

Causal abstractions of neural networks

Christopher Potts

Joint work with Atticus Geiger, Hanson Lu, Thomas Icard,
and Noah Goodman

Stanford Linguistics and the Stanford NLP Group

Amazon, August 31, 2021



My ACL talk engaging with 'NLP for Social Good'



My ACL talk engaging with 'NLP for Social Good'

Reliable characterizations of NLP systems as a
social responsibility

My ACL talk engaging with 'NLP for Social Good'

Reliable characterizations of NLP systems as a social responsibility

1. Benchmark datasets: Delimit responsible use
2. System assessment: Connect with real-world concerns

My ACL talk engaging with 'NLP for Social Good'

Reliable characterizations of NLP systems as a social responsibility

1. Benchmark datasets: Delimit responsible use
2. System assessment: Connect with real-world concerns

Do exactly what you said you would do.

My ACL talk engaging with 'NLP for Social Good'

Reliable characterizations of NLP systems as a social responsibility

1. Benchmark datasets: Delimit responsible use
2. System assessment: Connect with real-world concerns
3. Structural evaluation methods: Seek guarantees

Do exactly what you said you would do.



Overview: Structural evaluation methods

Overview: Structural evaluation methods

1. Motivations



Overview: Structural evaluation methods

1. Motivations
2. Probing



Overview: Structural evaluation methods

1. Motivations
2. Probing
3. Feature attribution



Overview: Structural evaluation methods

1. Motivations
2. Probing
3. Feature attribution
4. Causal abstraction

Overview: Structural evaluation methods

1. Motivations
2. Probing
3. Feature attribution
4. Causal abstraction
5. Case study: Monotonicity NLI

Motivations

Systematicity

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

1. Sandy loves the puppy.

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

1. Sandy loves the puppy.
2. The puppy loves Sandy.

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

1. Sandy loves the puppy.
2. The puppy loves Sandy.
3. the turtle ~ the puppy

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

1. Sandy loves the puppy.
2. The puppy loves Sandy.
3. the turtle ~ the puppy
4. The turtle loves the puppy.

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

1. Sandy loves the puppy.
2. The puppy loves Sandy.
3. the turtle ~ the puppy
4. The turtle loves the puppy.
5. The puppy loves the turtle.

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

1. Sandy loves the puppy.
2. The puppy loves Sandy.
3. the turtle ~ the puppy
4. The turtle loves the puppy.
5. The puppy loves the turtle.
6. The turtle loves Sandy.

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

1. Sandy loves the puppy.
2. The puppy loves Sandy.
3. the turtle ~ the puppy
4. The turtle loves the puppy.
5. The puppy loves the turtle.
6. The turtle loves Sandy.
7. ...

Systematicity

Fodor and Pylyshyn (1988:37):

“What we mean when we say that linguistic capacities are *systematic* is that the ability to produce/understand some sentences is *intrinsically* connected to the ability to produce/understand certain others.”

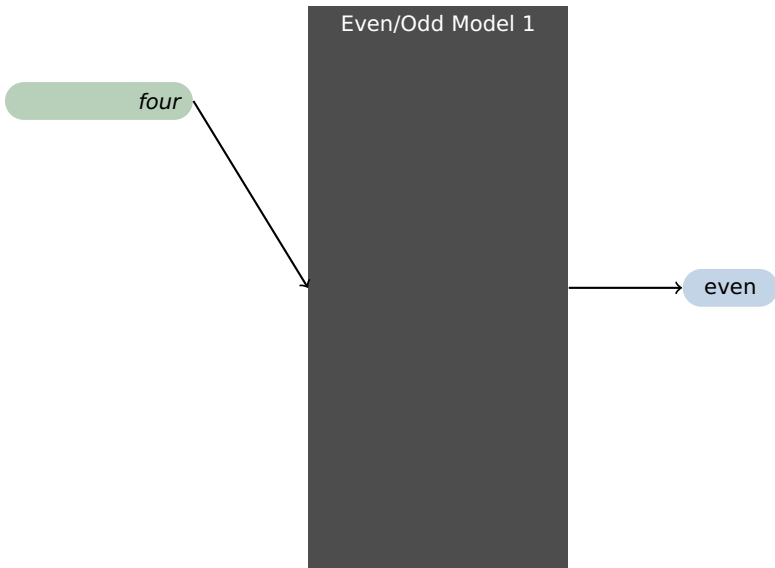
Example	Gold	Prediction
The bakery sells a mean apple pie.	pos	pos
They sell a mean apple pie.	pos	pos
She sells a mean apple pie.	pos	neg
He sells a mean apple pie.	pos	neg

Limits of behavioral testing

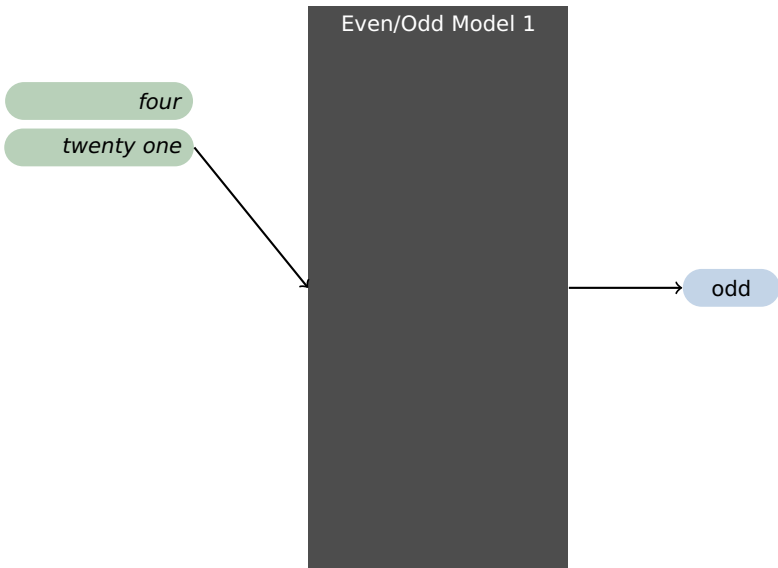
Even/Odd Model 1



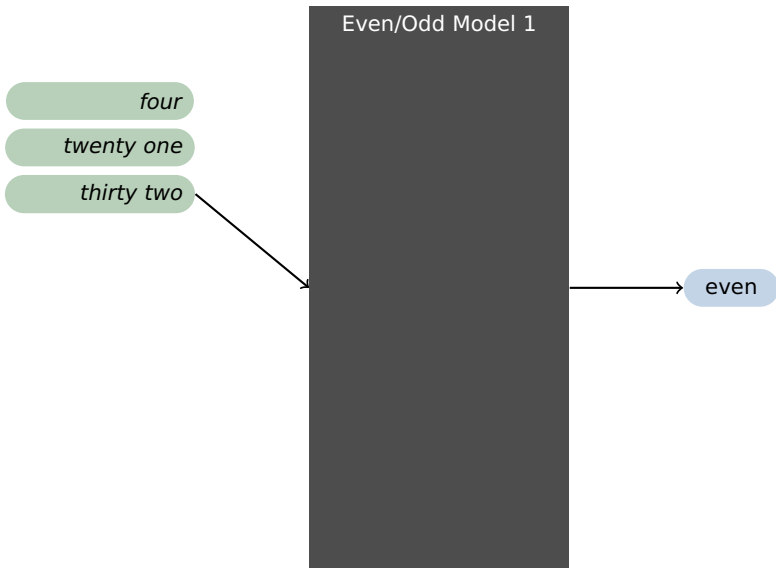
Limits of behavioral testing



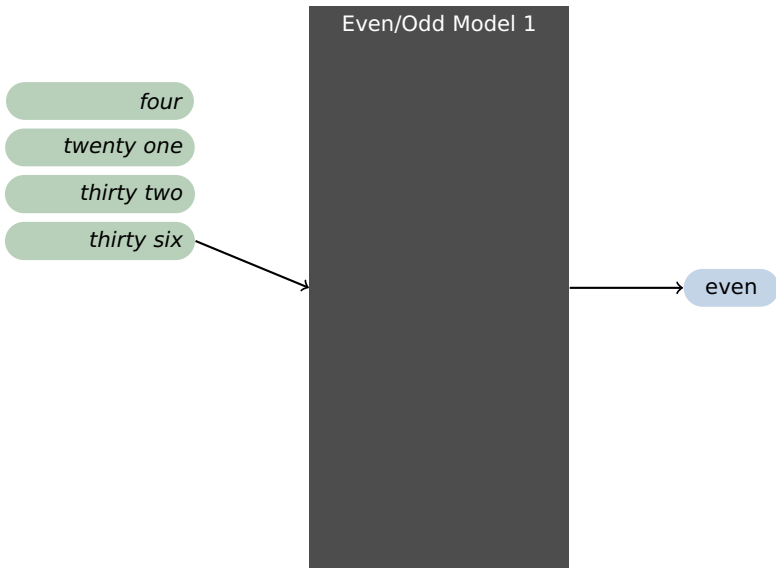
Limits of behavioral testing



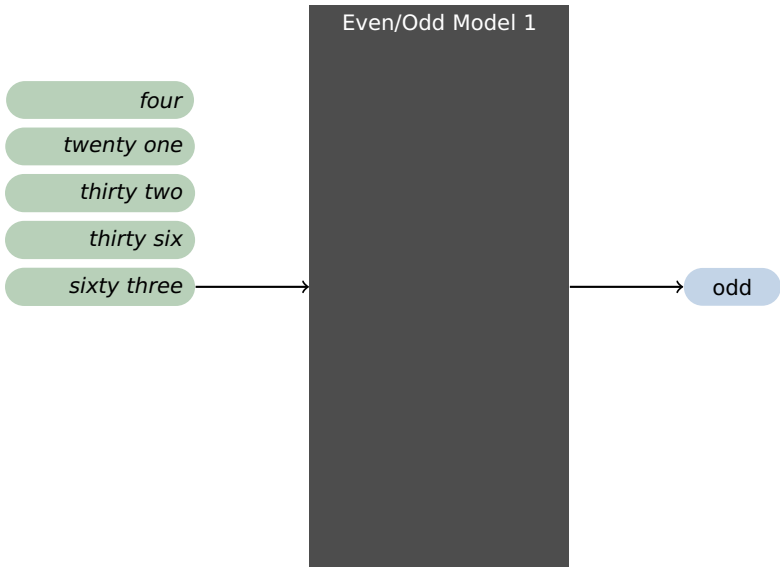
Limits of behavioral testing



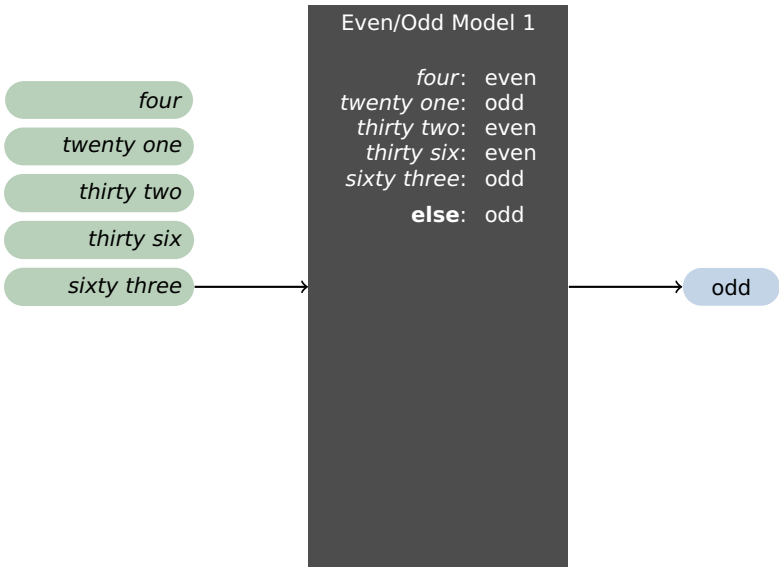
Limits of behavioral testing



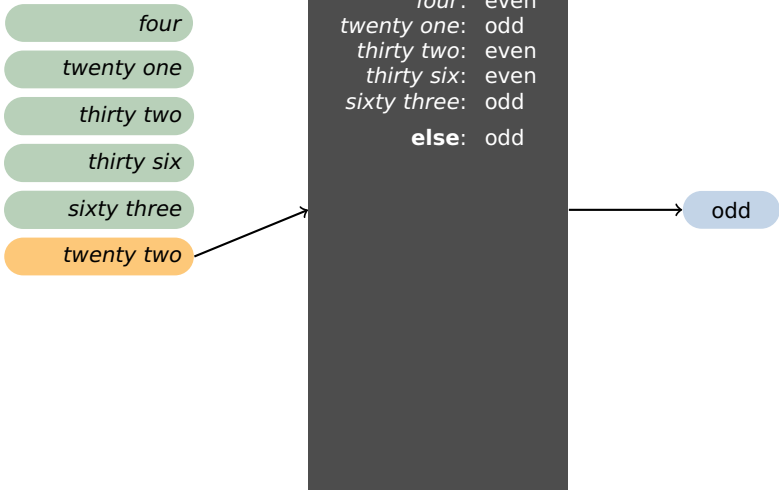
Limits of behavioral testing



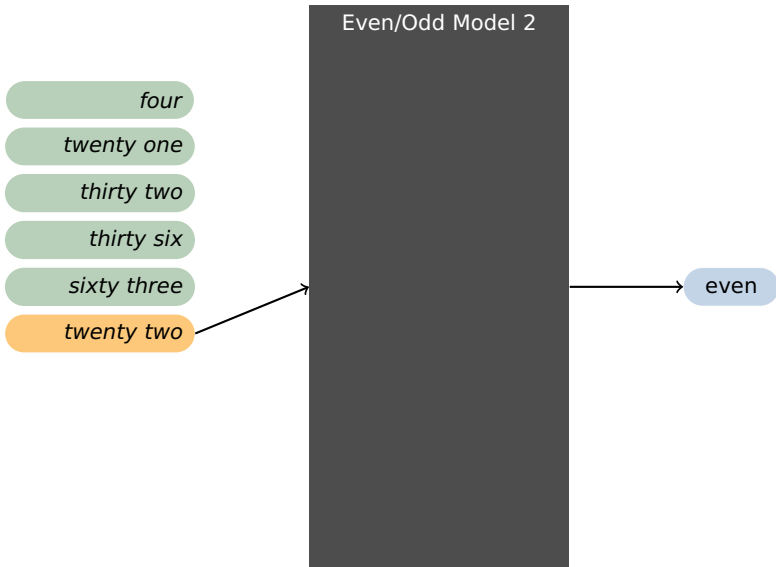
Limits of behavioral testing



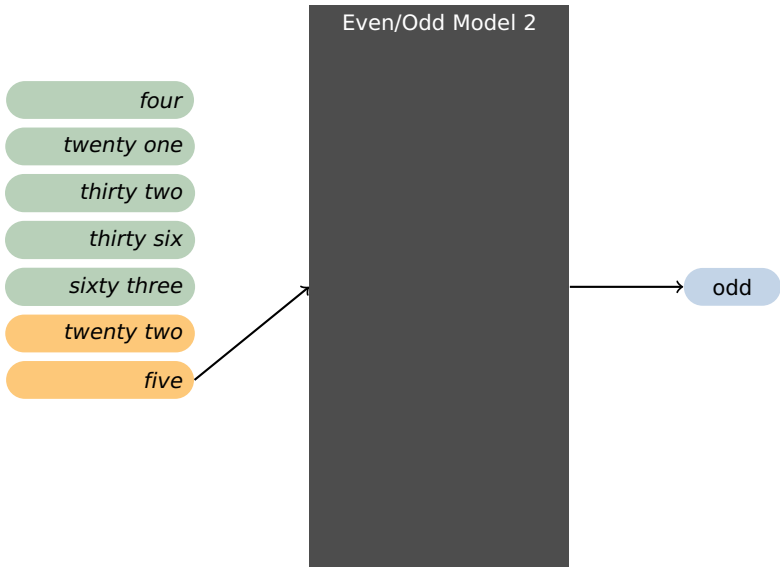
Limits of behavioral testing



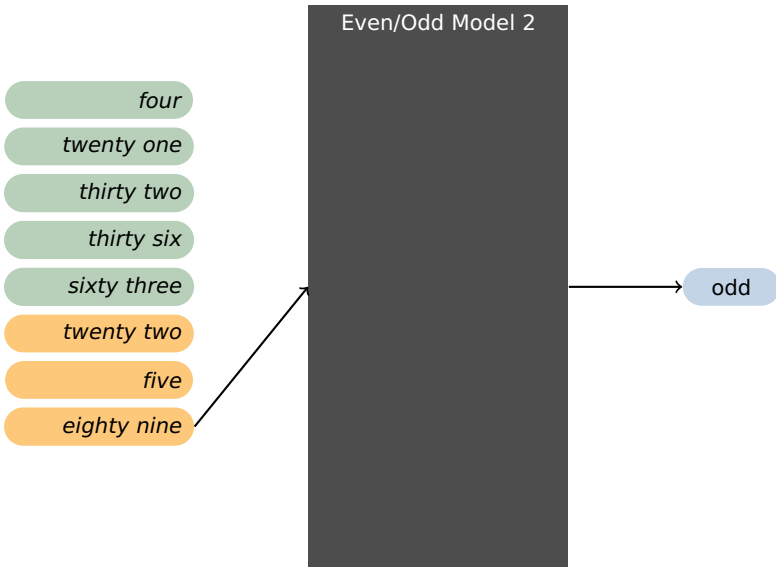
Limits of behavioral testing



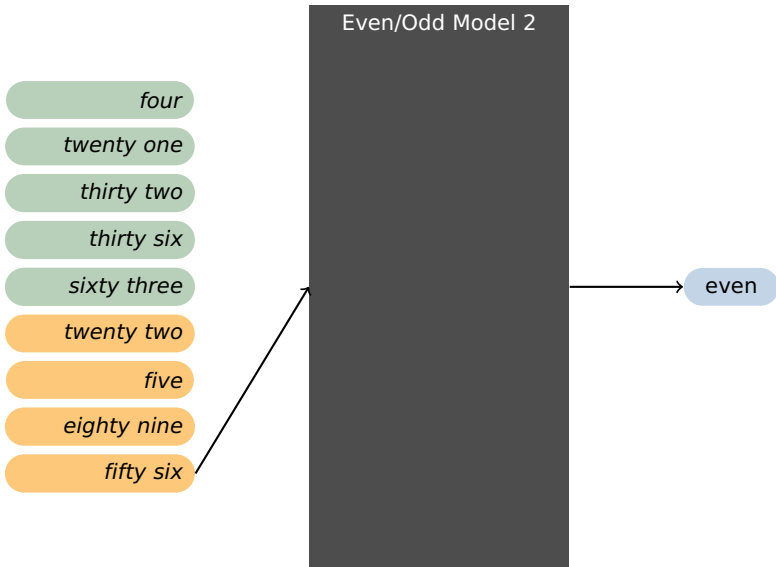
Limits of behavioral testing



Limits of behavioral testing



Limits of behavioral testing



Limits of behavioral testing

- four
- twenty one
- thirty two
- thirty six
- sixty three
- twenty two
- five
- eighty nine
- fifty six

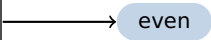
```

Even/Odd Model 2

d =
[
  one:  odd
  two:  even
  three: odd
  four: even
  five: odd
  six:  even
  seven: odd
  eight: even
  nine: odd
  else: odd
]

return
  d[input final token]

```



Limits of behavioral testing

- four
- twenty one
- thirty two
- thirty six
- sixty three
- twenty two
- five
- eighty nine
- fifty six
- sixteen

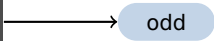
```

Even/Odd Model 2

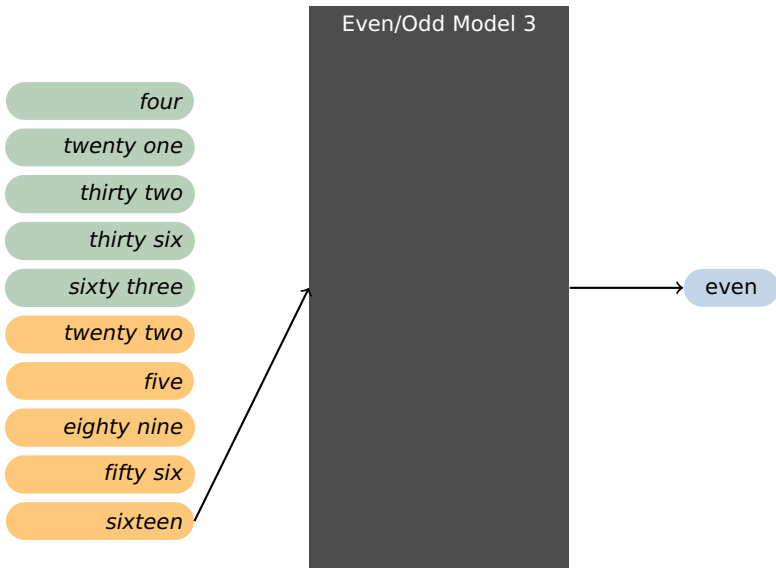
d =
[
  one:  odd
  two:  even
  three: odd
  four: even
  five: odd
  six:  even
  seven: odd
  eight: even
  nine: odd
  else: odd
]

return
d[input final token]

```



Limits of behavioral testing



Seeking generalization guarantees

Seeking generalization guarantees

- Goal: causal analysis of a model's structure, to obtain guarantees about how it will behave.

Seeking generalization guarantees

- Goal: causal analysis of a model's structure, to obtain guarantees about how it will behave.
- Further questions of
 - ▶ fairness
 - ▶ bias
 - ▶ reliability
 - ▶ robustnessare hard to address without such guarantees.

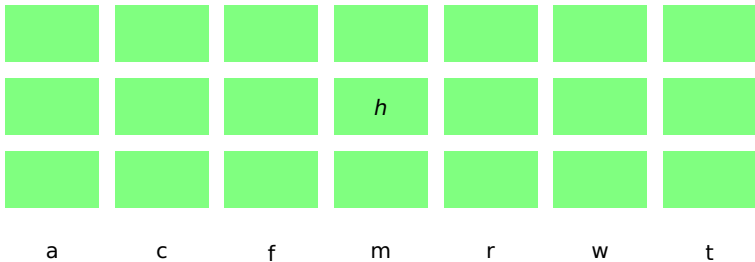
Probing

Core idea behind probing

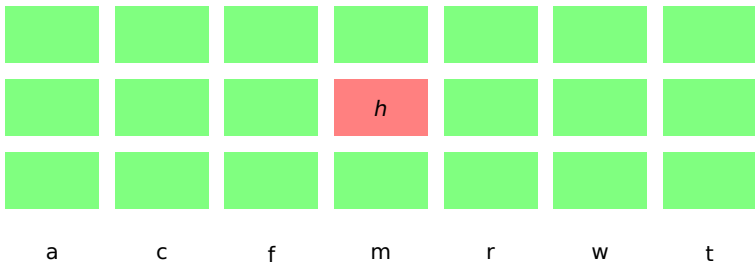
Use a supervised model (the probe) to determine what is latently encoded in the hidden representations of a target models.

Conneau et al. 2018; Tenney et al. 2019

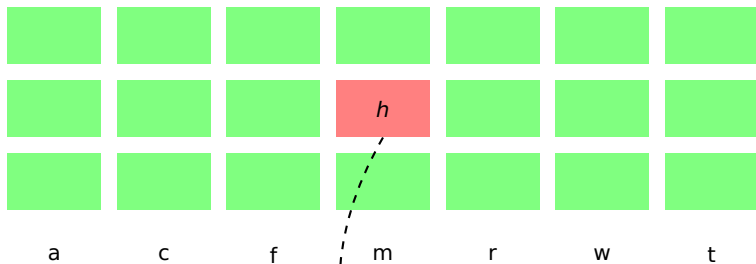
Core method



Core method

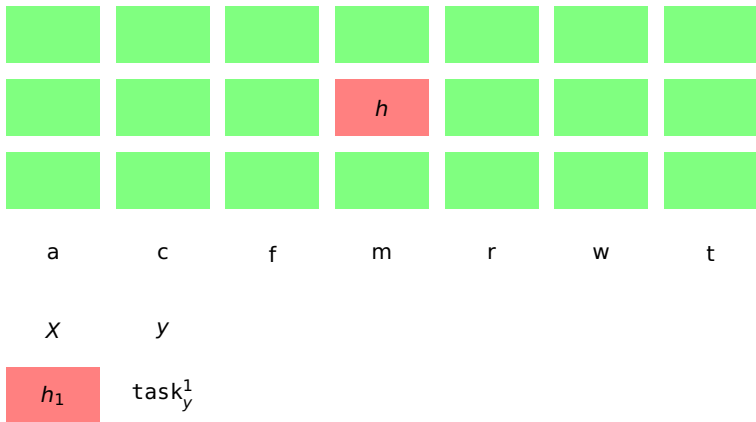


Core method

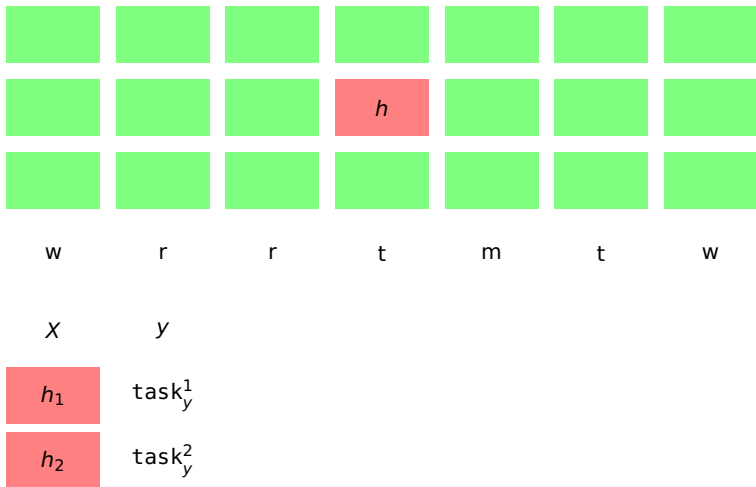


$\text{SmallLinearModel}(h) = \text{task}$

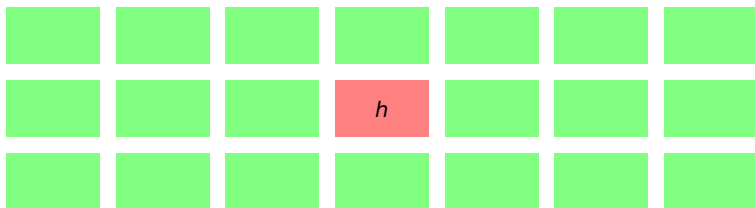
Core method



Core method



Core method



a b c t w w w

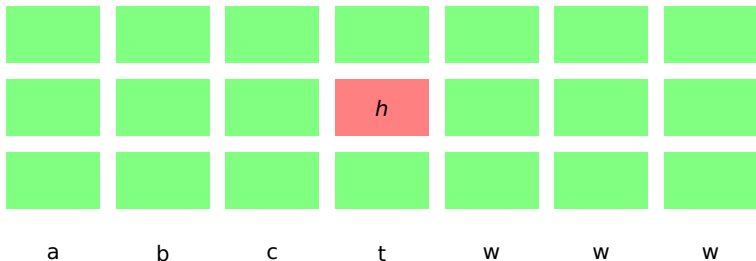
x y

h_1 task_y¹

h_2 task_y²

h_3 task_y³

Core method



X

y

h_1

task_y¹

h_2

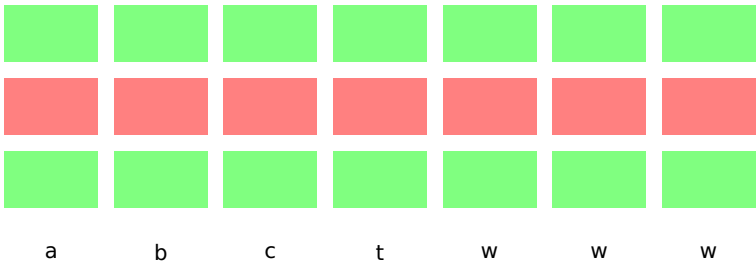
task_y²

h_3

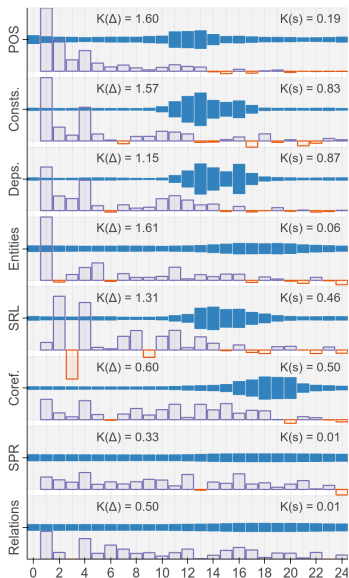
task_y³

SmallLinearModel(X, y)

Core method



Probing BERT



Tenney et al. 2019

Central limitations

Central limitations

Probing or learning a new model?

Central limitations

Probing or learning a new model?

1. A probe is a supervised model with a particular featurization choice.

Central limitations

Probing or learning a new model?

1. A probe is a supervised model with a particular featurization choice.
2. At least some of the information that we identify is likely to be stored in the probe model.

Central limitations

Probing or learning a new model?

1. A probe is a supervised model with a particular featurization choice.
2. At least some of the information that we identify is likely to be stored in the probe model.
3. Responses:
 - ▶ Unsupervised probes (Saphra and Lopez 2019; Clark et al. 2019; Hewitt and Manning 2019)
 - ▶ Control tasks (Hewitt and Liang 2019)

Central limitations

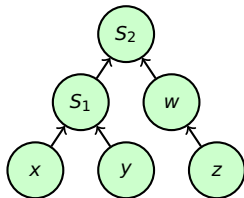
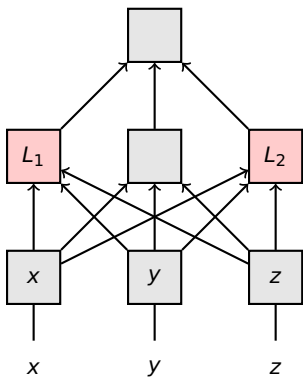
Probing or learning a new model?

1. A probe is a supervised model with a particular featurization choice.
2. At least some of the information that we identify is likely to be stored in the probe model.
3. Responses:
 - ▶ Unsupervised probes (Saphra and Lopez 2019; Clark et al. 2019; Hewitt and Manning 2019)
 - ▶ Control tasks (Hewitt and Liang 2019)

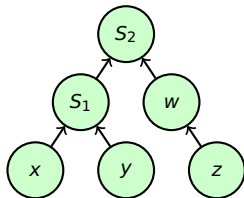
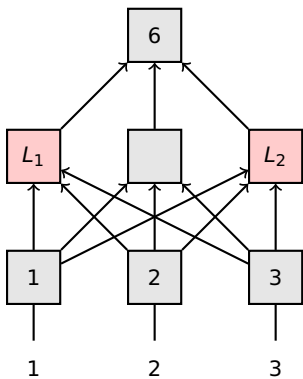
No causal inference

Probes cannot tell us about whether the information that we identify has any *causal* relationship with the target model's behavior (Belinkov and Glass 2019; Geiger et al. 2020, 2021).

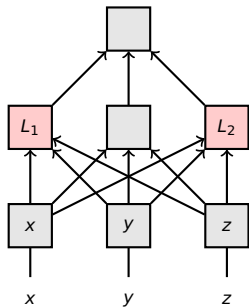
Simple running example



Simple running example

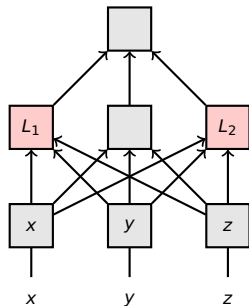


No causal inferences



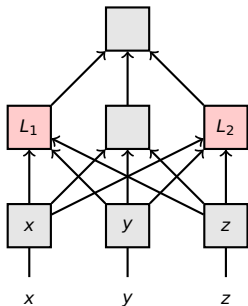
No causal inferences

1. Probe L_1 : it computes $x + y$



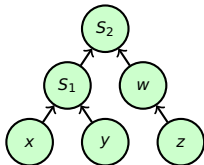
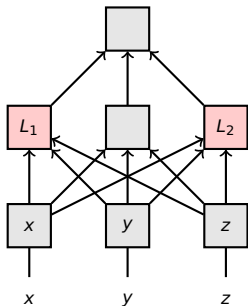
No causal inferences

1. Probe L_1 : it computes $x + y$
2. Probe L_2 : it computes z



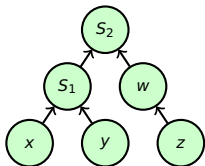
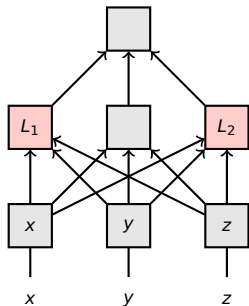
No causal inferences

1. Probe L_1 : it computes $x + y$
2. Probe L_2 : it computes z
3. Aha!



No causal inferences

1. Probe L_1 : it computes $x + y$
2. Probe L_2 : it computes z
3. Aha!



4. But neither has any impact on the output!

$$W_1 = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \quad W_2 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad W_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

$$\mathbf{w} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad (\mathbf{x}W_1; \mathbf{x}W_2; \mathbf{x}W_3) \mathbf{w}$$

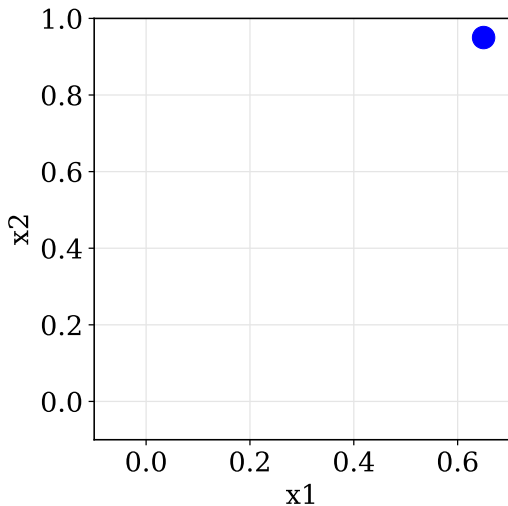
Feature attribution

captum.ai

1. Integrated gradients (Sundararajan et al. 2017)
2. Gradients
3. Saliency Maps (Simonyan et al. 2013)
4. DeepLift (Shrikumar et al. 2017)
5. Deconvolution (Zeiler and Fergus 2014)
6. LIME (Ribeiro et al. 2016)
7. Feature ablation
8. Feature permutation
9. ...

<https://captum.ai>

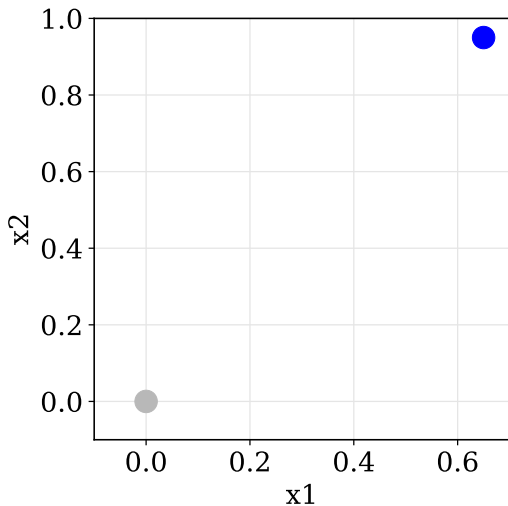
Integrated gradients: Intuition



Sundararajan et al. 2017;

[slide with IG definition](#)

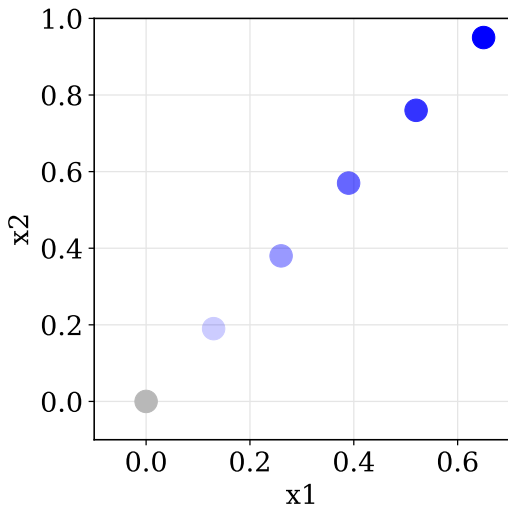
Integrated gradients: Intuition



Sundararajan et al. 2017;

[slide with IG definition](#)

Integrated gradients: Intuition



Sundararajan et al. 2017; [slide with IG definition](#)

Central properties

Central properties

Sensitivity

If two inputs x and x' differ only at dimension i and lead to different predictions, then feature f_i has non-zero attribution.

$$M([1, 0, 1]) = \text{positive}$$

$$M([1, 1, 1]) = \text{negative}$$

Central properties

Sensitivity

If two inputs x and x' differ only at dimension i and lead to different predictions, then feature f_i has non-zero attribution.

$$M([1, 0, 1]) = \text{positive}$$

$$M([1, 1, 1]) = \text{negative}$$

Completeness

For input x and baseline x' , the sum of attributions for x is equal to $M(x) - M(x')$.

Central properties

Sensitivity

If two inputs x and x' differ only at dimension i and lead to different predictions, then feature f_i has non-zero attribution.

$$M([1, 0, 1]) = \text{positive}$$

$$M([1, 1, 1]) = \text{negative}$$

Completeness

For input x and baseline x' , the sum of attributions for x is equal to $M(x) - M(x')$.

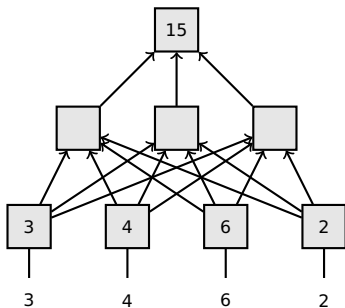
Implementation invariance

If two models M and M' have identical input/output behavior, then the attributions for M and M' are identical.

Reliable insights about causal structure

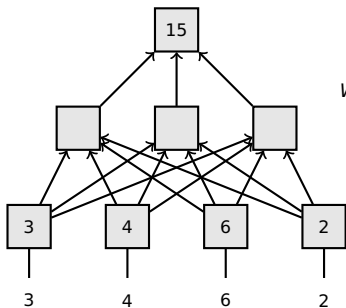
Sundararajan et al. 2017

Reliable insights about causal structure



Sundararajan et al. 2017

Reliable insights about causal structure

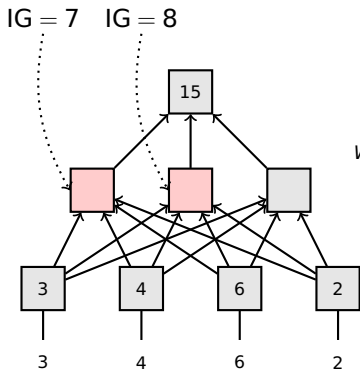


$$W_1 = \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} \quad W_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} \quad W_3 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\mathbf{w} = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \quad (\mathbf{x}W_1; \mathbf{x}W_2; \mathbf{x}W_3) \mathbf{w}$$

Sundararajan et al. 2017

Reliable insights about causal structure



$$W_1 = \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} \quad W_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} \quad W_3 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\mathbf{w} = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \quad (\mathbf{x}W_1; \mathbf{x}W_2; \mathbf{x}W_3) \mathbf{w}$$

Sundararajan et al. 2017

Causal abstraction

Recipe

Geiger et al. 2020, 2021

Recipe

1. State a hypothesis about (an aspect of) the target model's causal structure.

Geiger et al. 2020, 2021

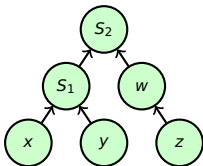
Recipe

1. State a hypothesis about (an aspect of) the target model's causal structure.
2. Search for an alignment between the causal model and target model.

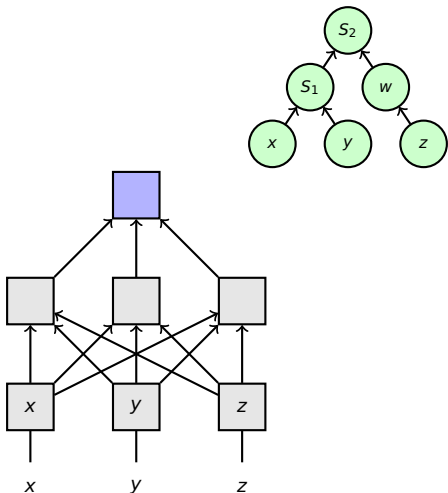
Recipe

1. State a hypothesis about (an aspect of) the target model's causal structure.
2. Search for an alignment between the causal model and target model.
3. Perform *interchange interventions*.

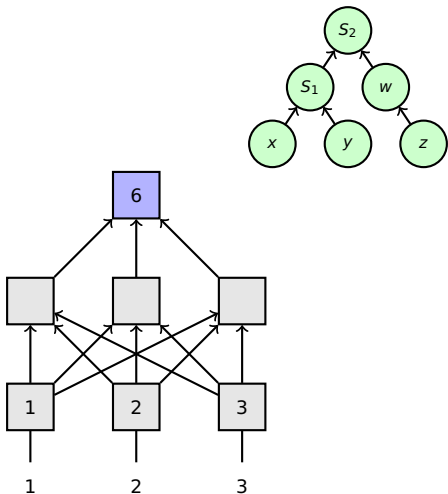
Interchange intervention analysis



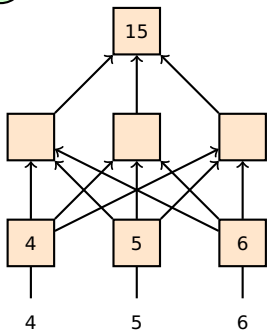
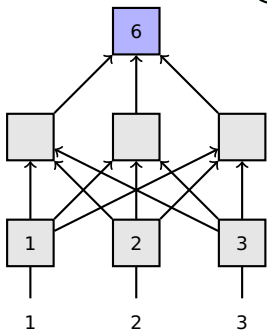
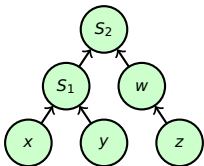
Interchange intervention analysis



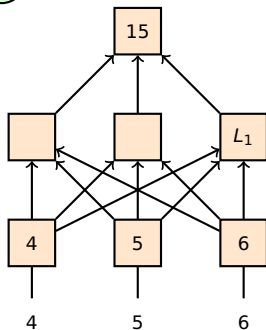
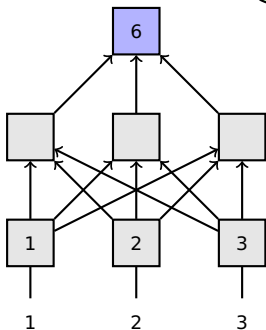
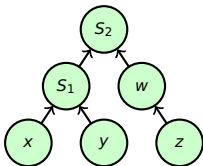
Interchange intervention analysis



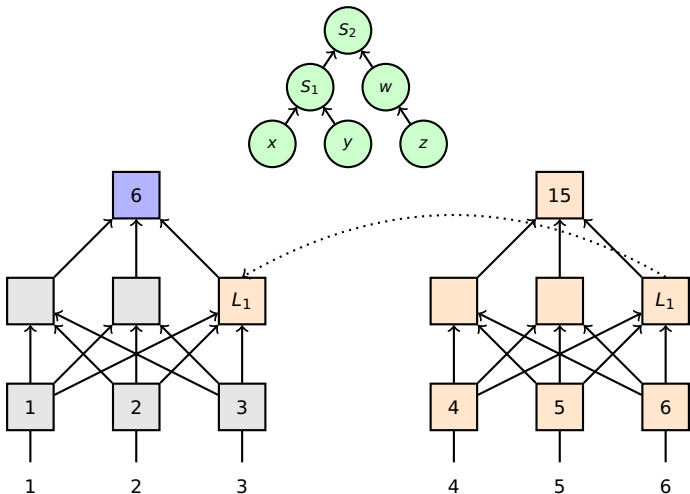
Interchange intervention analysis



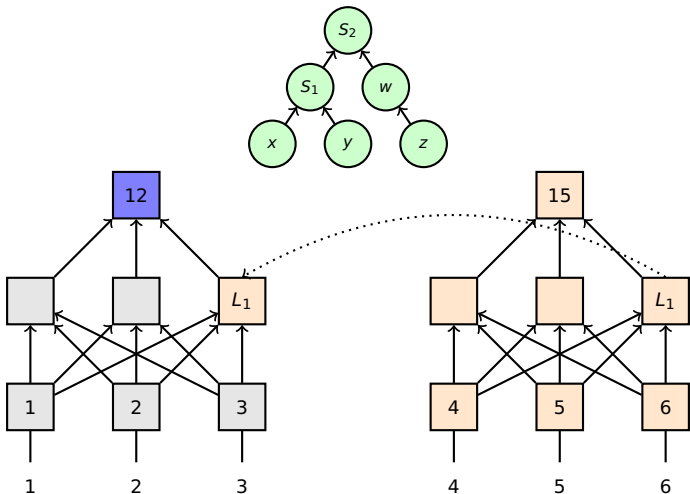
Interchange intervention analysis



Interchange intervention analysis



Interchange intervention analysis



Connections to the literature

- Constructive abstraction (Beckers et al. 2020)
- Causal mediation analysis (Vig et al. 2020)
- Role Learning Networks (Soulos et al. 2020)

Monotonicity NLI (MoNLI)

MoNLI dataset construction

MoNLI dataset construction

Positive MoNLI (PMoNLI; 1,476 examples)

MoNLI dataset construction

Positive MoNLI (PMoNLI; 1,476 examples)

SNLI hypothesis (A) Food was served.

MoNLI dataset construction

Positive MoNLI (PMoNLI; 1,476 examples)

SNLI hypothesis (A)
WordNet

Food was served.
pizza \sqsubset food

MoNLI dataset construction

Positive MoNLI (PMoNLI; 1,476 examples)

SNLI hypothesis (A)

WordNet

New example (B)

Food was served.

pizza \sqsubset food

Pizza was served.

MoNLI dataset construction

Positive MoNLI (PMoNLI; 1,476 examples)

SNLI hypothesis (A)

WordNet

New example (B)

Food was served.

pizza \sqsubset food

Pizza was served.

Positive MoNLI

(A) **neutral** (B)

Positive MoNLI

(B) **entailment** (A)

MoNLI dataset construction

Positive MoNLI (PMoNLI; 1,476 examples)

SNLI hypothesis (A)	Food was served.
WordNet	pizza \sqsubset food
New example (B)	Pizza was served.

Positive MoNLI	(A) neutral (B)
Positive MoNLI	(B) entailment (A)

Negative MoNLI (PMoNLI; 1,202 examples)

SNLI hypothesis (A)	The children are not holding plants.
WordNet	flowers \sqsubset plants
New example (B)	The children are not holding flowers.

Negative MoNLI	(A) entailment (B)
Negative MoNLI	(B) neutral (A)

MoNLI monotonicity algorithm

MoNLI monotonicity algorithm

Infer(*example*)

- 1 *lexrel* ← get-lexrel(*example*)
- 2 **if** contains-not(*example*)
- 3 **return** reverse(*lexrel*)
- 4 **return** *lexrel*

MoNLI monotonicity algorithm

Infer(*example*)

- 1 *lexrel* ← get-lexrel(*example*)
- 2 **if** contains-not(*example*)
- 3 **return** reverse(*lexrel*)
- 4 **return** *lexrel*

MoNLI
lexrel

Pizza was served.
Pizza

entailment
entailment

Food was served.
Food

MoNLI
lexrel

Pizza was not served.
Pizza

neutral
entailment
neutral

Food was not served.
Food

reverse(*lexrel*)

Models

Models

BiLSTM The bidirectional LSTM baseline from Williams et al. (2018).

ESIM The Enhanced Sequential Inference Model (Chen et al. 2016) is a hybrid TreeLSTM-based and biLSTM-based model that uses an inter-sentence attention mechanism to align words across sentences.

BERT A Transformer model trained to do masked language modeling and next-sentence prediction (Devlin et al. 2019).

MoNLI as challenge dataset

Model	Input pretrain	NLI train data	No MoNLI fine-tuning		
			SNLI	PMoNLI	NMoNLI
BiLSTM	GloVe	SNLI train	81.6	73.2	37.9

MoNLI as challenge dataset

Model	Input pretrain	NLI train data	No MoNLI fine-tuning		
			SNLI	PMoNLI	NMoNLI
BiLSTM	GloVe	SNLI train	81.6	73.2	37.9
ESIM	GloVe	SNLI train	87.9	86.6	39.4

MoNLI as challenge dataset

Model	Input pretrain	NLI train data	No MoNLI fine-tuning		
			SNLI	PMoNLI	NMoNLI
BiLSTM	GloVe	SNLI train	81.6	73.2	37.9
ESIM	GloVe	SNLI train	87.9	86.6	39.4
BERT	BERT	SNLI train	90.8	94.4	2.2

Model failure or dataset failure?

Liu et al. (2019)

“What should we conclude when a system fails on a challenge dataset? In some cases, a challenge might exploit blind spots in the design of the original dataset (*dataset weakness*). In others, the challenge might expose an inherent inability of a particular model family to handle certain natural language phenomena (*model weakness*). These are, of course, not mutually exclusive.”

Negation coverage in SNLI and MultiNLI

1. SNLI: Only 38 examples have negated premise and hypothesis.
2. MultiNLI: 18K examples ($\approx 4\%$) have negated premise and hypothesis, but few have the properties we are after.

A systematic generalization task

NMoNLI Train		NMoNLI Test	
person	198	dog	88
instrument	100	building	64
food	94	ball	28
machine	60	car	12
woman	58	mammal	4
music	52	animal	4
tree	52		
boat	46		
fruit	42		
produce	40		
fish	40		
plant	38		
jewelry	36		
anything	34		
hat	20		
man	20		
horse	16		
gun	12		
adult	10		
shirt	8		
shoe	6		
store	6		
cake	4		
individual	4		
clothe	2		
weapon	2		
creature	2		

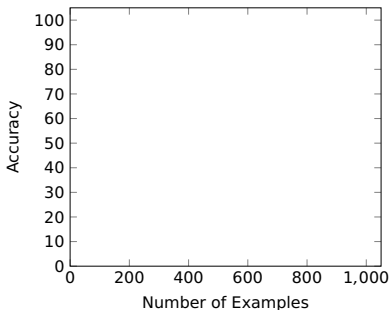
A systematic generalization task

NMoNLI Train		NMoNLI Test	
person	198	dog	88
instrument	100	building	64
food	94	ball	28
machine	60	car	12
woman	58	mammal	4
music	52	animal	4
tree	52		
boat	46		
fruit	42		
produce	40		
fish	40		
plant	38		
jewelry	36		
anything	34		
hat	20		
man	20		
horse	16		
gun	12		
adult	10		
shirt	8		
shoe	6		
store	6		
cake	4		
individual	4		
clothe	2		
weapon	2		
creature	2		

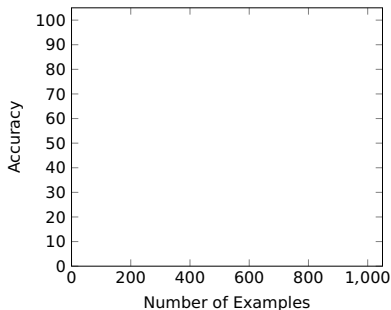
Our models know these lexical relations (high Positive MoNLI accuracy) and will be compelled to combine this knowledge with what they learn about negation during Negative MoNLI fine-tuning.

Fine-tuning on Negative MoNLI

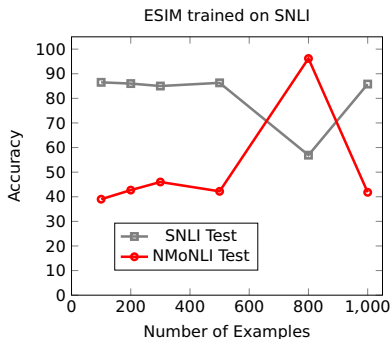
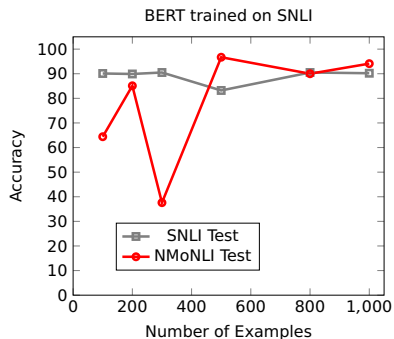
BERT trained on SNLI



ESIM trained on SNLI



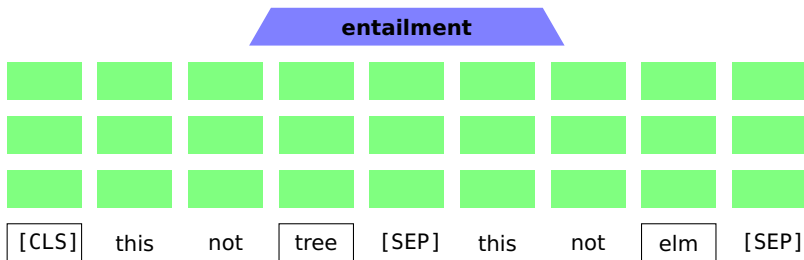
Fine-tuning on Negative MoNLI



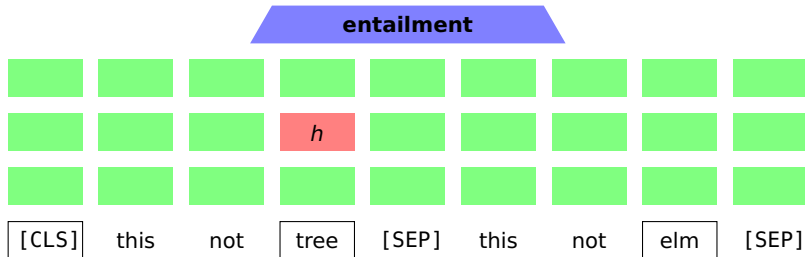
Fine-tuning results

Model	Input pretrain	NLI train data	No MoNLI fine-tuning			With NMoNLI fine-tuning	
			SNLI	PMoNLI	NMoNLI	SNLI	NMoNLI
BiLSTM	GloVe	SNLI train	81.6	73.2	37.9	74.6	93.5
ESIM	GloVe	SNLI train	87.9	86.6	39.4	56.9	96.2
BERT	BERT	SNLI train	90.8	94.4	2.2	90.5	90.0

Focusing on the BERT model

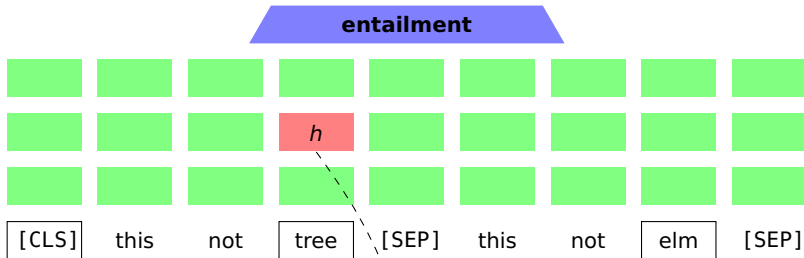


Probes



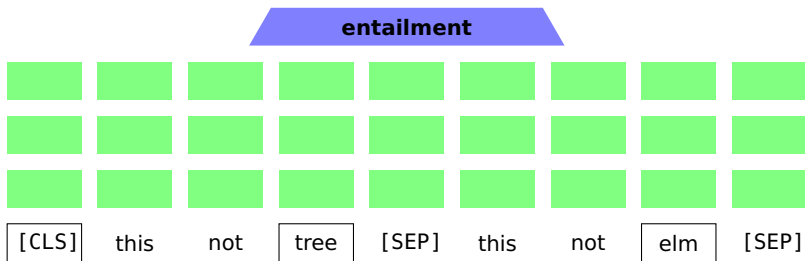
Hewitt and Liang 2019

Probes



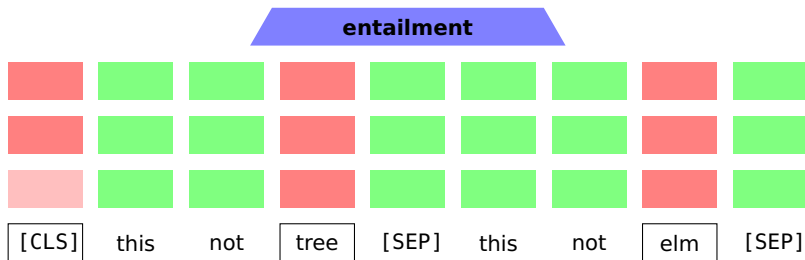
$$\text{SmallLinearModel}(h) = \text{get-lexrel}(\text{tree}, \text{elm})$$

Probe results for lexrel accuracy

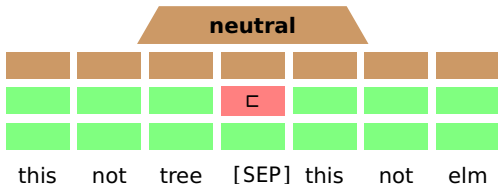
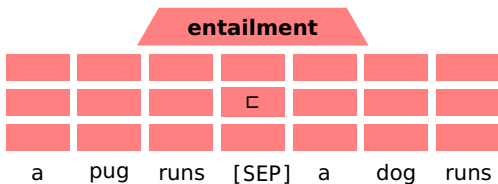
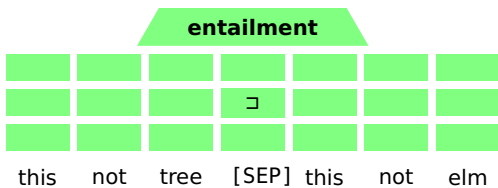


full probing results

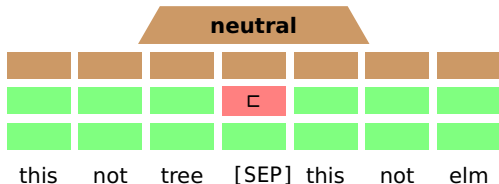
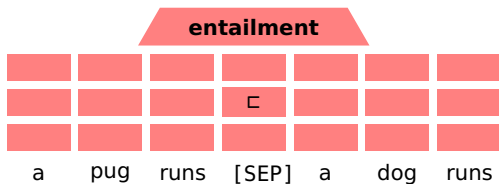
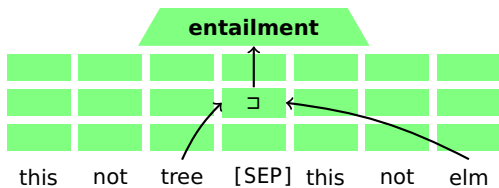
Probe results for lexrel accuracy

[full probing results](#)

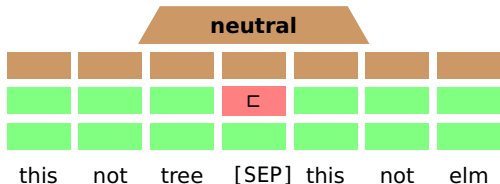
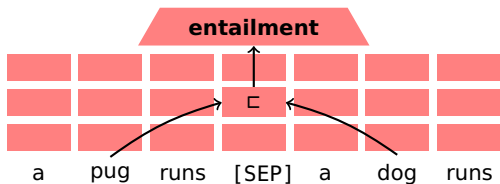
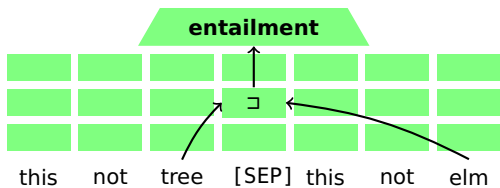
BERT NLI interventions



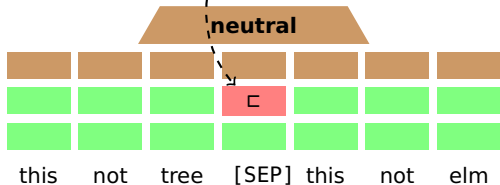
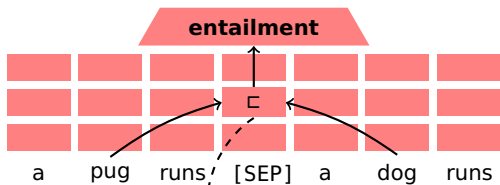
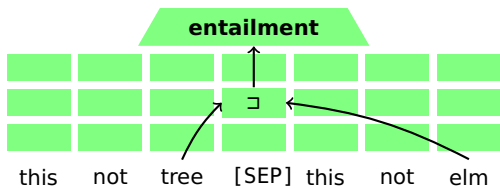
BERT NLI interventions



BERT NLI interventions



BERT NLI interventions



What it means for BERT to implement Infer

What it means for BERT to implement Infer

Infer(*example*)

- 1 *lexrel* ← get-lexrel(*example*)
- 2 **if** contains-not(*example*)
- 3 **return** reverse(*lexrel*)
- 4 **return** *lexrel*

What it means for BERT to implement Infer

Infer(*example*)

- 1 *lexrel* ← get-lexrel(*example*)
- 2 **if** contains-not(*example*)
- 3 **return** reverse(*lexrel*)
- 4 **return** *lexrel*

$$\text{Infer}_{\text{lexrel}(i) \rightarrow \text{lexrel}(j)}(i) = \begin{cases} \text{Infer}(i) & \text{lexrel}(i) = \text{lexrel}(j) \\ \text{reverse}(\text{Infer}(i)) & \text{lexrel}(i) \neq \text{lexrel}(j) \end{cases}$$

What it means for BERT to implement Infer

Infer(*example*)

- 1 *lexrel* ← get-lexrel(*example*)
- 2 **if** contains-not(*example*)
- 3 **return** reverse(*lexrel*)
- 4 **return** *lexrel*

$$\text{Infer}_{\text{lexrel}(i) \rightarrow \text{lexrel}(j)}(i) = \begin{cases} \text{Infer}(i) & \text{lexrel}(i) = \text{lexrel}(j) \\ \text{reverse}(\text{Infer}(i)) & \text{lexrel}(i) \neq \text{lexrel}(j) \end{cases}$$

$$\text{Infer}_{\text{lexrel}(i) \rightarrow \text{lexrel}(j)}(i) = \text{BERT}_{L(i) \rightarrow L(j)}(i)$$

Methods and findings

Methods and findings

1. Find a useful intervention point.

Methods and findings

1. Find a useful intervention point.
2. Interchange interventions for every pair of examples at that site.

Methods and findings

1. Find a useful intervention point.
2. Interchange interventions for every pair of examples at that site.
3. Find clusters of examples in which BERT mimics the causal dynamics of Infer.

Methods and findings

1. Find a useful intervention point.
2. Interchange interventions for every pair of examples at that site.
3. Find clusters of examples in which BERT mimics the causal dynamics of Infer.
4. The largest subsets we found 98, 63, 47, and 37.

Methods and findings

1. Find a useful intervention point.
2. Interchange interventions for every pair of examples at that site.
3. Find clusters of examples in which BERT mimics the causal dynamics of Infer.
4. The largest subsets we found 98, 63, 47, and 37.
 - a. For a random graph, the expected number of subsets larger than 20 is effectively 0.

Methods and findings

1. Find a useful intervention point.
2. Interchange interventions for every pair of examples at that site.
3. Find clusters of examples in which BERT mimics the causal dynamics of Infer.
4. The largest subsets we found 98, 63, 47, and 37.
 - a. For a random graph, the expected number of subsets larger than 20 is effectively 0.
 - b. If the site perfectly captured Infer, we would get a single huge cluster.

Largest exchangeable cluster

(cemetery,location)	(dogs,huskies)	
(house,location) (den,location)	(dog,husky) (dog,chiuahua)	(hood,thing)
(ghetto,location) (backyard,location) (park,location)	(dog,retriever) (dog,maltese)	(nut,thing) (capsule,thing)
(jungle,location) (meadow,location) (residence,location)	(dog,terrier) (dog,pomeranian)	(pouch,thing) (structure,thing)
(laboratory,location) (playground,location) (studio,location)	(beetle,insect)	(root,thing) (nugget,thing)
(slum,location) (station,location) (farm,location)	(grasshopper,insect) (bee,insect)	(tube,thing)
(lab,location) (campsite,location)	(wasp,insect) (fly,insect) (cricket,insect)	(box,object)
(town,location) (lawn,location)	(butterfly,insect) (bumblebee,insect)	(object,sweater) (hat,object)
(saxophone,instrument) (flute,instrument)	(flea,insect) (roach,insect) (moth,insect)	(object,jacket) (toy,object)
(bass,instrument) (piano,instrument)	(mosquito,insect)	(cane,object)
(violin,instrument) (tuba,instrument)	(person,vegetarian) (person,lunatic)	(water,rainwater)
(harmonica,instrument)	(person,republi can) (person,trooper)	(water,saltwater)
(liquid,whiskey)	(person,business) (person,navigator)	(sculptor,artist)
(liquid,margarita) (liquid,tequila)	(person,steward) (person,consultant)	(berry,blueberry)
(liquid,alcohol)	(person,farmer) (person,goalkeeper)	(tree,cypress)
(woman,granny)	(person,sophomore) (person,housekeeper)	(tree,magnolia)(trees,elms)
(woman,widow)	(person,cleanser) (person,physicist) (person,cop)	(tree,maple)
	(person,cambodian) (person,detective)	
	(person,genius) (person,sergeant) (person,californian)	
	(person,doctor) (person,runner)	

Which algorithm is BERT implementing then?

Which algorithm is BERT implementing then?

Infer(*example*)

- 1 *lexrel* ← get-lexrel(*example*)
- 2 **if** contains-not(*example*)
- 3 **return** reverse(*lexrel*)
- 4 **return** *lexrel*

Which algorithm is BERT implementing then?

Infer(*example*)

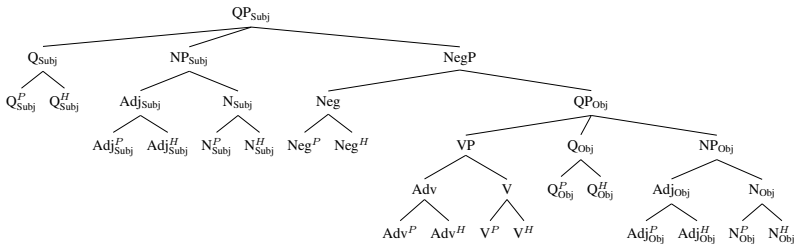
```
1 lexrel ← get-lexrel(example)
2 if contains-not(example)
3     return reverse(lexrel)
4 return lexrel
```

Infer(*example*)

```
1 if inCluster( $C_1$ , example)
2     lexrel1 ← get-lexrel(example)
3     if contains-not(example)
4         return reverse(lexrel1)
5     return lexrel1
6 if inCluster( $C_2$ , example)
7     lexrel2 ← get-lexrel(example)
8     if contains-not(example)
9         return reverse(lexrel2)
10    return lexrel2
11 if inCluster( $C_3$ , example)
12    lexrel3 ← get-lexrel(example)
13    if contains-not(example)
14        return reverse(lexrel3)
15    return lexrel3
16 ...
```

Conclusion

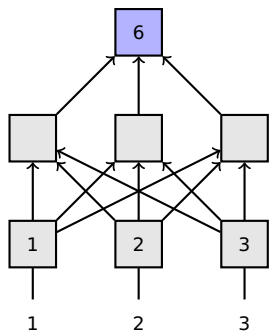
Compositional complexity



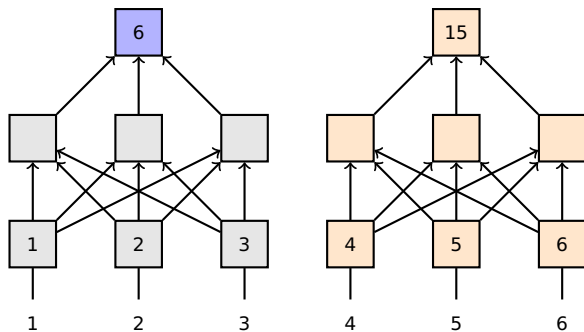
Geiger et al. 2021

Training models to conform to a hypothesis

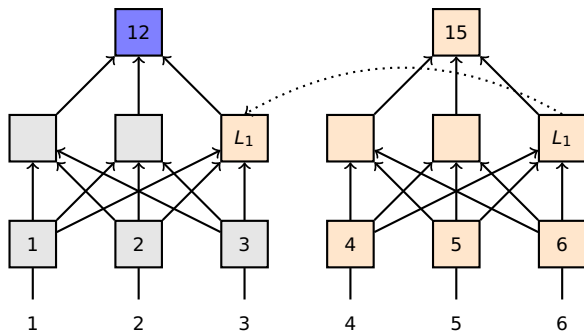
Training models to conform to a hypothesis



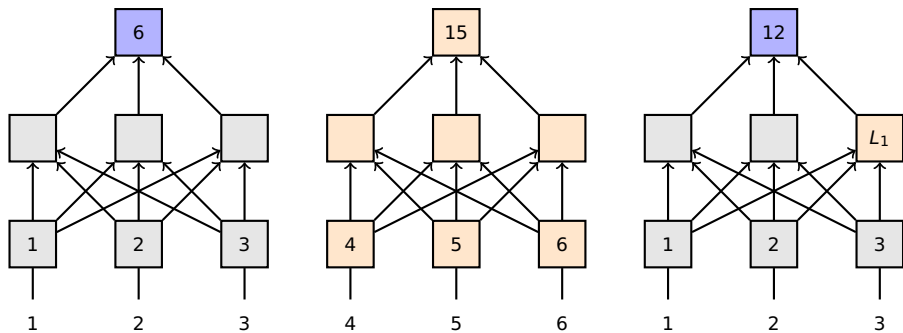
Training models to conform to a hypothesis



Training models to conform to a hypothesis



Training models to conform to a hypothesis



Open questions

Open questions

1. Can we more effectively leverage probes to find useful intervention points?

Open questions

1. Can we more effectively leverage probes to find useful intervention points?
2. What is the relationship between interchange interventions and integrated gradients?

Open questions

1. Can we more effectively leverage probes to find useful intervention points?
2. What is the relationship between interchange interventions and integrated gradients?
3. Can we characterize interchange interventions more generally so that they can be applied to more diverse models?

Open questions

1. Can we more effectively leverage probes to find useful intervention points?
2. What is the relationship between interchange interventions and integrated gradients?
3. Can we characterize interchange interventions more generally so that they can be applied to more diverse models?
4. Can interchanges be used to induce modularity during training?

Open questions

1. Can we more effectively leverage probes to find useful intervention points?
2. What is the relationship between interchange interventions and integrated gradients?
3. Can we characterize interchange interventions more generally so that they can be applied to more diverse models?
4. Can interchanges be used to induce modularity during training?

Thanks!

References I

- Sander Beckers, Frederick Eberhardt, and Joseph Y. Halpern. 2020. [Approximate causal abstractions](#). In *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, volume 115 of *Proceedings of Machine Learning Research*, pages 606–615. PMLR.
- Yonatan Belinkov and James Glass. 2019. [Analysis methods in neural language processing: A survey](#). *Transactions of the Association for Computational Linguistics*, 7:49–72.
- Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, and Hui Jiang. 2016. [Enhancing and combining sequential and tree LSTM for natural language inference](#). *CoRR*, abs/1609.06038.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. [What you can cram into a single \$\\$&!#*\$ vector: Probing sentence embeddings for linguistic properties](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jerry A. Fodor and Zenon W. Pylyshyn. 1988. [Connectionism and cognitive architecture: A critical analysis](#). *Cognition*, 28(1):3–71.
- Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. 2021. [Causal abstractions of neural networks](#). *ArXiv:2106.02997*.
- Atticus Geiger, Kyle Richardson, and Christopher Potts. 2020. [Neural natural language inference models partially embed theories of lexical entailment and negation](#). In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 163–173, Online. Association for Computational Linguistics.
- John Hewitt and Percy Liang. 2019. [Designing and interpreting probes with control tasks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.
- John Hewitt and Christopher D. Manning. 2019. [A structural probe for finding syntax in word representations](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.

References II

- Nelson F. Liu, Roy Schwartz, and Noah A. Smith. 2019. [Inoculation by fine-tuning: A method for analyzing challenge datasets](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2171–2179, Minneapolis, Minnesota. Association for Computational Linguistics.
- Marco Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. [“why should I trust you?”: Explaining the predictions of any classifier](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*, pages 97–101. Association for Computational Linguistics.
- Naomi Saphra and Adam Lopez. 2019. [Understanding learning dynamics of language models with SVCCA](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3257–3267, Minneapolis, Minnesota. Association for Computational Linguistics.
- Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3145–3153. JMLR. org.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*.
- Paul Soulos, R. Thomas McCoy, Tal Linzen, and Paul Smolensky. 2020. [Discovering the compositional structure of vector representations with role learning networks](#). In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 238–254, Online. Association for Computational Linguistics.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3319–3328. JMLR. org.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. [BERT rediscovered the classical NLP pipeline](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. [Causal mediation analysis for interpreting neural nlp: The case of gender bias](#).
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. [A broad-coverage challenge corpus for sentence understanding through inference](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.
- Matthew D Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer.

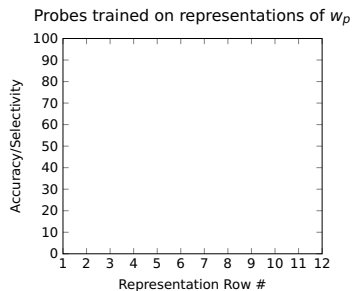
Integrated Gradients computation

$$IG_i(M, x, x') = \underbrace{(x_i - x'_i)}_5 \cdot \underbrace{\sum_{k=1}^4}_{4} \frac{\partial M(\underbrace{x' + \frac{k}{m} \cdot (x - x')}_2)}{\partial x_i} \cdot \underbrace{\frac{1}{m}}_4$$

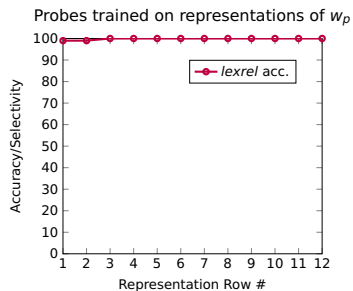
1. Generate $\alpha = [1, \dots, m]$
2. Interpolate inputs between baseline x' and actual input x
3. Compute gradients for each interpolated input
4. Integral approximation through averaging
5. Scaling to remain in the space region as the original

Adapted from the [TensorFlow integrated gradients tutorial](#)

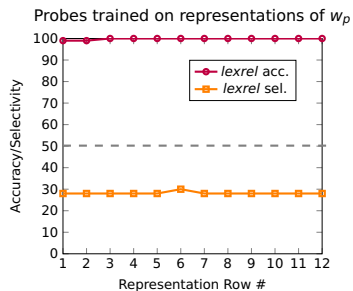
Probe results for lexrel accuracy



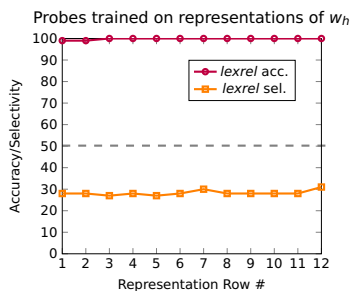
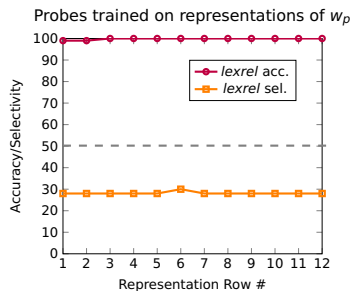
Probe results for lexrel accuracy



Probe results for lexrel accuracy

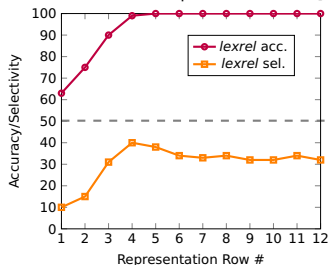


Probe results for lexrel accuracy

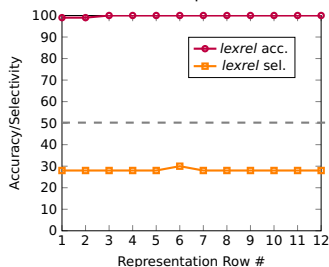


Probe results for lexrel accuracy

Probes trained on representations of [CLS]



Probes trained on representations of w_p



Probes trained on representations of w_h

