpics

xikun

Alon Talm

Ian Te

Google R

ifttenney@g



F

**RT:**  
**n Multilingual BERT**

ar<sup>1,2</sup> Yoav Goldberg<sup>1,2</sup>  
lan University  
elligence  
oav.goldberg}@gmail.com

*(But we can talk about this later!)*

# My Research - Today

Causal Attribution in  
Language Models

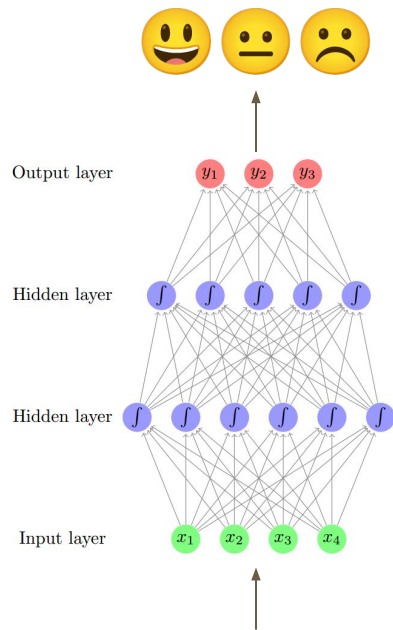
# Background

On NLP, Interpretation, Muppets, and  
Cramming a %&!\$ sentence into a single \$&!# vector



# The State of NLP (ML)

Output



Model

Input

*"Memories warm you up from the inside. But they also tear you apart."*

# The State of NLP (ML)

Output



Output layer

$y_1$   $y_2$   $y_3$

Hidden layer

$f$   $f$   $f$   $f$   $f$

Hidden layer

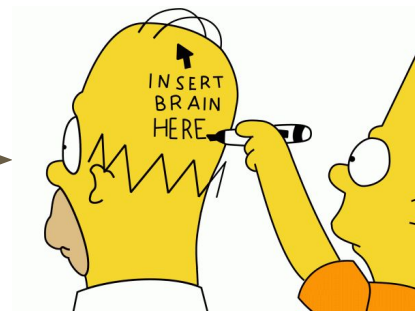
$f$   $f$   $f$   $f$   $f$   $f$

Input layer

$x_1$   $x_2$   $x_3$   $x_4$

Model

We train these models for some task



Input

*"Memories warm you up from the inside. But they also tear you apart."*

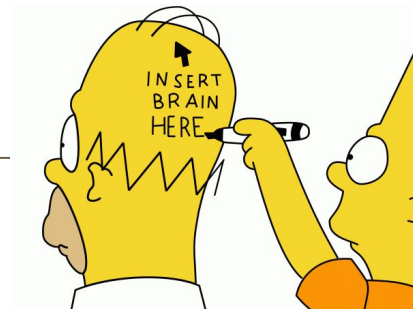
# The State of NLP (ML)

Output



And hope to get some  
"smart" model

Model



Input

*"Memories warm you up from the inside. But they also tear you apart."*

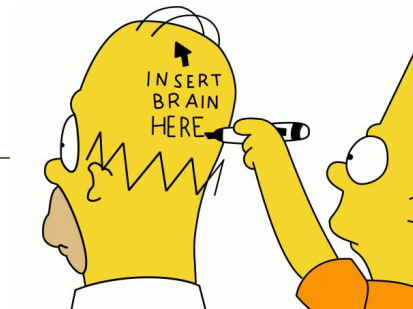
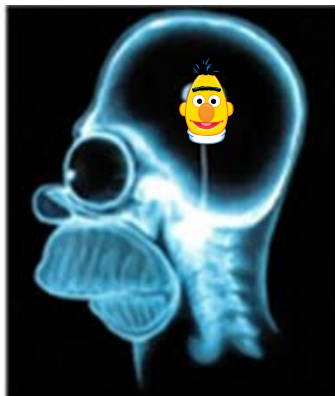
# The State of NLP (ML)

Output



And hope to get some  
"smart" model

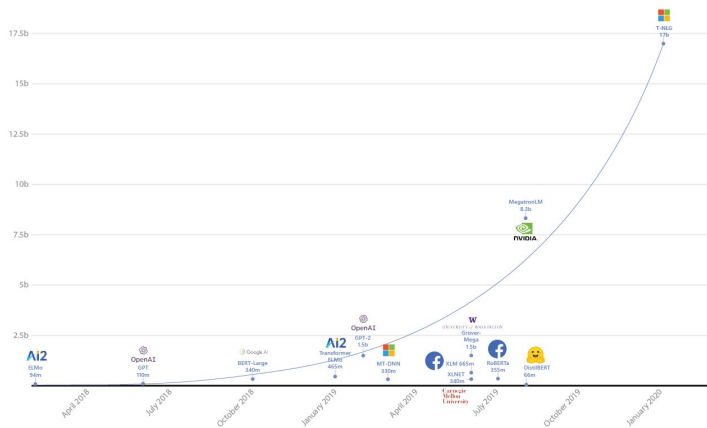
Model



Input

*"Memories warm you up from the inside. But they also tear you apart."*

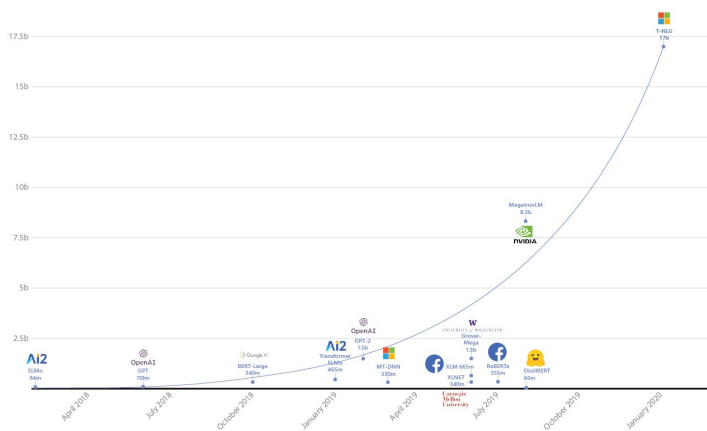
# The State of NLP: Sesame Street



Input

*"Memories warm you up from the inside. But they also tear you apart."*

# The State of NLP: Inside Sesame Street



Input

↑  
*"Memories warm you up from the inside. But they also tear you apart."*

# The State of NLP: Inside Sesame Street



*Input*

*"Memories warm you up from the inside. But they also tear you apart."*

# Opening the BlackBox

*you cannot cram the meaning of a whole sentence into a single vector*

*-- Ray Mooney*



# Opening the BlackBox

*you cannot cram the meaning of a whole sentence into a single vector*

*-- Ray Mooney*

- So what can be crammed into that?

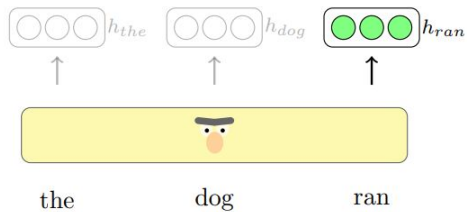
# Probing

Popular approach

# Probing

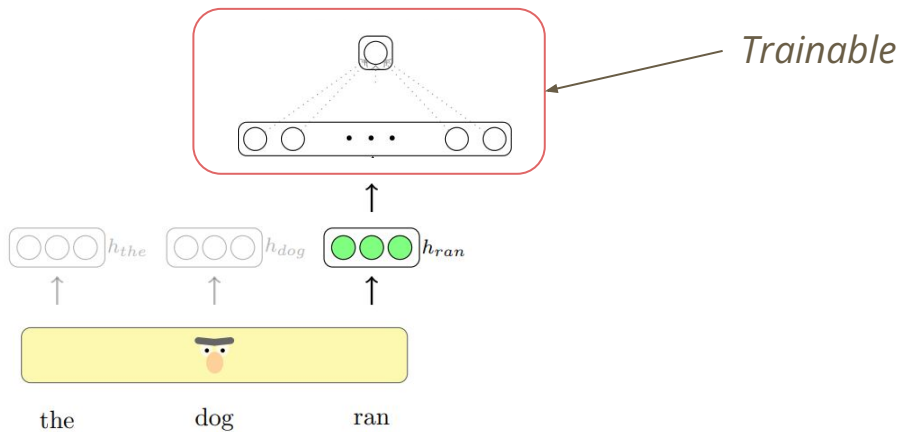
Popular approach

- Encode some text and retrieve its representation



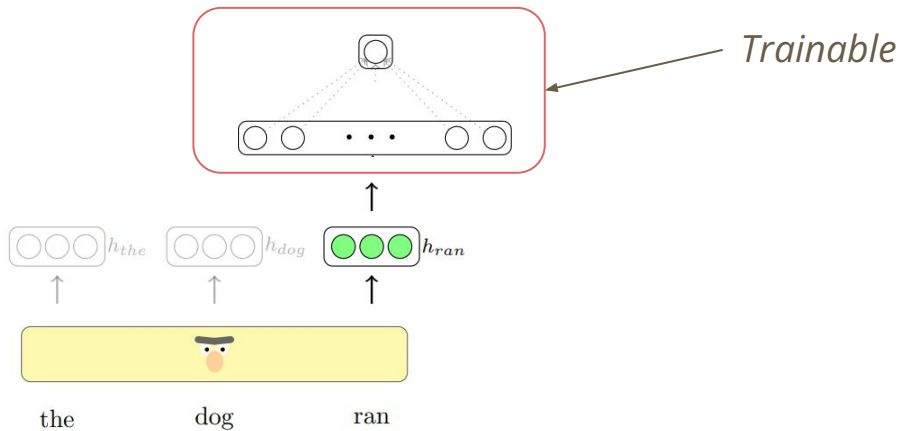
# Probing

- Encode some text and retrieve its representation
- Train a classifier to predict a property of interest



# Probing

- Encode some text and retrieve its representation
- Train a classifier to predict a property of interest
- High performance is interpreted as the encoding of the property



# People Probe for...

- Sentence Length
- Word Order
- Tense

*Adi et al., 2016, Conneau et al., 2018, Hewitt and Manning, 2019, Tenney et al., 2019, Chi et al., 2020*

# People Probe for...

- Sentence Length
- Word Order
- Tense
- POS
- Tree depth
- Entities
- Coref.
- ...

*Adi et al., 2016, Conneau et al., 2018, Hewitt and Manning, 2019, Tenney et al., 2019, Chi et al., 2020*

# What's Wrong with Probing?



# Probing - The Problem

Probing answers:

“What is *encoded* in the representation?”

But the interesting question is:

“What is *being used* for prediction?”



# Probing - The Problem

Probing answers:

“What is *encoded* in the representation?”

But the interesting question is:

“What is *being used* for prediction?”

Which are very **different** questions!



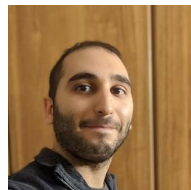
# Part I

## Amnesic Probing: Behavioral Explanation with Amnesic Counterfactuals

**Yanai Elazar<sup>1,2</sup> Shauli Ravfogel<sup>1,2</sup> Alon Jacovi<sup>1</sup> Yoav Goldberg<sup>1,2</sup>**

<sup>1</sup>Computer Science Department, Bar Ilan University

<sup>2</sup>Allen Institute for Artificial Intelligence



**Our Solution: *Amnesic Probing*, A Behavioral Probe**



# Amnesic Probing: A Behavioral Probe

- Interpretability tool, which allows to:
  - Answer scientific questions (e.g. *does an LM use POS information?*)
  - Answer applicative questions (e.g. *does the model use gender for making a decision?*)

Probing answers:

“What is *encoded* in the representation?”

But the interesting question is:

“What is *being used* for prediction?”



Probing answers:

“What is *encoded* in the representation?”

Probing

But the interesting question is:

“What is *being used* for prediction?”



Probing answers:

“What is *encoded* in the representation?”

Probing

But the interesting question is:

“What is *being used* for prediction?”

Amnesic Probing





# The Intuition: Counterfactuals

**What would the model predict *without* a given concept?**

# Amnesic Probing: The Intuition

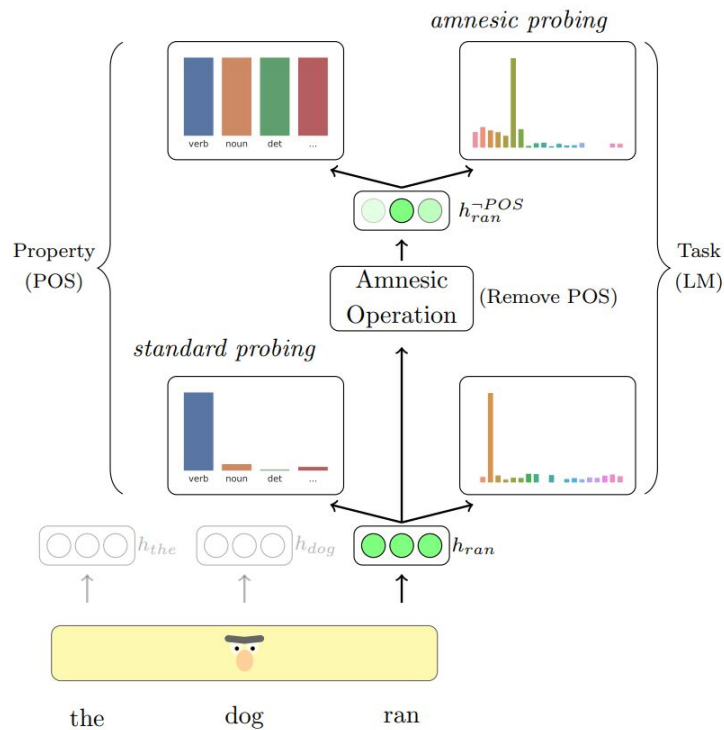
- Counterfactuals (or *ablation* on a trained model):
  - Remove a certain component, property
  - Measure how it affects the results
- Since it is hard to intervene on the input text...  
**...we intervene on the representation**

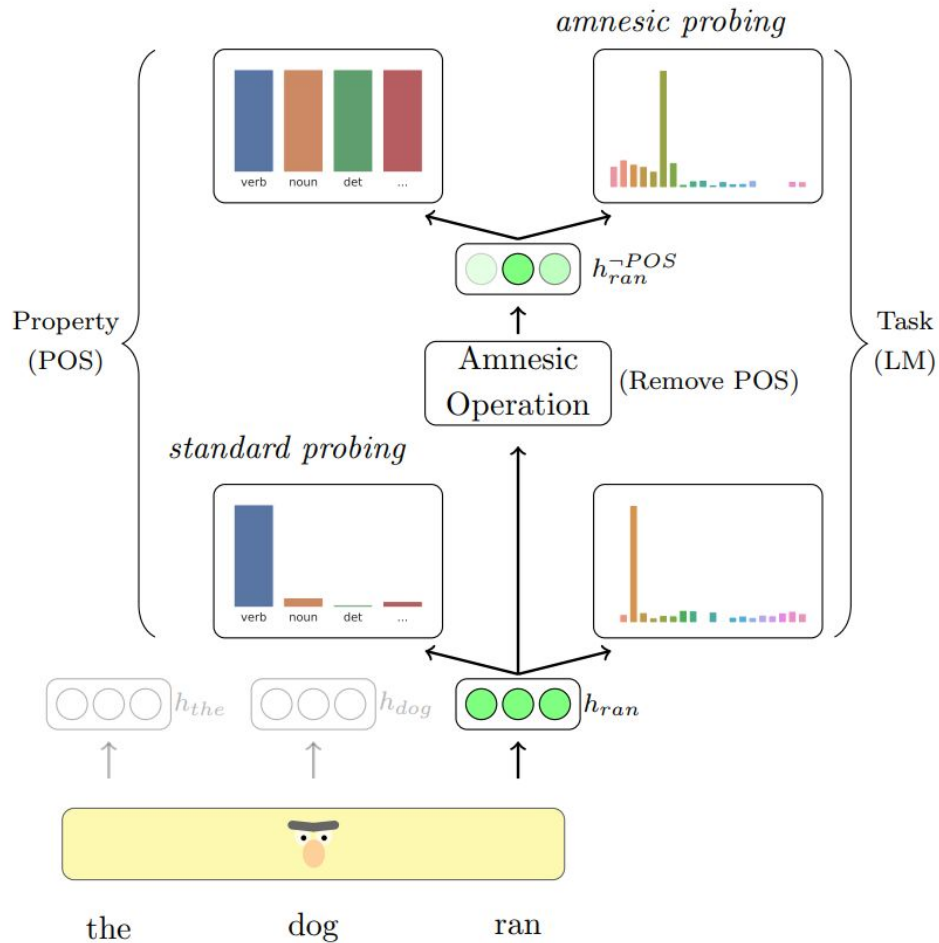


# Amnesic Probing: The Intuition

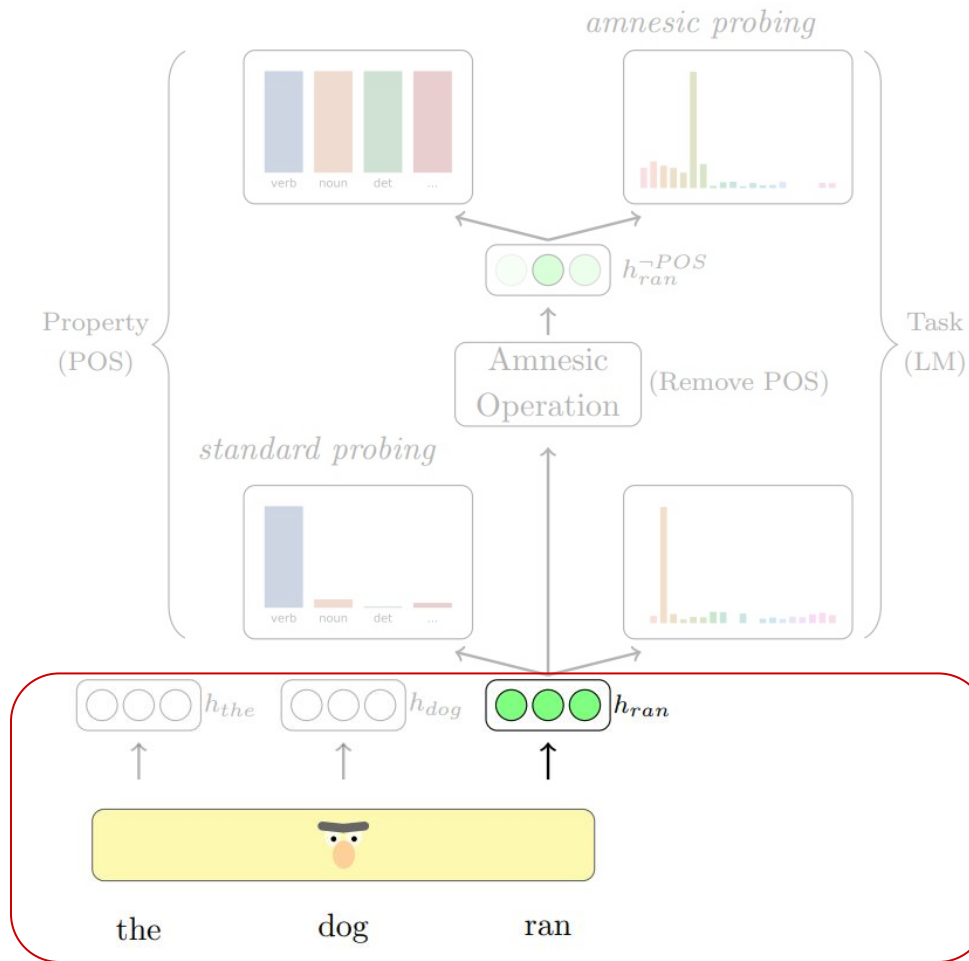
- We remove a feature from the representation (e.g. **remove POS information**)
- Does the model change its behavior?
  
- Yes:
  - The model **uses** this information for its predictions
- No:
  - The model **does not use** this information for its predictions

# Amnesic Probing: Overview

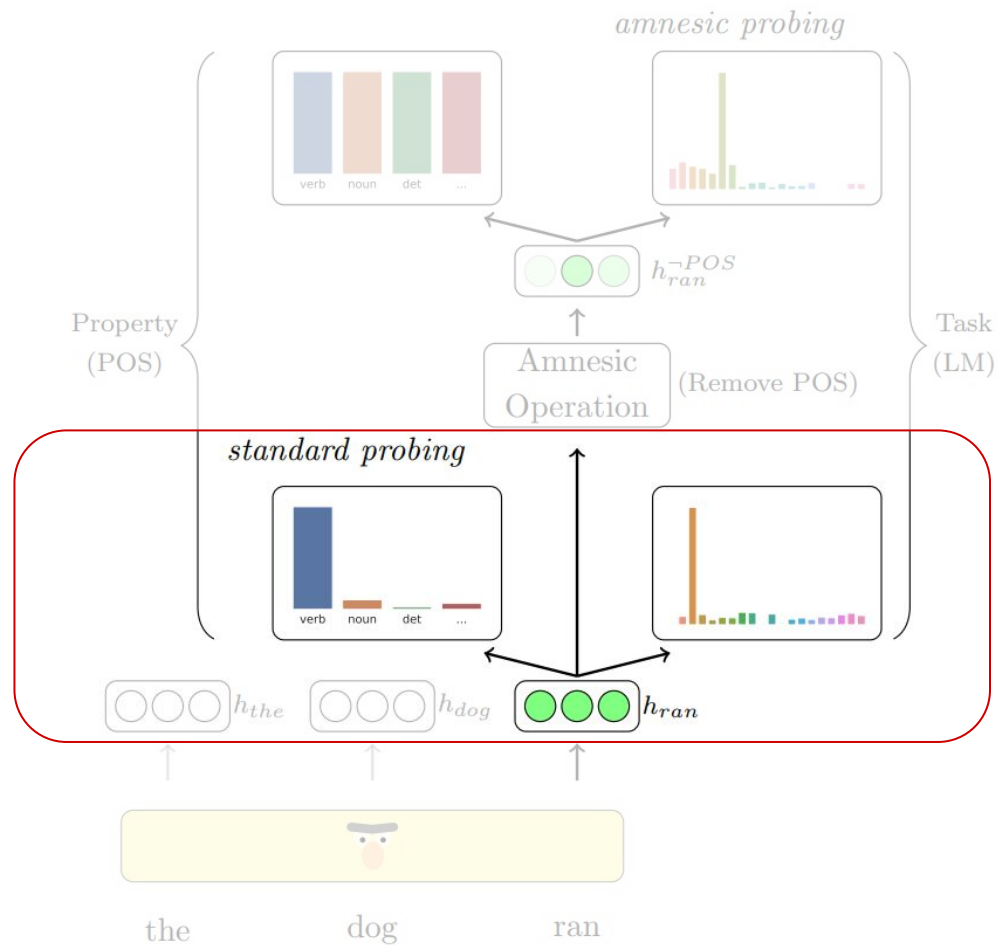




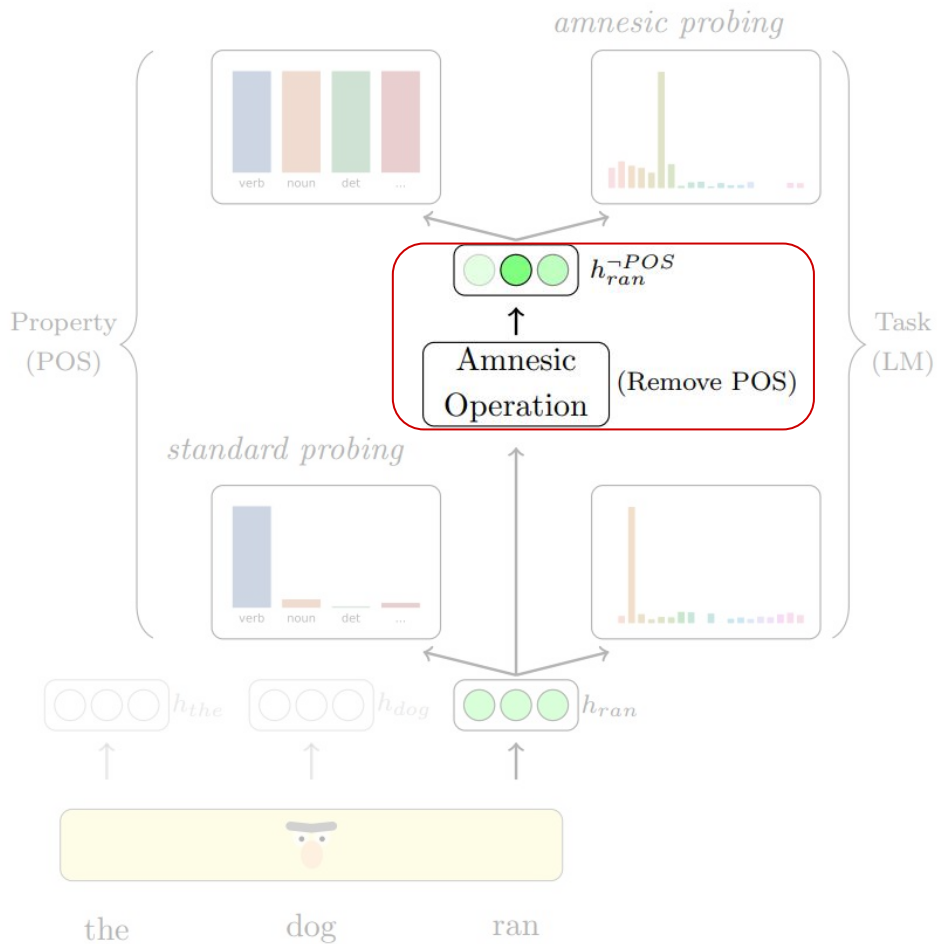
# 1. Encode



1. Encode
2. Probe

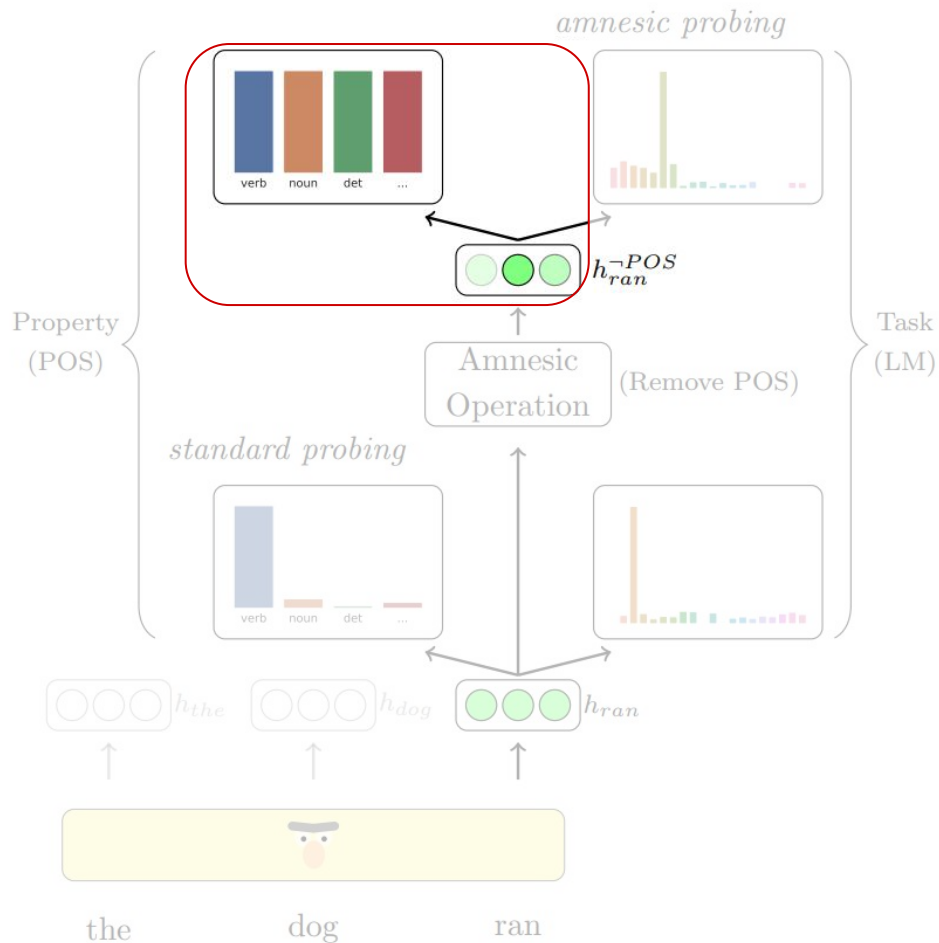


1. Encode
2. Probe
3. Amnesia

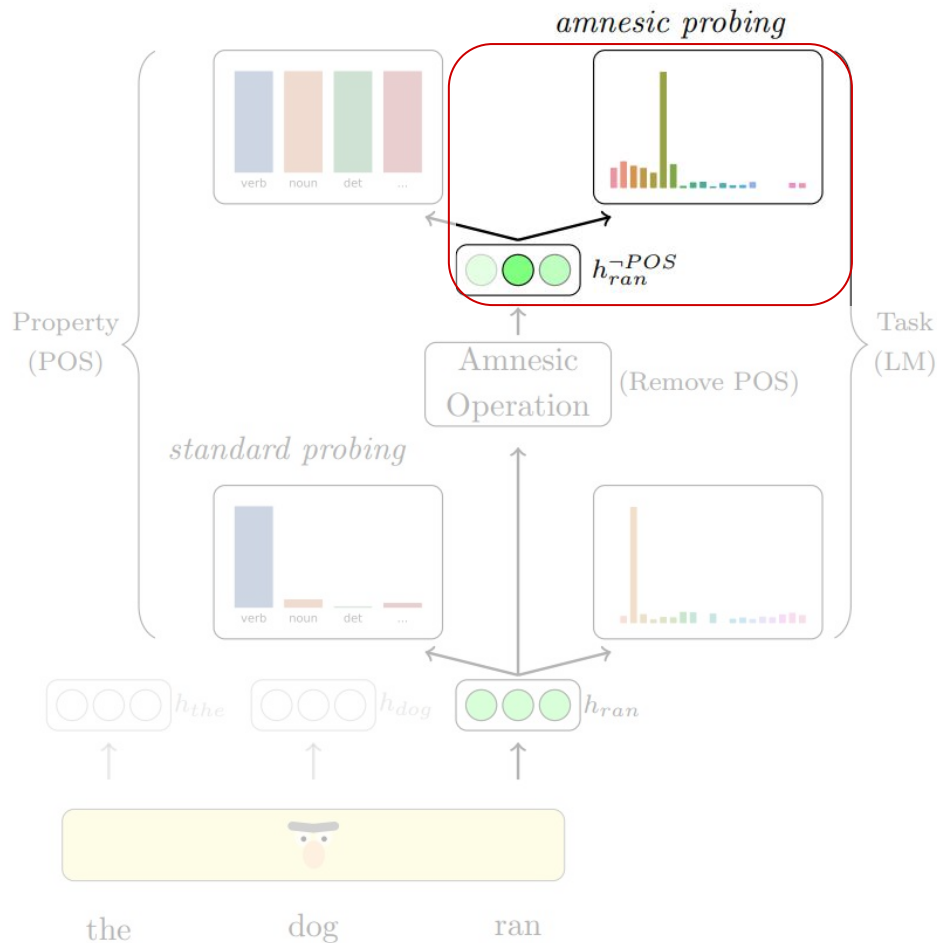




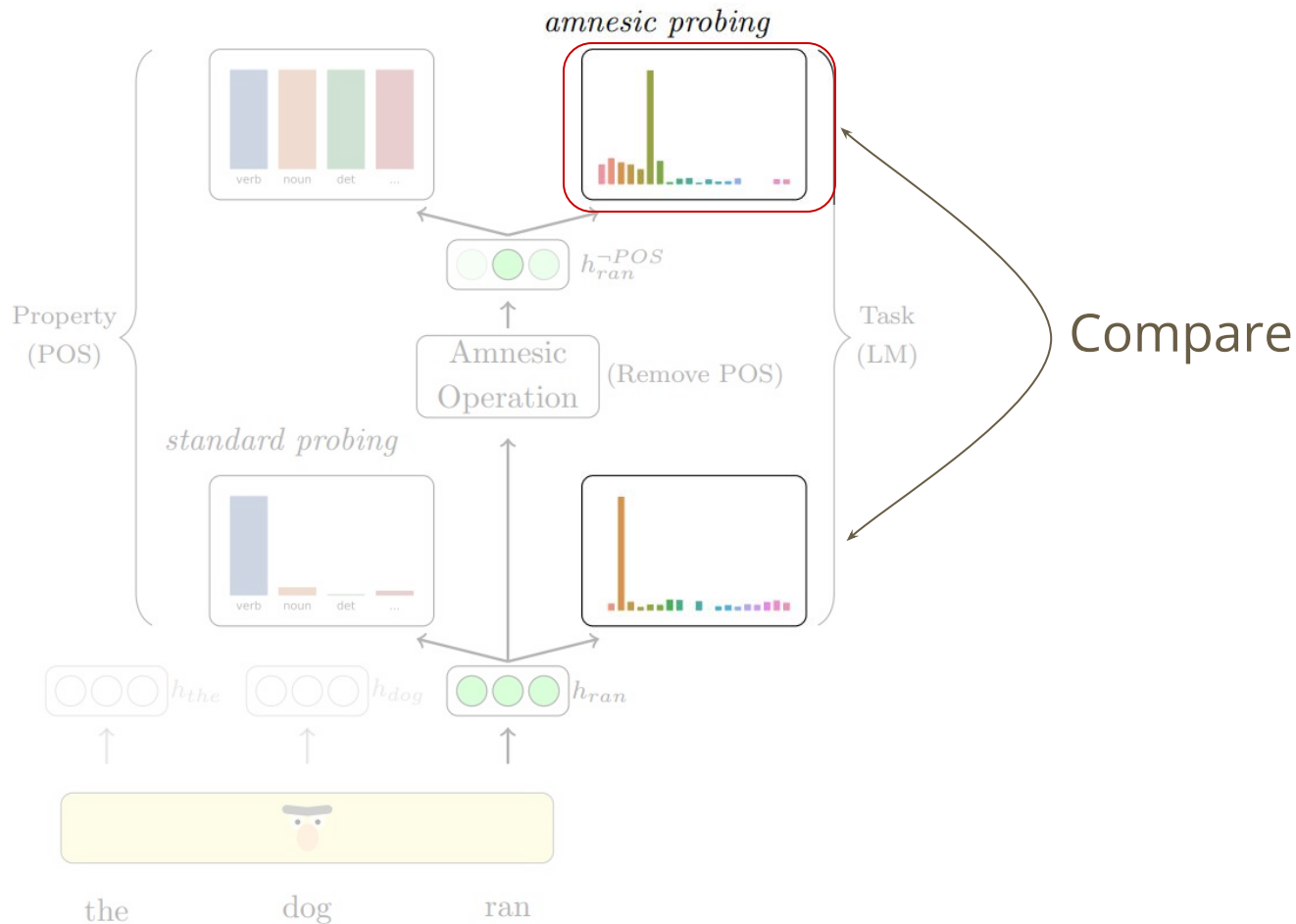
1. Encode
  2. Probe
  3. Amnesia
- 3.1. Verify



1. Encode
2. Probe
3. Amnesia
  - 3.1. Verify
4. *Amnesic Probing*



1. Encode
2. Probe
3. Amnesia
  - 3.1. Verify
4. *Amnesic Probing*



# The Amnesic Operation

# Amnesic Probing: The Amnesia

One option: Adversarial Training

## **Adversarial Removal of Demographic Attributes from Text Data**

**Yanai Elazar<sup>†</sup>** and **Yoav Goldberg<sup>†\*</sup>**

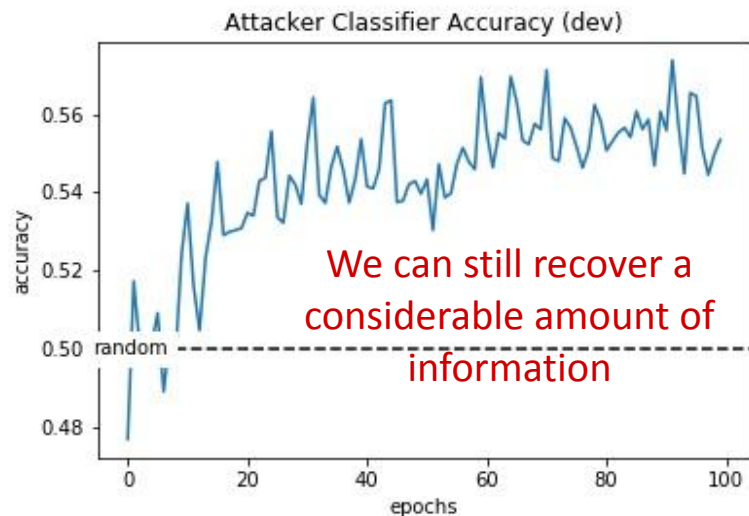
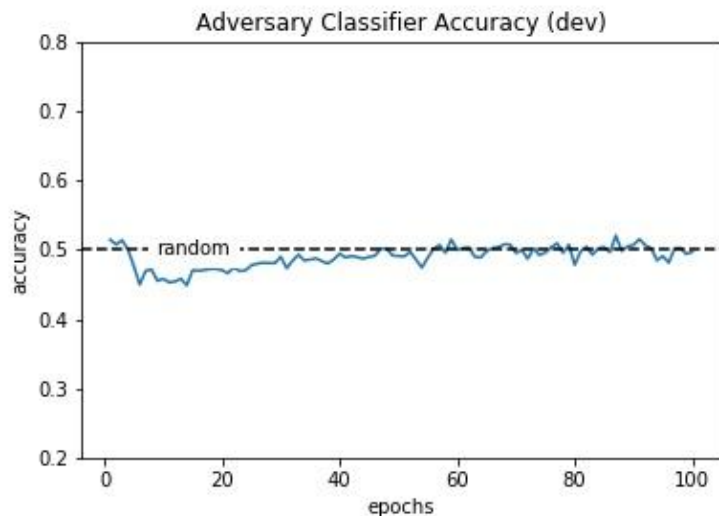
<sup>†</sup>Computer Science Department, Bar-Ilan University, Israel

\*Allen Institute for Artificial Intelligence

{yanaiela, yoav.goldberg}@gmail.com

# Amnesic Probing: The Amnesia

One option: Adversarial Training

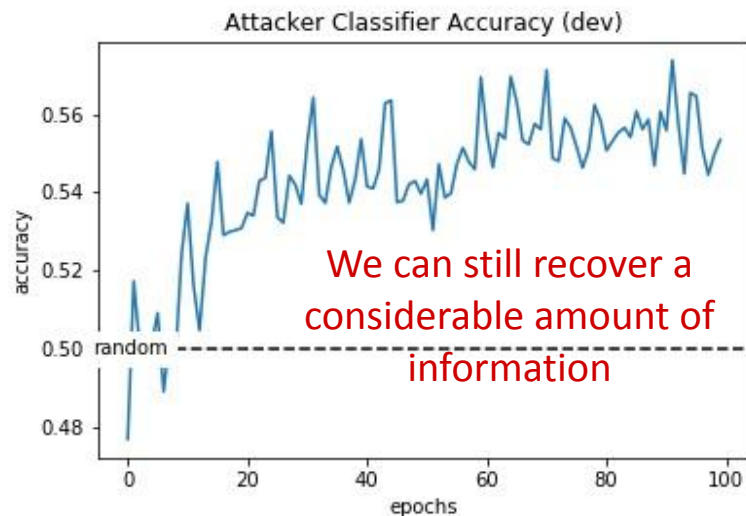


# Amnesic Probing: The Amnesia

One option: Adversarial Training

But also:

- Slow & unstable training
- Is it the same model afterwards?



# Amnesic Probing: The Amnesia

## Null It Out: Guarding Protected Attributes by Iterative Nullspace Projection

**Shauli Ravfogel<sup>1,2</sup>**

**Yanai Elazar<sup>1,2</sup>**

**Hila Gonen<sup>1</sup>**

**Michael Twiton<sup>3</sup>**

**Yoav Goldberg<sup>1,2</sup>**

<sup>1</sup>Computer Science Department, Bar Ilan University

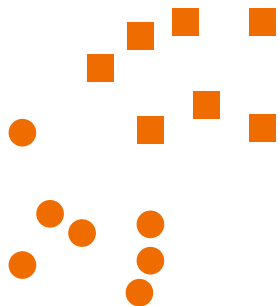
<sup>2</sup>Allen Institute for Artificial Intelligence

<sup>3</sup>Independent researcher



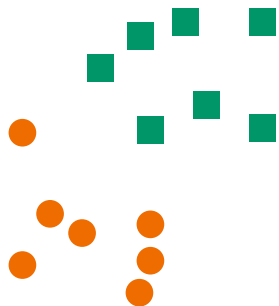
# Amnesic Operation: Using INLP

- An algorithm for removing linear information from deep networks



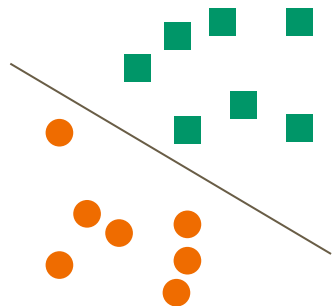
# Amnesic Operation: Using INLP

- An algorithm for removing linear information from deep networks
- Receives representations and labels, and returns a function



# Amnesic Operation: Using INLP

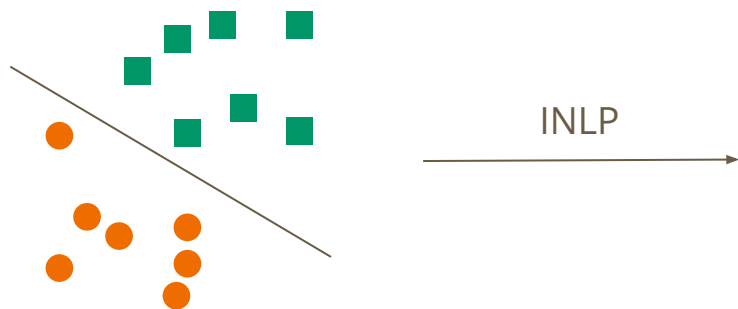
- An algorithm for removing linear information from deep networks
- Receives representations and labels, and returns a function



*Ravfogel et al., 2020*

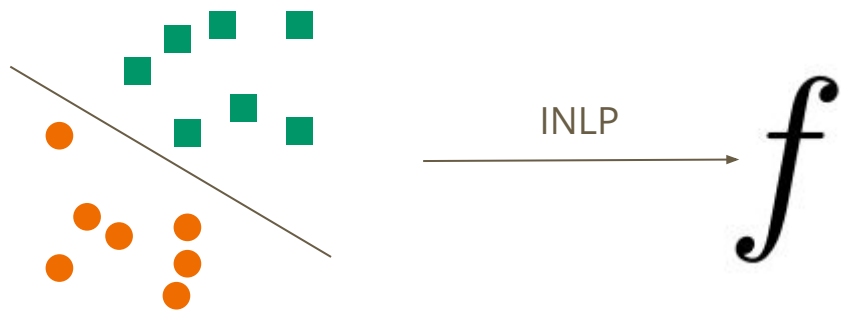
# Amnesic Operation: Using INLP

- An algorithm for removing linear information from deep networks
- Receives representations and labels, and returns a function



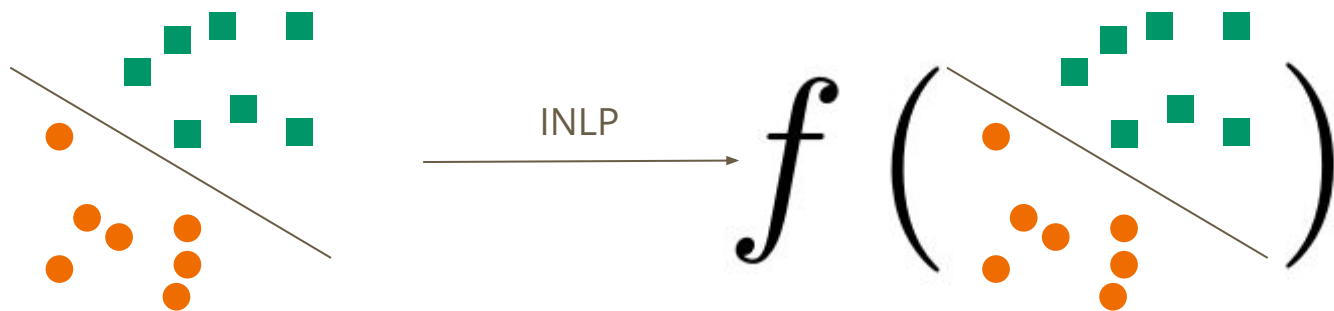
# Amnesic Operation: Using INLP

- An algorithm for removing linear information from deep networks
- Receives representations and labels, and returns a function



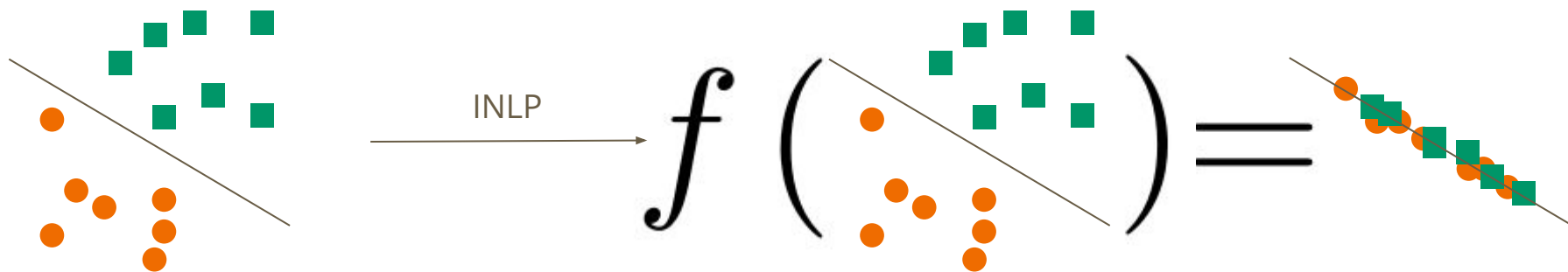
# Amnesic Operation: Using INLP

- An algorithm for removing linear information from deep networks
- Receives representations and labels, and returns a function
- When applied to vectors, any linear model cannot predict the labels



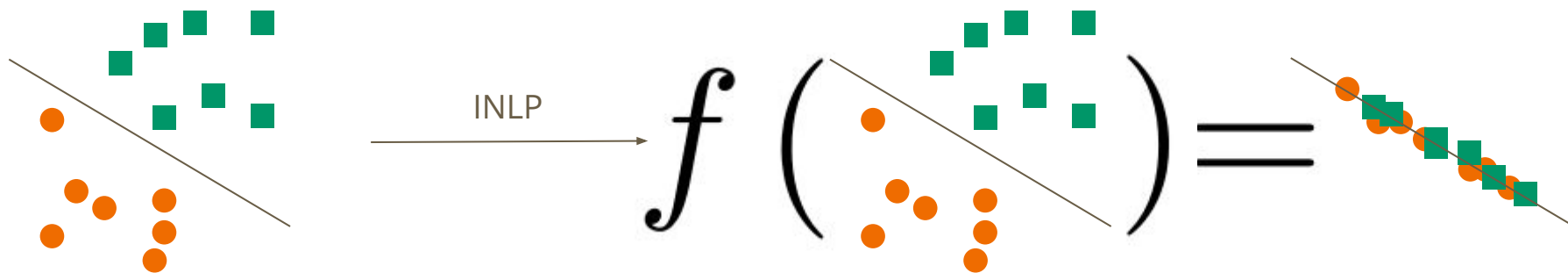
# Amnesic Operation: Using INLP

- An algorithm for removing linear information from deep networks
- Receives representations and labels, and returns a function
- When applied to vectors, any linear model cannot predict the labels



# Amnesic Operation: Using INLP

- An algorithm for removing linear information from deep networks
- Receives representations and labels, and returns a function
- When applied to vectors, any linear model cannot predict the labels



*Ravfogel et al., 2020*

(\*) We use INLP in this work, but this is a component that can be replaced with a future (non-linear) alternative

*Feder et al. 2021*



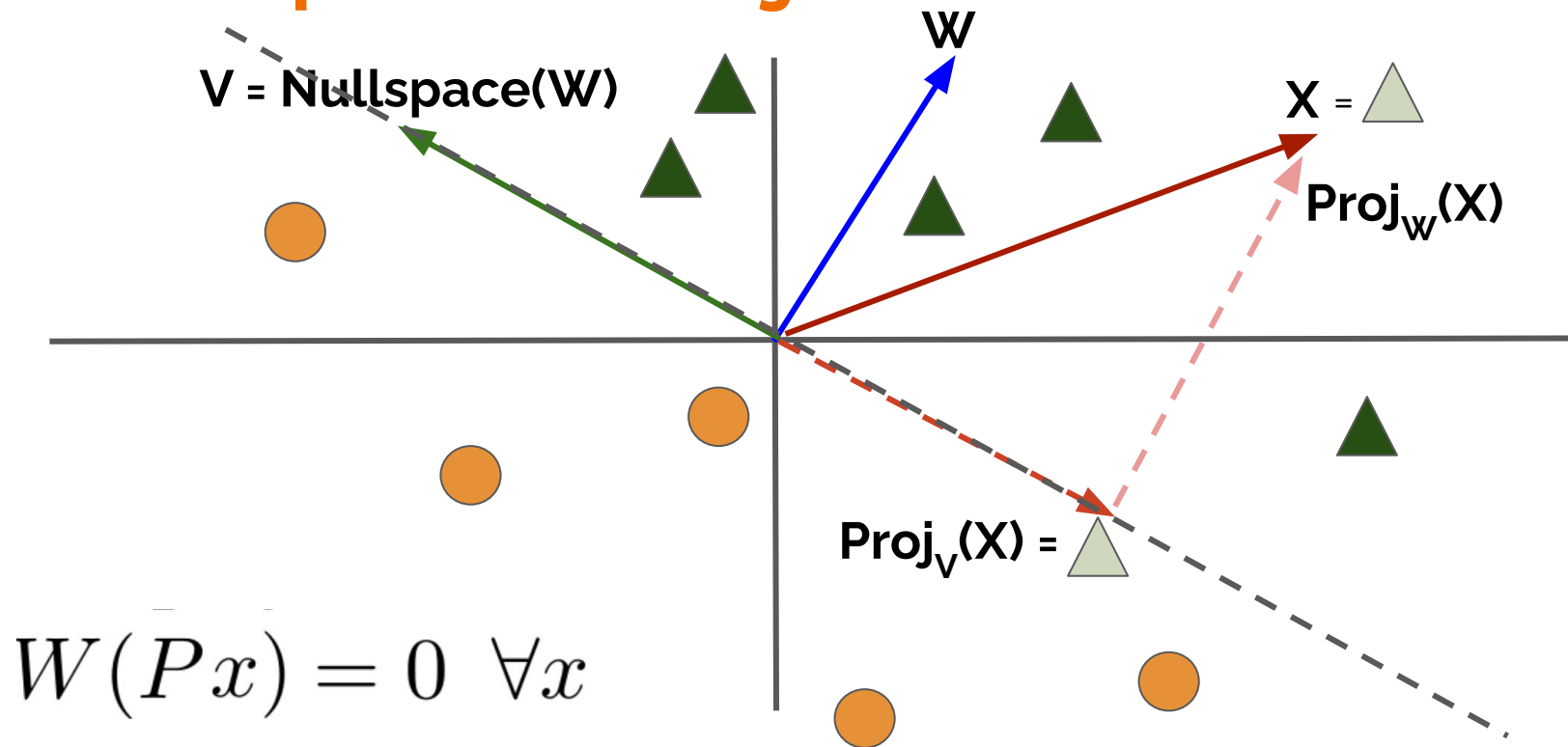
# Amnesic Operation: Using INLP

**INLP:** Iterative Nullspace Projection

- Find a projection matrix  $P$ , which projects into the nullspace

$$N(W) = \{x \mid Wx = 0\}$$

# Amnesic Operation: Using INLP



# Amnesic Operation: Using INLP

- Each projection only removes a single direction
- Therefore the “iterative” part:
- We repeat this process until convergence

# Amnesic Operation: Using INLP

- Debiasing applications (Ravfogel et al., 2020)

Check it out!

		BoW	FastText	BERT
Accuracy (profession)	Original	78.2	78.1	80.9
	+INLP	80.1	73.0	75.2
$GAP_{male}^{TPR,RMS}$	Original	0.203	0.184	0.184
	+INLP	0.124	0.089	0.095

Table 2: Fair classification on the Biographies corpus.

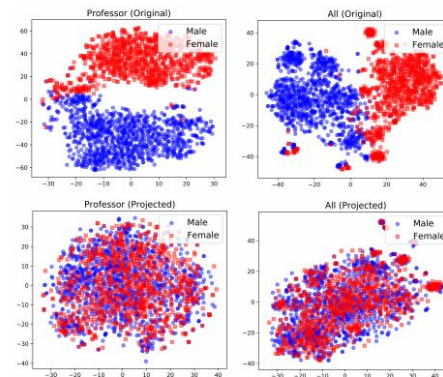


Figure 3: t-SNE projection of BERT representations for the profession “professor” (left) and for a random sample of all professions (right), before and after the projection.

# Amnesic Probing: Setup

- Start with a trained model
- Encode and obtain the representations
- Choose properties/features of interest
- Remove them
- Measure the difference (behavioral!), via:
  - Accuracy (of predicting the “right” label)

# Verifying that the Amnesic Operation Works

# Amnesic Probing: Controls

- Did the amnesic operation remove too little?
- Did the amnesic operation remove too much?



**TOO  
LITTLE**



**TOO MUCH**



**JUST  
RIGHT**

# Amnesic Probing: Controls

- Did the amnesic operation remove too little?
- Did the amnesic operation remove too much?
- Control over Information
  - Removing random features



**TOO  
LITTLE**



**TOO MUCH**



**JUST  
RIGHT**

imgflip.com



# Amnesic Probing: Controls

- Did the amnesic operation remove too little?
- Did the amnesic operation remove too much?
  
- Control over Information
  - Removing random features
- Control over Selectivity
  - Add back the “real” features, and retrain



**TOO  
LITTLE**



**TOO MUCH**



**JUST  
RIGHT**



# Amnesic Probing: Controls

- Did the amnesic operation remove too little?
- Did the amnesic operation remove too much?
  
- Control over Information
  - Removing random features
- Control over Selectivity
  - Add back the “real” features, and retrain
- Hopefully we’ll be here →



**TOO  
LITTLE**

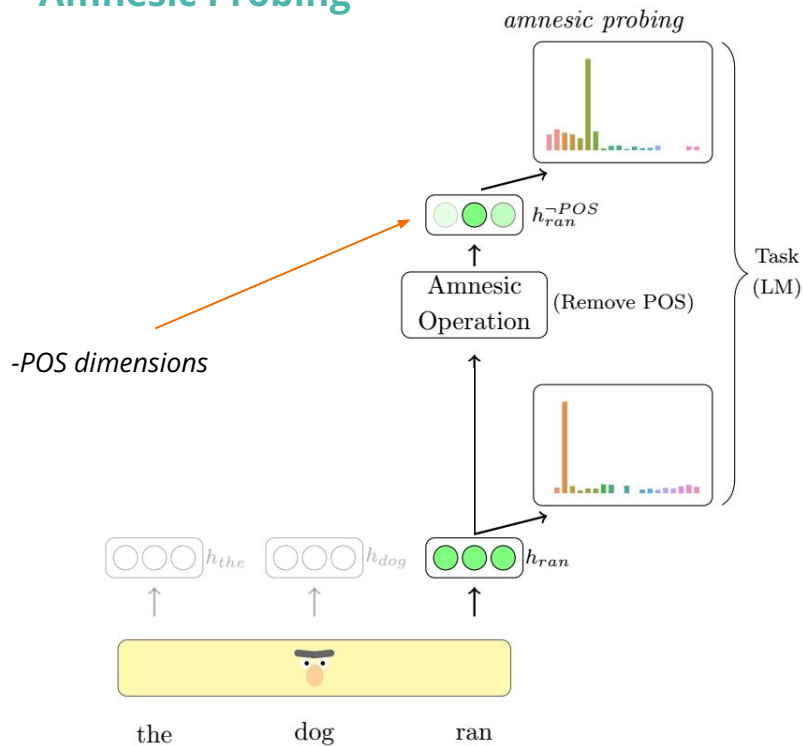


**TOO MUCH**

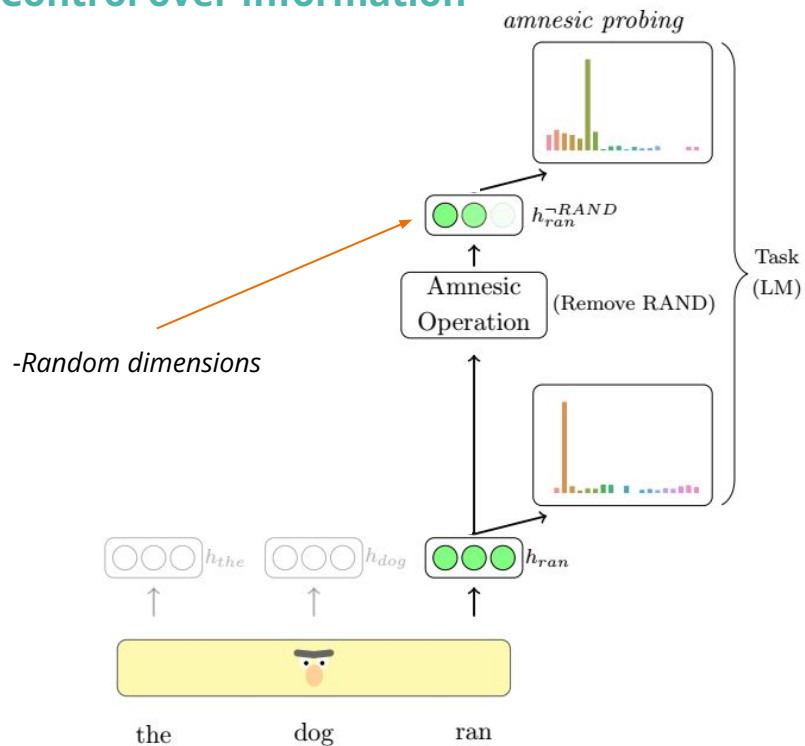


**JUST  
RIGHT**

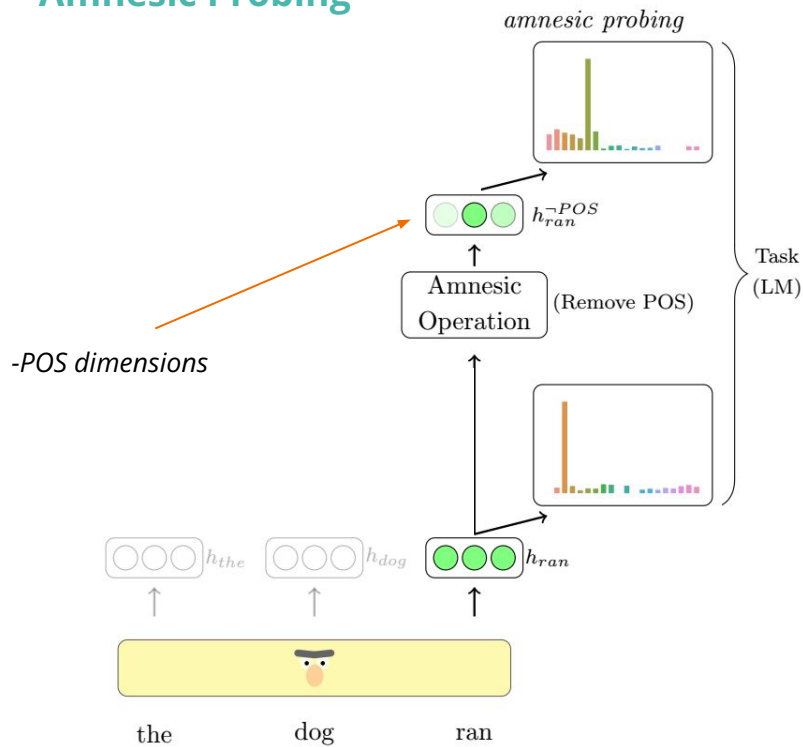
## Amnesic Probing



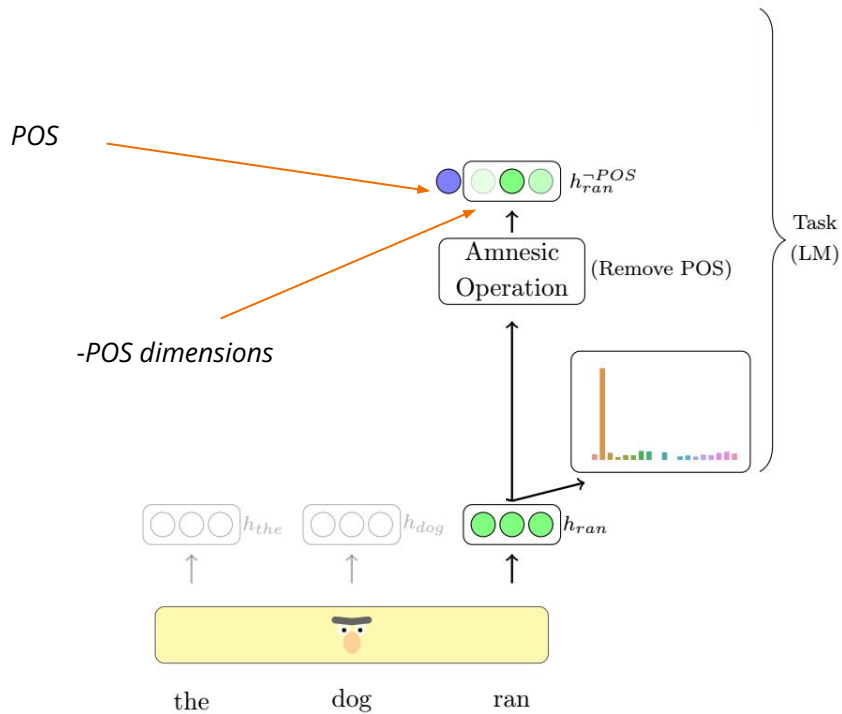
## Control over Information



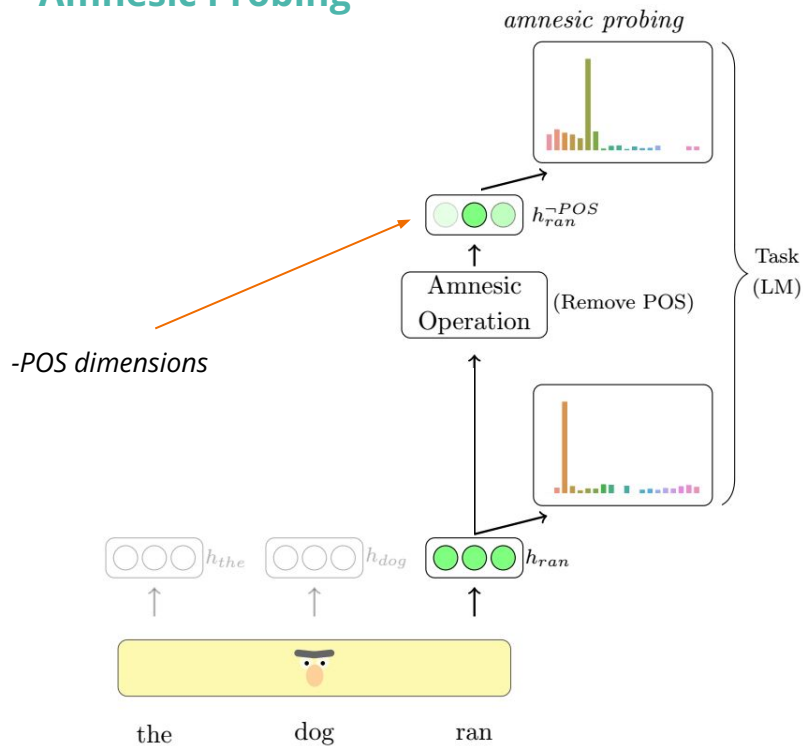
## Amnesic Probing



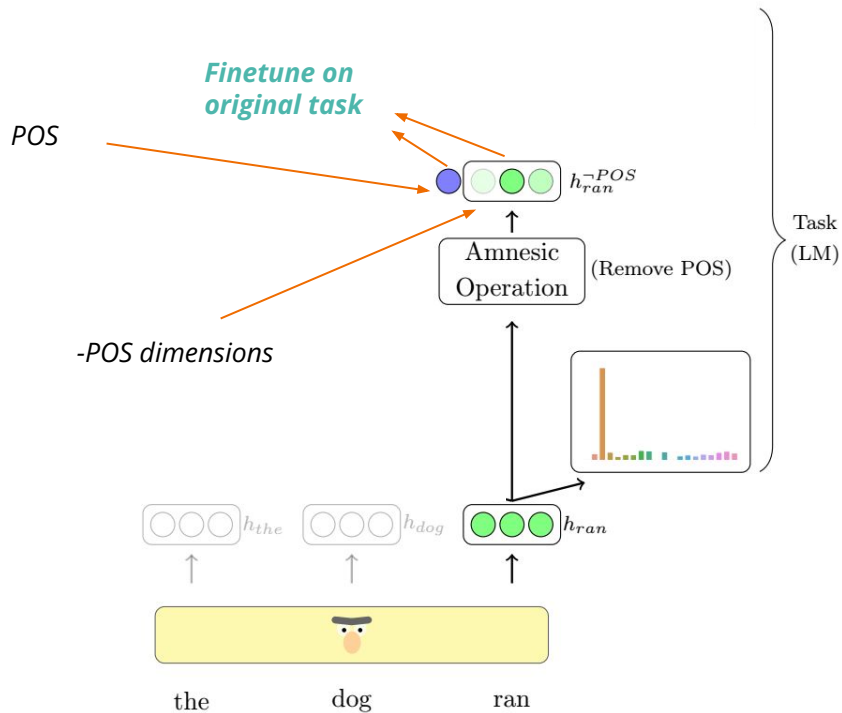
## Control over Selectivity



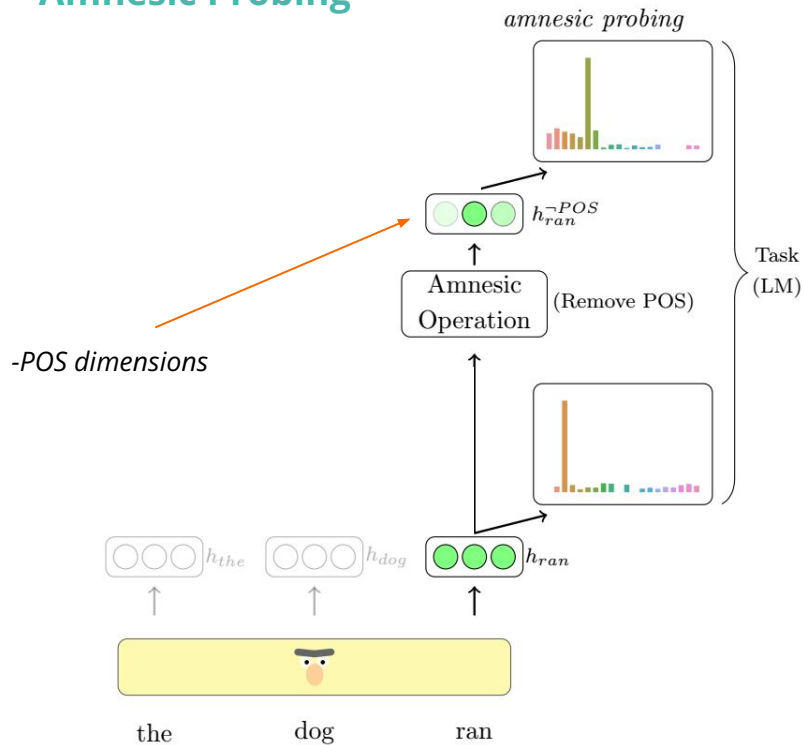
## Amnesic Probing



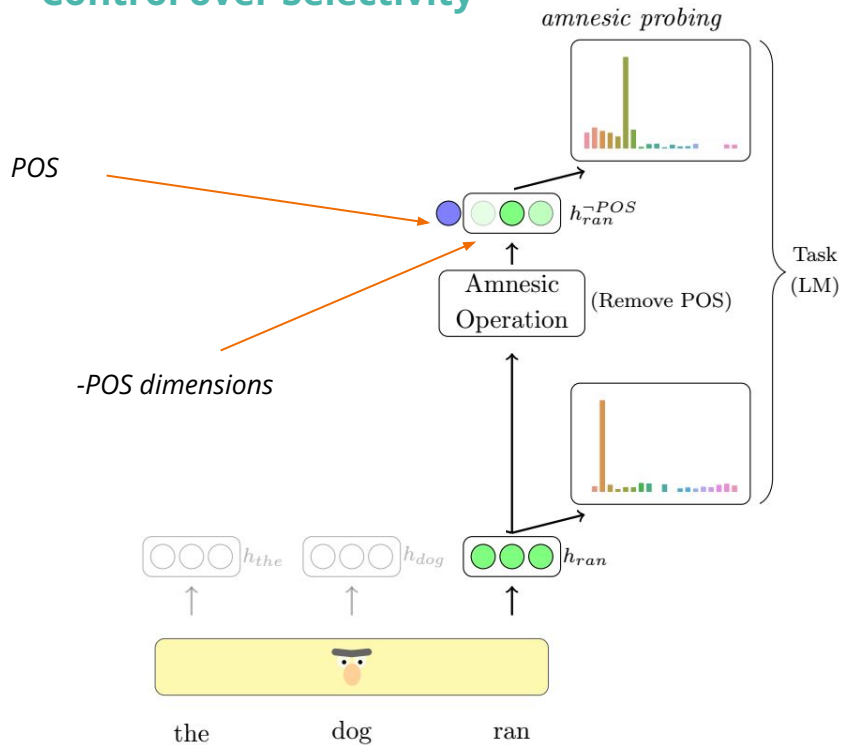
## Control over Selectivity



## Amnesic Probing



## Control over Selectivity



# Case Study: Pre-trained BERT



What linguistic properties are encoded used in BERT

# Amnesic Probing: Setup

- The model: BERT-base





# Amnesic Probing: Setup

- The model: BERT-base
- Properties:
  - POS



Does Bert make use of linguistic information ?

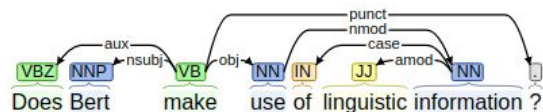
VBZ NNP VB NN IN JJ NN

# Amnesic Probing: Setup

- The model: BERT-base
- Properties:
  - POS
  - Dependency edges



Does Bert make use of linguistic information ?

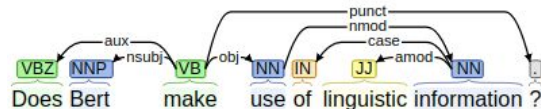


# Amnesic Probing: Setup

- The model: BERT-base
- Properties:
  - POS
  - Dependency edges
  - NER



Does Bert make use of linguistic information ?



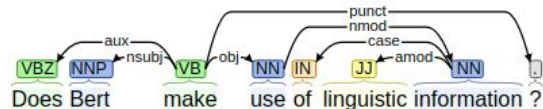
Does Bert make use of linguistic information ?

# Amnesic Probing: Setup

- The model: BERT-base
- Properties:
  - POS
  - Dependency edges
  - NER
  - Constituency boundaries



Does Bert make use of linguistic information ?



Does Bert make use of (linguistic information) ?

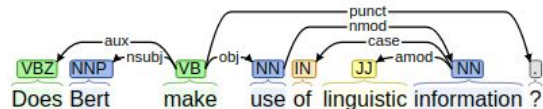
# Amnesic Probing: Setup

- The model: BERT-base
- Properties:
  - POS
  - D

Does BERT make use of POS, Dep-edge, NER and Const-boundaries when predicting words?

Does Bert use linguistic information ?

Does Bert make use of (linguistic information) ?



# Amnesic Probing: Results

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM- $D_{KL}$	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

# Amnesic Probing: Results

*Linguistic Properties*



		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

# Amnesic Probing: Results

Standard Probing

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
<b>Probing</b>	<b>Vanilla</b>	<b>76.00</b>	<b>89.50</b>	<b>92.34</b>	<b>93.53</b>	<b>85.12</b>	<b>83.09</b>
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM- $D_{KL}$	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01



# Amnesic Probing: Results

LM Accuracy Results

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

# Amnesic Probing: Results

Amnesic Comparison

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

# Amnesic Probing: Results

Amnesic Comparison

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

# Amnesic Probing: Results

Comparison to Control:  
**Information**

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

# Amnesic Probing: Results

Comparison to Control:  
**Information**

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

# Amnesic Probing: Results

Comparison to Control:  
**Selectivity**

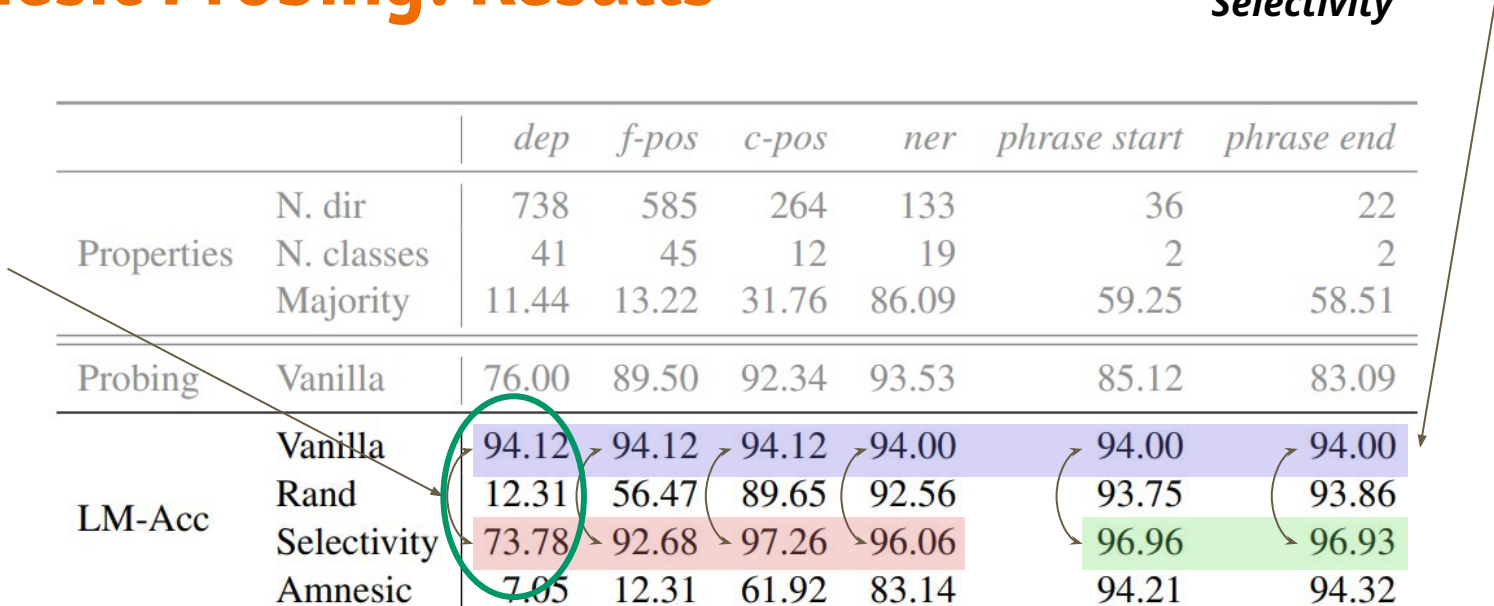
		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

# Amnesic Probing: Results

Comparison to Control:  
**Selectivity**

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

Doesn't  
Recover

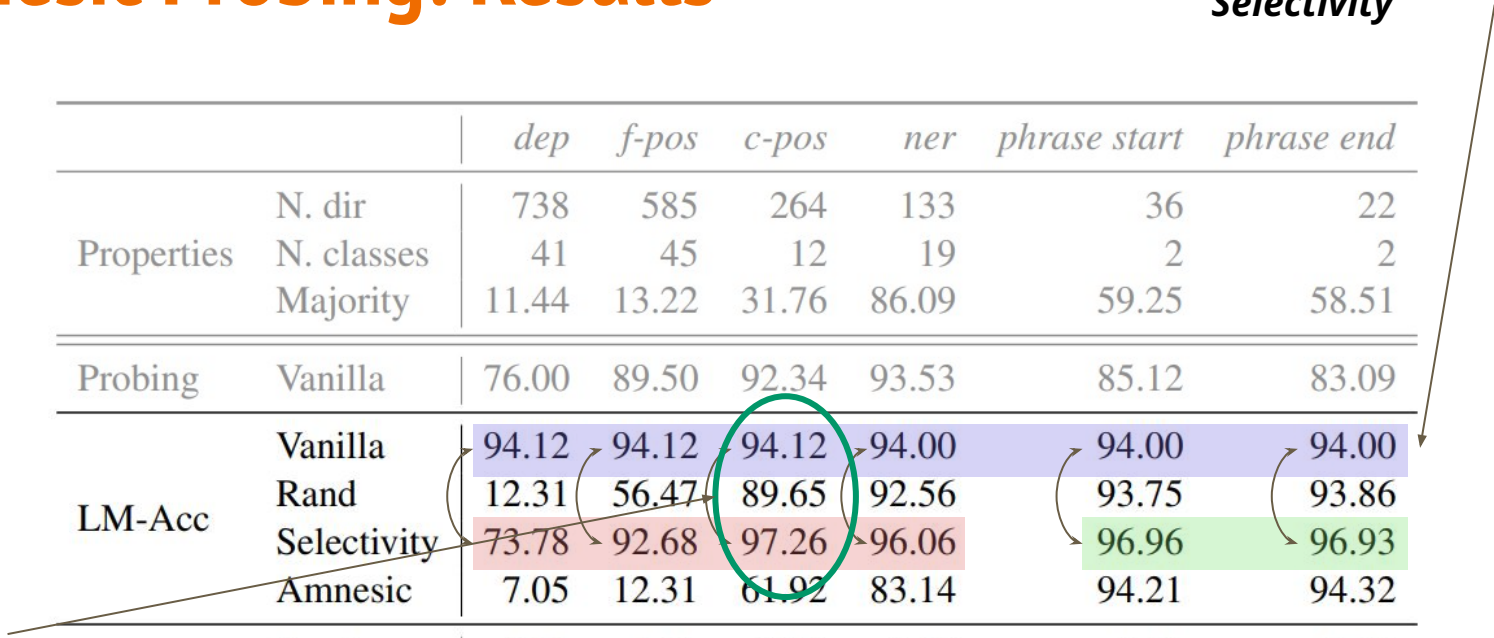


# Amnesic Probing: Results

Comparison to Control:  
**Selectivity**

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

Does  
Recover





# Amnesic Probing: Results

Phrase markers **are not** being used

Conclusions from all this:

		<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>		
		264	133	36	22		
		12	19	2	2		
		31.76	86.09	59.25	58.51		
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

POS and NER are being **used** by the model

# Amnesic Probing: Results

		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM- $D_{KL}$	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

*DKL Results*



# Amnesic Probing: Results

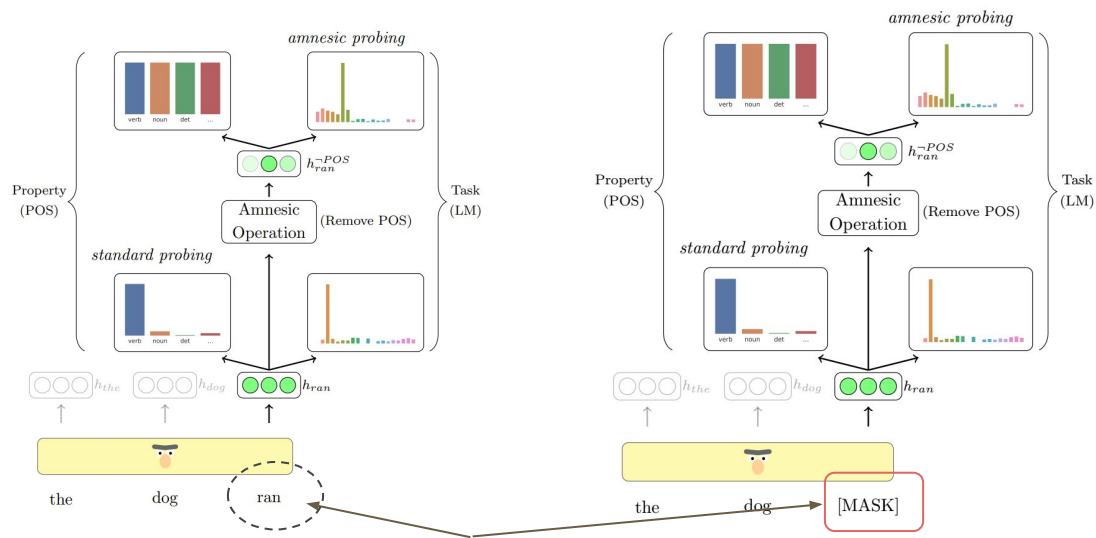
		<i>dep</i>	<i>f-pos</i>	<i>c-pos</i>	<i>ner</i>	<i>phrase start</i>	<i>phrase end</i>
Properties	N. dir	738	585	264	133	36	22
	N. classes	41	45	12	19	2	2
	Majority	11.44	13.22	31.76	86.09	59.25	58.51
Probing	Vanilla	76.00	89.50	92.34	93.53	85.12	83.09
LM-Acc	Vanilla	94.12	94.12	94.12	94.00	94.00	94.00
	Rand	12.31	56.47	89.65	92.56	93.75	93.86
	Selectivity	73.78	92.68	97.26	96.06	96.96	96.93
	Amnesic	7.05	12.31	61.92	83.14	94.21	94.32
LM-D <sub>KL</sub>	Rand	8.11	4.61	0.36	0.08	0.01	0.01
	Amnesic	8.53	7.63	3.21	1.24	0.01	0.01

*DKL Results*



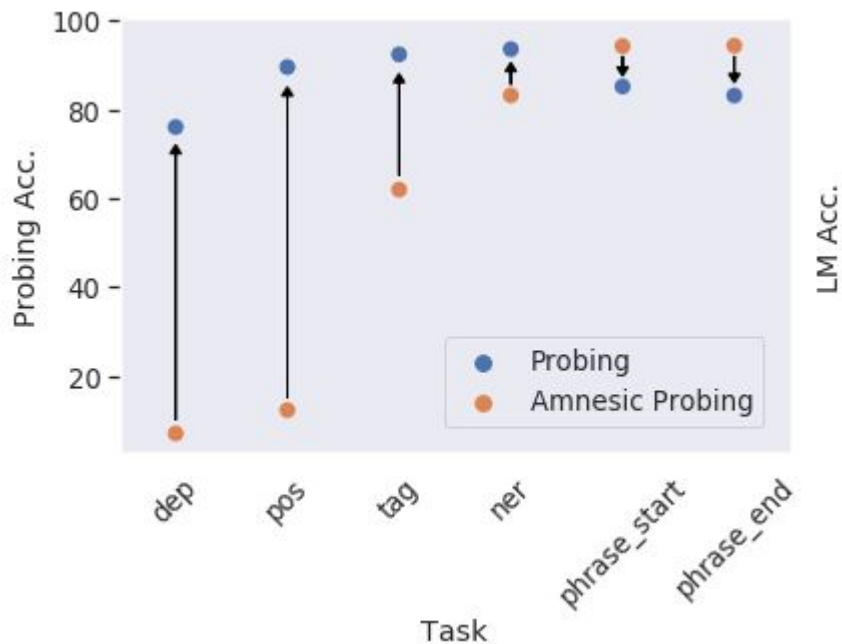
# Amnesic Probing: Results

- We perform the same experiments on another setup, where the words are masked
  - (Similar results, will elaborate if time permits)



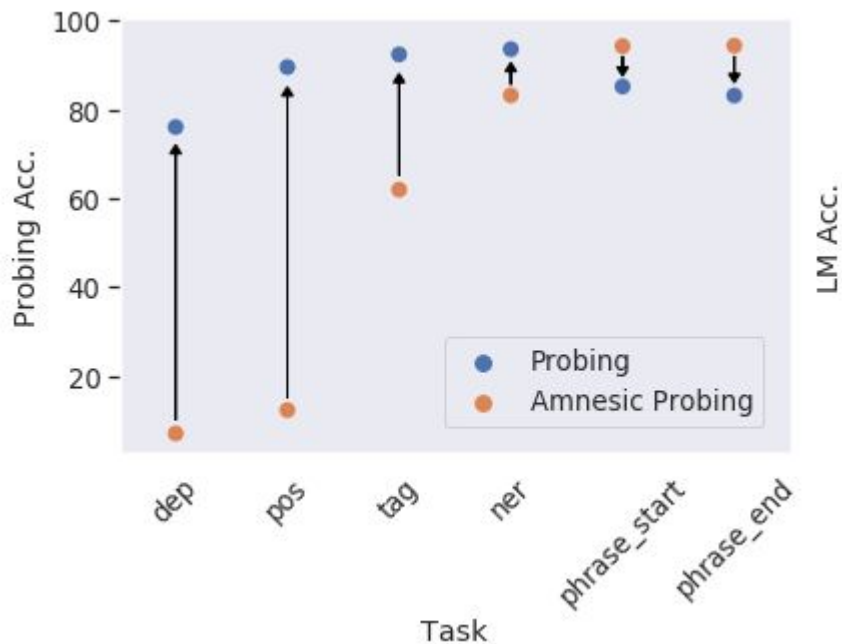
# Amnesic Probing vs. Standard Probing

- We plot the probing extractability performance vs. *amnesic probing*
- We observe no correlation between the two metrics



# Amnesic Probing vs. Standard Probing

- We plot the probing extractability performance vs. *amnesic probing*
- We observe no correlation between the two metrics
- **Can't make behavioural conclusions from standard probing results**



# Amnesic Probing: Diving In

# Amnesic Probing Fine Grained

- How individuals POS are affected by the removal of POS information?
- Open vs. Closed vocabulary

Large changes

Small changes

<i>c-pos</i>	Vanilla	Rand	Amnesic	$\Delta$
verb	46.72	44.85	34.99	11.73
noun	42.91	38.94	34.26	8.65
adposition	73.80	72.21	37.86	35.93
determiner	82.29	83.53	16.64	65.66
numeral	40.32	40.19	33.41	6.91
punctuation	80.71	81.02	47.03	33.68
particle	96.40	95.71	18.74	77.66
conjunction	78.01	72.94	4.28	73.73
adverb	39.84	34.11	23.71	16.14
pronoun	70.29	61.93	33.23	37.06
adjective	46.41	42.63	34.56	11.85
other	70.59	76.47	52.94	17.65



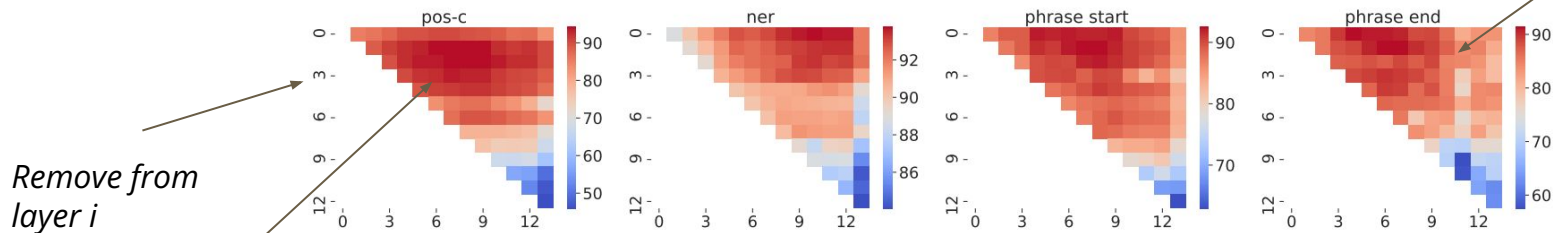
# Amnesic Probing: Inside The Model

# The Inner Layers

- Until now, querying the last layer
  - INLP removes linear information, last layer is only multiplied by a matrix
- We perform the same analysis on the Inner layers
- Standard Probe (after the amnesic operation)
- Behavioral Probe

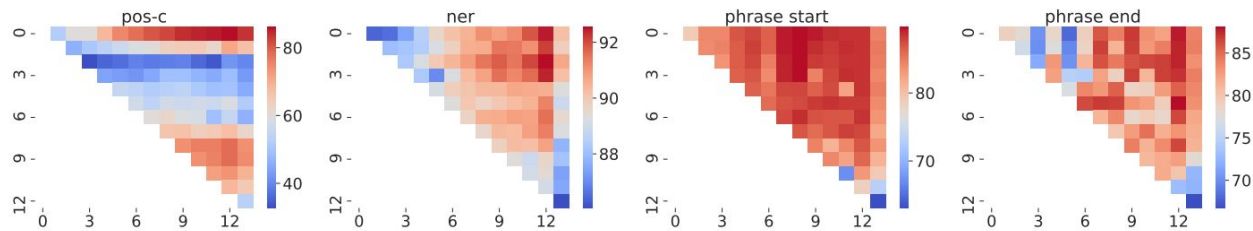
# The Inner Layers: Probing

- Removing information from layer  $i$ , and probing in layer  $j$



(a) Non-Masked version

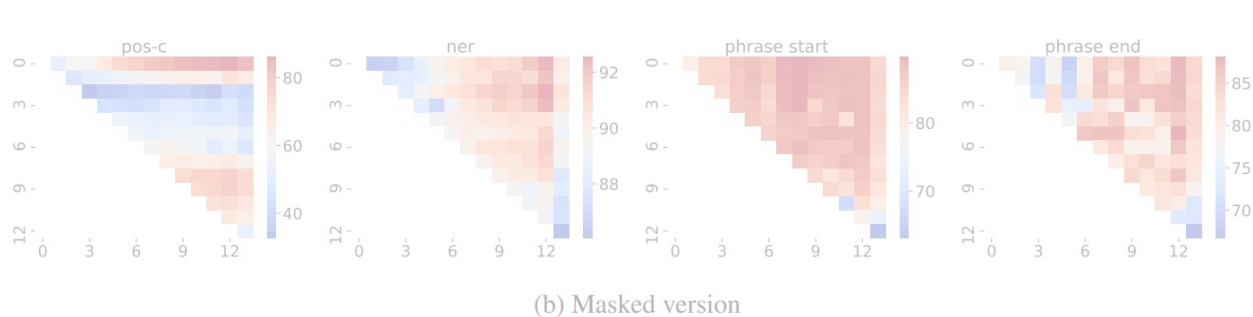
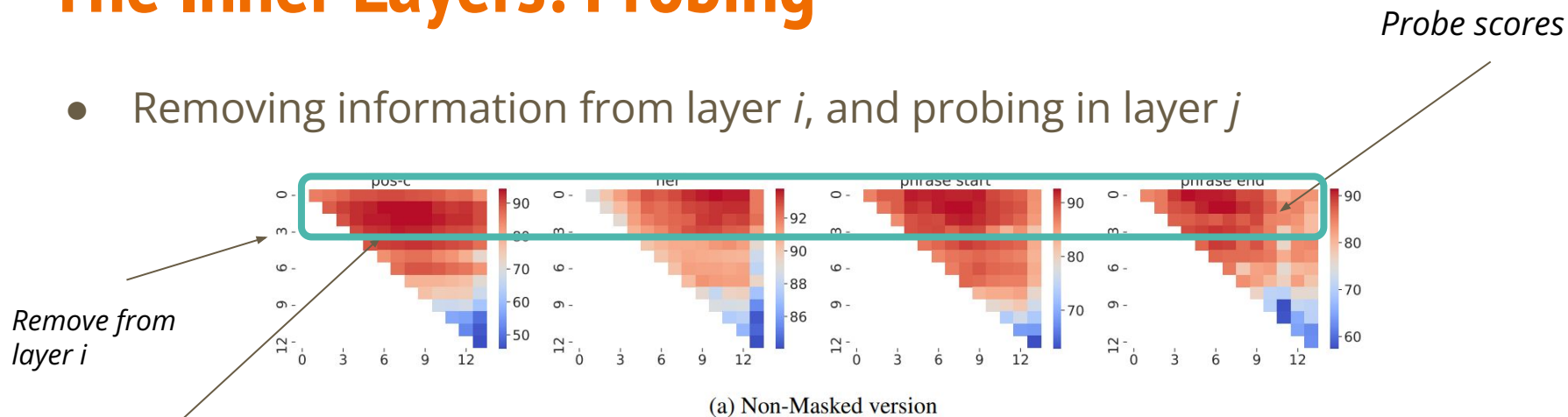
Probe layer  $j$



(b) Masked version

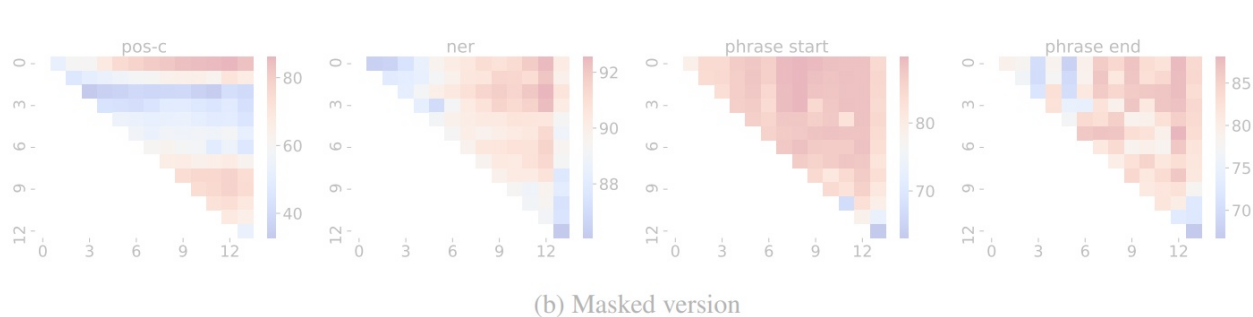
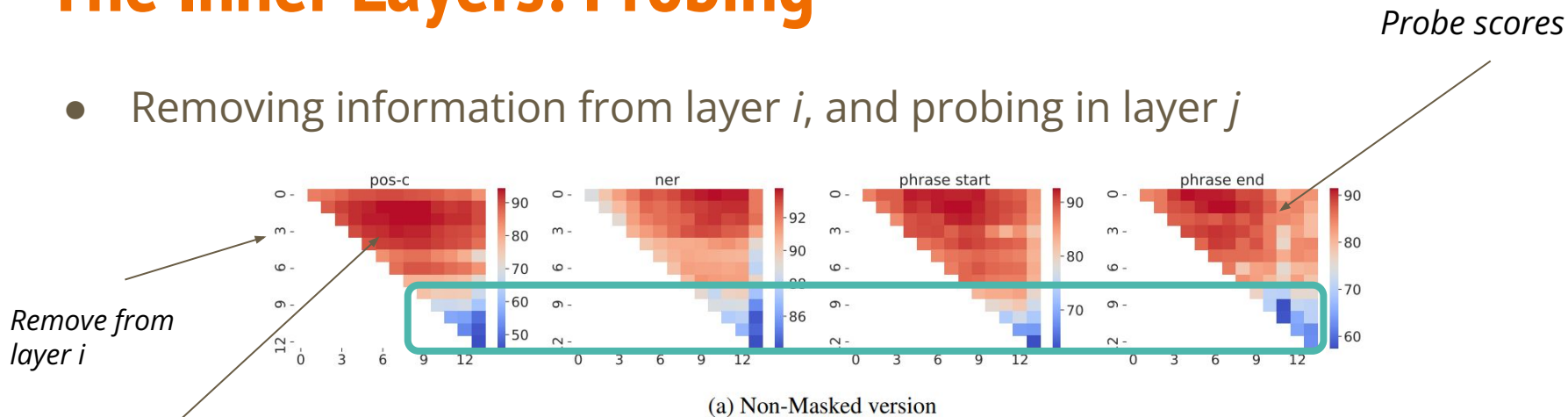
# The Inner Layers: Probing

- Removing information from layer  $i$ , and probing in layer  $j$



# The Inner Layers: Probing

- Removing information from layer  $i$ , and probing in layer  $j$



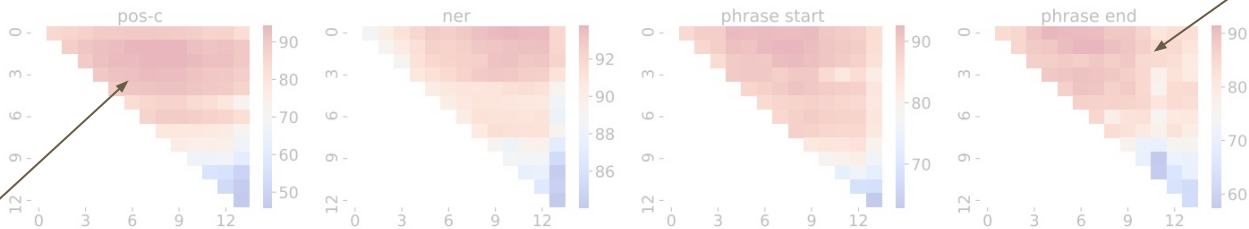
# The Inner Layers: Probing

- Removing information from layer  $i$ , and probing in layer  $j$

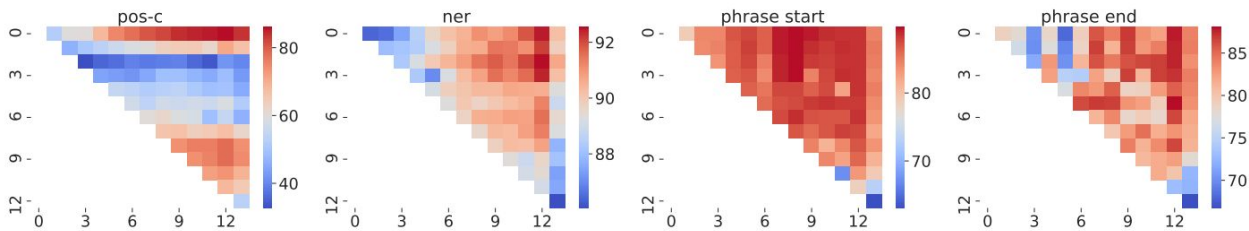
Probe scores

Remove from layer  $i$

Probe layer  $j$



(a) Non-Masked version



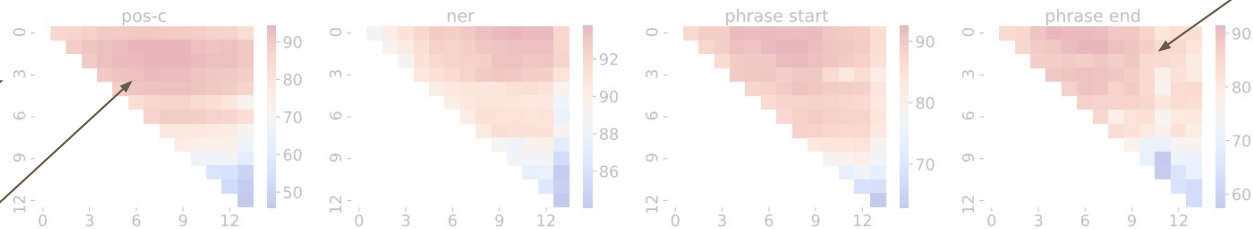
(b) Masked version

# The Inner Layers: Probing

- Removing information from layer  $i$ , and probing in layer  $j$

Probe scores

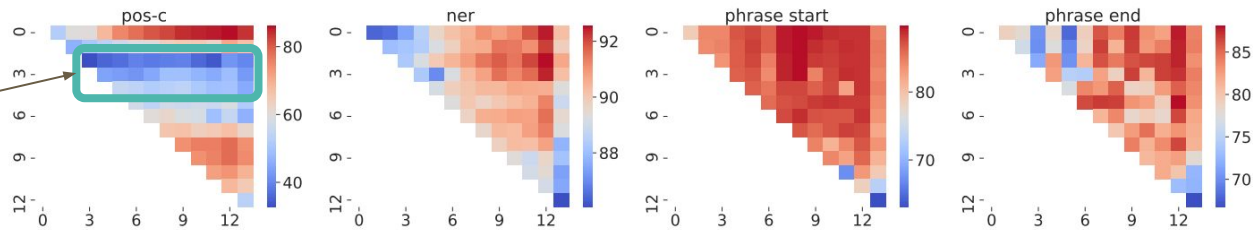
Remove from layer  $i$



(a) Non-Masked version

Probe layer  $j$

Irreversible removal



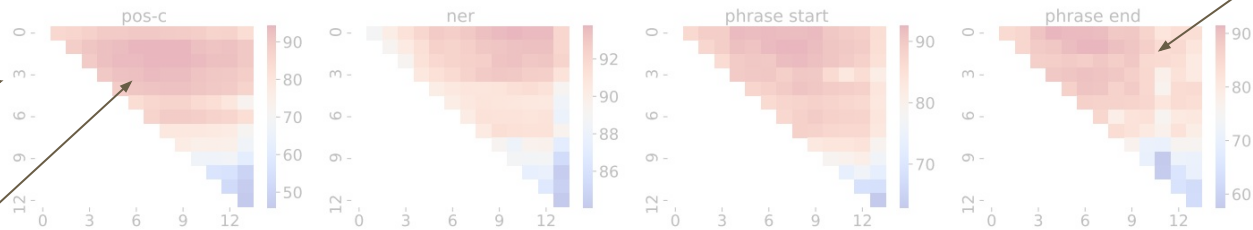
(b) Masked version

# The Inner Layers: Probing

- Removing information from layer  $i$ , and probing in layer  $j$

Probe scores

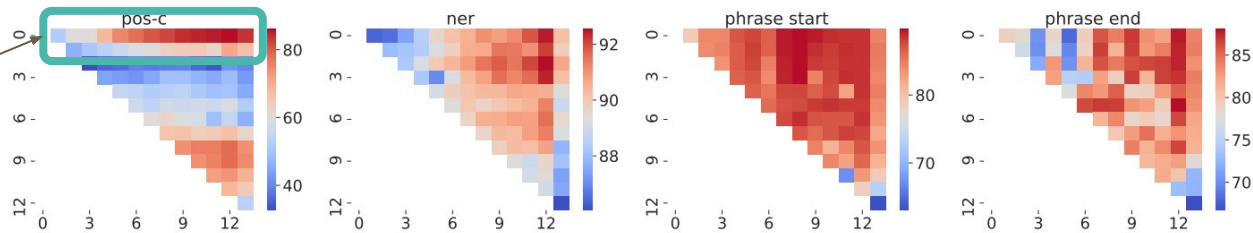
Remove from layer  $i$



(a) Non-Masked version

Probe layer  $j$

Reversible removal

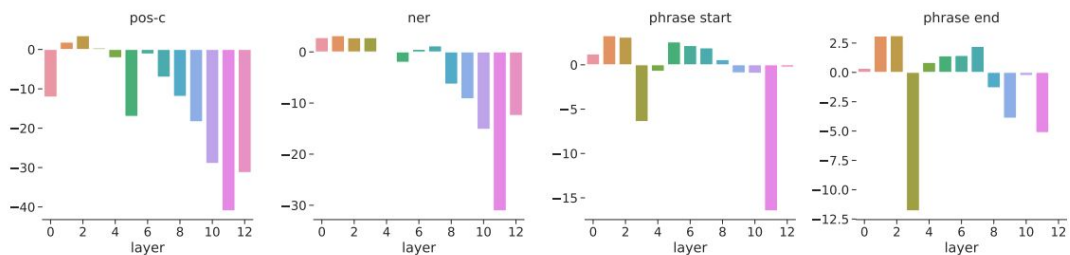


(b) Masked version

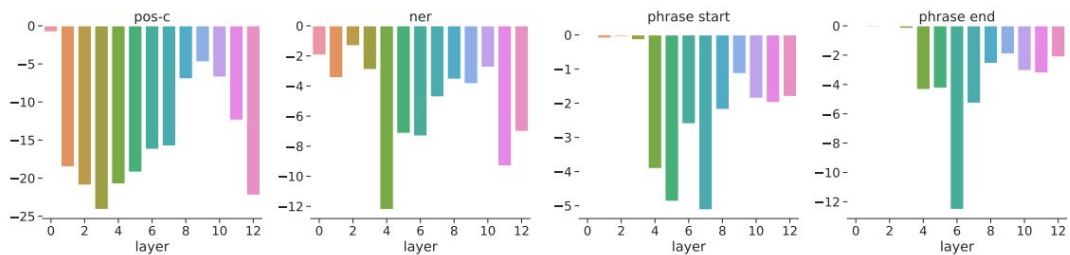


# The Inner Layers: Amnesic Probing

- Removing information from layer  $i$ , and inspecting the model's predictions



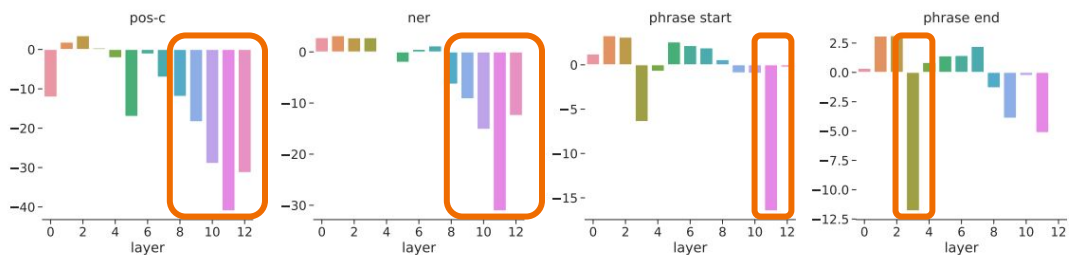
(a) Non-Masked version



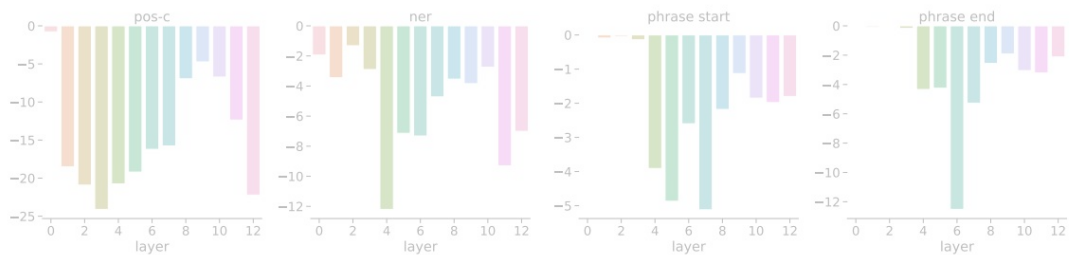
(b) Masked version

# The Inner Layers: Amnesic Probing

- Removing information from layer  $i$ , and inspecting the model's predictions



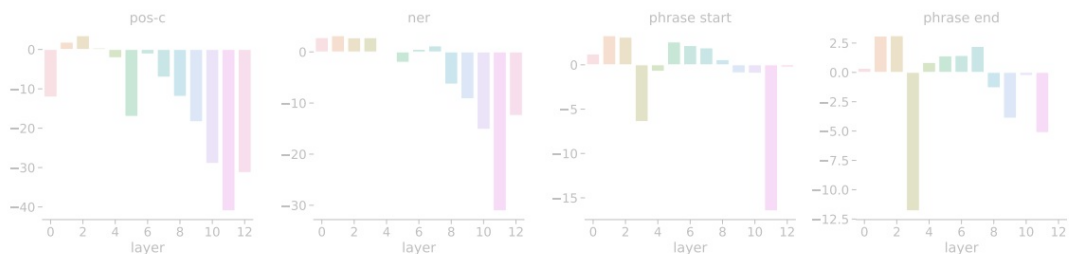
(a) Non-Masked version



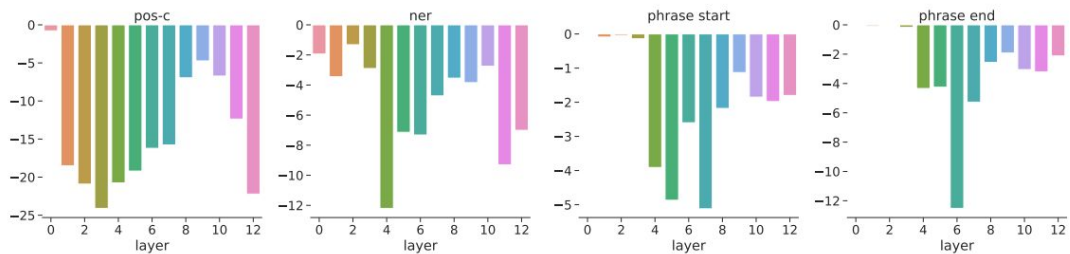
(b) Masked version

# The Inner Layers: Amnesic Probing

- Removing information from layer  $i$ , and inspecting the model's predictions



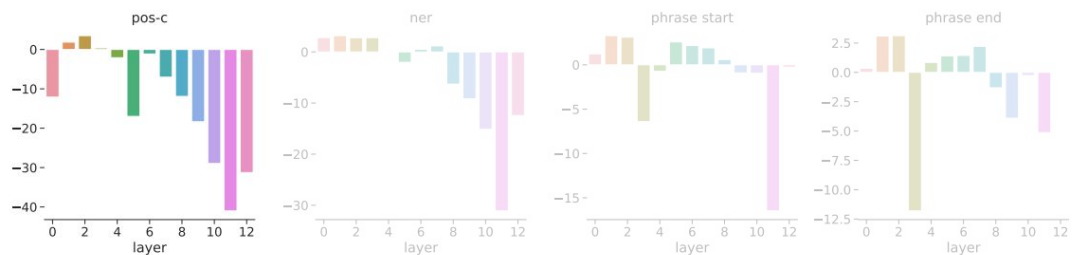
(a) Non-Masked version



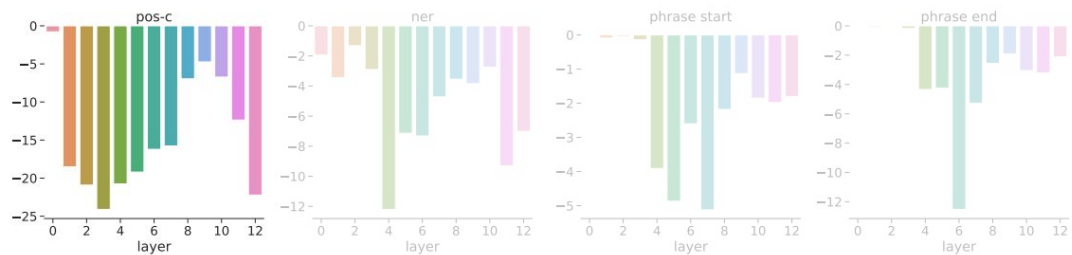
(b) Masked version

# The Inner Layers: Amnesic Probing

- Removing information from layer  $i$ , and inspecting the model's predictions



(a) Non-Masked version

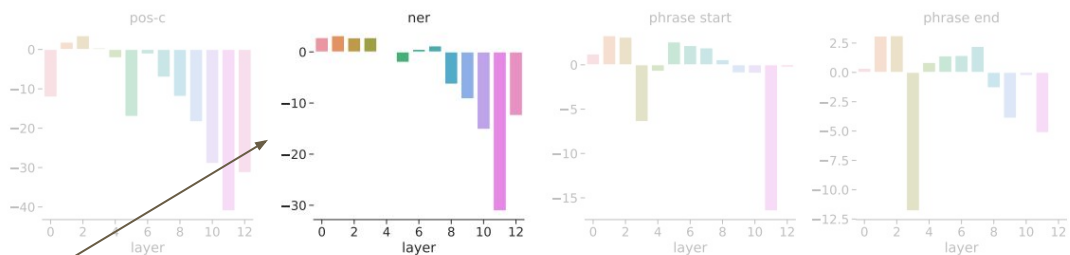


(b) Masked version

*Masked vs  
Non-Masked*

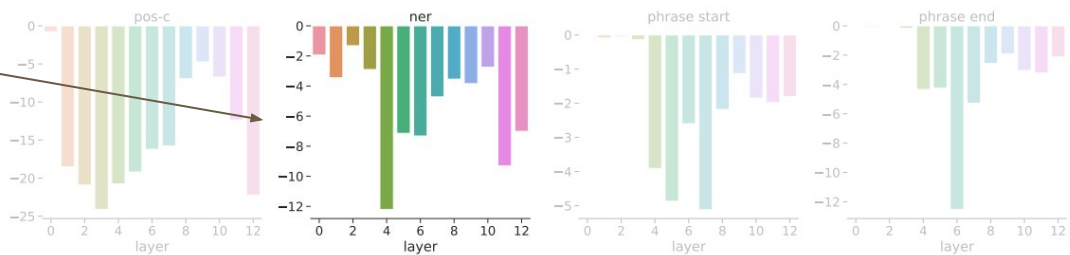
# The Inner Layers: Amnesic Probing

- Removing information from layer  $i$ , and inspecting the model's predictions



(a) Non-Masked version

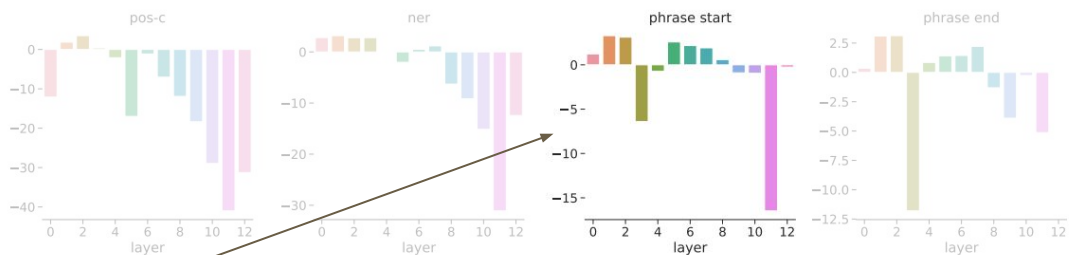
*Masked vs  
Non-Masked*



(b) Masked version

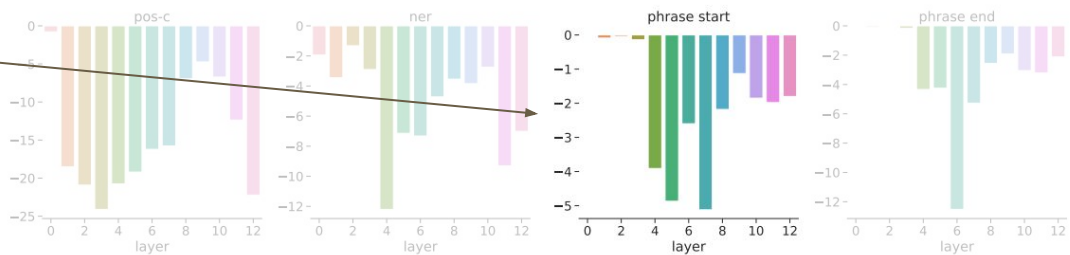
# The Inner Layers: Amnesic Probing

- Removing information from layer  $i$ , and inspecting the model's predictions



(a) Non-Masked version

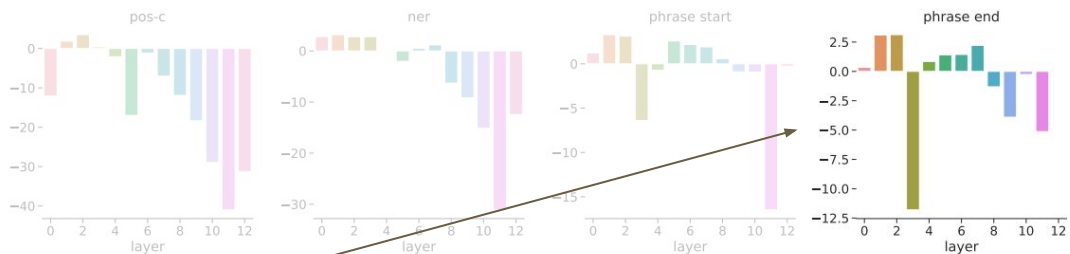
*Masked vs  
Non-Masked*



(b) Masked version

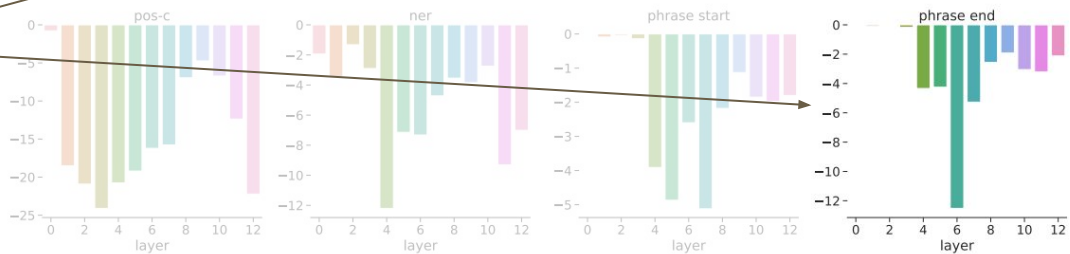
# The Inner Layers: Amnesic Probing

- Removing information from layer  $i$ , and inspecting the model's predictions



(a) Non-Masked version

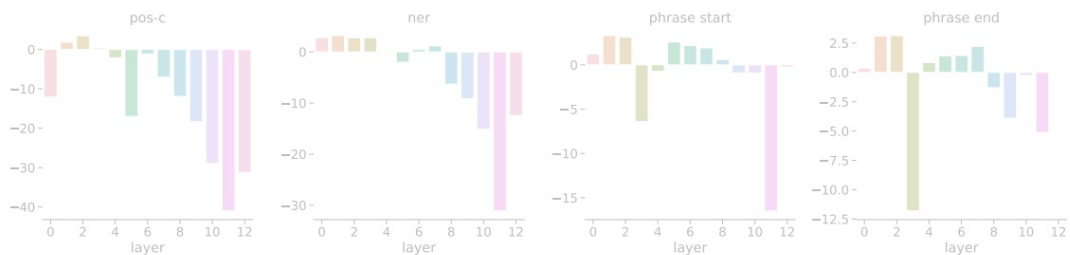
*Masked vs  
Non-Masked*



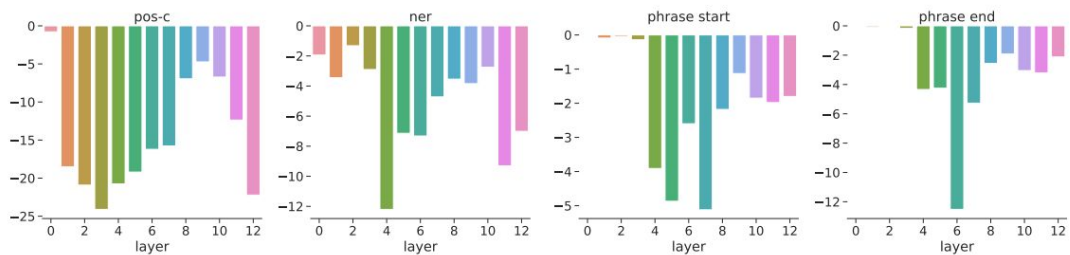
(b) Masked version

# The Inner Layers: Amnesic Probing

- Removing information from layer  $i$ , and inspecting the model's predictions



(a) Non-Masked version



(b) Masked version

*Strong impact  
in the first few  
layers!!*



# To conclude

- Probing answers the question of “**what/how properties are encoded?**”
- We are often interested in a **different** question: “**what is being used?**”
- We propose to ask the causal question and **offer a method** to answer it:  
***Amnesic Probing***
- **We encourage you to use it!**



# Going Forward

- What **does** it mean that some information is extractable?
- ... or, why is it there from the first place?
- Algorithms that remove also non-linear information

# Part II

## Measuring and Improving Consistency in Pretrained Language Models

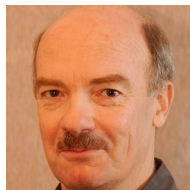
**Yanai Elazar<sup>1,2</sup> Nora Kassner<sup>3</sup> Shauli Ravfogel<sup>1,2</sup> Abhilasha Ravichander<sup>4</sup>  
Eduard Hovy<sup>4</sup> Hinrich Schütze<sup>3</sup> Yoav Goldberg<sup>1,2</sup>**

<sup>1</sup>Computer Science Department, Bar Ilan University

<sup>2</sup>Allen Institute for Artificial Intelligence

<sup>3</sup>Center for Information and Language Processing (CIS), LMU Munich

<sup>4</sup>Language Technologies Institute, Carnegie Mellon University



*TACL 2021*

# Model's Failure Mode



How many birds?	<b>A: 1</b>
Is there 1 bird?	<b>A: no</b>
Are there 2 birds?	<b>A: yes</b>
Are there any birds?	<b>A: no</b>

# Model's Failure Mode

**Context:** 826 Doctor Who instalments have been televised since 1963 ... Starting with the 2009 special "Planet of the Dead", the series was filmed in 1080i for HDTV ...

**Q1:** In what year did Doctor Who begin being shown in HDTV? **A:** 2009



# Model's Failure Mode

**Context:** 826 Doctor Who instalments have been televised since 1963 ... Starting with the 2009 special "Planet of the Dead", the series was filmed in 1080i for HDTV ...

**Q1:** In what year did Doctor Who begin being shown in HDTV? **A:** 2009



*Inconsistent*

**Q2:** Since what year has Doctor Who been televised in HDTV? **A:** 1963

# Model's Failure Mode

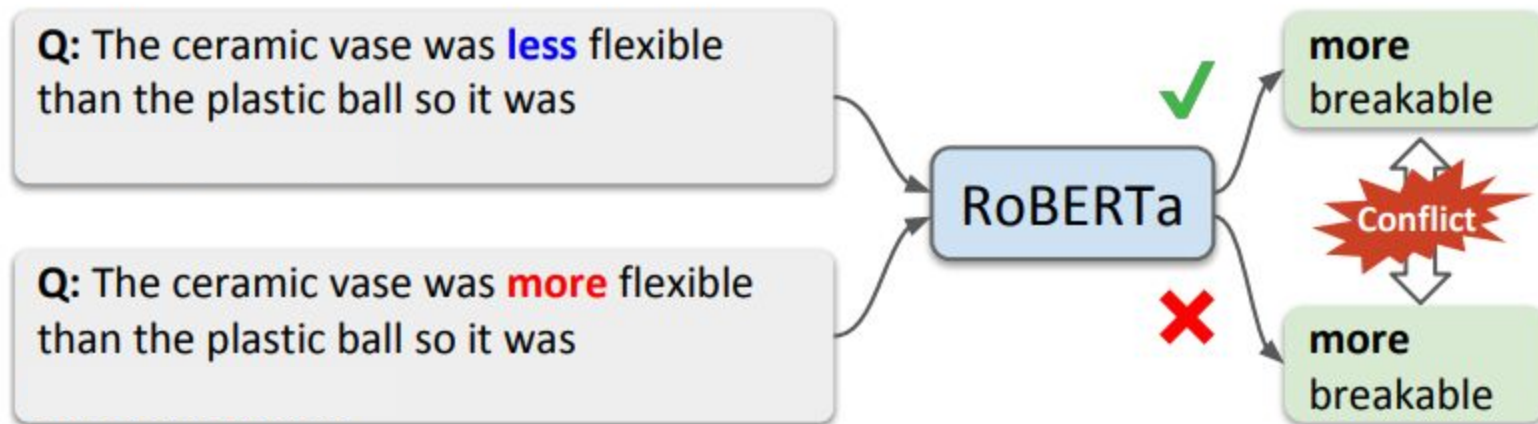
Kublai originally named his eldest son, Zhenjin, as the Crown Prince, but he died before Kublai in 1285.

(c) Excerpt from an input paragraph, **SQuAD dataset**.

**Q:** When did Zhenjin die?      **A:** 1285

**Q:** Who died in 1285?      **A:** Kublai

# Model's Failure Mode

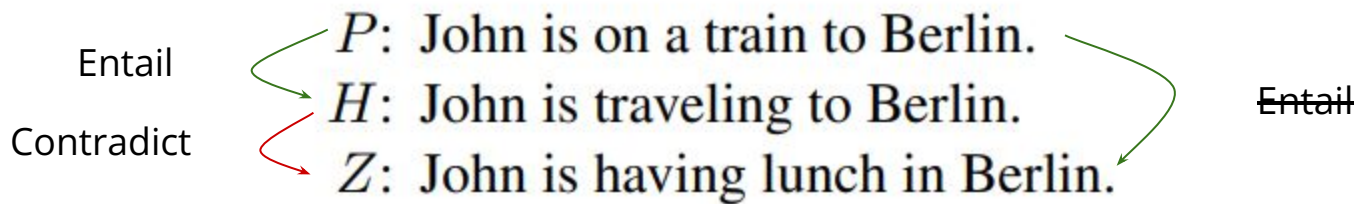




# Model's Failure Mode

Context	Match
<i>A robin is a ___</i>	<i>bird</i>
<i>A robin is not a ___</i>	<i>bird</i>

# Model's Failure Mode



# Consistency in Models

- End-task models suffer from inconsistency
- Today's standard pipeline is: Pretrain -> Finetune
- **In this work:** we show that *Inconsistency starts in the PLM itself*

# Consistency in Humans

1:1s Advance Sign-Up Sheet - Yanai Elazar

File Edit View Insert Format Data Tools Extensions Help

100% \$ % .0 .00 123 Arial 12 B I S A

A1 AI2 TALK PRESENTER

	A	B	C
1	<b>AI2 TALK PRESENTER</b>		
2	Yanai Elazar		
3			
4	<b>TITLE:</b>		
5	Causal Attributions in Language Models		
6			
7	<b>DATE</b>		
8	Tuesday, November 23		
9			
10	<b>1:1 TIME SLOT (30 mins ea)</b>	<b>NAME</b>	<b>LOCATION</b>
11	11:00	Noah Smith	<a href="https://meet.google.com/rxg-dvmv-sdy?authuser=0">https://meet.google.com/rxg-dvmv-sdy?authuser=0</a>
12	11:30	Jungo	"
13	12:00	Pete Clark/Lunch Break (45 mins)	"
14	12:45	Yejin	"
15	1:15	KyleL (happy to switch if need)	"
16	1:45	Ronan	"
17	2:15	<b>BREAK</b>	
18	2:30	Yuling Gu	"
19	3:00	Nishant Subramani	"
20	3:30	Alexis Ross	"
21			
22	<b>ABSTRACT</b>		
	The outstanding results of enormous language models are largely unexplained, and different methods in interpretability aim to and analyze these models to understand their working mechanisms. Probing, one of these tools suggests that accurately predicted properties from models' representations are likely to explain some of the features or concepts that these models utilize in their predictions.		
	In the first part of this talk, I'll propose Amnesic Probing, a new interpretability method that takes inspiration from counterfactual world have been the prediction if the model had not accessed certain information?" Amnesic Probing is a more suitable method asking causal questions about how attributes are used by models.		
23	In the second part, I'll talk about a different kind of probing that treats the model as a black box and uses cloze patterns to query model for world knowledge under the LAMA framework.		

+

Emma Strubell

# Consistency in Humans

**Yanai Elazar**  
Bar Ilan University **AI2**  
yanaiela@gmail.com

**Yanai Elazar**<sup>1,3</sup> **Yoav Goldberg**<sup>1,3</sup>  
**<sup>1</sup>The Allen Institute for AI**

**er**<sup>1,2</sup> **Yanai Elazar**<sup>3,4</sup> **Benoît Sagot**<sup>1</sup>  
Paris, France <sup>2</sup>Sorbonne Université, Paris  
Computer Science Department, Bar Ilan Unive  
**<sup>4</sup>Allen Institute for Artificial Intelligence**

# Consistency in Models: This Part


1. Language Models as Knowledge Bases
2. Why is consistency crucial?

} *background*

# Consistency in Models: This Part

1. Language Models as Knowledge Bases
  2. Why is consistency crucial?
- } *background*
3. ParaRel 🙌: a new resource that enables us to measure consistency

# Consistency in Models: This Part

1. Language Models as Knowledge Bases
  2. Why is consistency crucial?
- 
- background*
3. ParaRel 🙌: a new resource that enables us to measure consistency
  4. A framework for measuring (In)Consistency in Language Models
    - In the context of factual knowledge



# Consistency in Models: This Part

1. Language Models as Knowledge Bases
  2. Why is consistency crucial?
- } *background*
3. ParaRel 🙌: a new resource that enables us to measure consistency
  4. A framework for measuring (In)Consistency in Language Models
    - In the context of factual knowledge
  5. A proposal to improve consistency in LMs.

# Consistency in Models: This Part

1. Language Models as Knowledge Bases
2. Why is consistency crucial?

} *background*

3. ParaRel 🙌: a new resource that enables us to measure consistency
4. A framework for measuring (In)Consistency in Language Models
  - In the context of factual knowledge

} *novelty*

5. A proposal to improve consistency in LMs.

# Setup: LMs as Knowledge Bases

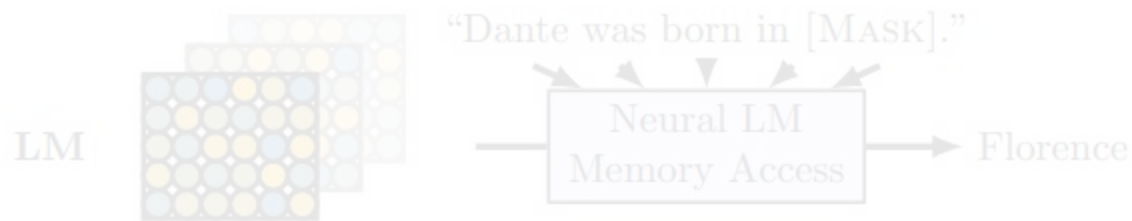
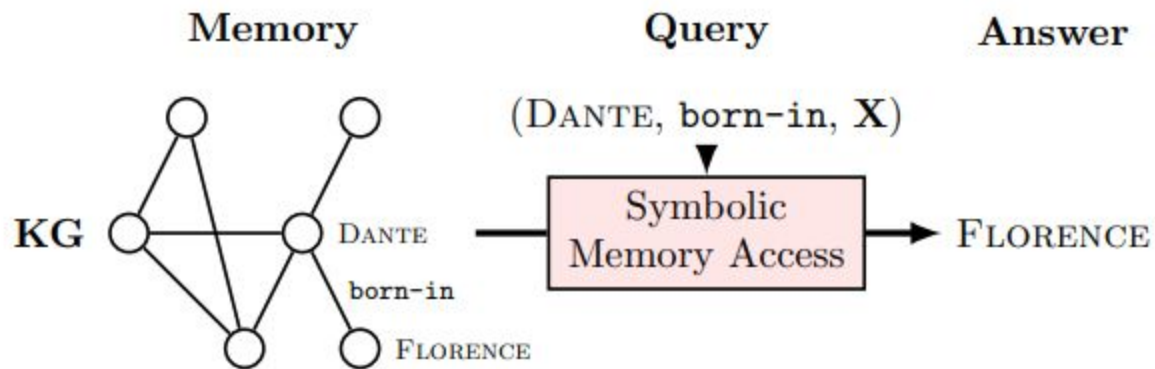
# Language Models as Knowledge Bases?

**Fabio Petroni<sup>1</sup> Tim Rocktäschel<sup>1,2</sup> Patrick Lewis<sup>1,2</sup> Anton Bakhtin<sup>1</sup>  
Yuxiang Wu<sup>1,2</sup> Alexander H. Miller<sup>1</sup> Sebastian Riedel<sup>1,2</sup>**

<sup>1</sup>Facebook AI Research

<sup>2</sup>University College London

{fabiopetroni, rockt, plewis, yolo, yuxiangwu, ahm, sriedel}@fb.com



*e.g.* ELMo/BERT

Figure 1: Querying knowledge bases (KB) and language models (LM) for factual knowledge.

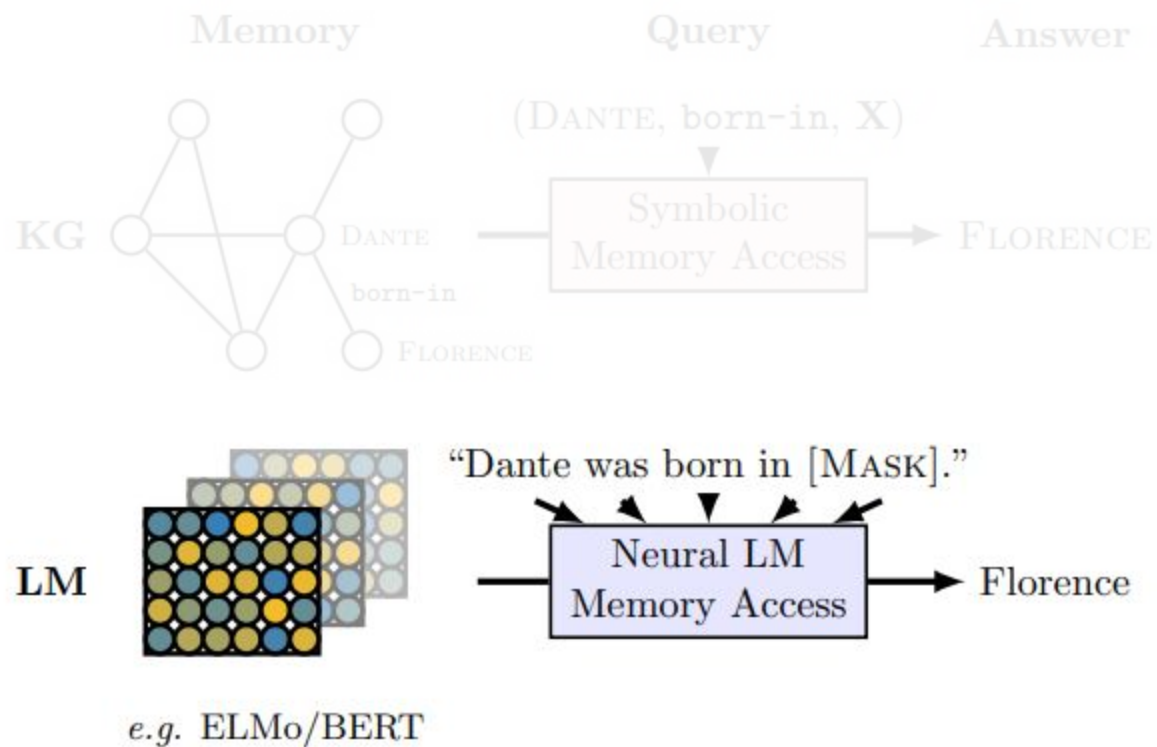


Figure 1: Querying knowledge bases (KB) and language models (LM) for factual knowledge.

# Using Patterns to Query LMs

- Born-In: “[X] was born in [Y] .”
  - *Barack Obama was born in [MASK].*
- Broadcasting Channel: “[X] was originally aired on [Y] .”
  - *Lost was originally aired on [MASK].*
- ...

# Language Models as KBs - Setup

- The data is of the form <subject, relation, object>
  - E.g. <“Barack Obama”, “born-in”, “Hawaii”>
- To query an LM, we write a ‘pattern’ that expresses a relation
  - E.g. “[X] was born in [Y]”
- Given the subject and relation, the task is to predict the object
  - E.g. <“Barack Obama”, born-in> -> “Hawaii”
  - In Petroni et al., 2019, used 1 pattern for every relation



Corpus	Relation	Statistics		Baselines		KB		LM					Bl
		#Facts	#Rel	Freq	DrQA	RE <sub>n</sub>	RE <sub>o</sub>	Fs	Txl	Eb	E5B	Bb	
Google-RE	birth-place	2937	1	4.6	-	3.5	13.8	4.4	2.7	5.5	7.5	14.9	<b>16.1</b>
	birth-date	1825	1	1.9	-	0.0	<b>1.9</b>	0.3	1.1	0.1	0.1	1.5	1.4
	death-place	765	1	6.8	-	0.1	7.2	3.0	0.9	0.3	1.3	13.1	<b>14.0</b>
	Total	5527	3	4.4	-	1.2	7.6	2.6	1.6	2.0	3.0	9.8	<b>10.5</b>
T-REx	1-1	937	2	1.78	-	0.6	10.0	17.0	36.5	10.1	13.1	68.0	<b>74.5</b>
	<i>N-1</i>	20006	23	23.85	-	5.4	<b>33.8</b>	6.1	18.0	3.6	6.5	32.4	34.2
	<i>N-M</i>	13096	16	21.95	-	7.7	<b>36.7</b>	12.0	16.5	5.7	7.4	24.7	24.3
	Total	34039	41	22.03	-	6.1	<b>33.8</b>	8.9	18.3	4.7	7.1	31.1	32.3
ConceptNet	Total	11458	16	4.8	-	-	-	3.6	5.7	6.1	6.2	15.6	<b>19.2</b>
SQuAD	Total	305	-	-	<b>37.5</b>	-	-	3.6	3.9	1.6	4.3	14.1	17.4

# Language Models as KBs

- LMs were trained on large sources of knowledge (e.g. Wikipedia)
- Can capture (memorize) some of these facts as part of the pretraining objective

# Pretraining a Language Model

## Background

---

### Early life of Barack Obama

*Main articles: [Early life and career of Barack Obama](#) and [Ann Dunham](#)*

People who express doubts about Obama's eligibility or reject details about his early life are often informally called "birthers", a term that parallels<sup>[23]</sup> the nickname "truthers" for adherents of [9/11 conspiracy theories](#).<sup>[24][25]</sup> These [conspiracy theorists](#) reject at least some of the following facts about his early life:

Barack Obama was born on August 4, 1961, at Kapi'olani Maternity & Gynecological Hospital (now called [Kapi'olani Medical Center for Women & Children](#)) in Honolulu, Hawaii,<sup>[26][27][28][29]</sup> to [Ann Dunham](#),<sup>[30]</sup> from [Wichita, Kansas](#),<sup>[31]</sup> and her husband [Barack Obama Sr.](#), a Luo from [Nyang'oma Kogelo, Nyanza Province](#) (in what was then the [Colony and Protectorate of Kenya](#)), who was attending the University of Hawaii. Birth notices for Barack Obama were published in [The Honolulu Advertiser](#) on August 13 and the [Honolulu Star-Bulletin](#) on August 14, 1961.<sup>[26][31]</sup> Obama's father's immigration file also clearly states Barack Obama was born in Hawaii.<sup>[32]</sup> One of his high school teachers, who was acquainted with his mother at the time, remembered hearing about the day of his birth.<sup>[30]</sup>

# Pretraining a Language Model

And it actually works!

## LM predictions

#1 mask: Tel Aviv is located in **[MASK]**.

bert_large_cased	
0	Israel
1	Jerusalem
2	Palestine
3	Haifa
4	Egypt
5	Europe
6	Ukraine
7	Lebanon
8	Jordan
9	Germany

# Pretraining a Language Model

Well, sometimes...

## LM predictions

#1 mask: Barack Obama was born in **[MASK]**.

```
bert_large_cased
```

0	Chicago
1	Philadelphia
2	Detroit
3	Houston
4	Atlanta
5	Georgia
6	Boston
7	Texas
8	Paris
9	Dallas

# Pretraining a Language Model

Well, sometimes...

LM predictions

#1 mask: Barack Obama was born in [MASK].

	bert_large_cased
0	Chicago
1	Philadelphia
2	Detroit
3	Houston
4	Atlanta
5	Georgia
6	Boston
7	Texas
8	Paris
9	Dallas

**Teaser:**

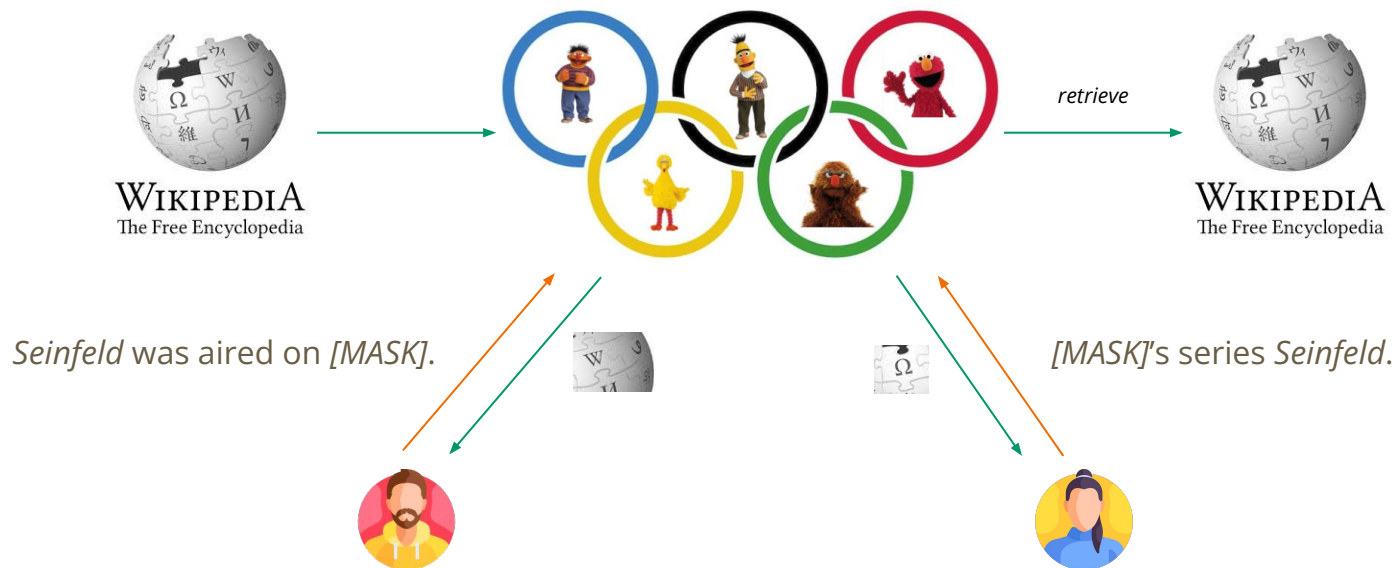
*we'll get back to the reason  
behind this prediction in Part III*

# Language Models as KBs - Setup



- Restricting to MLM predictions: single token objects
- Restricting to the possible objects for a specific relation

# Language Models as KBs





# Language Models as KBs

So the real question is



Does It Generalize?

# Language Models as KBs - Consistency?

We'd like that an LM would make the same prediction across paraphrases

E.g.:

*"Seinfeld* was aired on [Y]."

-  *"Seinfeld, that was aired on [Y],"*
-  *"[Y]'s series Seinfeld,"*

# Language Models as KBs - Consistency?

We'd like that an LM would make the same prediction across paraphrases

E.g.:

"*Seinfeld* was aired on [Y]."

- ↔ "*Seinfeld*, that was aired on [Y],"

- ↔ "[Y]'s series *Seinfeld*,"



*Consistent*



*Inconsistent*

# Measuring Consistency:

ParaRel 🤘

# Language Models as KBs - ParaRel 🤘

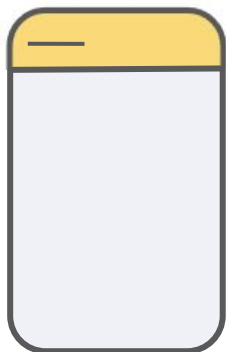
But where can we get these patterns?

We build a new resource:

ParaRel 🤘 (**Par**aphrase **Rel**ations)

# ParaRel 🙌 - Creation

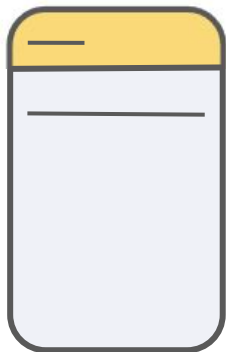
- For every relation, we manually build a set of patterns that are paraphrases of each other, in 4 steps:



# ParaRel 🙌 - Creation

(a) A single pattern

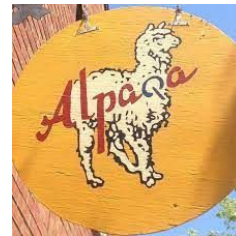
- For every relation, we manually build a set of patterns that are paraphrases of each other, in 4 steps:
  - a. Starting with the LAMA patterns (*Petroni et al., 2019*)



# ParaRel 🙌 - Creation

- For every relation, we manually build a set of patterns that are paraphrases of each other, in 4 steps:
  - Starting with the LAMA patterns (*Petroni et al., 2019*)
  - Augmenting with LPAQA patterns (*Jiang et al., 2020*)

(a) A single pattern



(b) Multiple patterns,  
noisy

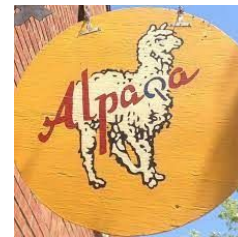




# ParaRel 🙌 - Creation

- For every relation, we manually build a set of patterns that are paraphrases of each other, in 4 steps:
  - Starting with the LAMA patterns (*Petroni et al., 2019*)
  - Augmenting with LPAQA patterns (*Jiang et al., 2020*)
  - Searching for patterns in wikipedia using SPIKE (*Shlain et al., 2020*)

(a) A single pattern



(b) Multiple patterns, noisy

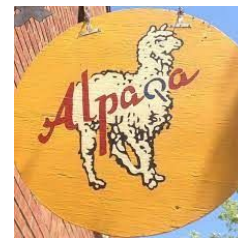
(c) Searching for syntactic patterns



# ParaRel 🙌 - Creation

- For every relation, we manually build a set of patterns that are paraphrases of each other, in 4 steps:
  - Starting with the LAMA patterns (*Petroni et al., 2019*)
  - Augmenting with LPAQA patterns (*Jiang et al., 2020*)
  - Searching for patterns in wikipedia using SPIKE (*Shlain et al., 2020*)
  - Additional patterns using linguistic expertise

(a) A single pattern



(b) Multiple patterns, noisy

(d) linguistic expertise, expanding previous patterns



(c) Searching for syntactic patterns



# ParaRel 🙌 - Summary

aired-on

[X] was aired on [Y].  
[X], that was aired on [Y].  
[Y]'s series [X]

instrument

[X] plays [Y].  
[Y] player [X].  
[X] is a [Y] player.

# Relations	38
# Patterns	328
Min # patterns	2
Max # patterns	20
Avg # patterns	8.63

employer

[X] used to work in [Y].  
[X] found employment in [Y].  
[X] took up work in [Y].

twin-cities

[X] and [Y] are twin cities.  
[Y] and [X] are twin cities.  
[X] is a twin city of [Y].

# ParaRel 🤘 - Creation

- For every relation, we manually build a set of patterns that are paraphrases of each other, in 4 steps:
  - Starting with the single pattern from LAMA (Petroni et al., 2019)
  - Augmenting with automatically extracted patterns from LPAQA (Jiang et al., 2020)
  - Searching for patterns in wikipedia using SPIKE (Shlain et al., 2020)
  - Additional patterns using linguistic expertise of the authors

# ParaRel 🙌 - Verification

- Was collected manually by the authors of this paper
- 2 additional authors verified the quality, while engaging in discussion to reach an agreement (discarding otherwise)
- Human Eval: Sampled 156 pairs, and ask NLP grad students to annotate. Reaching **95.5%** agreement (and later fixed the errors)

# Setup & Evaluation

# Consistency - Setup

Data Pairs ( $D$ )

$D_1$  (*Lou Reed, Brooklyn*)  
(*Masako Natsume, Tokyo*)  
...  
...  
(*Seinfeld, NBC*)  
 $D_i$  (*Homeland, Showtime*)  
...  
...

$(D_1, r_1, P_1), \dots, (D_i, r_i, P_i), \dots, (D_n, r_n, P_n)$



$r_i = \text{originally-aired-on}$

(*Homeland* originally aired on [MASK]  
*Homeland* premiered on [MASK]  
...  
*Seinfeld* originally aired on [MASK]  
*Seinfeld* premiered on [MASK])

Patterns ( $P$ )

( $X$  was born in  $Y$ )  
( $X$  is native to  $Y$ )  $P_1$   
...  
...  
( $X$  originally aired in  $Y$ )  
( $X$  premiered on  $Y$ )  $P_i$   
...  
...

# Consistency - Setup

## Data Pairs ( $D$ )

(*Lou Reed, Brooklyn*)

$D_1$  (*Masako Natsume, Tokyo*)

...

...

(*Seinfeld, NBC*)

$D_i$  (*Homeland, Showtime*)

...

...

$(D_1, r_1, P_1), \dots, (D_i, r_i, P_i), \dots, (D_n, r_n, P_n)$



$r_i = \text{originally-aired-on}$

(*Homeland* originally aired on [MASK])

*Homeland* premiered on [MASK]

...

(*Seinfeld* originally aired on [MASK])

*Seinfeld* premiered on [MASK]

## Patterns ( $P$ )

( $X$  was born in  $Y$ )

( $X$  is native to  $Y$ )  $P_1$

...

...

( $X$  originally aired in  $Y$ )

( $X$  premiered on  $Y$ )  $P_i$

...

...



# Consistency - Setup

Data Pairs ( $D$ )

$D_1$  (*Lou Reed, Brooklyn*)  
(*Masako Natsume, Tokyo*)  
...  
...  
(*Seinfeld, NBC*)  
 $D_i$  (*Homeland, Showtime*)  
...  
...

$(D_1, r_1, P_1), \dots, (D_i, r_i, P_i), \dots, (D_n, r_n, P_n)$



$r_i = \text{originally-aired-on}$

(*Homeland* originally aired on [MASK])  
(*Homeland* premiered on [MASK])  
...  
(*Seinfeld* originally aired on [MASK])  
(*Seinfeld* premiered on [MASK])

Patterns ( $P$ )

( $X$  was born in  $Y$ )  
( $X$  is native to  $Y$ )  $P_1$   
...  
...  
( $X$  originally aired in  $Y$ )  
( $X$  premiered on  $Y$ )  $P_i$   
...  
...

# Consistency - Setup

Data Pairs ( $D$ )

$D_1$  (*Lou Reed, Brooklyn*)  
(*Masako Natsume, Tokyo*)  
...  
...  
(*Seinfeld, NBC*)  
 $D_i$  (*Homeland, Showtime*)  
...  
...

$(D_1, r_1, P_1), \dots, (D_i, r_i, P_i), \dots, (D_n, r_n, P_n)$



$r_i = \text{originally-aired-on}$

(*Homeland* originally aired on [MASK])  
(*Homeland* premiered on [MASK])  
...  
(*Seinfeld* originally aired on [MASK])  
(*Seinfeld* premiered on [MASK])

Patterns ( $P$ )

( $X$  was born in  $Y$ )  
( $X$  is native to  $Y$ )  $P_1$   
...  
...  
( $X$  originally aired in  $Y$ )  
( $X$  premiered on  $Y$ )  $P_i$   
...  
...

# Consistency - Setup

Data Pairs ( $D$ )

$D_1$  (*Lou Reed, Brooklyn*)  
(*Masako Natsume, Tokyo*)  
...  
...  
(*Seinfeld, NBC*)  
 $D_i$  (*Homeland, Showtime*)  
...  
...

$(D_1, r_1, P_1), \dots, (D_i, r_i, P_i), \dots, (D_n, r_n, P_n)$



$r_i = \text{originally-aired-on}$

(*Homeland* originally aired on [MASK]  
*Homeland* premiered on [MASK]  
...  
*Seinfeld* originally aired on [MASK]  
*Seinfeld* premiered on [MASK])

Patterns ( $P$ )

( $X$  was born in  $Y$ )  
( $X$  is native to  $Y$ )  $P_1$   
...  
...  
( $X$  originally aired in  $Y$ )  
( $X$  premiered on  $Y$ )  $P_i$   
...  
...

# Consistency - Models

- BERT
- BERT Whole-Word-Masking
- RoBERTa
- ALBERT

And a Baseline:

- Most common object (consistent by definition)

# Consistency - Evaluation

- **Accuracy:** Accurate prediction of the LAMA pattern
- **Consistency:** For each relation and tuple, compute all paraphrases pairs, and test if the predictions are equal:  $n(n-1)/2$  pairs
- **Consistent-Acc:** Consistent and accurate prediction of all paraphrases

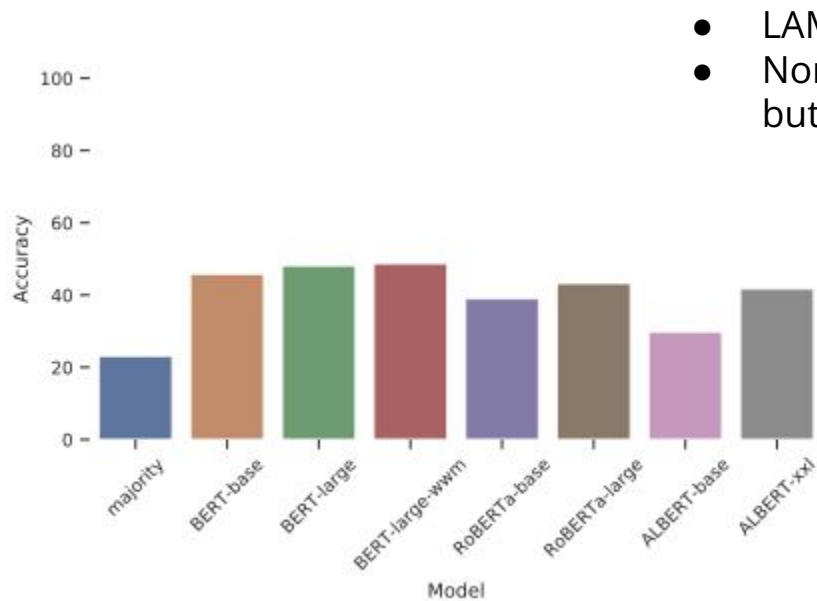
# Results

# Consistency - Results

Are LMs Consistent?

No!

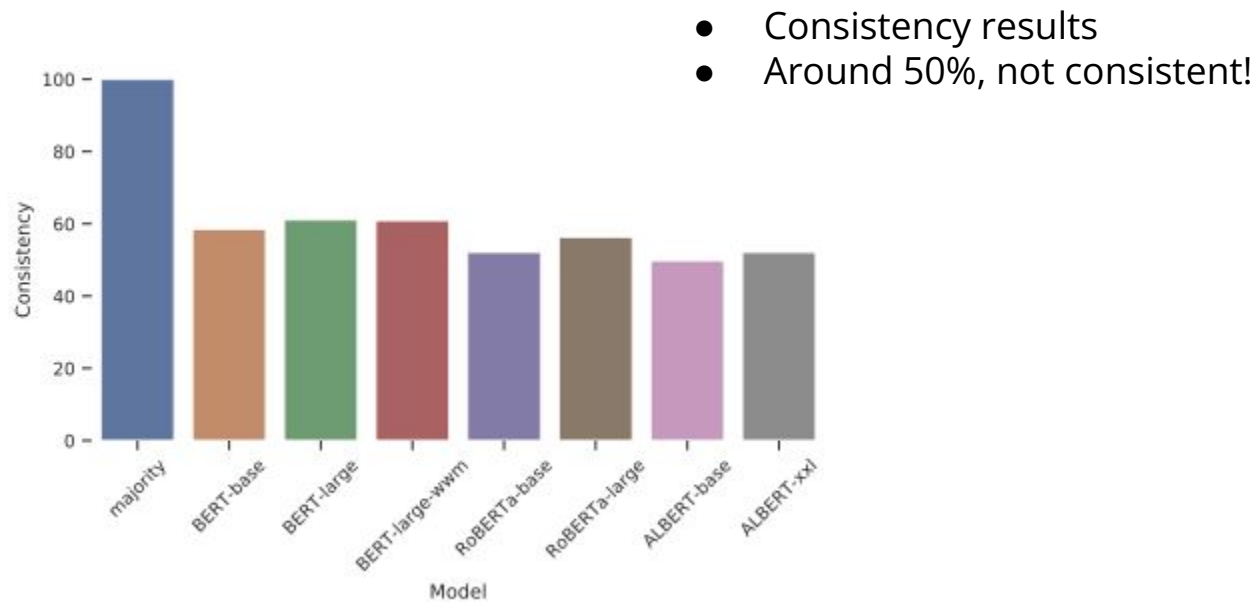
# Consistency - Results



- LAMA accuracy performance
- Non-trivial retrieval abilities (~40%), but not good in any way

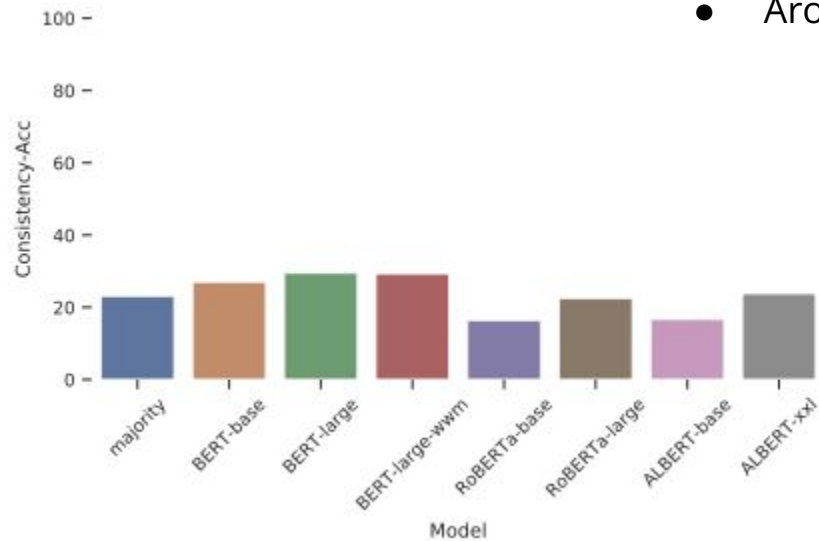


# Consistency - Results

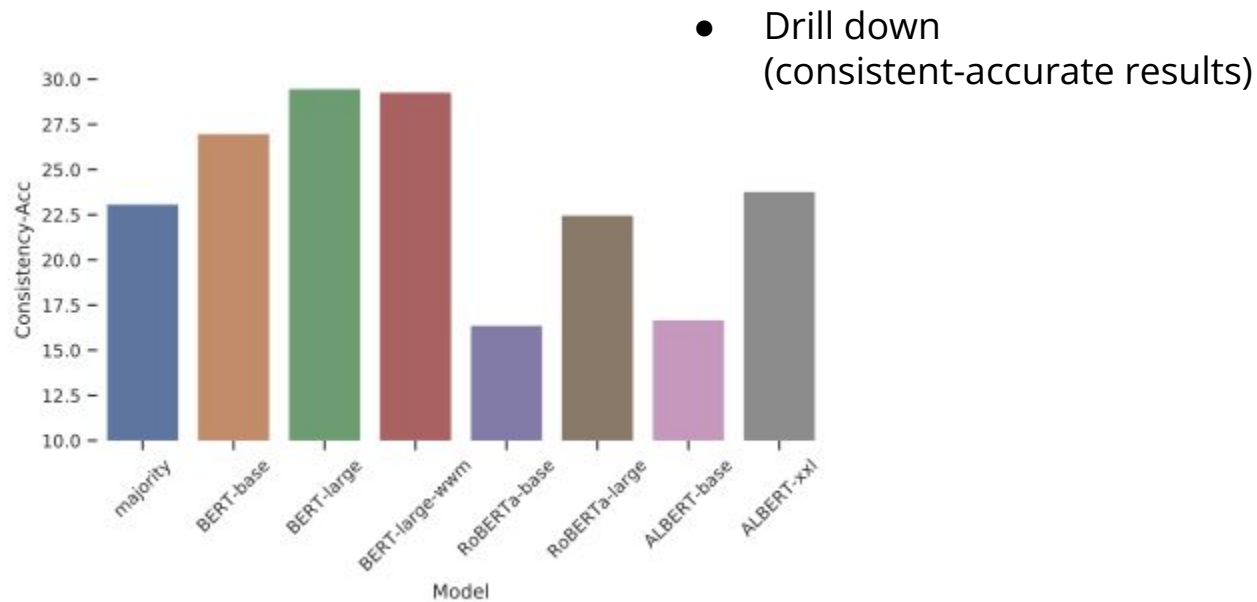


# Consistency - Results

- Consistent-accurate results
- Around 20-30%, much worse!

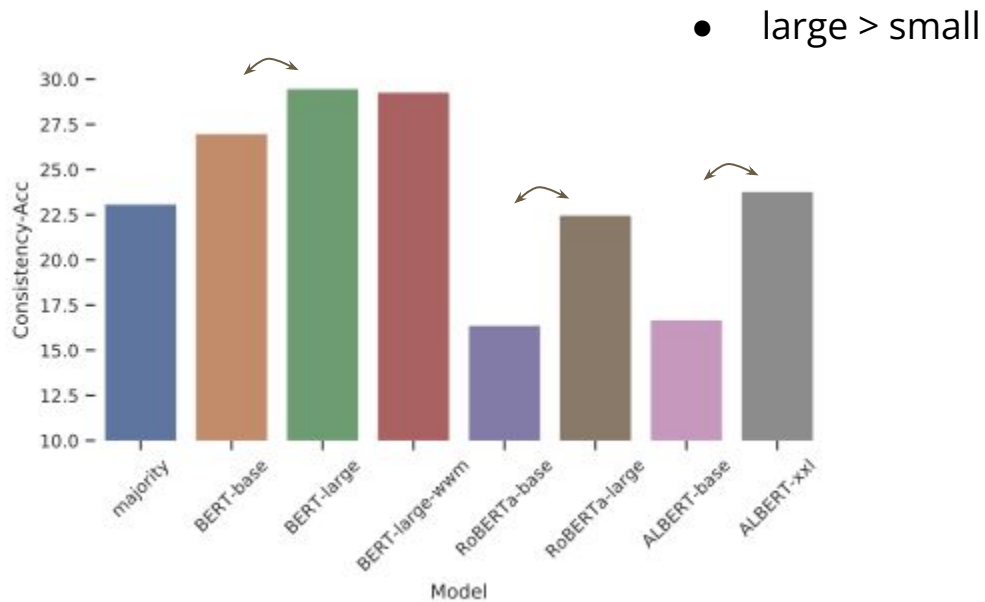


# Consistency - Results



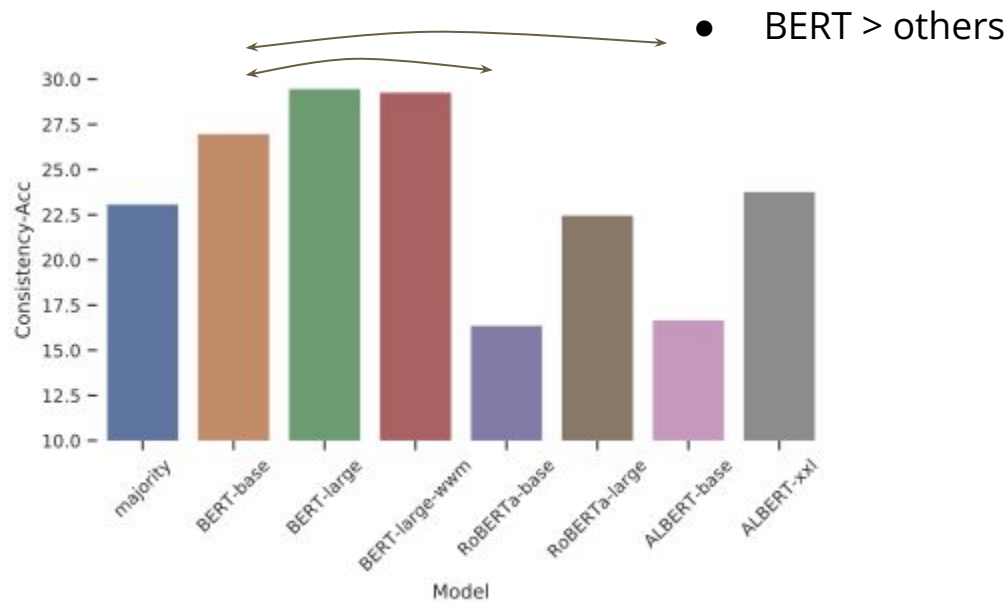
# Consistency - Results

Interesting trends: *base vs. large*



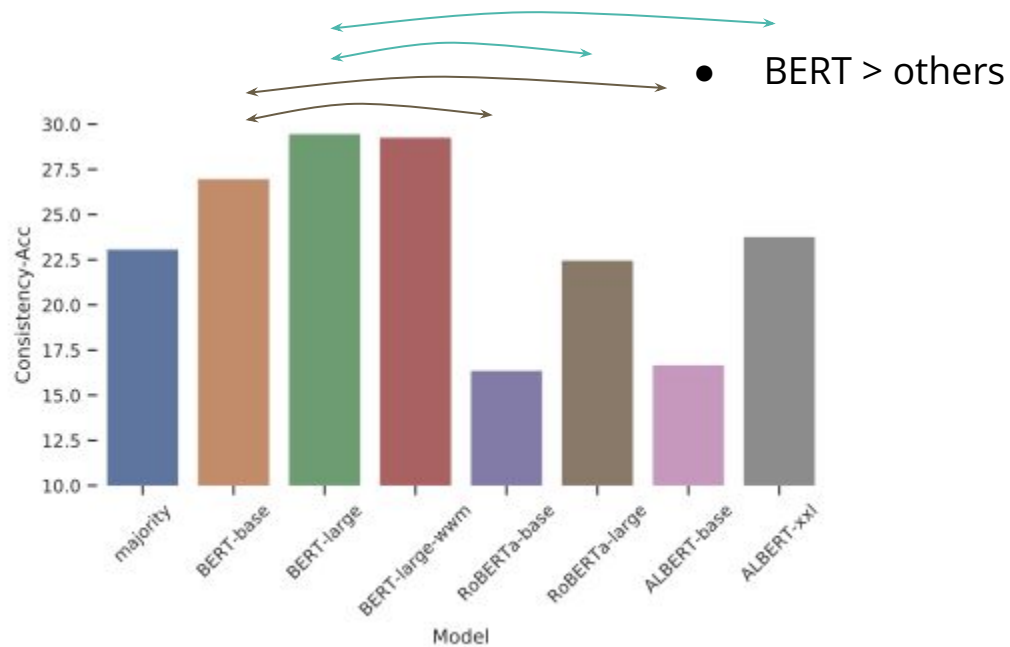
# Consistency - Results

Interesting trends: *BERT vs. others*



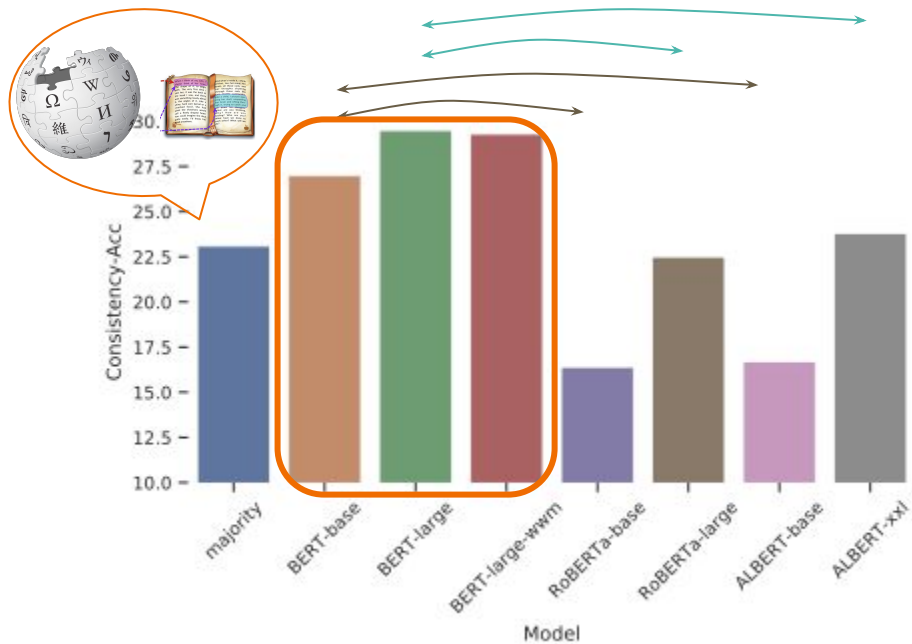
# Consistency - Results

Interesting trends: *BERT vs. others*

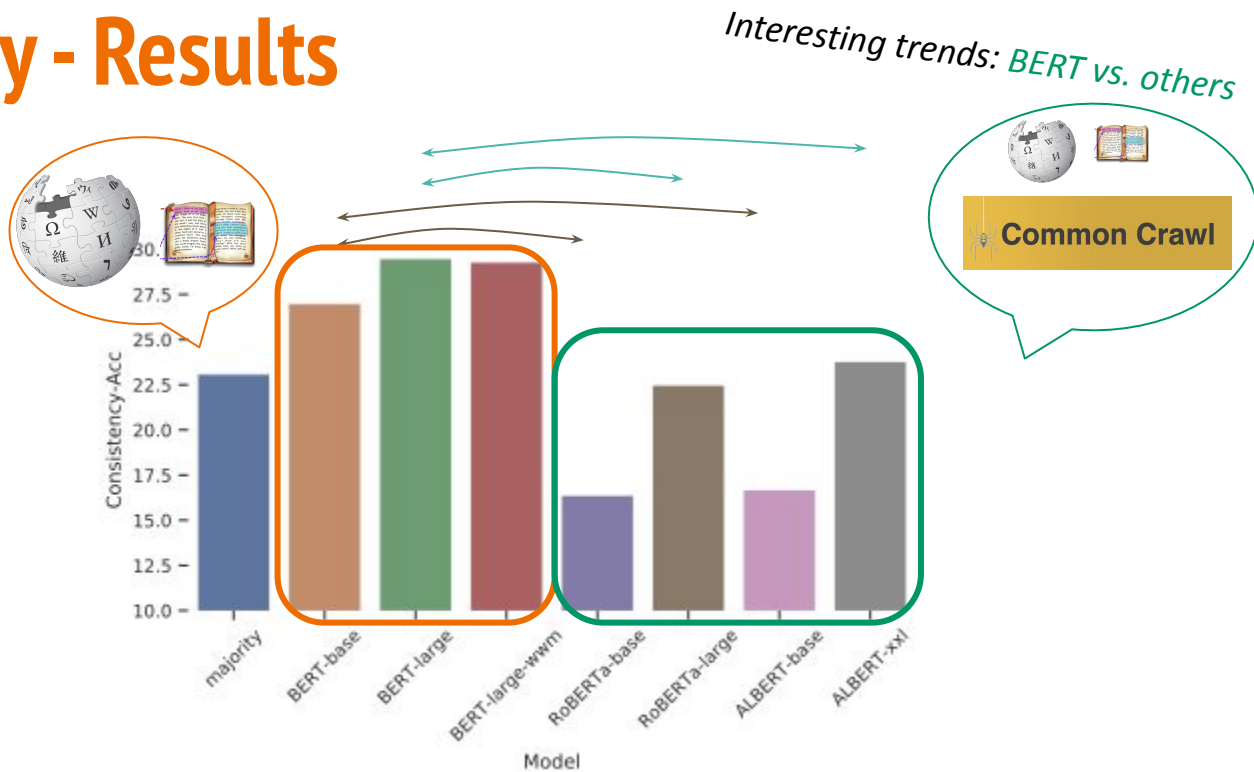


# Consistency - Results

Interesting trends: *BERT vs. others*



# Consistency - Results





# Consistency - Summary

We have shown that:

- Some relations are more consistent than others
- Some models are more consistent than others

But overall, **models are inconsistent!**

*Much more analysis and experiments in the paper!!*

# Improving Consistency

# Improving Consistency

*More details in the paper!*

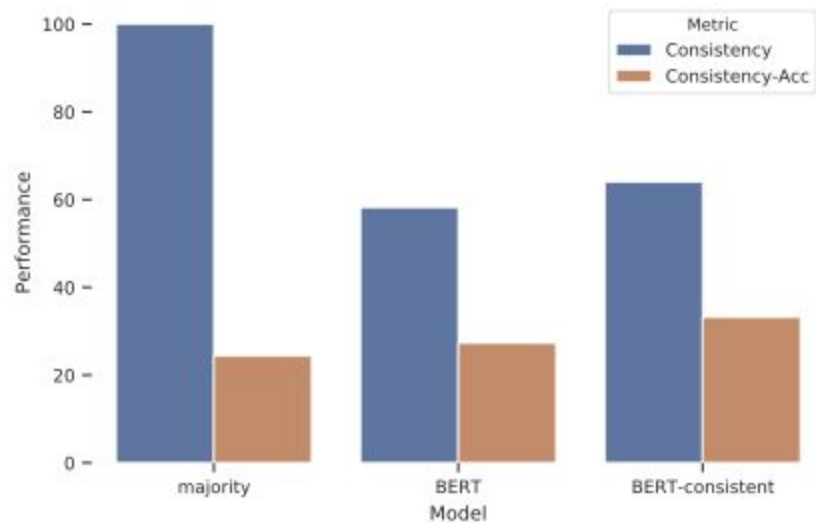
- Can we improve the consistency of PLMs?
- We want predictions from paraphrases to be equal
- We try to make the distributions alike

$$Q_n = \text{softmax}(f_\theta(P_n))$$

$$\mathcal{L}_c = \sum_{n=1}^k \sum_{m=n+1}^k D_{KL}(Q_n^{r_i} || Q_m^{r_i}) + D_{KL}(Q_m^{r_i} || Q_n^{r_i})$$

$$\mathcal{L} = \lambda \mathcal{L}_c + \mathcal{L}_{MLM}$$

# Improved Consistency



## Part III

# Explaining The “Knowledge”

Work In Progress

# Explaining Knowledge in PLMs

## LM predictions

#1 mask: Barack Obama was born in [MASK].

bert_large_cased	
0	Chicago
1	Philadelphia
2	Detroit
3	Houston
4	Atlanta
5	Georgia
6	Boston
7	Texas
8	Paris
9	Dallas

roberta_large	
0	Kenya
1	Hawaii
2	1961
3	1964
4	Chicago
5	Honolulu
6	Indonesia
7	1965
8	1969
9	1963

WHY???

# Explaining Knowledge in PLMs

## LM predictions

#1 mask: Barack Obama was born in [MASK].

	bert_large_cased
0	Chicago
1	Philadelphia
2	Detroit
3	Houston
4	Atlanta
5	Georgia
6	Boston
7	Texas
8	Paris
9	Dallas

	roberta_large
0	Kenya
1	Hawaii
2	1961
3	1964
4	Chicago
5	Honolulu
6	Indonesia
7	1965
8	1969
9	1963

### University of Chicago Law School and civil rights attorney

In 1991, Obama accepted a two-year position as Visiting Law and Government Fellow at the University of Chicago Law School to work on his first book<sup>[120][121]</sup> He then taught constitutional law at the University of Chicago Law School for twelve years, first as a lecturer from 1992 to 1996, and then as a senior lecturer from 1996 to 2004.<sup>[122]</sup>

From April to October 1992, Obama directed Illinois's Project Vote, a voter registration campaign with ten staffers and seven hundred volunteer registrars; it achieved its goal of registering 150,000 of 400,000 unregistered African Americans in the state, leading *Crain's Chicago Business* to name Obama to its 1993 list of "40 under Forty" powers to be.<sup>[123]</sup>

He joined Davis, Miner, Barnhill & Galland, a 13-attorney law firm specializing in civil rights litigation and neighborhood economic development, where he was an associate for three years from 1993 to 1996, then of counsel from 1996 to 2004. In 1994, he was listed as one of the lawyers in *Buycks-Roberson v. Citibank Fed. Sav. Bank*, 94 C 4094 (N.D. Ill.).<sup>[124]</sup> This class action lawsuit was filed in 1994 with Selma Buycks-Roberson as lead plaintiff and alleged that Citibank Federal Savings Bank had engaged in practices forbidden under the Equal Credit Opportunity Act and the Fair Housing Act.<sup>[125]</sup>

The case was settled out of court.<sup>[126]</sup> Final judgment was issued on May 13, 1998, with Citibank Federal Savings Bank agreeing to pay attorney fees.<sup>[127]</sup> His law license became inactive in 2007.<sup>[128][129]</sup>

From 1994 to 2002, Obama served on the boards of directors of the Woods Fund of Chicago—which in 1985 had been the first foundation to fund the Developing Communities Project—and of the Joyce Foundation.<sup>[56]</sup> He served on the board of directors of the Chicago Annenberg Challenge from 1995 to 2002, as founding president and chairman of the board of directors from 1995 to 1999.<sup>[56]</sup>



Search

0 Obama Born In Kenya? His Grandmother Says Yes

Posted by u/Slipgrid 12 years ago

0 Obama Born In Kenya? His Grandmother Says Yes

israelenews.com/view.a... [🔗](#)

# Explaining Knowledge in PLMs

Data as a source of explanations



# Explaining Knowledge in PLMs

- Taking another look at the data: Wikipedia

## Background

### Early life of Barack Obama

*Main articles: [Early life and career of Barack Obama](#) and [Ann Dunham](#)*

People who express doubts about Obama's eligibility or reject details about his early life are often informally called "birthers", a term that parallels<sup>[23]</sup> the nickname "truthers" for adherents of [9/11 conspiracy theories](#).<sup>[24][25]</sup> These [conspiracy theorists](#) reject at least some of the following facts about his early life:

Barack Obama was born on August 4, 1961, at Kapi'olani Maternity & Gynecological Hospital (now called [Kapi'olani Medical Center for Women & Children](#)) in Honolulu, Hawaii,<sup>[26][27][28][29]</sup> to [Ann Dunham](#),<sup>[30]</sup> from [Wichita, Kansas](#),<sup>[31]</sup> and her husband [Barack Obama Sr.](#), a Luo from [Nyang'oma Kogelo, Nyanza Province](#) (in what was then the [Colony and Protectorate of Kenya](#)), who was attending the University of Hawaii. Birth notices for Barack Obama were published in *The Honolulu Advertiser* on August 13 and the *Honolulu Star-Bulletin* on August 14, 1961.<sup>[26][31]</sup> Obama's father's immigration file also clearly states Barack Obama was born in Hawaii.<sup>[32]</sup> One of his high school teachers, who was acquainted with his mother at the time, remembered hearing about the day of his birth.<sup>[30]</sup>

### LM predictions

#1 mask:Barack Obama was born in [MASK].

bert_large_cased	
0	Chicago
1	Philadelphia
2	Detroit
3	Houston
4	Atlanta
5	Georgia
6	Boston
7	Texas
8	Paris
9	Dallas

# Explaining Knowledge in PLMs

- Taking another look at the data: Wikipedia

## Background

### Early life of Barack Obama

Main articles: *Early life and career of Barack Obama* and *Ann Dunham*

People who express doubts about Obama's eligibility or reject details about his early life are often informally called "birthers", a term that parallels<sup>[23]</sup> the nickname "truthers".

Barack Obama was born in Honolulu (in what is now Hawaii).



WIKIPEDIA  
The Free Encyclopedia

chicago

1/97



hawaii

1/16



## Barack Obama

From Wikipedia, the free encyclopedia

### LM predictions

#1 mask: Barack Obama was born in [MASK].

bert_large_cased	
0	Chicago
1	Philadelphia
2	Detroit
3	Houston
4	Atlanta
5	Georgia
6	Boston
7	Texas
8	Paris
9	Dallas

WHY???

# Explaining Knowledge in PLMs

- Maybe these models rely on co-occurrences?

What else?

(How to predict a word, given a cloze sentence such as:  
“Barack Obama was born in [MASK].”)

# Pitfalls of LMs as KBs

# Explaining Knowledge in PLMs

- We inspect 3 pitfalls (or heuristics) in LMs with respect to knowledge extraction

# Explaining Knowledge in PLMs

- We inspect 3 pitfalls (or heuristics) in LMs with respect to knowledge extraction
  - Model ignores the subject, and only uses the pattern to make prediction
- Example:
  - Barack Obama was born in [MASK]. (*born-in*)

# Explaining Knowledge in PLMs

- We inspect 3 pitfalls (or heuristics) in LMs with respect to knowledge extraction
  - Model ignores the subject, and only uses the pattern to make prediction
  - Model ignores the pattern, and only uses the subject to make prediction
- Example:
  - Barack Obama was born in [MASK]. (*born-in*)
  - Barack Obama ~~was born in~~ [MASK]. (~~*born-in*~~)

# Explaining Knowledge in PLMs

- We inspect 3 pitfalls (or heuristics) in LMs with respect to knowledge extraction
  - Model ignores the subject, and only uses the pattern to make prediction
  - Model ignores the pattern, and only uses the subject to make prediction
  - Model ignores the abstract relation, and uses the subject+pattern to make prediction
- Example:
  - ~~Barack Obama~~ was born in [MASK]. (*born-in*)
  - Barack Obama ~~was born in~~ [MASK]. (~~*born-in*~~)
  - Barack Obama was born in [MASK]. (~~*born-in*~~)



# Explaining Knowledge in PLMs

- We inspect 3 pitfalls (or heuristics) in LMs with respect to knowledge extraction
  - Model ignores the subject, and only uses the pattern to make prediction
  - Model ignores the pattern, and only uses the subject to make prediction
  - Model ignores the abstract relation, and uses the subject+pattern to make prediction
- Example:
  - Barack Obama was born in [MASK]. (*born-in*) ← *Pattern's preference*
  - Barack Obama ~~was born in~~ [MASK]. (~~*born-in*~~) ← *Subj-obj cooccurrences*
  - Barack Obama was born in [MASK]. (~~*born-in*~~) ← *Memorization*

# Explaining Knowledge in PLMs

- Example:
  - ~~Barack Obama~~ was born in [MASK]. (*born-in*)
  - Barack Obama ~~was born in~~ [MASK]. (*born-in*)
  - Barack Obama was born in [MASK]. (~~born-in~~)
  
- Tests using the model:
  - Default Behavior
    - was born in [MASK]
  - Entities association
    - Barack Obama died in [MASK]
  - Consistency
    - Barack Obama, born in [MASK].

# Explaining Knowledge in PLMs

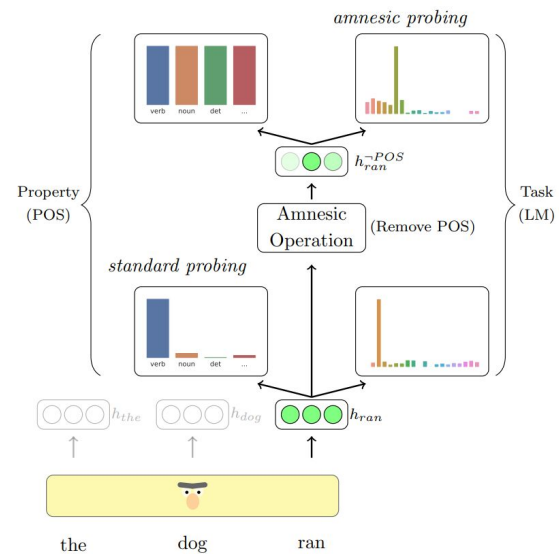
- Example:
  - ~~Barack Obama~~ was born in [MASK]. (*born-in*)
  - Barack Obama ~~was born in~~ [MASK]. (*born-in*)
  - Barack Obama was born in [MASK]. (~~*born-in*~~)
- Explaining the Model through the Data:
  - Occurrences of pattern+object (count in wiki: “was born in Hawaii”)
  - Entities Co-occurrence (count <“Barack Obama, Hawaii”>, <“Barack Obama, Chicago”>, ...)
  - Memorization (count “Barack Obama was born in Hawaii”)

# Explaining Knowledge in PLMs

- Entities Association:
  - Probability that BERT predicts the most co-occurred entity when the pattern describes a correct relation is 40%, compared to 35% when the relation doesn't hold
- Similar trends for the memorization
- But this is not a causal attribution!

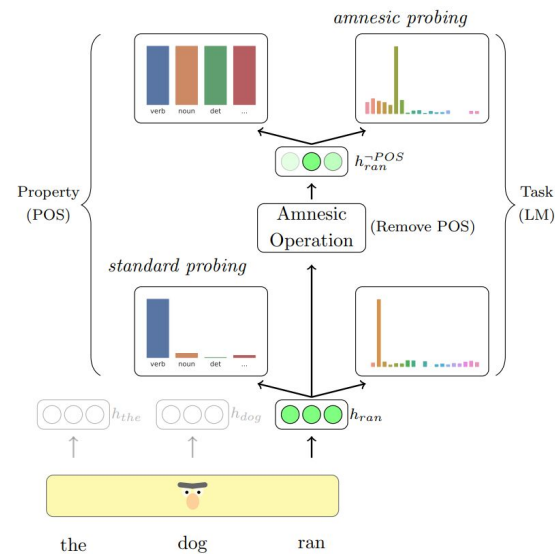
# Causal Explanation through the Data

- Can't use *amnesic probing*
  - Concepts aren't clear



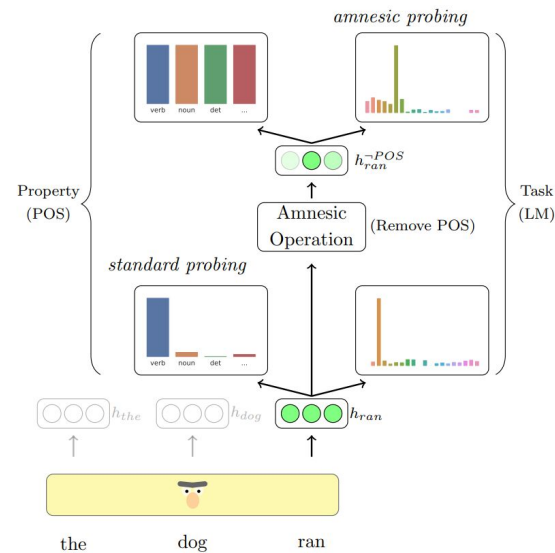
# Causal Explanation through the Data

- Can't use *amnesic probing*
  - Concepts aren't clear
- Can't perform intervention on the data
  - Retraining BERT (on each combination) is expensive



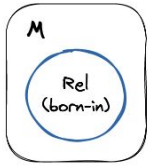
# Causal Explanation through the Data

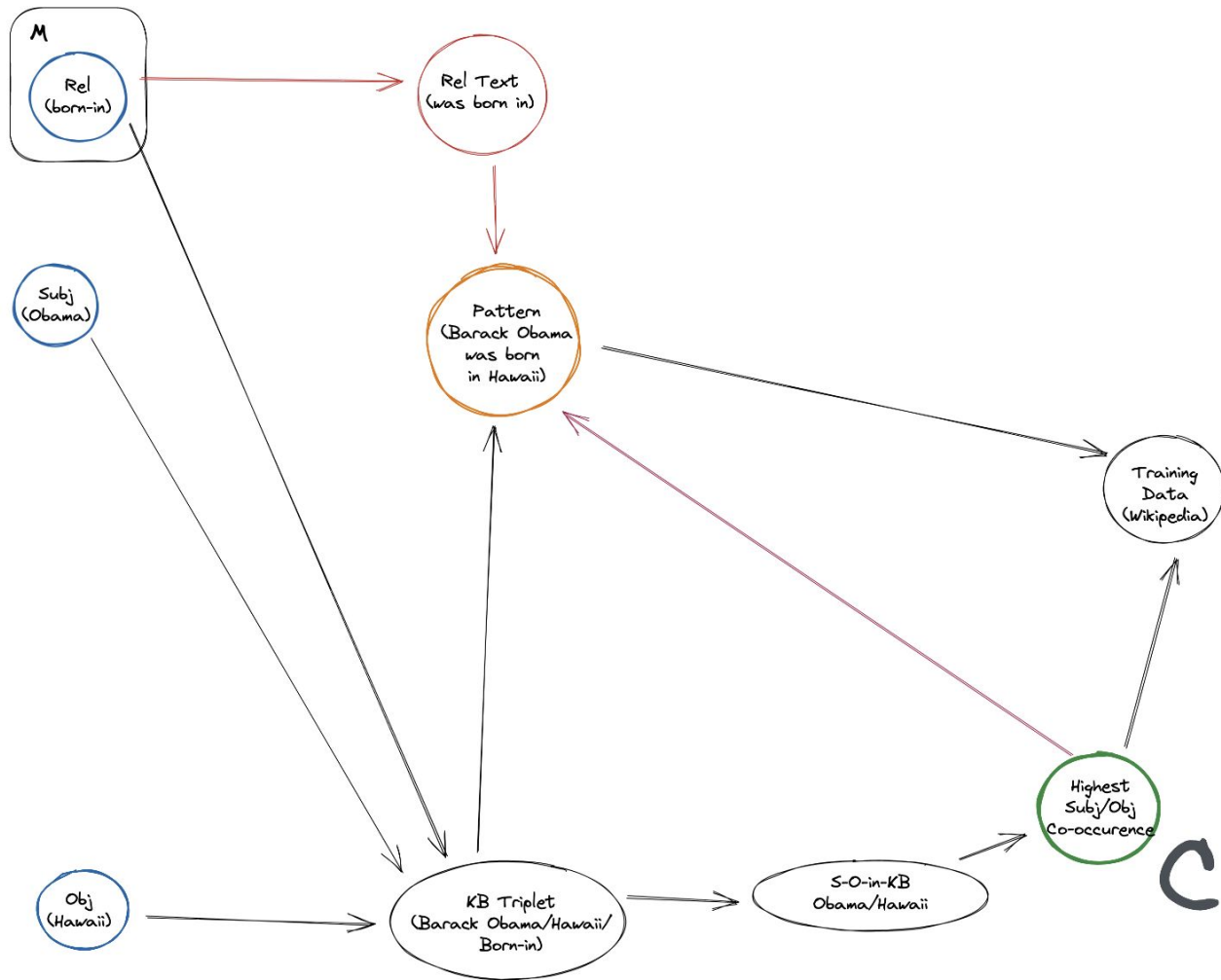
- Can't use *amnesic probing*
  - Concepts aren't clear
- Can't perform intervention on the data
  - Retraining BERT (on each combination) is expensive
- Solution: Measure **Average Treatment Effect (ATE)** using observational data
  - Assuming we can observe the measurable variables

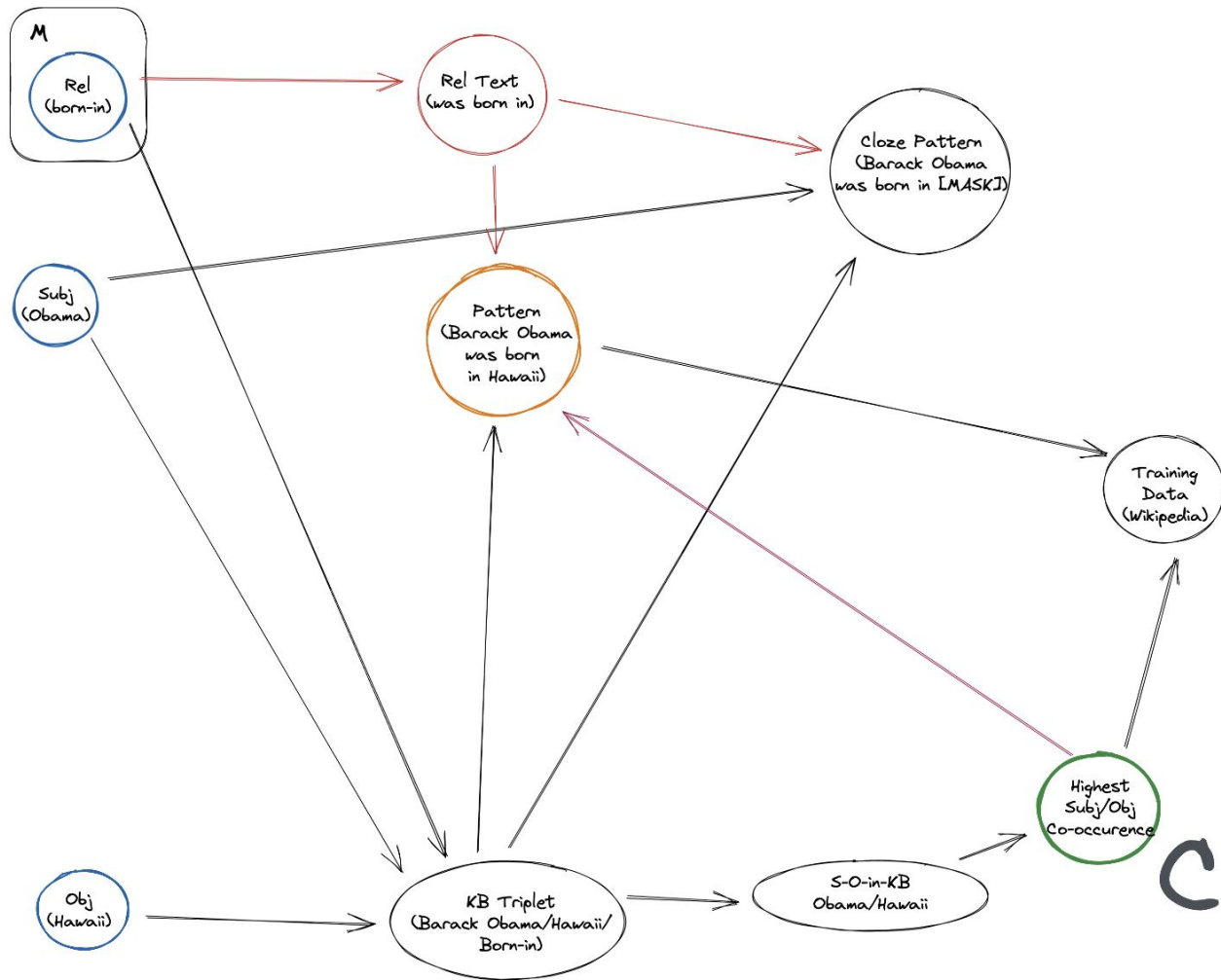


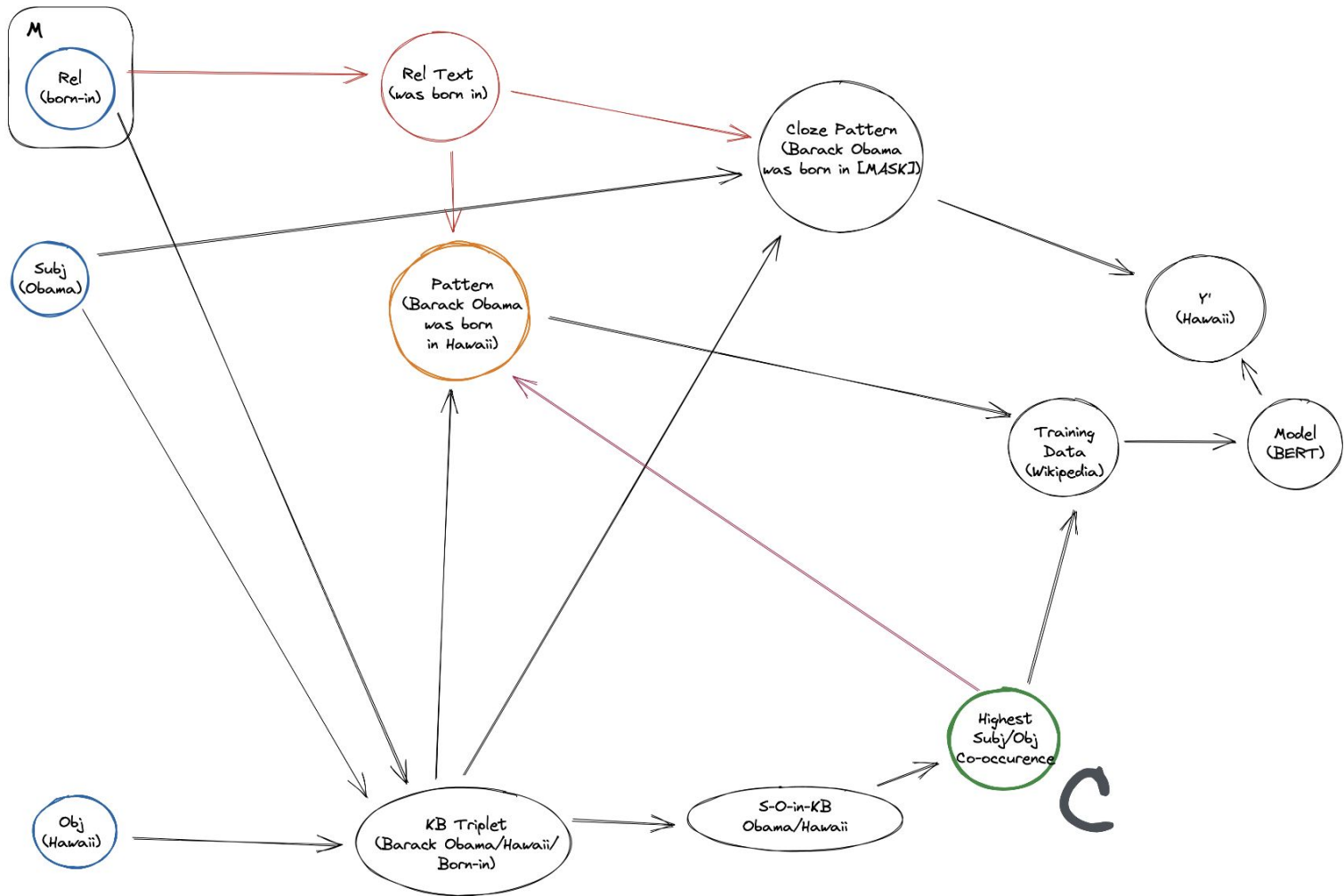
# Causal Diagram

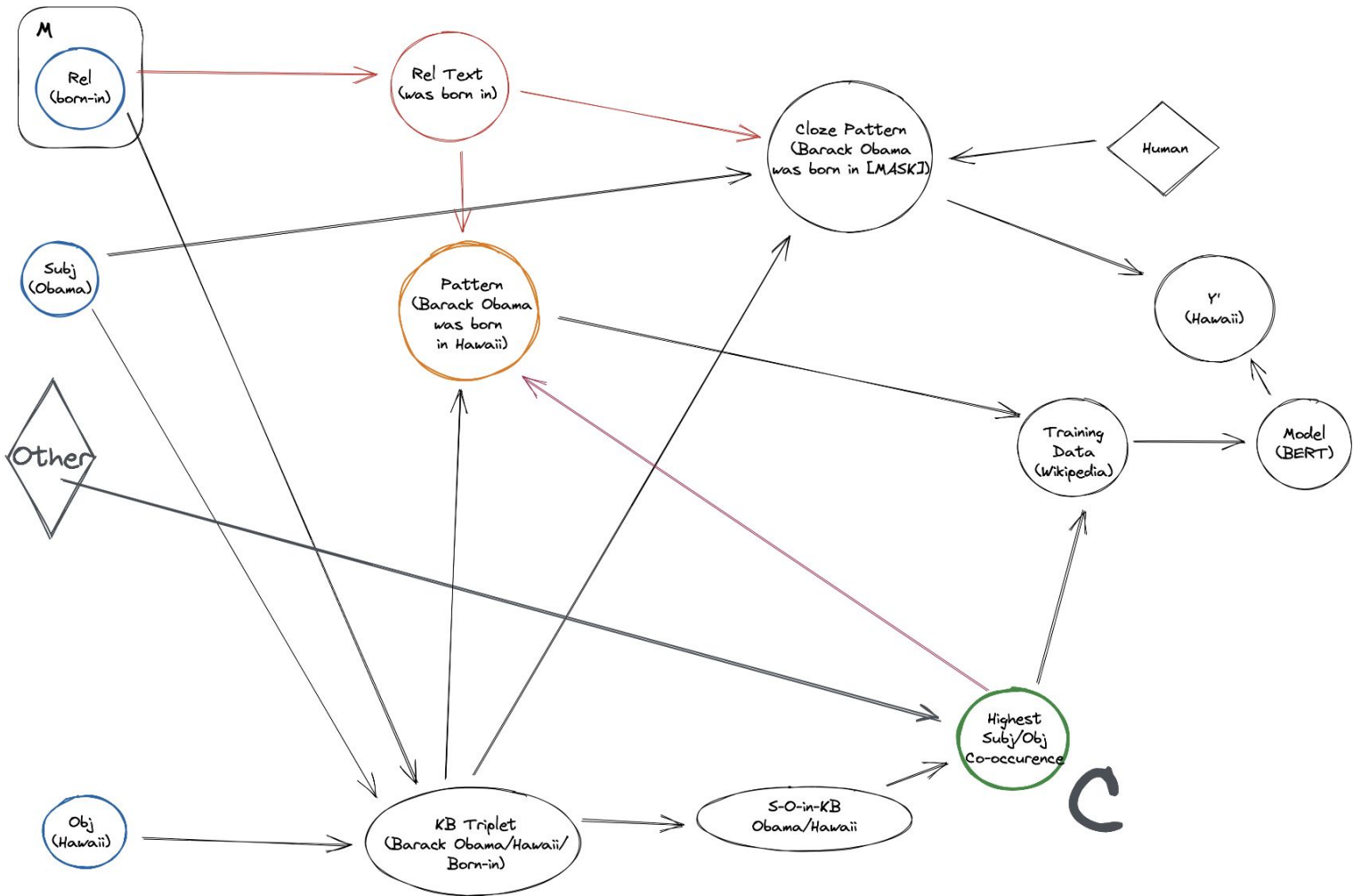


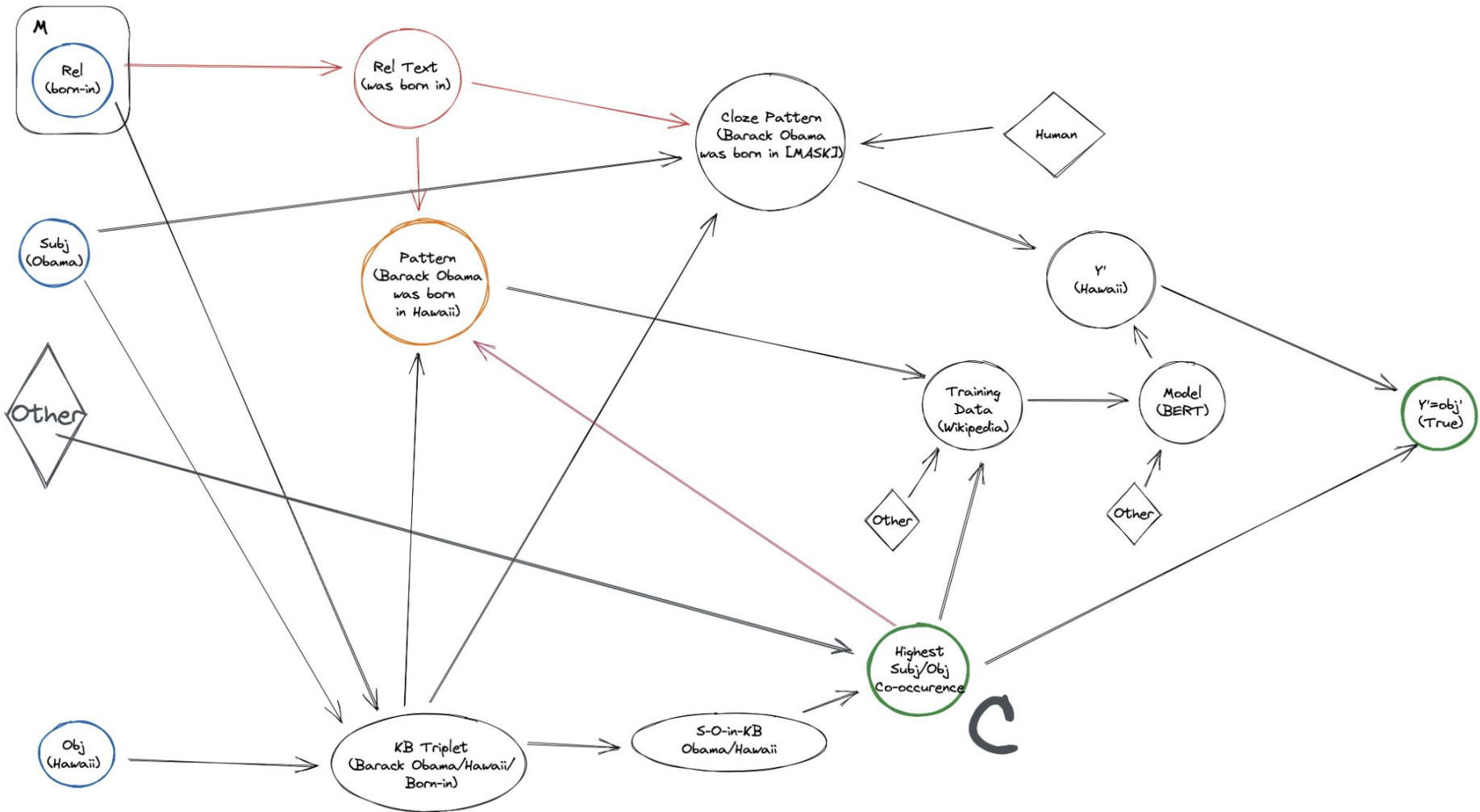












# Explaining Knowledge - Causal Explanation

- Given that we believe this graph accurately describes the world...
- ... and we find the relevant back/front door criterion to control for confounding variables
- We can measure the effect of the heuristics on the models' predictions

NEAT! AND A STRONG RESULT

# Explaining Knowledge - Causal Explanation

- If the effect is strong, what does it tell us about this model?
  - The model memorize, and uses correlations for making predictions
  - It has a limited understanding of linguistic relations
  - More?



# Results

- Example:

- ~~Barack Obama~~ was born in [MASK]. (*born-in*) ← *Pattern's preference*
- Barack Obama ~~was born in~~ [MASK]. (~~*born in*~~) ← *Subj-obj cooccurrences*
- Barack Obama was born in [MASK]. (~~*born in*~~) ← *Memorization*

Hypothesis	ATE
Pattern's Preference	4.1
Subj-Obj Co-occurrence	19.0
Memorization	10.4

# Data as a Source of Explanations

# Summary

- **Amnesic Probing:** a method that answers a causal question: “what is being used?”
- **Consistency** of PLMs knowledge is limited
- **Data as Explanation:** A graph describing causal relations
  - Allows to ask how concepts/heuristics **associated with training data** are used by models

**Thanks!**  
**Questions?**

Yanai Elazar

@yanaiela 

[yanaiela.github.io](https://yanaiela.github.io)

---