



Principles and Interactive Tools for Evaluating and Improving the Behavior of NLP models

Tongshuang (Sherry) Wu

@tongshuangwu / wtshuang@cs.washington.edu

University of Washington

How do I check if my model works?



Should I replace my doctor with OSCAR?

Should we use OSCAR in our products?



If not, what do I need to fix?

OSCAR: Pre-training of Neural Networks Directly on Human Brains

Gnome Chompsky

Arcadia Research

chompsky@arcadia.com

Waltolomew Strickler

Arcadia Oaks High

stricklander@aoh.edu

Abstract

We train neural networks on human brains and achieve SOTA in everything.

1 Introduction

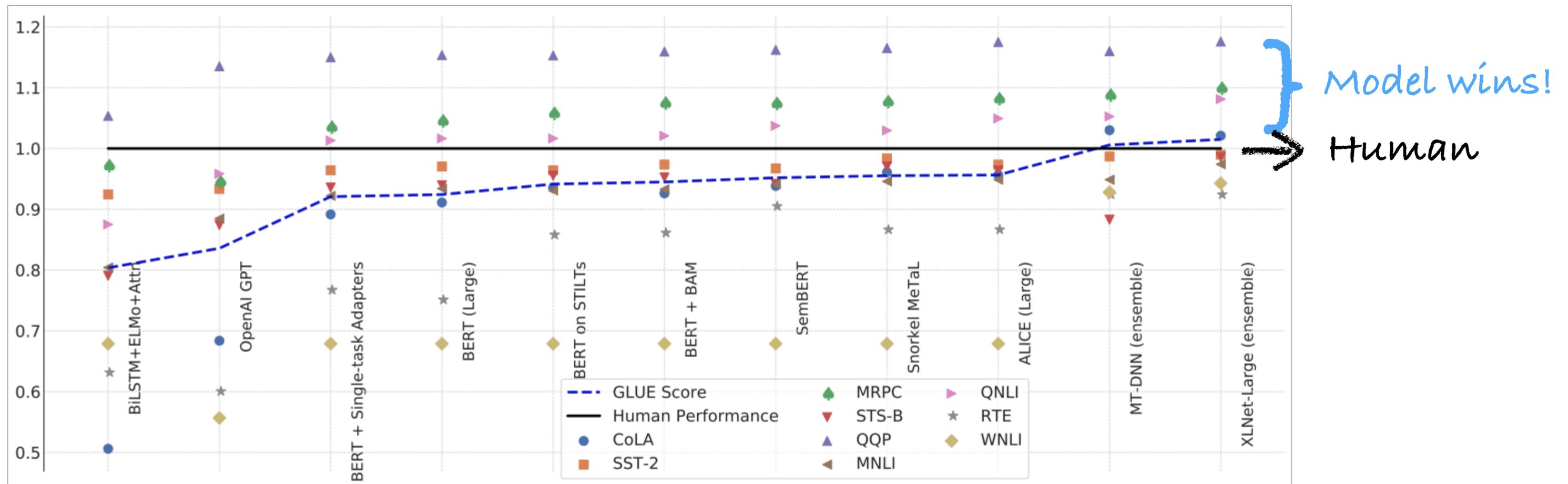
3 Model and Architecture

Our basic human brain training approach is similar in spirit to the process described in [XYZ] et al, with the relative straightforward difference that we train directly on brains rather than on text.

We train a transformer model with a CNN on top

Accuracy seems a good solution?

GLUE: "performance on the benchmark has recently come close to the level of non-expert humans, suggesting limited headroom for further research."

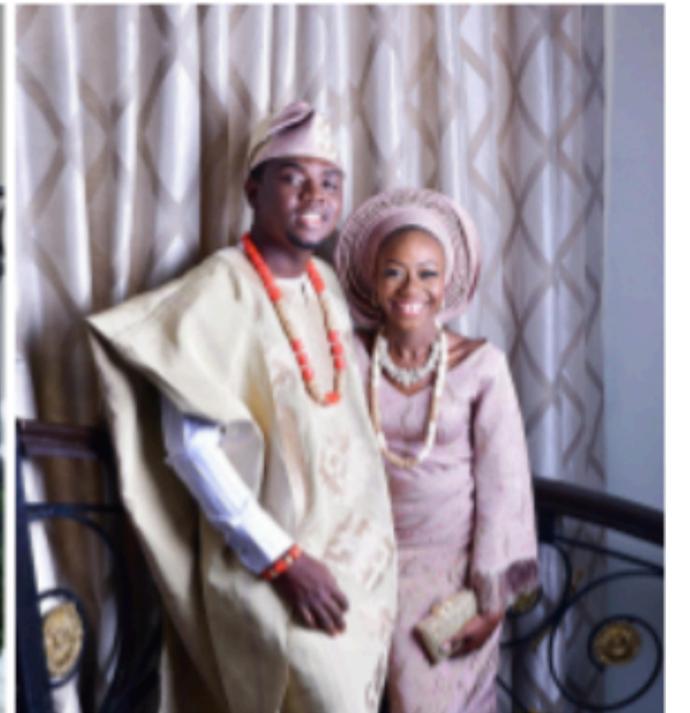
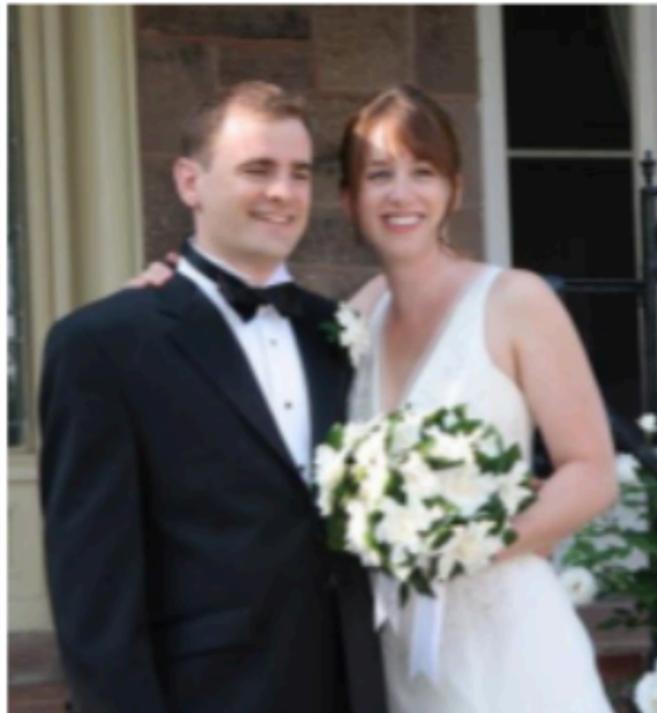
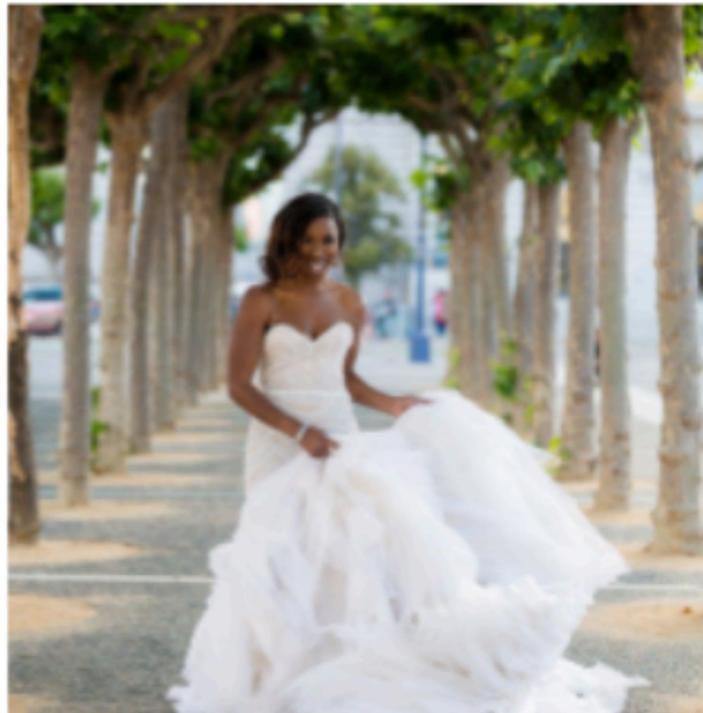
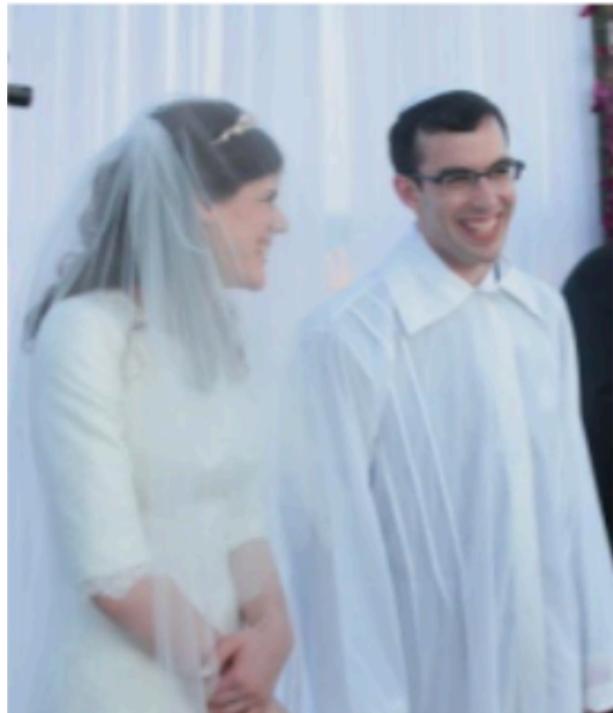


What could go wrong?

Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., ... & Bowman, S. (2019). Superglue: A stickier benchmark for general-purpose language understanding systems. In Advances in Neural Information Processing Systems (pp. 3266-3280).

Missing critical data slices (bias, fairness)

What's in the figure?



> ceremony, wedding, bride, groom, dress

> person, people

High accuracy \neq Model succeeding.

How can we make sure our model can handle particular data slices?

Shortcuts/right for wrong reasons



What is the moustache made of?

> Banana

What are the *eyes* made of?

> Banana

What is?

> Banana

What?

> Banana



High accuracy \neq Model succeeding.

How can we make sure our model can handle particular data slices?

Correct prediction \neq correct reasoning.

How can practitioners ensure the model learns important features & avoid spurious correlations?

Analyzing structured data is easy.

How good is our model on records with different city entries?



What happens if I change the city column to `New York`?

Give me 0.001 seconds to run a SQL script!



Analyzing ~~structured data~~ text is ~~easy~~ hard.

How good is our model on passive sentences on cities?



What happens if I change the passive voice to positive?

Ugh...POS? Named entities? clustering?



"State-of-the-art"

we sampled 200 question answer pairs and manually analyzed their properties.

Joshi et al.
ACL'17

Chen et al.
ACL'16

We randomly select 50 incorrect questions and categorize them into 6 classes.

We sample 100 incorrect predictions and try to find common error categories.

Wadhwa et al.
ACL'18

Joshi, Mandar, et al. "Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension." arXiv preprint arXiv:1705.03551 (2017).
Chen, Danqi, Jason Bolton, and Christopher D. Manning. "A thorough examination of the cnn/daily mail reading comprehension task." arXiv preprint arXiv:1606.02858 (2016).
.Wadhwa, Soumya, Khyathi Raghavi Chandu, and Eric Nyberg. "Comparative analysis of neural qa models on squad." arXiv preprint arXiv:1806.06972 (2018).

"State-of-the-art"

*we sampled 200 question answer pairs and manually
analyze*

Joshi et al.
ACL'17

**"We randomly select 50-100 instances and
roughly label them into N error groups."**

Chen et al.
ACL'16

them into 6 classes.

orize

**Under-representative, subjective, high
variance, low reproducibility
(TOCHI 19, ACL 19, CHI 21)**

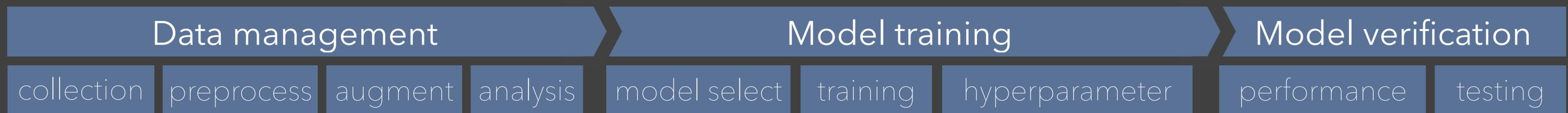
*We sam
comm*

Wadhwa et al.
ACL'18

Joshi, Mandar, et al. "Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension." arXiv preprint arXiv:1705.03551 (2017).
Chen, Danqi, Jason Bolton, and Christopher D. Manning. "A thorough examination of the cnn/daily mail reading comprehension task." arXiv preprint arXiv:1606.02858 (2016).
Wadhwa, Soumya, Khyathi Raghavi Chandu, and Eric Nyberg. "Comparative analysis of neural qa models on squad." arXiv preprint arXiv:1806.06972 (2018).

I **uncover pitfalls in the status-quo analysis** process, and help NLP practitioners **inspect the inputs and outputs** of their models, such that they can gain **more systematic insights** into their models' behaviors.

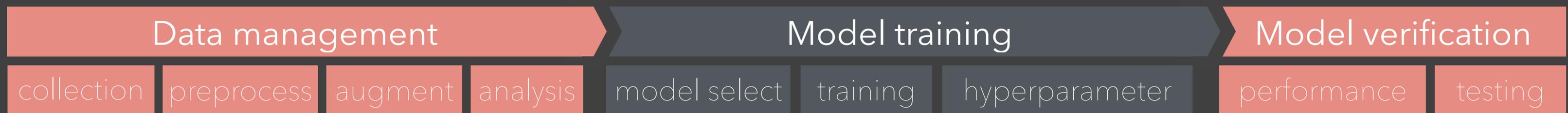
I will focus on this in the talk!



(Paleyes et al. 2020)

I **uncover pitfalls in the status-quo analysis** process, and help NLP practitioners **inspect the inputs and outputs** of their models, such that they can gain **more systematic insights** into their models' behaviors.

I will focus on this in the talk!



To improve ↗

(Paleyes et al. 2020)

Building blocks

High accuracy \neq Model succeeding.

How can we make sure our model can handle particular data slices?

Correct prediction \neq correct reasoning.

How can practitioners ensure the model learns important features & avoid superficial correlations?

Building blocks

Quantitative grouping

Inspect similar instances,
semantically & syntactically

```
ENT(g) != ""  
and count(token(c, pattern=ENT(g))) >  
    count(token(g, pattern=ENT(g)))  
and ENT(g) == ENT(p(m))  
and f1(m) == 0
```

Correct prediction \neq correct reasoning.

How can practitioners ensure the model learns important features & avoid superficial correlations?

Building blocks

Quantitative grouping

Inspect similar instances,
semantically & syntactically

```
ENT(g) != ""  
and count(token(c, pattern=ENT(g))) >  
    count(token(g, pattern=ENT(g)))  
and ENT(g) == ENT(p(m))  
and f1(m) == 0
```

Counterfactual perturbation

Isolate important components
targeted minimal rewrites

This is a good movie.

Lexical

This is ~~good~~ great a movie.

This is a good ~~movie~~ film.

Negation

This **would be** a good movie.

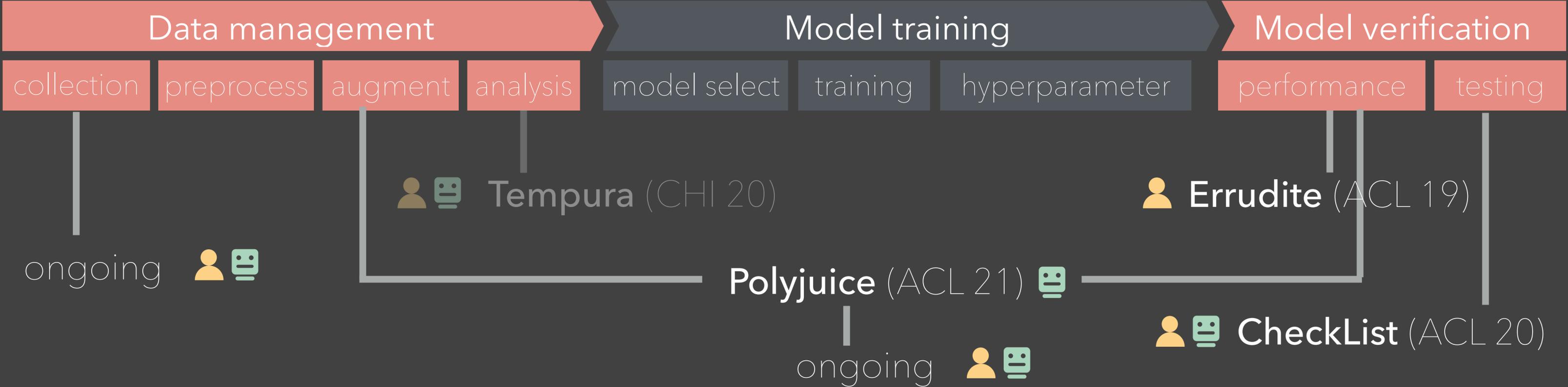
This is **not** a good movie.

Building blocks

Quantitative grouping

Counterfactual perturbation

Model development stages (Paleyes et al. 2020)



ACL 2019

Errudite: Scalable, Reproducible, and Testable Error Analysis

Tongshuang (Sherry) Wu @tongshuangwu
University of Washington

Marco Tulio Ribeiro
Microsoft Research

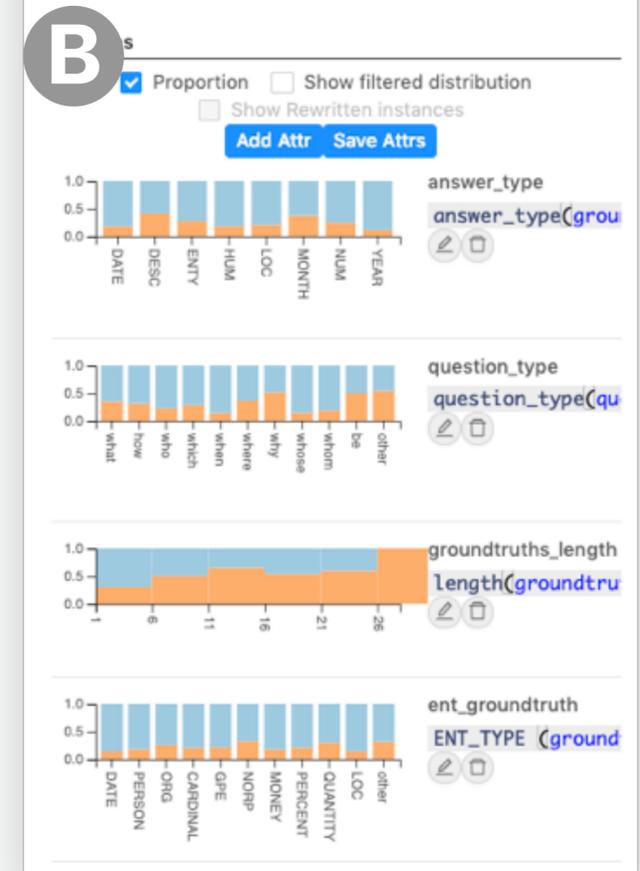
Jeffrey Heer @jeffrey_heer
Daniel S. Weld @dsweld
University of Washington



A Overview

Model Comparison

Model	em	f1	sent	precision	recall
bidaf	0.68	0.77	0.91	0.78	0.81



Errudite: An Interactive Tool for Scalable and Reproducible Error Analysis

Load Undo Query Redo Query

C Instances to explore (without edits)

Get Instances Sample 10 instances randomly that are in and not in

Filter CMD ENT (groundtruth) == ""
 Preview the filter on 10570 instances

Filtered instances: NaN (0.0% of total), Error: undefined (NaN% of slice, NaN% of total, NaN% of all errors)

Record the Group Get samples

D Selected instances (answer encoding:groundtruth,prediction by bidaf(correct,incorrect),model prediction distributions)

Who created the 2005 theme for Doctor Who?

A different arrangement was recorded by Peter Howell for season 18 (1980), which was in turn replaced by Dominic Glynn's arrangement for the season-long serial The Trial of a Time Lord in season 23 (1986). Keff McCulloch provided the new arrangement for the Seventh Doctor's era which lasted from season 24 (1987) until the series' suspension in 1989. American composer John Debney created a new arrangement of Ron Grainer's original theme for Doctor Who in 1996. For the return of the series in 2005, **Murray Gold** provided a new arrangement which featured samples from the 1963 original with further elements added; in the 2005 Christmas episode "The Christmas Invasion", Gold introduced a modified closing credits arrangement that was used up until the conclusion of the 2007 series.[citation needed]

DID YOU MEAN TO FILTER INSTANCES THAT ARE... Close Now

- starts_with(prediction(model="bidaf"), pattern="NNP")
- starts_with(prediction(model="bidaf"), pattern="PERSON")
- attr:answer_type == answer_type(prediction(model="bidaf"))
- exact_match(model="bidaf") == 0
- is_correct_sent(prediction(model="bidaf")) == 0
- overlap(question, sentence(prediction(model="bidaf"))) > overlap(question, sentence(groundtruths))

Prev page Next page
 Displaying #0-4 samples.

E Groups

Proportion Show filtered distribution
 Export the Groups Compare models

all_instances		10570	32%	68%
is_entity	Length (ENT (groundtri	4240	20%	80%
has_distractor	Length (ENT (groundtri	3495	21%	79%
correct_type	Length (attr:ent_g) >	2988	12%	88%
is_distracted	Length (attr:ent_g) >	192	100%	

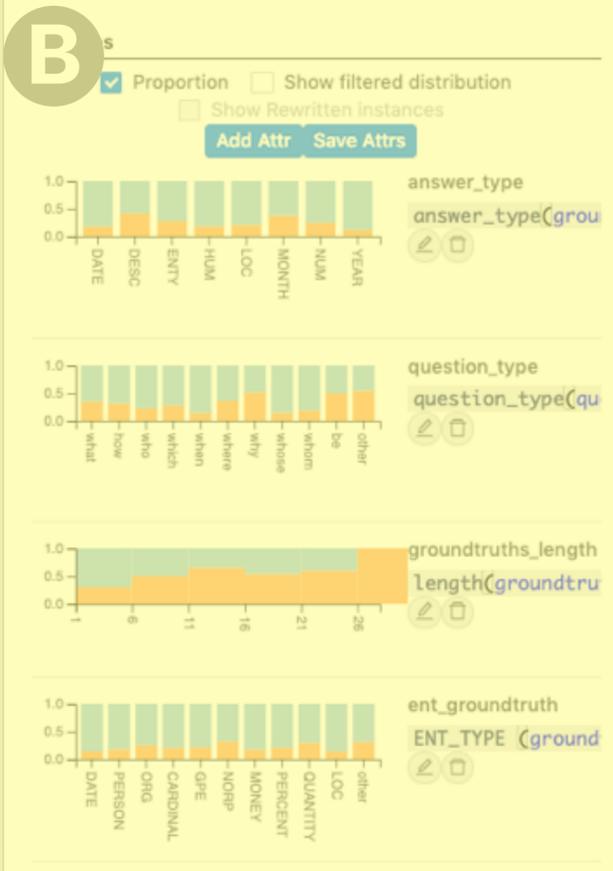
F Selected Re-write Rules

Proportion Show filtered distribution
 Add a rule Save rules

STRING (prediction (model="ANCHOR"))	→ ‡	202	46%	51%
keep_correct_sentence		92	100%	86%
remove_clues		17	82%	82%
resolve_coref		5	100%	

A Overview Model Comparison

Model	em	f1	sent	precision	recall
bidaf	0.68	0.77	0.91	0.78	0.81



C Instances to explore (without edits)

Get Instances Sample 10 instances randomly that are in and not in

Filter CMD ENT (groundtruth) == ""

Preview the filter on 10570 instances

Filtered instances: NaN (0.0% of total), Error: undefined (NaN% of slice, NaN% of total, NaN% of all errors)

Record the Group Get samples

D Selected instances (answer encoding:groundtruth,prediction by bidaf(correct,incorrect),model prediction distributions)

Who created the 2005 theme for Doctor Who?

A different arrangement was recorded by Peter Howell for season 18 (1980), which was in turn replaced by Dominic Glynn's arrangement for the season-long serial The Trial of a Time Lord in season 23 (1986). Keff McCulloch provided the new arrangement for the Seventh Doctor's era which lasted from season 24 (1987) until the series' suspension in 1989.

American composer John Debney created a new arrangement of Ron Grainer's original theme for Doctor Who in 1996.

For the return of the series in 2005, **Murray Gold** provided a new arrangement which featured samples from the 1963 original with further elements added; in the 2005 Christmas episode "The Christmas Invasion", Gold introduced a modified closing credits arrangement that was used up until the conclusion of the 2007 series.[citation needed]

DID YOU MEAN TO FILTER INSTANCES THAT ARE... Close Now

- starts_with(prediction(model="bidaf"), pattern="NNP")
- starts_with(prediction(model="bidaf"), pattern="PERSON")
- attr:answer_type == answer_type(prediction(model="bidaf"))
- exact_match(model="bidaf") == 0
- is_correct_sent(prediction(model="bidaf")) == 0
- overlap(question, sentence(prediction(model="bidaf"))) > overlap(question, sentence(groundtruths))

Prev page Next page
Displaying #0-4 samples.

E Groups

Proportion Show filtered distribution

Export the Groups Compare models

Attribute	Length	Correct	Incorrect
all_instances	10570	32%	68%
is_entity	length (ENT (groundtri	4240	20% 80%
has_distractor	length (ENT (groundtri	3495	21% 79%
correct_type	length (attr:ent_g) >	2988	12% 88%
is_distracted	length (attr:ent_g) >	192	100%

F Selected Re-write Rules

Proportion Show filtered distribution

Add a rule Save rules

Rule	Count	Correct	Incorrect
STRING (prediction (model="ANCHOR")) > #	202	46%	54%
keep_correct_sentence	92	100%	86%
remove_clues	17	82%	82%
resolve_coref	5	100%	

Building blocks

Quantitative grouping

Inspect similar instances,
semantically & syntactically

Counterfactual perturbation

Isolate important components
targeted minimal rewrites

Errudite

Precise & reproducible hypotheses

+

Scale up to the entire dev set

+

Cover **errors & correct** instances

+

Test via counterfactual analysis

Scenario: distractor hypothesis

Common error hypothesis in question answering, on **BiDAF** (Seo et al., 2016), with **SQuAD** (10570 instances; Rajpurkar et al., 2016)

Common belief: BiDAF...

- ✓ **Matches entity types**
Knows to find a *PERSON*
- ✗ **Finds the exact answer spans**
Distracted by other *PERSON* spans

Who created the 2005 theme for Doctor Who?

...**John Debney** created a new arrangement of Ron Grainer's original theme for Doctor Who in 1996. For the return of the series in 2005, **Murray Gold** provided a new arrangement... sampled from the 1963 original.

Precise DSL (Domain Specific Language)

Attribute Extractor + Target + Operators

`length(q) > 20`

The screenshot displays the Errudite tool interface with several key components:

- A Overview:** A table showing model performance metrics for 'em', 'f1', 'sent', 'precision', and 'recall' across different models.
- B Instance Attribute:** A panel with a bar chart showing the distribution of 'answer_type' (DATE, DESC, ENTRY, HUM, LOC, MONTH, NUM, YEAR) and a list of attributes like 'answer_type(grou...)'.
- C Filter:** A panel showing a filter command: `ENT (groundtruth) == ""` and a list of filter rules such as `starts_with(prediction(model="bidaf"), pattern="NNP")`.
- D Instance Groups:** A panel showing a sample instance: "Who created the 2005 theme for Doctor Who?" with a detailed answer paragraph.
- E Groups:** A panel showing a table of instance groups with columns for attribute name, value, and percentage. The table includes rows for 'all_instances', 'is_entity', 'has_distractor', 'correct_type', and 'is_distracted'.

Build distractor groups with DSL

The screenshot shows the Errudite tool interface. At the top, the title bar reads "Errudite: An Interactive Tool for Scalable and Reproducible Error Analysis" with buttons for "Load", "Undo Query", and "Redo Query".

The main panel is titled "Instances to explore (without edits)". It contains a "Get Instances" section with a dropdown menu set to "randomly" and two "Select groups" input fields. Below this is a "Filter CMD" section with the query: `ENT (groundtruth) == ""`. A preview indicates the filter is applied to 10570 instances. Below the filter is a status message: "Filtered instances: NaN (0.0% of total), Error: undefined (NaN% of slice, NaN% of total, NaN% of all errors)". At the bottom of this section are two buttons: "Record the Group" and "Get samples".

A second section, titled "Selected instances (answer exceeding groundtruth prediction by bidirectional model prediction distributions)", shows a list of instances. A large grey box highlights a DSL query:

```
1 ENT (g) != ""
2 and count (token (c, pattern=ENT (g))) >
3 count (token (g, pattern=ENT (g)))
4 and ENT (g) == ENT (p (m))
5 and f1 (m) == 0
```

The background shows a list of instances with columns for "question", "context", "groundtruth", and "ent_g". The first instance is "Who created the...". The second instance is "A different arrangement...". The third instance is "Keff McCloskey's suspension in...". The fourth instance is "American...". The fifth instance is "Who in 1996...".

On the right side, there is a "GROUPS" panel with a list of groups: "all_inst...", "is_entity", "has_dis...", "correct...", "is_distr...", and "REWR...".

Build distractor groups **with DSL**

ENT (Murray Gold) == PERSON

```
1 ENT (g) != ""  
2   and count(token(c, pattern=ENT(g))) >  
3     count(token(g, pattern=ENT(g)))  
4   and ENT(g) == ENT(p(m))  
5   and f1(m) == 0
```

`is_entity`

"The **g**roundtruth is an **ENT**ity."

Build distractor groups with DSL

```
1 ENT(g) != ""
2 and count(token(c, pattern=ENT(g))) >
3   count(token(g, pattern=ENT(g)))
4 and ENT(g) == ENT(p(m))
5 and f1(m) == 0
```

count(PERSON : Murray Gold, John Dubney, Ron Grainer) == 1

is_entity
has_distractor

count(PERSON : Murray Gold) == 1

“There are more tokens matching the ground truth entity type (**ENT(g)**) in the whole **c**ontext than in the **g**roundtruth.”

Build distractor groups **with DSL**

```
1 ENT(g) != ""
2 and count(token(c, pattern=ENT(g))) >
3   count(token(g, pattern=ENT(g)))
4 and ENT(g) == ENT(p(m))
5 and f1(m) == 0 → ENT(John Debney) == PERSON
```

is_entity
has_distractor
correct_type

“The **m**odel **p**rediction **ENT**ity type matches the
groundtruth **ENT**ity type.”

Build distractor groups **with DSL**

```
1 ENT(g) != ""
2   and count(token(c, pattern=ENT(g))) >
3     count(token(g, pattern=ENT(g)))
4   and ENT(g) == ENT(p(m))
5   and f1(m) == 0
```

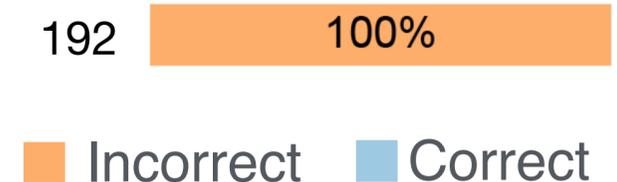
```
is_entity
has_distractor
correct_type
is_distracted
```

"The **m**odel prediction is *incorrect*."

Build distractor groups with DSL

```
1 ENT(g) != ""
2 and count(token(c, pattern=ENT(g))) >
3   count(token(g, pattern=ENT(g)))
4 and ENT(g) == ENT(p(m))
5 and f1(m) == 0
```

```
is_entity
has_distractor
correct_type
is_distracted
```



5.7% of all BiDAF errors:

The distractor hypothesis seems **correct**!

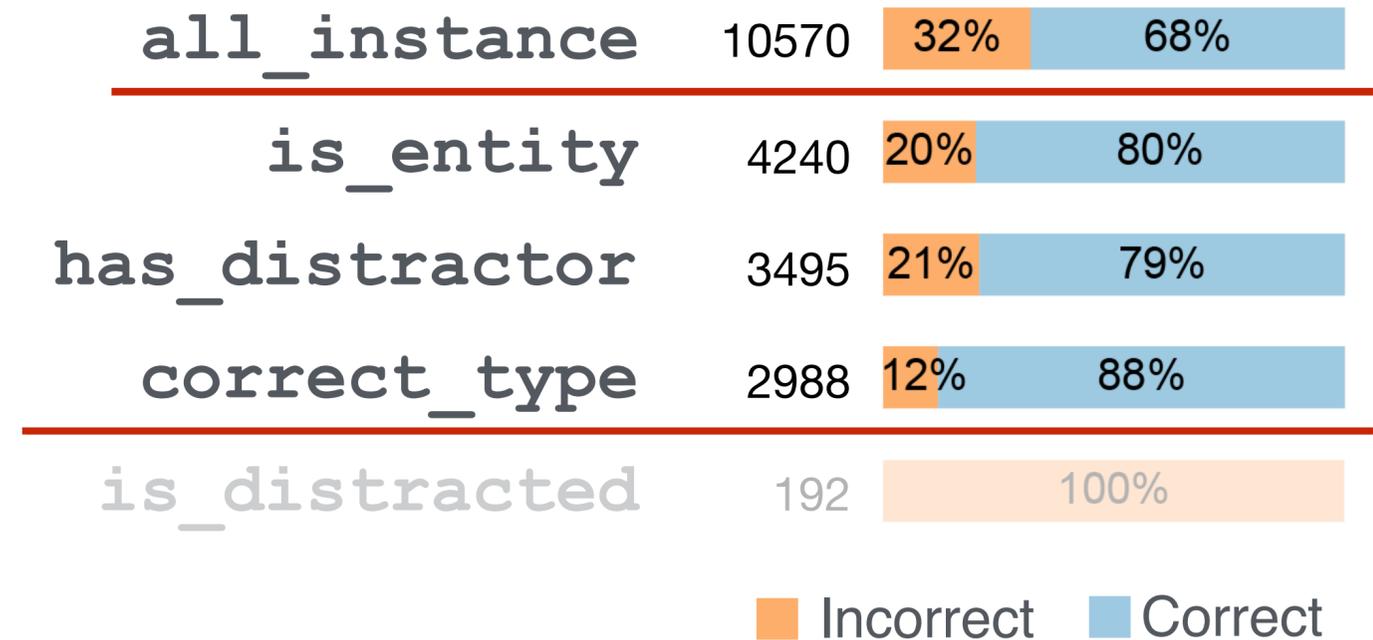
Build distractor groups with DSL

```
1 ENT (g) != ""
2 and count(token(c, pattern=ENT(g))) >
3   count(token(g, pattern=ENT(g)))
4 and ENT(g) == ENT(p(m))
5 and f1(m) == 0
```

88% EM > 68% EM:

BiDAF performs **better** when have distractors & entity type is matched, than overall.

Reject / revise the hypothesis!



Building blocks

Quantitative grouping

Inspect similar instances,
semantically & syntactically

Counterfactual perturbation

Isolate important components
targeted minimal rewrites

Errudite

Precise & reproducible hypotheses

+

Scale up to the entire dev set

+

Cover **errors & correct** instances

+

Test via counterfactual analysis

Scenario: distractor hypothesis

is_distracted 192  100%

HAS distractor prediction

≠

IS WRONG due to distractor prediction

Who created the 2005 theme for **Doctor Who**?

...**John Debney** created a new arrangement of Ron Grainer's original theme for Doctor Who in 1996. For the return of the **series** in 2005, **Murray Gold** provided a new arrangement... sampled from the 1963 original.

Distractor entity?

Multi-sentence reasoning?



C Instances to explore (without edits)

Get Instances Sample 10 instances randomly that are in Select groups and not in Select groups

Filter CMD ENT (groundtruth) == ""

Preview the filter on 10570 instances

Filtered instances: NaN (0.0% of total), error: undefined (NaN% of slice, NaN% of total, NaN% of all errors)

Record the Group Get samples

D Selected instances (answer encoding: groundtruth, prediction by bidaf (correct, incorrect), model prediction distributions)

Who created the 2005 theme for Doctor Who?

A different arrangement was recorded by Peter Howell for season 18 (1980), which was in turn replaced by another arrangement for the season-long serial The Trial of a Time Lord in season 23 (1986). Keff McCulloch provided the new arrangement for the Seventh Doctor's era which lasted from season 24 to season 26, followed by a suspension in 1989. An Irish composer John Debney created a new arrangement of Ron Grainer's original theme for Doctor Who in 1996. For the return of the series in 2005, Murray Gold provided a new arrangement which featured samples from the original theme; in the 2005 Christmas episode "The Christmas Invasion", Gold introduced a modified arrangement that was used up until the conclusion of the 2007 series. [citation needed]

DISQUIN... LER IN... T...
 starts_with(prediction(model="bidaf"), pattern="NNP")
 starts_with(prediction(model="bidaf"), pattern="PERSON")
 attr:answer_type == answer_type(prediction(model="bidaf"))
 exact_match(model="bidaf") == 0
 is_correct_sent(prediction(model="bidaf")) == 0
 overlap(question, sentence(prediction(model="bidaf"))) > overlap(question, sentence(groundtruth))

Prev page Next page
 Displaying #0-4 samples.



F Add Re-write Rules

Proportion Show filtered distribution Add a rule Save rules

Details for a rewrite rule

APPLY THE CHANGE TO context

CHANGE FROM PATTERN STRING (prediction (model="bidaf"))

CHANGE TO PATTERN #

Are the 192 instances really wrong because of the distractor?

Would BiDAF work perfectly if we remove the distractors?

Answer what-if questions with counterfactual analysis!

Counterfactual Analysis with Rewrite Rules

`rewrite (target , from → to)`

Re-write the **target** part of an instance
by replacing **from**
with **to**

Counterfactual Analysis with Rewrite Rules

Would BiDAF work perfectly if we **remove** the distractors?

`rewrite (target , from → to)`

Re-write the **target** part of an instance
by replacing **from**
with **to**

Counterfactual Analysis with Rewrite Rules

Would BiDAF work perfectly if we **remove** the distractors?

```
rewrite ( c , from → to )
```

Re-write the **context** part of an instance
by replacing **from**
with **to**

Counterfactual Analysis with Rewrite Rules

Would BiDAF work perfectly if we **remove** the distractors?

```
rewrite ( c , string ( p ( m ) ) → to )
```

Re-write the **context** part of an instance
by replacing the **model predicted distractor string**
with **to**

Counterfactual Analysis with Rewrite Rules

Would BiDAF work perfectly if we **remove** the distractors?

```
rewrite( c , string(p(m)) → "#" )
```

Re-write the **context** part of an instance
by replacing the **model predicted distractor string**
with a placeholder token **"#"**

Counterfactual Analysis with Rewrite Rules

Would BiDAF work perfectly if we **remove** the distractors?

```
rewrite ( c , string ( p ( m ) ) → "#" )
```

Incorrect
Incorrect

Q: Who created the 2005 theme for Doctor Who?

C: ...~~John Dobney~~ # created a new arrangement of Ron Grainer's ...
Murray Gold provided a new arrangement...

Counterfactual Analysis with Rewrite Rules

Would BiDAF work perfectly if we **remove** the distractors?

```
rewrite ( c , string ( p ( m ) ) → "#" )
```

↪ Incorrect
↪ Incorrect

☑ Another distractor is still confusing the model!

Counterfactual Analysis with Rewrite Rules

```
rewrite( c , string(p(m)) → "#" )
```

p(m) for the 192 rewritten `is_distracted` instances are...

Incorrect
Incorrect

✔ Another distractor is still confusing the model!

Counterfactual Analysis with Rewrite Rules

```
rewrite ( c , string(p(m)) → "#" )
```

`p(m)` for the 192 rewritten `is_distracted` instances are...

↪ Incorrect
↪ Incorrect

29%

✓ Another distractor is still confusing the model!

Incorrect
Correct

48%

✓ The distractor was fooling the model!

age of 18, ~~10.5%~~ # from 18 to 24...

Unchanged

23%

✗ Other factors are at play!

Building blocks

Quantitative grouping

Inspect similar instances,
semantically & syntactically

Counterfactual perturbation

Isolate important components
targeted minimal rewrites

Errudite

Precise & reproducible hypotheses

+

Scale up to the entire dev set

+

Cover **errors & correct** instances

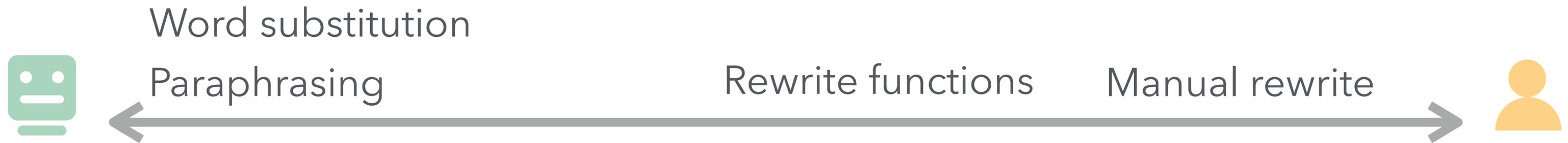
+

Test via counterfactual analysis

Grouping



Counterfactual rewrite



Errudite

✓ scalable, interpretable

Grouping

rules + *guiding matrix*

Model testing



Clustering

Filtering rules

Manual categorize



x Prone-to human error, not representative

Counterfactual rewrite

error analysis



Word substitution

Paraphrasing

- Precise
- Reproducible
- Scalable
- Testable

- will then...**
- Uncover bugs
- Improve the state-of-art
- Safeguard deployments

ACL 2020 Best Paper Award

Beyond Accuracy: Behavioral Testing of NLP Models with **Checklist**

Marco Tulio Ribeiro
Microsoft Research

Tongshuang (Sherry) Wu @tongshuangwu
University of Washington

Jeffrey Heer @jeffrey_heer
Daniel S. Weld @dsweld
University of Washington



CheckList – Framework + Tooling

Applying the principles for Software Engineering testing to NLP

Software engineering → NLP

Capabilities	Descriptions
Vocab/POS	important words or word types for the task.
Named entities	appropriately understanding named entities.
Nagation	understand the negation words.
Taxonomy	synonyms, antonyms, etc.
Robustness	to typos, irrelevant changes, etc.
Coreference	resolve ambiguous pronouns, etc.
Fairness	not biasing towards certain gender/race groups.
Semantic Role Labeling	understanding roles such as agent, object, etc.
Logic	handle symmetry, consistency, and conjunctions.
Temporal	understand order of events.

Principle: test small units



What to test: capabilities

Why do we have the universal list?

Models' **required capabilities** are task-independent.

Models' **expected behaviors** w.r.t capabilities are task-dependent.

This is **not** an exhaustive list!

Software engineering → NLP

Capabilities

Vocab/POS

Named entities

Nagation

...

Behavioral testing: decouple tests from implementation



Decouple tests from training

Meets users' needs
Works with black box models

Software engineering → NLP

Capabilities			
Vocab/POS			
Named entities			
Nagation			
...			

Behavioral testing: decouple tests from implementation



Decouple tests from training

How to test:

Test behaviors with different test types!

Illustrating task: **sentiment analysis**
with **Google Cloud's Natural Language**



Software engineering → NLP

Capabilities	MFT		
Vocab/POS			
Named entities			
Nagation			
...			

Unit tests: known in-/out-puts



Minimum Functionality Test

Software engineering → NLP

Capabilities	MFT		
Vocab/POS			
Named entities			
Nagation			
...			

Unit tests: known in-/out-puts



Minimum Functionality Test

Expectation: Exact labels

This was a great flight. (positive)

I hated this seat. (negative)



A group of n=500 test cases

Software engineering → NLP

Capabilities	MFT		
Vocab/POS	Pos/Neg: 15%	←	1 test, with failure rate
Named entities			
Nagation			
...			

Unit tests: known in-/out-puts



Minimum Functionality Test

Expectation: Exact labels

This was a great flight. (positive)

I hated this seat. (negative)



A group of $n=500$ test cases

Software engineering → NLP

Capabilities	MFT		
Vocab/POS	Pos/Neg: 15%		
Named entities			
Nagation			
...			

Unit tests: known in-/out-puts



Minimum Functionality Test

Expectation: Exact labels

This was a great flight. (positive)
I hated this seat. (negative)

Expectation: Exact labels

This is a commercial flight. (neutral)
I flew to Indiana yesterday. (neutral)

Software engineering → NLP

Capabilities	MFT		
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%	← multiple tests per cell	
Named entities			
Nagation			
...			

Unit tests: known in-/out-puts



Minimum Functionality Test

Expectation: Exact labels

This was a great flight. (positive)
I hated this seat. (negative)

Expectation: Exact labels

This is a commercial flight. (neutral)
I flew to Indiana yesterday. (neutral)

Software engineering → NLP

Capabilities	MFT		
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		
Named entities			
Nagation			
...			

Unit tests: known in-/out-puts



Minimum Functionality Test

Expectation: Exact labels

The cabin crew was not great. (negative)

I can't say I enjoyed the food. (negative)

Software engineering → NLP

Capabilities	MFT		
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		
Named entities			
Nagation	Easy: 49.2%		
...			

Unit tests: known in-/out-puts



Minimum Functionality Test

Expectation: Exact labels

The cabin crew was not great. (negative)

I can't say I enjoyed the food. (negative)

Software engineering → NLP

Capabilities	MFT		
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		
Named entities			
Nagation	Easy: 49.2%		
...			

Metamorphic (perturbations)
& property-based testing

~~Start from scratch~~ → Perturb existing ones

~~Expect exact label~~ → Expect predictions to (not) change

Software engineering → NLP

Capabilities	MFT	INV	
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		
Named entities			
Nagation	Easy: 49.2%		
...			

Metamorphic (perturbations)
& property-based testing



INVariance Tests

Software engineering → NLP

Capabilities	MFT	INV	
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		
Named entities			
Nagation	Easy: 49.2%		
...			

Metamorphic (perturbations)
& property-based testing



INVariance Tests

*No need to specify
the exact prediction!*

Expectation: Same prediction after the change.

@AmericanAir thank you we got on a different flight to ~~Chicago~~ Dallas.

@VirginAmerica I can't lose my luggage, moving to ~~Brazil~~ Turkey soon.

Software engineering → NLP

Capabilities	MFT	INV	
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		
Named entities		LOC: 21%	
Nagation	Easy: 49.2%		
...			

Metamorphic (perturbations)
& property-based testing



INVariance Tests

*No need to specify
the exact prediction!*

Expectation: Same prediction after the change.

@AmericanAir thank you we got on a different flight to ~~Chicago~~ Dallas.

@VirginAmerica I can't lose my luggage, moving to ~~Brazil~~ Turkey soon.

Software engineering → NLP

Capabilities	MFT	INV	DIR
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		
Named entities		LOC: 21%	
Nagation	Easy: 49.2%		
...			

Metamorphic (perturbations)
& property-based testing



INVariance Tests

DIRectional Expectation Tests

Software engineering → NLP

Capabilities	MFT	INV	DIR
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		
Named entities		LOC: 21%	
Nagation	Easy: 49.2%		
...			

Metamorphic (perturbations)
& property-based testing



INVariance Tests

DIRectional Expectation Tests

Expectation: Sentiment monotonic decreasing (↓)

@AmericanAir service wasn't great. *You are lame.*

@JetBlue why won't YOU help them?! Ugh. *I dread you.*

*expectation on
probability!*

Software engineering → NLP

Capabilities	MFT	INV	DIR
Vocab/POS	Pos/Neg: 15% Neutral: 7.6%		Add neg: 34.6%
Named entities		LOC: 21%	
Nagation	Easy: 49.2%		
...			

Metamorphic (perturbations)
& property-based testing



INVariance Tests

DIRectional Expectation Tests

Expectation: Sentiment monotonic decreasing (↓)

@AmericanAir service wasn't great. *You are lame.*

@JetBlue why won't YOU help them?! Ugh. *I dread you.*

*expectation on
probability!*

NLP testing in a nutshell: fill in the matrix

Tests are **grouped by** (capability, test type, expectation).

how?

what?

Capabilities	MFT	INV	DIR
Vocab/POS	✓	✗	✗
Named entities	✓	✓	✗
Nagation	✗	✓	✗
...			

Find a cell of (cap, test type)

Define (maybe ≥ 1) tests

test = test case + expectation

Run the model, get passes/fails

Form a test suite – reuse for other models!

CheckList – Framework + Tooling

Abstractions that ease the pain of the test generation, increase coverage.

CheckList as a tool

Templates

RoBERTa suggestions

Lexicons

Perturbation library

Expectation functions

Test inspecting/sharing

Visualization

CheckList as a tool

Templates

RoBERTa suggestions

Lexicons

Perturbation library

Expectation functions

Test inspecting/sharing

Visualization

More in our repo!

<https://github.com/marcotcr/checklist>

CheckList as a tool

```
In [27]: ▶ editor.visual_suggest('This is {a:mask} movie.')
```

> This is **a:mask** movie .

FILL IN WITH...

- Check All
- a good
- an amazing
- an excellent
- an awful

Preview

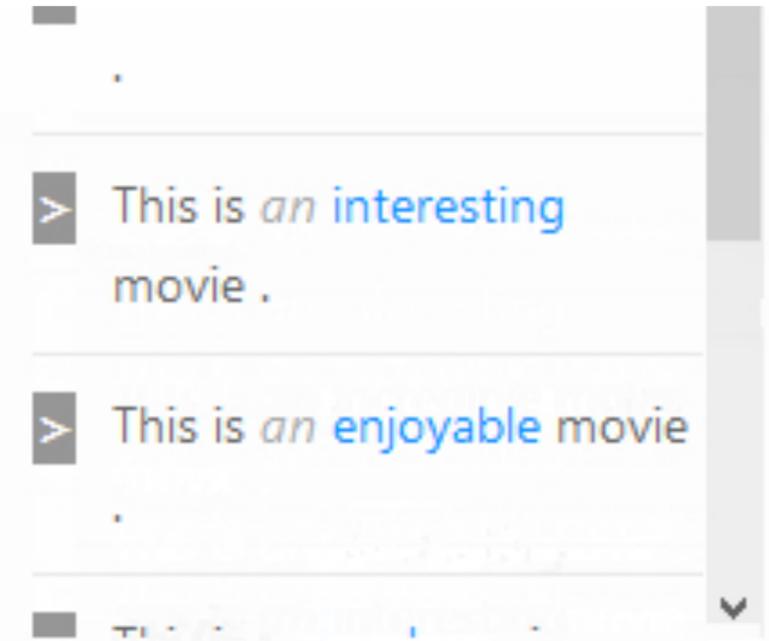
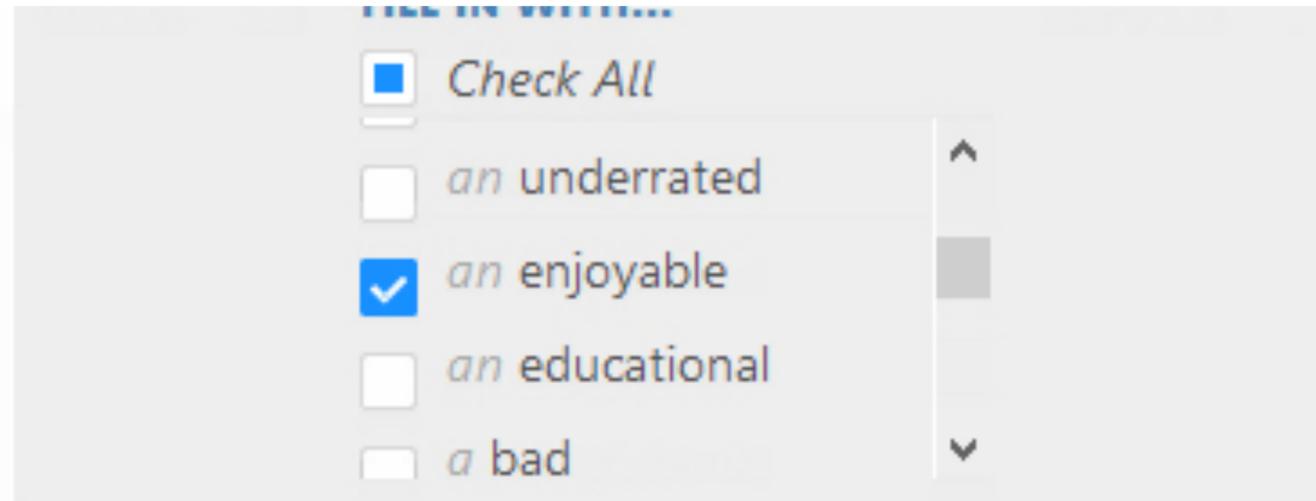


No Data

```
In [26]: ▶ editor.selected_suggestions
```

Wordnet

CheckList as a tool



```
In [28]: ▶ editor.selected_suggestions
```

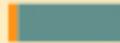
```
Out[28]: ['good',  
          'amazing',  
          'excellent',  
          'interesting',  
          'awesome',  
          'incredible',  
          'great',  
          'enjoyable']
```

Capabilities	Minimum Functionality Test <i>failure rate % (over N tests)</i>	INVariance Test <i>failure rate % (over N tests)</i>	DIRectional Expectation Test <i>failure rate % (over N tests)</i>
+ Vocabulary	100.0% (5)	10.2% (1)	0.8% (4)
+ Robustness		11.4% (5)	
+ NER		7.6% (3)	
+ Fairness		96.4% (4)	
+ Temporal	18.8% (1)		100.0% (1)
+ Negation	99.8% (9)		
+ SRL	100.0% (5)		

- ▶
- ▶
- ▶
- ▶
- ▶
- ▶
- ▶
- ▶

+ NER		7.6% (3)	
+ Fairness		96.4% (4)	
+ Temporal	18.8% (1)		100.0% (1)
- Negation	99.8% (9)		

MINIMUM FUNCTIONALITY TEST

	test name	failure rate
+	simple negations: negative	42 / 500 = 8.4% 
+	simple negations: not negative	66 / 500 = 13.2% 
+	simple negations: not neutral is still neutral	492 / 500 = 98.4% 
+	simple negations: I thought x was positive, but it was not (should be negative)	11 / 500 = 2.2% 
+	simple negations: I thought x was negative, but it was not (should be neutral or positive)	424 / 500 = 84.8% 
+	simple negations: but it was not (neutral) should still be neutral	493 / 500 = 98.6% 
+	Hard: Negation of positive with neutral stuff in the middle (should be negative)	370 / 500 = 74.0% 
+	Hard: Negation of negative with neutral stuff in the middle	

+ simple negations: not negative	66 / 500 = 13.2%	
+ simple negations: not neutral is still neutral	492 / 500 = 98.4%	
+ simple negations: I thought x was positive, but it was not (should be negative)	11 / 500 = 2.2%	
+ simple negations: I thought x was negative, but it was not (should be neutral or positive)	424 / 500 = 84.8%	
+ simple negations: but it was not (neutral) should still be neutral	493 / 500 = 98.6%	
- Hard: Negation of positive with neutral stuff in the middle (should be negative)	370 / 500 = 74.0%	

Test Summary

Test [MFT] on [NEGATION]

Hard: Negation of positive with neutral stuff in the middle (should be negative)

Result FAILURE RATE ON ALL CASES

370/500=74.0%

FILTER TEST CASES

Input free text and enter

Examples Failed cases only

- > I would n't say , given that I am from Brazil , that this food was extraordinary . Expect: 0 | Pred: 2 (1.00)
- > I would n't say , given it 's a Tuesday , that that is a beautiful aircraft . Expect: 0 | Pred: 2 (1.00)
- > I would n't say , given that I am from Brazil , that the service is wonderful . Expect: 0 | Pred: 2 (1.00)
- > I ca n't say , given all that I 've seen . Expect: 0 | Pred: 2 (1.00)

+ Hard: Negation of negative with neutral stuff in the middle

400 / 500 = 80.0%



+ simple negations: not negative	66 / 500 = 13.2%	
+ simple negations: not neutral is still neutral	492 / 500 = 98.4%	
+ simple negations: I thought x was positive, but it was not (should be negative)	11 / 500 = 2.2%	
+ simple negations: I thought x was negative, but it was not (should be neutral or positive)	424 / 500 = 84.8%	
+ simple negations: but it was not (neutral) should still be neutral	493 / 500 = 98.6%	
- Hard: Negation of positive with neutral stuff in the middle (should be negative)	370 / 500 = 74.0%	

Test Summary

Test [MFT] on [NEGATION]

Hard: Negation of positive with neutral stuff in the middle (should be negative)

Result FAILURE RATE ON ALL CASES

370/500=74.0%

FILTER TEST CASES

Input free text and enter

Examples Failed cases only

- > I would n't say , given that I am from Brazil , that this food was extraordinary . Expect: 0 | Pred: 2 (1.00)
- > I would n't say , given it 's a Tuesday , that that is a beautiful aircraft . Expect: 0 | Pred: 2 (1.00)
- > I would n't say , given that I am from Brazil , that the service is wonderful . Expect: 0 | Pred: 2 (1.00)
- > I ca n't say , given all that I 've seen . Expect: 0 | Pred: 2 (1.00)

+ Hard: Negation of negative with neutral stuff in the middle

400 / 500 = 80.0%



+ simple negations: not negative	66 / 500 = 13.2%	
+ simple negations: not neutral is still neutral	492 / 500 = 98.4%	
+ simple negations: I thought x was positive, but it was not (should be negative)	11 / 500 = 2.2%	
+ simple negations: I thought x was negative, but it was not (should be neutral or positive)	424 / 500 = 84.8%	
+ simple negations: but it was not (neutral) should still be neutral	493 / 500 = 98.6%	
- Hard: Negation of positive with neutral stuff in the middle (should be negative)	370 / 500 = 74.0%	

Test Summary

Test [MFT] on [NEGATION]

Hard: Negation of positive with neutral stuff in the middle (should be negative)

Result FAILURE RATE ON ALL CASES

370/500=74.0%

FILTER TEST CASES

Input free text and enter

Examples Failed cases only

- > I would n't say , given that I am from Brazil , that this food was extraordinary . Expect: 0 | Pred: 2 (1.00)
- > I would n't say , given it 's a Tuesday , that that is a beautiful aircraft . Expect: 0 | Pred: 2 (1.00)
- > I would n't say , given that I am from Brazil , that the service is wonderful . Expect: 0 | Pred: 2 (1.00)
- > I ca n't say , given all that I 've seen ... Expect: 0 | Pred: 2 (1.00)

+ Hard: Negation of negative with neutral stuff in the middle

400 / 500 = 80.0%



+ simple negations: not negative	66 / 500 = 13.2%	
+ simple negations: not neutral is still neutral	492 / 500 = 98.4%	
+ simple negations: I thought x was positive, but it was not (should be negative)	11 / 500 = 2.2%	
+ simple negations: I thought x was negative, but it was not (should be neutral or positive)	424 / 500 = 84.8%	
+ simple negations: but it was not (neutral) should still be neutral	493 / 500 = 98.6%	
- Hard: Negation of positive with neutral stuff in the middle (should be negative)	370 / 500 = 74.0%	

Test Summary

Test [MFT] on [NEGATION]

Hard: Negation of positive with neutral stuff in the middle (should be negative)

Result FAILURE RATE ON ALL CASES

370/500=74.0%

FILTER TEST CASES

Input free text and enter

Examples Failed cases only

- > I would n't say , given that I am from Brazil , that this food was extraordinary . Expect: 0 | Pred: 2 (1.00)
- > I would n't say , given it 's a Tuesday , that that is a beautiful aircraft . Expect: 0 | Pred: 2 (1.00)
- > I would n't say , given that I am from Brazil , that the service is wonderful . Expect: 0 | Pred: 2 (1.00)
- > I ca n't say , given all that I 've seen . Expect: 0 | Pred: 2 (1.00)

+ Hard: Negation of negative with neutral stuff in the middle

400 / 500 = 80.0%





This is too simple, you won't find any bugs.

Case Study & User Study

*Let's test some SOTA models (that some people **consider solved**)!
sentiment analysis, QQP, QA*



Sentiment analysis

Task Twitter sentiment analysis

@AmericanAir thank you for a delightful flight to Chicago!
(positive)

Claimed to be a use case by all commercial models!

Models

Commercial models

Microsoft's Text Analytics

Google Cloud's Natural Language

Amazon's Comprehend

Research models

BERT (trained on SST-2)

RoBERTa (trained on SST-2)

[.https://azure.microsoft.com/en-us/services/cognitive-services/text-analytics/](https://azure.microsoft.com/en-us/services/cognitive-services/text-analytics/)

<https://cloud.google.com/natural-language>

<https://aws.amazon.com/cn/comprehend/>

Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A. Y., & Potts, C. (2013, October). Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing* (pp. 1631-1642).

Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS			
Taxonomy			
Robustness			
NER			
Fairness			
Temporal			
Nagation			
Coreference			
SRL			
Logic			
...			

Replace neutral words with BERT

Inputs (n=500) & expectations

~~the~~ ~~our~~ nightmare continues **(INV)**

@Virgin should I be concerned ~~that~~ ~~when~~ I'm about to fly... **(INV)**

Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS		X	
Taxonomy			
Robustness			
NER			
Fairness			
Temporal			
Nagation			
Coreference			
SRL			
Logic			
...			

Replace neutral words with BERT

Inputs (n=500) & expectations

~~the~~ our nightmare continues (INV)

@Virgin should I be concerned ~~that~~ when I'm about to fly... (INV)

				RoBERTa
9.4	16.2	12.4	10.2	10.2

Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS		X	
Taxonomy			
Robustness			
NER			
Fairness			
Temporal			
Nagation			
Coreference			
SRL			
Logic			
...			

Add negative phrases

Inputs (n=500) & expectations

@SouthwestAir ok, gotcha! I abhor you (↓)

Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS		✗	✗
Taxonomy			
Robustness			
NER			
Fairness			
Temporal			
Nagation			
Coreference			
SRL			
Logic			
...			

Add negative phrases

Inputs (n=500) & expectations

@SouthwestAir ok, gotcha! I abhor you (↓)



Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS		X	X
Taxonomy			
Robustness			
NER			
Fairness			
Temporal			
Nagation			
Coreference			
SRL			
Logic			
...			

Add random url or @

Inputs (n=500) & expectations

@JetBlue that selfie was extreme. @pi9QDK (INV)

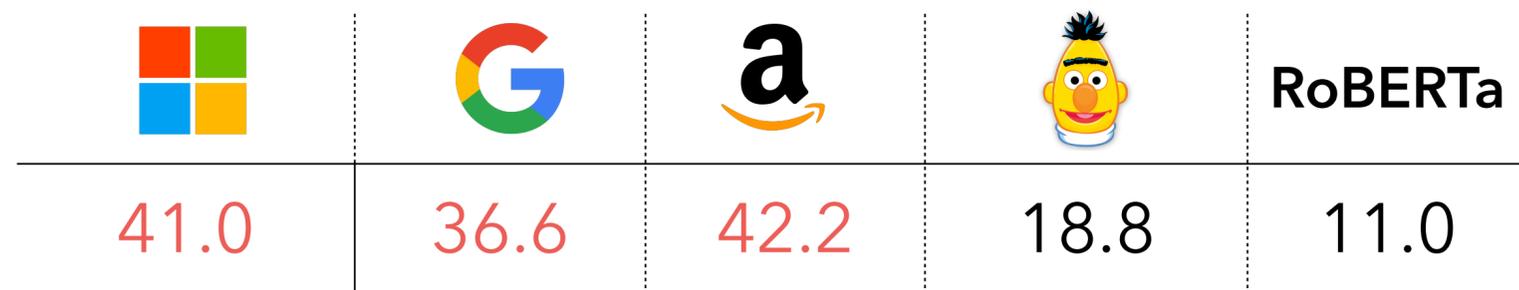
				RoBERTa
9.6	13.4	24.8	11.4	7.4

Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS		X	X
Taxonomy			
Robustness		X	
NER			
Fairness			
Temporal			
Nagation			
Coreference			
SRL			
Logic			
...			

Temporal change

Inputs (n=500) & expectations
 I used to hate this airline, although now I like it (Pos)

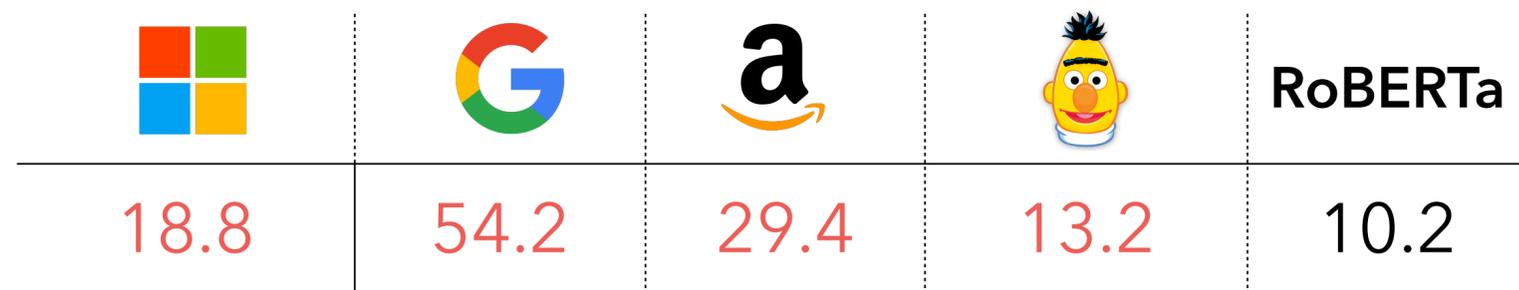


Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS		X	X
Taxonomy			
Robustness		X	
NER			
Fairness			
Temporal	X		
Nagation			
Coreference			
SRL			
Logic			
...			

Negated negation

Inputs (n=500) & expectations
 It wasn't a lousy customer service (Pos or Neutral)



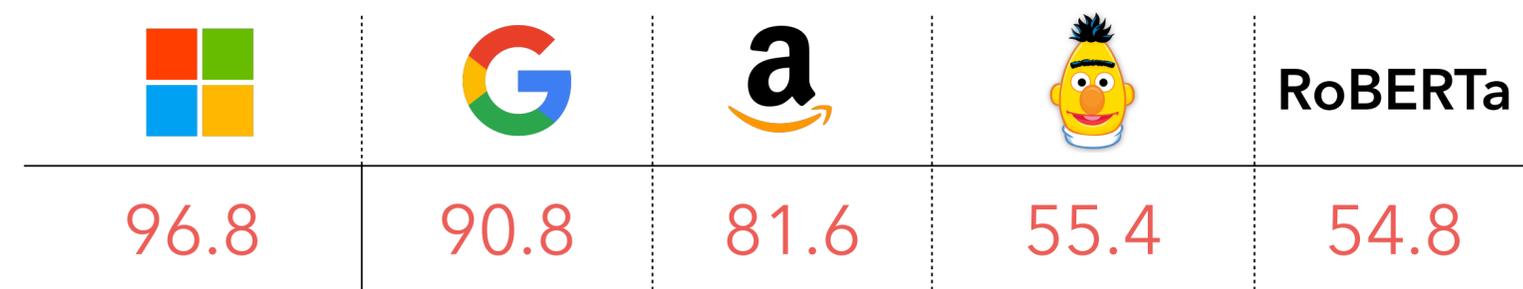
Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS		X	X
Taxonomy			
Robustness		X	
NER			
Fairness			
Temporal	X		
Nagation	X		
Coreference			
SRL			
Logic			
...			

Q&A form

Inputs (n=500) & expectations

Do I think this company is bad? No (**Pos or Neutral**)



Sentiment analysis

Capabilities	MFT	INV	DIR
Vocab/POS		X	X
Taxonomy			
Robustness		X	
NER			
Fairness			
Temporal	X		
Nagation	X		
Coreference			
SRL	X		
Logic			
...			

Q&A form

Inputs (n=500) & expectations

Do I think this company is bad? No (**Pos or Neutral**)

				RoBERTa
96.8	90.8	81.6	55.4	54.8

Case study: Microsoft Sentiment Analysis



Model already stress tested, continue to improve

Public benchmarks

In-house benchmarks (e.g. negation)

User complaint benchmarks

CheckList: 5 hour session

Find many new bugs

Test new capabilities

Test old capabilities better

Same process, different tasks & models

Sentiment analysis

Question Pair Detection

Question Answering

Label: duplicate ≡, or non-duplicate ≠; INV: sam		Failure Rate		Test TYPE and Description	Example Test cases (with expected behavior and prediction)		
Test TYPE and Description	Failure Rate	RoB	Failure Rate (%)				
Vocab.	MFT: Modifiers changes question intent	78.4	78.0	{ Is Mark Wr	MFT: comparisons	20.0	C: Victoria is younger than Dylan. Q: Who is less young? A: Dylan → Victoria
Taxonomy	MFT: Synonyms in simple templates	22.8	39.2	{ How can I t	MFT: intensifiers to superlative: most/least	91.3	C: Anna is worried about the project. Matthew is extremely worried about the project. Q: Who is least worried about the project? A: Anna → Matthew
	INV: Replace words with synonyms in real pairs	13.1	12.7	Is it necessa Is it necessa	MFT: match properties to categories	82.4	C: There is a tiny purple box in the room. Q: What size is the box? A: tiny → purple
	MFT: More X = Less antonym(X)	69.4	100.0	{ How can I t	MFT: nationality vs job	49.4	C: Stephanie is an Indian accountant. Q: What is Stephanie's job? A: accountant → Indian accountant
Robust.	INV: Swap one character with its neighbor (typo)	18.2	12.0	{ Why am I	MFT: animal vs vehicles	26.2	C: Jonathan bought a truck. Isabella bought a hamster. Q: Who bought an animal? A: Isabella → Jonathan
	DIR: Paraphrase of question should be duplicate	69.0	25.0	Can I gain w Can I → Do y	MFT: comparison to antonym	67.3	C: Jacob is shorter than Kimberly. Q: Who is taller? A: Kimberly → Jacob
NER	INV: Change the same name in both questions	11.8	9.4	Why isn't H Is Hillary C	MFT: more/less in context, more/less antonym in question	100.0	C: Jeremy is more optimistic than Taylor. Q: Who is more pessimistic? A: Taylor → Jeremy
	DIR: Change names in one question, expect ≠	35.1	30.1	What does I What India t	INV: Swap adjacent characters in Q (typo)	11.6	C: ...Newcomen designs had a duty of about 7 million, but most were closer to 5 million.... Q: What was the ideal duty → uddy of a Newcomen engine? A: INV → 7 million + 5 million
Temporal	DIR: Keep first word and entities of a question, fill in the gaps with RoBERTa; expect ≠	30.0	32.8	Will it be di Will the US	INV: add irrelevant sentence to C	9.8	(no example)
	MFT: Is ≠ used to be, non-duplicate	61.8	96.8	{ Is Jordan Pe	MFT: change in one person only	41.5	C: Both Luke and Abigail were writers, but there was a change in Abigail, who is now a model. Q: Who is a model? A: Abigail → Abigail were writers, but there was a change in Abigail
	MFT: before ≠ after, non-duplicate	98.0	34.4	{ Is it unhealt	MFT: Understanding before/after, last/first	82.9	C: Logan became a farmer before Danielle did. Q: Who became a farmer last? A: Danielle → Logan
Negation	MFT: before becoming ≠ after becoming	100.0	0.0	What was D What was D	MFT: Context has negation	67.5	C: Aaron is not a writer. Rebecca is. Q: Who is a writer? A: Rebecca → Aaron
	MFT: simple negation, non-duplicate	18.6	0.0	{ How can I t	MFT: Q has negation, C does not	100.0	C: Aaron is an editor. Mark is an actor. Q: Who is not an actor? A: Aaron → Mark
Coref	MFT: negation of antonym, should be duplicate	81.6	88.6	{ How can I t	MFT: Simple coreference, he/she.	100.0	C: Melissa and Antonio are friends. He is a journalist, and she is an adviser. Q: Who is a journalist? A: Antonio → Melissa
	MFT: Simple coreference: he ≠ she	79.0	96.6	If Joshua an If Joshua an	MFT: Simple coreference, his/her.	100.0	C: Victoria and Alex are friends. Her mom is an agent Q: Whose mom is an agent? A: Victoria → Alex
SRL	MFT: Simple resolved coreference, his and her	99.6	100.0	If Jack and I If Jack and I	MFT: former/latter	100.0	C: Kimberly and Jennifer are friends. The former is a teacher Q: Who is a teacher? A: Kimberly → Jennifer
	MFT: Order is irrelevant for comparisons	99.6	100.0	{ Are tigers h	MFT: subject/object distinction	60.8	C: Richard bothers Elizabeth. Q: Who is bothered? A: Elizabeth → Richard
	MFT: Orders is irrelevant in symmetric relations	81.8	100.0	{ Is Nicole re	MFT: subj/obj distinction with 3 agents	95.7	C: Jose hates Lisa. Kevin is hated by Lisa. Q: Who hates Kevin? A: Lisa → Jose
Logic	MFT: Order is relevant for asymmetric relations	71.4	100.0	{ Is Sean hurt			
	MFT: Active / passive swap, same semantics	65.8	98.6	{ Does Anna			
	MFT: Active / passive swap, different semantics	97.4	100.0	{ Does Danie			
	INV: Symmetry: pred(a, b) = pred(b, a)	4.4	2.2	{ (q1, q2) (q			
	DIR: Implications, eg. (a=b) ∧ (a=c) ⇒ (b=c)	9.7	8.5	no example			

See our paper!

Building blocks

Quantitative grouping

Inspect similar instances,
semantically & syntactically

Counterfactual perturbation

Isolate important components
targeted minimal rewrites

Checklist

What to test

Capabilities, shared across tasks

+ How to test

Simple examples (MFTs),
perturbations (INVs, DIRs)

+ Tooling

BERT fill-ins, visualizations, lexicons,...

How to fix bugs found in CheckList?

Perturbations as feedback to model training, dataset augmentation, etc.

Data manage.

augment

Model verification

perform - explanation

Submitted to ACL 2021

Polyjuice: Generating Counterfactuals for Explaining, Evaluating, and Improving models

Tongshuang (Sherry) Wu @tongshuangwu
University of Washington

Marco Tulio Ribeiro
Microsoft Research

Jeffrey Heer @jeffrey_heer
Daniel S. Weld @dsweld
University of Washington





Most work!

Word substitution

Paraphrasing

× limited patterns

✓ Scalable

Errudite

CheckList

Snorkel

Challenge sets

Heuristic rules

Rewrite functions

✓ Middle ground!

× Domain expertise

× Tedious

did → didn't

would never → would

Counterfactual data aug.

Contrast Set

Manual rewrite

✓ Diverse

× Costly to scale

Counterfactual rewriting as an NLG task

Prompt: I **did** like the movie.

Generate: I **didn't** like the movie

Language models can **complete paragraphs** → **be finetuned for rewriting.**

More **diverse** patterns from models like GPT-2

Design prompts to teach the model "**how**" and "**where**" to change

did → **didn't**, **would never** → **would**



"add negation modifiers to aux"

how

where

“How to change”: Control codes

Control codes

Definition, color: delete → insert

“How to change”: Control codes

Constraints
concrete

Control codes	Definition, color: delete → insert
negation	A dog is not embraced by the woman.
quantifier	A dog is → Three dogs are embraced by the woman.
shuffle	<i>To move (or swap) key phrases or entities around the sentence.</i> A dog → woman is embraced by the woman → dog .

"How to change": Control codes

Constraints

concrete → loose

Control codes	Definition, color: delete → insert
negation	A dog is not embraced by the woman.
quantifier	A dog is → Three dogs are embraced by the woman.
shuffle	<i>To move (or swap) key phrases or entities around the sentence.</i> A dog → woman is embraced by the woman → dog .
lexical	<i>Changing just one word or noun chunks without breaking the POS tags.</i> A dog is embraced → attacked by the woman.
resemantic	<i>To replace short phrases or clauses without affecting the parsing tree.</i> A dog is embraced by the woman → wrapped in a blanket .

“How to change”: Control codes

Constraints

concrete → loose

Group by

syntactic > semantic

Control codes	Definition, color: delete → insert
negation	A dog is not embraced by the woman.
quantifier	A dog is → Three dogs are embraced by the woman.
shuffle	<i>To move (or swap) key phrases or entities around the sentence.</i> A dog → woman is embraced by the woman → dog .
lexical	<i>Changing just one word or noun chunks without breaking the POS tags.</i> A dog is embraced → attacked by the woman.
resemantic	<i>To replace short phrases or clauses without affecting the parsing tree.</i> A dog is embraced by the woman → wrapped in a blanket .
insert	<i>To add constraints without affecting the parsing structure of other parts.</i> A dog is embraced by the little woman.
delete	<i>To remove constraints without affecting the parsing structure of other parts.</i> A dog is embraced by the woman .
restructure	<i>To alter the dependency tree structure, e.g. changing from passive to positive.</i> A dog is embraced by → hugging he woman.

"How to change": Control codes

No training data?

Merge existing ones!

5 datasets →

191,415 sentence pairs

Control codes
negation
quantifier
shuffle
lexical
resemantic
insert
delete
restructure

PAWS, high lexical overlap, but non-paraphrasing sentences.
Can a **bad** → **good** person be **good** → **bad**?

WinoGrande, commonsense on lexical
The lions ate the zebras because they are **predators** → **meaty**.

HANS, challenge set for Natural Language Inference (NLI).
The banker **near the judge** saw the actor.

“Where to change”: Fill-in-the-blank (Donahue, ACL'20)

(It is great for kids., it is not great for children.)

It is **not** great for **kids** → **children**.

```
It is great for kids. <|perturb|> [negation]
```

```
It is [BLANK] great for [BLANK]. [SEP] not [ANSWER] children [ANSWER]
```

```
<|endoftext|>
```

“Where to change”: Fill-in-the-blank (Donahue, ACL'20)

(It is great for kids., it is not great for children.)

It is **not** great for **kids** → **children**.

It is great for kids. <|perturb|> [negation]

It is [BLANK] great for [BLANK]. [SEP] not [ANSWER] children [ANSWER]

→ It [BLANK] great [BLANK]. [SEP] is not [ANSWER] children [ANSWER]

It is [BLANK]. [SEP] not great for children [ANSWER]

→ [BLANK] [SEP] It is not great for children. [ANSWER]

<|endoftext|>

191,415 sentence pairs → 657,144 training prompts

Filter based on desired tasks

It is **not** great for **kids** → **children**.
It is great for **anyone but** kids.
It is **great** → **disastrous** for kids.
It is **great** → **good** for kids.
It is **great** → **unnecessary** for kids.

auto-augment



semantic dist

analyze "great"



It is **great** → **good** for kids.

It is **great** → **disastrous** for kids.

It is **great** → **good** for kids.

It is **great** → **unnecessary** for kids.

Counterfactual data augmentation

Does the model get more robust afterwards?

Similar to Kaushik et al.

Counterfactual data aug: crowd labeling

Ranking: diversity!

prefer

It is great for anyone but kids.
It is great → disastrous for kids.

over

It is great → good for kids.
It is great → disastrous for kids.

Crowd labeling on MTurk

Validity? Class label?

Three tasks

Sentiment analysis

Natural Language Inference (NLI)

Duplicate question detection (QQP)

Reference Example

Old Q1 Who likes more sex , men or women ?

Old Q2 Is sex more pleasurable for men or for women ?

Label Duplicate

Label the following! [Review the instructions!](#)

The green color highlights new words added in New Q2 , compared to Old Q2 in the Reference example above. • indicates something is deleted.

Old Q1 Who likes more sex , men or women ?

New Q2 Is sex more pleasurable for women or for men ?

Valid? Invalid Valid

Label Non-duplicate Duplicate

Old Q1 Who likes more sex , men or women ?

New Q2 Is sex too pleasurable for men or for women ?

Valid? Invalid Valid

Label Non-duplicate Duplicate

Counterfactual data aug: Training results

Sentiment analysis (on Stanford Sentiment Treebank)

I have not been this **disappointed by** → **excited to see** a movie in a long time. (**negative** → **positive**)
We just **don't** really care about this love story. (**negative** → **positive**)

polyjuice: 4k original + 2k counterfactuals (6k total)

v.s. **baseline** : 4k original + 2k extra, compensating original (6k total)

Model	SST-2	Senti140	SemEval	Amzbook	Yelp	IMDB	IMDB-Cont.	IMDB-CAD
m-baseline	92.9 ± 0.2	88.9 ± 0.3	84.8 ± 0.5	85.1 ± 0.4	90.0 ± 0.3	90.8 ± 0.5	92.2 ± 0.6	86.5 ± 0.2
m-polyjuice	92.7 ± 0.2	90.7 ± 0.4	86.4 ± 0.1	85.6 ± 0.8	90.1 ± 0.0	90.6 ± 0.3	94.0 ± 0.3	89.7 ± 0.5

human generated counterfactuals from prior papers

Improve on Twitter data (Senti140, SemEval), and contrast sets (IMDB-CDA, IMDB-Cont.)

Maintain accuracies on in domain, and other reviews data.

Counterfactual data aug: Training results

Natural Language Inference (on SNLI)

Premise: An airborne man on a surfboard.

Hypothesis: A man in mid air while his surfboard is **beneath** → **lying on** him. (**entail** → **contradict**)

polyjuice: 20k training + 1.5k perturbations, automatically gen. from polyjuice

v.s. **CAD** : 20k training + 1.5k perturbations, manually gen. from Kaushik et al.

v.s. **baseline** : 20k training + 1.5k extra, compensating original

Model	SNLI	MNLI-m	MNLI-mm	SNLI-CAD	break	DNC	stress	diagnostic
m-baseline	85.7 ± 0.4	86.1 ± 0.2	86.6 ± 0.2	72.8 ± 0.3	86.4 ± 1.5	54.5 ± 0.6	65.1 ± 0.6	56.0 ± 0.8
m-CAD	85.8 ± 0.6	86.6 ± 0.1	85.6 ± 0.3	73.8 ± 0.2	89.4 ± 2.9	55.8 ± 0.9	65.5 ± 0.5	56.4 ± 0.4
m-polyjuice	85.3 ± 0.3	86.0 ± 0.1	86.4 ± 0.0	73.6 ± 0.2	89.1 ± 1.2	57.7 ± 0.3	65.1 ± 0.2	57.5 ± 0.5

Similarly, improvements on multiple **contrast/challenge sets**, even better than CAD

Counterfactual data aug: Training results

Natural Language Inference (on SNLI)

Premise: An airborne man on a surfboard.

Hypothesis: A man in mid air while his surfboard is **beneath** → **lying on** him. (**entail** → **contradict**)

polyjuice: 20k training + 1.5k perturbations, automatically gen. from polyjuice

v.s. **CAD** : 20k training + 1.5k perturbations, manually gen. from Kaushik et al.

v.s. **baseline** : 20k training + 1.5k extra, compensating original

Model	SNLI	MNLI-m	MNLI-mm	SNLI-CAD	break	DNC	stress	diagnostic
m-baseline	85.7 ± 0.4	86.1 ± 0.2	86.6 ± 0.2	72.8 ± 0.3	86.4 ± 1.5	54.5 ± 0.6	65.1 ± 0.6	56.0 ± 0.8
m-CAD	85.8 ± 0.6	86.6 ± 0.1	85.6 ± 0.3	73.8 ± 0.2	89.4 ± 2.9	55.8 ± 0.9	65.5 ± 0.5	56.4 ± 0.4
m-polyjuice	85.3 ± 0.3	86.0 ± 0.1	86.4 ± 0.0	73.6 ± 0.2	89.1 ± 1.2	57.7 ± 0.3	65.1 ± 0.2	57.5 ± 0.5
m-polyjuice-rand	85.7 ± 0.4	86.1 ± 0.1	86.2 ± 0.1	73.4 ± 0.5	87.2 ± 0.6	54.7 ± 0.3	64.6 ± 0.6	56.9 ± 0.8

Similarly, improvements on multiple **contrast/challenge sets**, even better than CAD

But only when we do targeted augmentation!

Model already knows **woman** ≠ **man** well, **need to focus on its error cases, e.g., preposition, quantifier!**

v.s. Manual creation: more effective

Kaushik et al.
ICLR'20

Workers spent roughly 5 minutes per revised review, and 4 minutes per revised sentence (for NLI).

Gardner et al.
EMNLP 20

30 seconds to change one image caption sentence.

Verify the machine generated ones, **cheaper than** writing new ones

Focus on **perturbed** parts, **cheaper than** parsing full sentences

→ 30s per round (3 perturbations)

Counterfactual explanation & analysis

Do counterfactuals help with model understanding?

Counterfactual explanation: Complement *Status-quo*

Q1: How can I help a friend experiencing serious depression?

Q2: How do I help a friend who is in depression?

Predict $f(x)$: = *Duplicate* (98.2% confident)

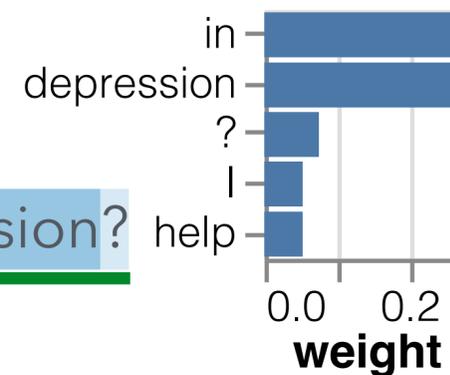
Counterfactual explanation: Complement *Status-quo*

Evaluate with popular saliency score (e.g. SHAP)

Q1: How can I help a friend experiencing serious depression?

Q2: How do I help a friend who is in depression?

Predict $f(x)$: = *Duplicate* (98.2% confident)



Counterfactual explanation: Complement *Status-quo*

Evaluate with popular saliency score (e.g. SHAP)

Q1: How can I help a friend experiencing serious depression?
A serious depression?

Q2: How do I help a friend who is in depression?
Predict $f(x)$: = Duplicate (98.2% confident)

Word	Weight
in	~0.25
depression	~0.25
?	~0.05
I	~0.05
help	~0.05

But...

\hat{x} , perturbed Q2

$f(\hat{x})$

Q2: How do I find a friend who is in depression?

=

✗ salience score is too abstract

Counterfactual explanation: Complement *Status-quo*

Evaluate with popular saliency score (e.g. SHAP)

Q1: How can I help a friend experiencing serious depression?
A

Q2: How do I help a friend who is in depression?
Predict $f(x)$: = Duplicate (98.2% confident)

Word	Weight
in	0.25
depression	0.25
?	0.05
I	0.05
help	0.05

But...

\hat{x} , perturbed Q2

$f(\hat{x})$

Q2: How do I ● find a friend who is in depression?

=

× salience score is too abstract

Q2: How do I help a ● woman who is in depression?

≠

× they miss error regions

Counterfactual explanation: Complement *Status-quo*

Evaluate with popular saliency score (e.g. SHAP)

Q1: How can I help a friend experiencing serious depression?
A

Q2: How do I help a friend who is in depression?
Predict $f(x)$: = Duplicate (98.2% confident)

Word	Weight
in	0.25
depression	0.25
?	0.05
I	0.02
help	0.01

But...

\hat{x} , perturbed Q2

$f(\hat{x})$

Q2: How do I ●find a friend who is in depression? =

Q2: How do I help a ●woman who is in depression? ≠

Q2: How do I help a friend who is ●suicidal? =

× salience score is too abstract

× they miss error regions

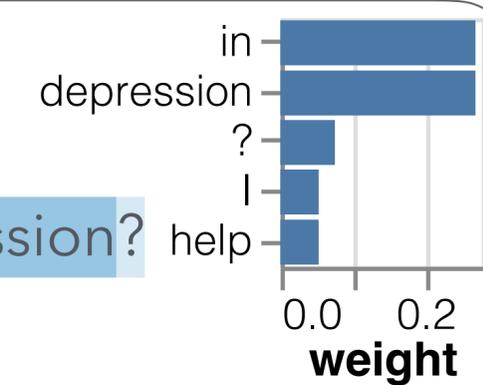
Counterfactual explanation: Complement *Status-quo*

Evaluate with popular saliency score (e.g. SHAP)

Q1: How can I help a friend experiencing serious depression?

A Q2: How do I help a friend who is in depression?

Predict $f(x)$: = *Duplicate* (98.2% confident)



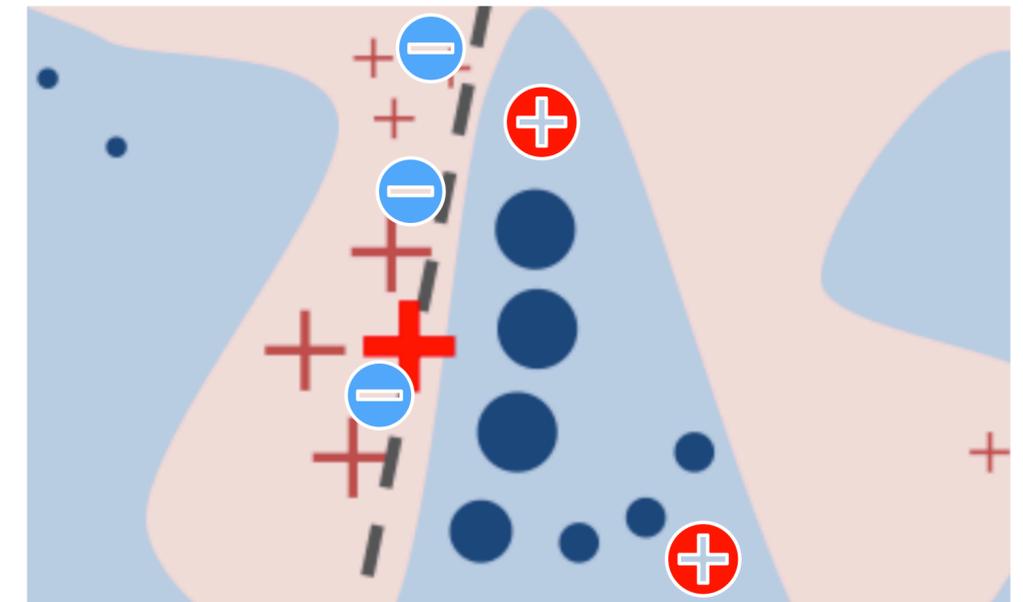
Word	Weight
in	~0.25
depression	~0.25
?	~0.10
I	~0.05
help	~0.05

But...

\hat{x} , perturbed Q2

$f(\hat{x})$

- | | |
|---|---|
| Q2: How do I find a friend who is in depression? | = |
| Q2: How do I help a woman who is in depression? | ≠ |
| Q2: How do I help a friend who is suicidal ? | = |



SHAP: Estimate scores by masking
How <mask> help <mask> <mask> ...

Not reflect model behavior on natural counterfactuals...

Complement SHAP with surprising (violate expectation) counterfactuals

Counterfactual explanation: Complement *Status-quo*

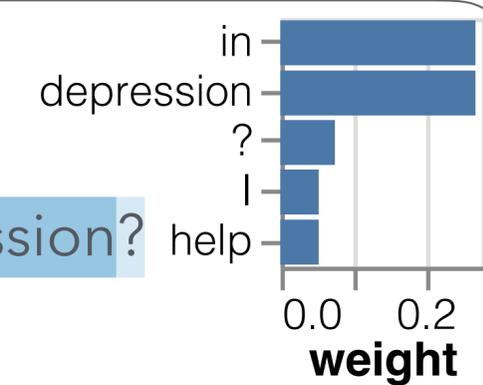
Evaluate with popular saliency score (e.g. SHAP)

Q1: How can I help a friend experiencing serious depression?

A

Q2: How do I help a friend who is in depression?

Predict $f(x)$: = *Duplicate* (98.2% confident)



Word	Weight
in	~0.25
depression	~0.25
?	~0.05
I	~0.02
help	~0.05

But...

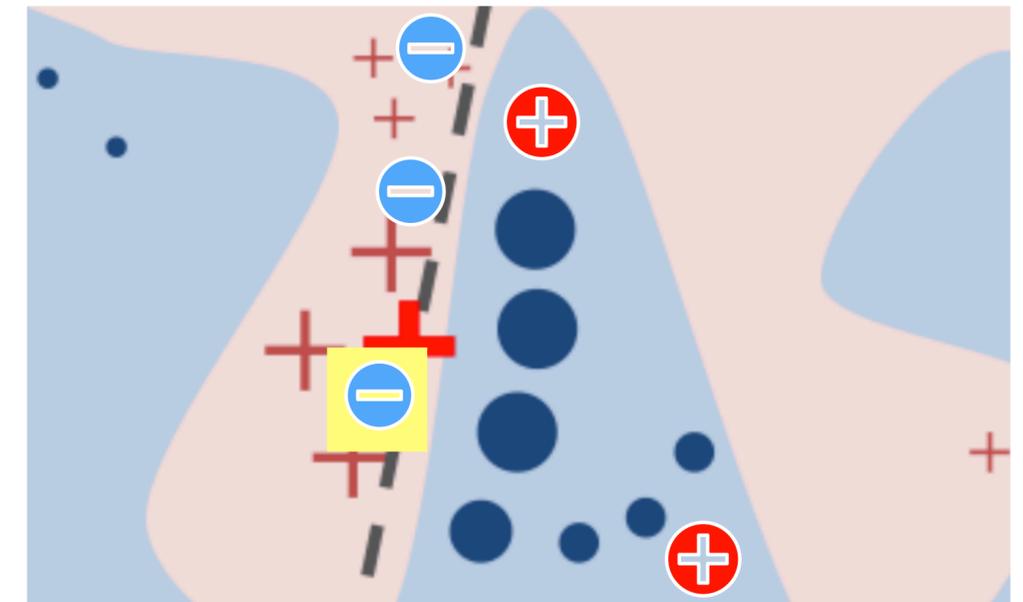
\hat{x} , perturbed Q2

$f(\hat{x})$

Q2: How do I **find** a friend who is in depression? =

Q2: How do I help a **woman** who is in depression? **≠**

Q2: How do I help a friend who is **suicidal**? =



SHAP: Estimate scores by masking
How <mask> help <mask> <mask> ...

Not reflect model behavior on natural counterfactuals...

Complement SHAP with surprising (violate expectation) counterfactuals

Counterfactual explanation: Complement *Status-quo*

Evaluate with popular saliency score (e.g. SHAP)

Q1: How can I help a friend experiencing serious depression?

Q2: How do I help a friend who is in depression?

Predict $f(x)$: = *Duplicate* (98.2% confident)

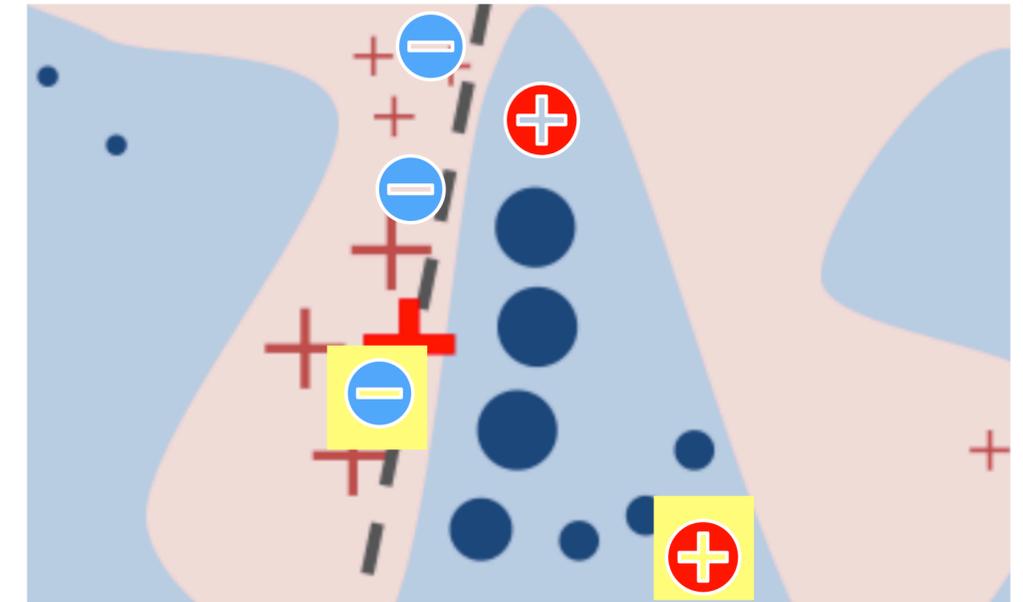
Word	Weight
in	~0.2
depression	~0.2
?	~0.05
I	~0.05
help	~0.2

But...

\hat{x} , perturbed Q2

$f(\hat{x})$

Q2: How do I find a friend who is in depression?	=
Q2: How do I help a woman who is in depression?	≠
Q2: How do I help a friend who is suicidal ?	=



SHAP: Estimate scores by masking
How <mask> help <mask> <mask> ...

Not reflect model behavior on natural counterfactuals...

Complement SHAP with surprising (violate expectation) counterfactuals

User study: Surprising points bring additional insights

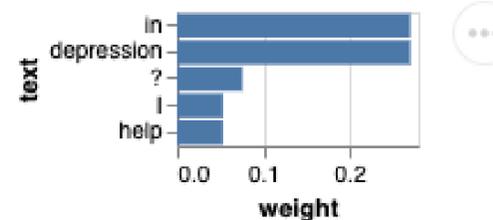
REFERENCE EXAMPLE AND WORD IMPORTANCE

Old Q1 How can I help a friend experiencing serious depression ?

Old Q2 How do I help a friend who is in depression ?

Model predicts Duplicate (97.3% confident)

Model correct? Correct



ASK QUESTIONS!

You can edit Q2 and get model's prediction on variations of this example. Use it wisely, and try to come up with edits that will best help you understand **the model behavior around** the instance!

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I play with a friend who is in depression ?

Get prediction

QUERY RESULTS WILL BE DISPLAYED BELOW

The green color highlights new words added in New Q2, compared to Old Q2 in the Reference example above. ● indicates something is deleted.

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I play with a friend who is in depression ?

Model predicts Non-duplicate (67.4% confident)

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I help a friend who is in frustrated ?

Model predicts Non-duplicate (95.1% confident)

GIVE Original examples, SHAP

Q2 in the Reference example above. ● indicates something is deleted.

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I ask a friend who is in depression ?

Model will predict Non-duplicate Duplicate

GIVE Manual counterfactual analysis

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I help a friend who is in a relationship ?

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I help someone who is depressed ?

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do ● a friend help me who is in depression ?

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I help a friend ● ?

Model will predict Non-duplicate Duplicate

User study: Surprising points bring additional insights

REFERENCE EXAMPLE AND WORD IMPORTANCE

Old Q1 How can I help a friend experiencing serious depression ?

Old Q2 How do I help a friend who is in depression ?

Model predicts Duplicate (97.2% confident)



The bar chart shows the importance of words in the reference example. The y-axis is labeled 'text' and lists 'depression' and 'help'. The x-axis is labeled 'weight' and ranges from 0 to 0.2. The bar for 'depression' is significantly higher than the bar for 'help'.

GET Simulate model's behavior on Polyjuice surprising counterfactuals

13 NLP PhD students did only slightly better than random guessing. (57% correct)

The green color highlights new words added in New Q2, compared to Old Q2 in the Reference example above. ● indicates something is deleted.

these examples are surprising even after seeing explanations and creating manual counterfactuals!

Model predicts Non-duplicate (95.1% confident)

LABEL THE FOLLOWING!

The green color highlights new words added in New Q2, compared to Old Q2 in the Reference example above. ● indicates something is deleted.

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I ask a friend who is in depression ?

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I help a woman who is in depression ?

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I help a friend who is in a relationship ?

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I help someone who is depressed ?

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do ● a friend help me who is in depression ?

Model will predict Non-duplicate Duplicate

Old Q1 How can I help a friend experiencing serious depression ?

New Q2 How do I help a friend ● ?

Model will predict Non-duplicate Duplicate

Error analysis: Extend back to grouping

x	$f(x)$
P: A woman is holding a baby by a window. H: This woman is looking out the window.	<u>Neutral</u>
\hat{x} , perturbed H through [negation]	$f(\hat{x})$
H: No woman is looking out the window.	Contradiction
H: This woman isn't looking out the window.	Contradiction
H: This woman is not looking out the window.	<u>Neutral</u>

Neg correlates with Contradiction in NLI?
Has more nuances!

Inconsistency between "n't" and "not"!

$x \rightarrow f(\hat{x})$	Template	Coverage (%N→C)
...is not looking...	AUX → AUX not	412 (42.3%)
...aren't playing...	* → * not	
The → No girls like...	* → * n't	434 (43.5%)
A → No man in...	* → * PART	180 (92.8%)
	DET → No	

DET → No flips model prediction
much more frequently!

Building blocks

Polyjuice

Quantitative grouping
Inspect similar instances,
semantically & syntactically

Counterfactual perturbation
Isolate important components
targeted minimal rewrites

Supports



General-purpose generator with control, for various applications

Data management

collection preprocess augment analysis

Model training

model select training hyperparameter

Model verification

performance testing

grouping

Counterfactual

Polyjuice (ACL 21) 

Structured data collection

Enhanced counterfactual generator
Human-generator collaboration

Polyjuice 2.0, with structured labels & content keywords

Polyjuice is nice, but not bias free.

- ✗ Control codes are predefined and not exhaustive
- ✗ Paired data makes certain perturbations more or less likely.
- ✗ Does not accept any context-control.

→ Design generators that rely on

Multiple blanks,

linguistic structural labels (SRL), for choosing which blank to fill in + rough content
content signal, for more concrete fill-in control

Polyjuice 2.0, with structured labels & content keywords

→ Design generators that rely on

Multiple blanks,

linguistic structural labels (SRL), for choosing which blank to fill in + rough content
content signal, for more concrete fill-in control

MANNER	*
TEMPORAL	*
CAUSE	*
LOCATIVE	*

And [BLANK] he [BLANK] nearly always bought and sold [BLANK].

[EMPTY] [EMPTY]

like a good guy.

[EMPTY] through a 17-month period

[EMPTY]

moreover because he had always known him,

[EMPTY]

[EMPTY] [EMPTY]

in the United States

Polyjuice 2.0, with structured labels & content keywords

→ Design generators that rely on

Multiple blanks,

linguistic structural labels (SRL), for choosing which blank to fill in + rough content
content signal, for more concrete fill-in control

MANNER	*
TEMPORAL	*
CAUSE	*
LOCATIVE	*
LOCATIVE	profit

And [BLANK] he [BLANK] nearly always bought and sold [BLANK].

[EMPTY] [EMPTY]

like a good guy.

[EMPTY] through a 17-month period

[EMPTY]

moreover because he had always known him,

[EMPTY]

[EMPTY] [EMPTY]

in the United States

[EMPTY] [EMPTY]

at a profit

Interactive workflow: Human-generator interaction



Scalable

Possibly have more vocabulary



Can pick most/least impact changing places

Can be more creative after seeing the machine generation

saw **is** → **isn't**, can think of **is** → **could have been**



Humans assist

Machine generates/selects

Human only verify



Humans lead

Human provide seed perturbations

Human BLANKs change points

Machine follows up to broaden vocab.



Take turns

Machine starts

Human corrects imperfect ones

Human supplies missing cases

Extend to labeling: Get structures at collection

Improve data collection, instead of apply structures “after the fact”

Infer the structure from labeled counterfactual pairs.

Labels

+

Has-negation

negative



positive

I didn't → did like the movie.

I don't → do think it is a good book, and → even though it is expensive.

You would never → always care about this movie if you hadn't seen the previous one.

Building blocks

Quantitative grouping

Counterfactual perturbation

Model development stages (Paleyes et al. 2020)

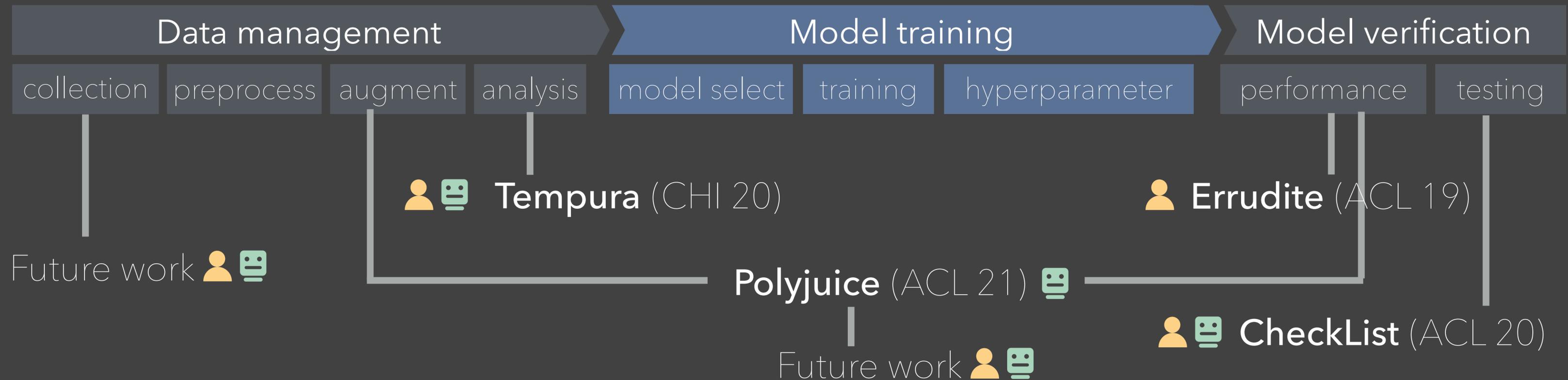


Building blocks

Quantitative grouping

Counterfactual perturbation

Model development stages (Paleyes et al. 2020)



Thanks!

And I'm happy to take questions :)

Backup slides

ACL 2019

Errudite: Scalable, Reproducible, and Testable Error Analysis

Tongshuang (Sherry) Wu @tongshuangwu
University of Washington

Marco Tulio Ribeiro
Microsoft Research

Jeffrey Heer @jeffrey_heer
Daniel S. Weld @dsweld
University of Washington



Deliveries: Precise + Reproducible + Re-applicable

Attribute

```
ENT (g)
```

Groups

```
all_instance  
is_entity  
has_distractor  
correct_type  
is_distracted
```

Rewrite rule

```
rewrite (  
  c,  
  string (p (m) ) → "#")
```

Deliveries: **Precise + Reproducible + Re-applicable**

Attribute + Groups + Rewrite rule



applied to...

BiDAF is ...

not particularly **bad** at distractors.

Seemingly distractor errors can be due to **other factors**.

Deliveries: **Precise + Reproducible + Re-applicable**

Attribute + Groups + Rewrite rule



Re-applied to...

Other datasets & **Other** models...

? at handling distractor.

User study: What is imprecise answer boundaries?

“The model is making predictions with missing or additional words...?”

D1 No exact match, but high overlap

```
exact_match(p(m)) == 0  
and f1(p(m)) > 0.7
```

D2 Off by at most 2 tokens both on the left and right

```
exact_match(p(m)) == 0  
and abs(answer_offset(p(m), "left")) <= 2  
and abs(answer_offset(p(m), "right")) <= 2
```

User study: What is imprecise answer boundaries?

“The model is making predictions with missing or additional words...?”

D1 No exact match, but high overlap

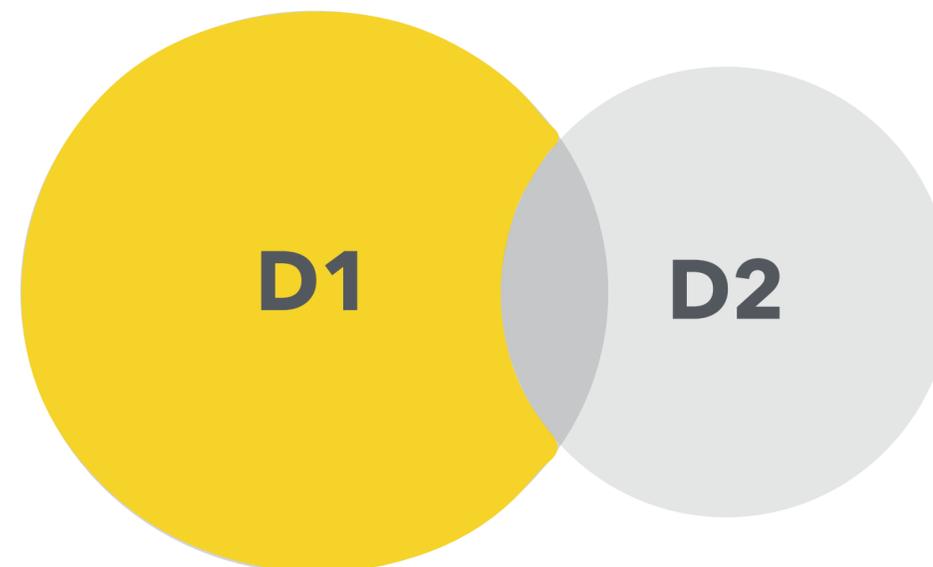
```
exact_match(p(m)) == 0  
and f1(p(m)) > 0.7
```

D2 Off by at most 2 tokens both on the left and right

```
exact_match(p(m)) == 0  
and abs(answer_offset(p(m), "left")) <= 2  
and abs(answer_offset(p(m), "right")) <= 2
```

User study: **What is imprecise answer boundaries?**

- D1** No exact match, but **high overlap** **D2** Off by at most **2** tokens both on the **left** and **right**



↑ **groundtruth**
...the **polynomial time hierarchy** collapses.
...believed that the polynomial hierarchy does..
prediction

ACL 2020 Best Paper Award

Beyond Accuracy: Behavioral Testing of NLP Models with **Checklist**

Marco Tulio Ribeiro
Microsoft Research

Tongshuang (Sherry) Wu @tongshuangwu
University of Washington

Jeffrey Heer @jeffrey_heer
Daniel S. Weld @dsweld
University of Washington



Discussion: translate failure rate to **success** / **failure**?

"passed" if failures are on rare tokens

Capabilities	MFT
Vocab/POS	Pos/Nec: 15% Neutral: 7.6%
Named entities	
Nagation	Easy: 49.2%
...	

Affected by the test cases selected

Abs. value is not as interesting as "high enough"

Can be subjective & case-to-case

The failure is ~50%!

Discussion: Cautious on what to claim!

Failing a test \neq failing what the test name indicates.

Linguistic capabilities are more intertwined. Should try to further isolate compounds through INV tests. And should fix the pattern anyways!

Passing a test \neq model working.

Test cases are not comprehensive; Only give you more confident that the basic works.

Discussion: who are the users?

Model developers, Experts on model evaluation & task
Common and intuitive tests that are crucial for deployment

Researchers, Experts on model evaluation
Investigate into sophisticated tests (that may worth a paper)

Customers, Experts on the specific data/application
Tests specific to the dataset (e.g., NER tests on medical terms)

Ultimate goal: Have a shared test suite for each NLP task

user study: people test the same model/capability with different test cases!

Discussion: why do we care?

It's true ...

Some of the failures are by design and are not surprising.

e.g. MFT tests are usually out of distribution; SQuAD dataset do not have very short paragraphs.

But!

It is annotation artifact.

Dataset collection does not reflect the real world what we care about.

The training data will never be comprehensive.

Language is high dimension and selection bias is unavoidable.

The training data will keep getting more biased.

Concept drift caused by the deployed model interacting with the world.

Discussion: why do we care?

It's true ...

The testing does not necessarily point to the source of bug / a fix.
NER-INV failure is due to contextual embedding, not my model/data.

But!

We should first find the bug, and then try to isolate the source.

Detecting bugs is paramount for evaluation, and a prerequisite for further exploration of what caused them.

It's true ...

Testing sophisticated capabilities can be hard.

Test cases for sarcasm require more effort than simple negation.

But!

We can start with the simple ones as demo-ed!

Test models with the basics, & write tests close to models' capability.
Make sure your model pass level 1 MFTs before you reach level 3!

Data manage.

augment

Model verification

perform - explanation

Submitted to ACL 2021

Polyjuice: Generating Counterfactuals for Explaining, Evaluating, and Improving models

Tongshuang (Sherry) Wu @tongshuangwu
University of Washington

Marco Tulio Ribeiro
Microsoft Research

Jeffrey Heer @jeffrey_heer
Daniel S. Weld @dsweld
University of Washington



Have we selected criticizing points?

CLAIM

The "surprises" compensate SHAP.

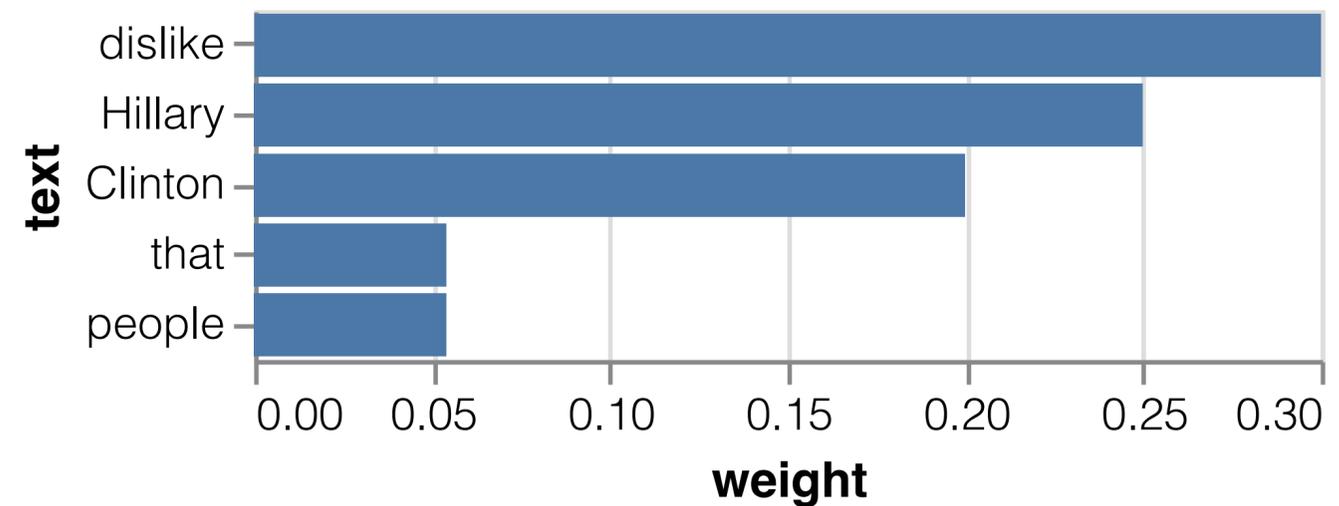
HOW? for each round... (20 in total)

GIVE Provide original examples, SHAP

Q1: Why do so many people hate Hilary Clinton?

Q2: What are the reasons that people dislike Hillary Clinton?

Predict: Duplicate (99.3% confident)



Have we selected criticizing points?

CLAIM

The “surprises” compensate SHAP.

HOW? for each round... (20 in total)

GIVE Provide original examples, SHAP

GIVE Allow asking models questions

MEASURE:

unsuccessful simulation

ASK QUESTIONS!

You can edit **Q2** and get model's prediction on **variations** of this example. Use it wisely, and try to come up with edits that will best help you understand **the model behavior around** the instance!

You can do **10** more model query! If you think you have learned enough, please **start the task**.

Old Q1 Why do so many people hate Hilary Clinton ?

New Q2

Get prediction

Have we selected criticizing points?

CLAIM

The “surprises” compensate SHAP.

HOW? for each round... (20 in total)

GIVE Provide original examples, SHAP

GIVE Allow asking models questions

GET Simulate model’s behavior

on perturb. (criticism, random, human)

MEASURE:

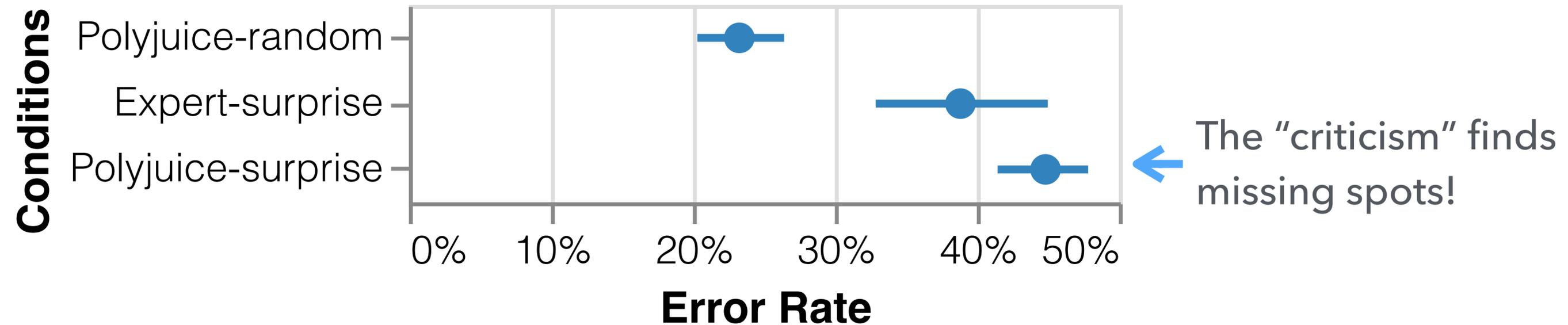
Error rate: # unsuccessful simulation

The more they miss (V.S. random/human),
the more information we add!

Old Q1	Why do so many people hate Hilary Clinton ?
New Q2	What are the reasons that people dislike like Hillary Clinton ?
Model will predict	<input type="radio"/> Non-duplicate <input type="radio"/> Duplicate
Old Q1	Why do so many people hate Hilary Clinton ?
New Q2	What are the reasons that people dislike Hillary Clinton Trump ?
Model will predict	<input type="radio"/> Non-duplicate <input type="radio"/> Duplicate
Old Q1	Why do so many people hate Hilary Clinton ?
New Q2	What are the reasons that not so many people dislike Hillary Clinton ?
Model will predict	<input type="radio"/> Non-duplicate <input type="radio"/> Duplicate

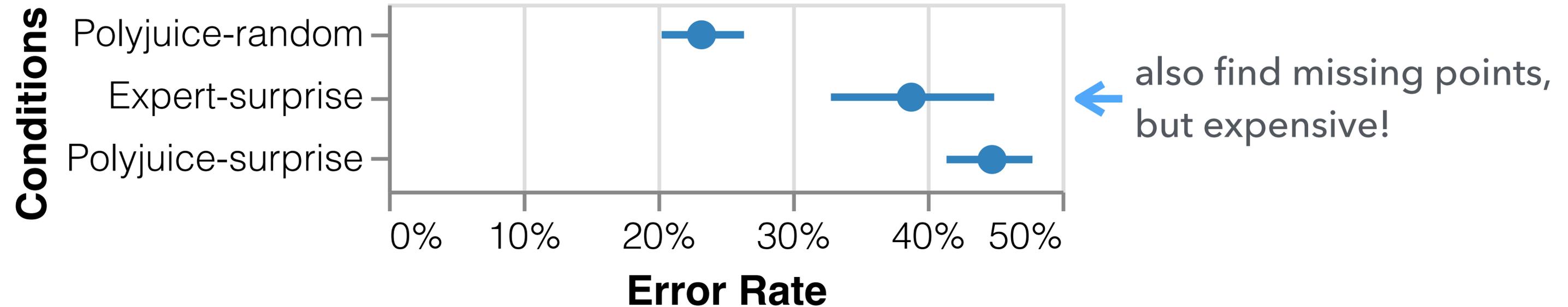
Results

13 people, 20 rounds, each round contains 2 questions × 3 selection conditions



Results

13 people, 20 rounds, each round contains 2 questions \times 3 selection conditions

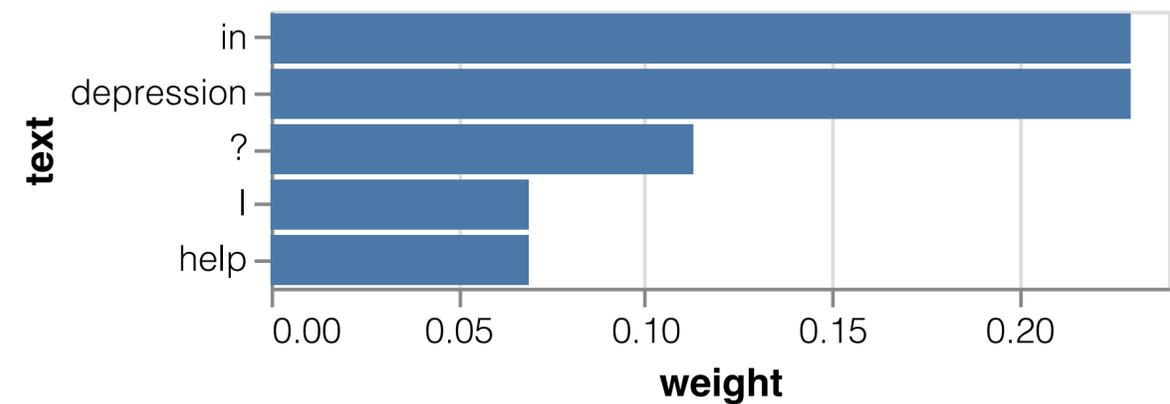


Why did people miss the surprises?

Q1: How can I help a friend who is experiencing serious depression?

Q2: How do I help a friend who is in depression?

Predict: Duplicate (99.3% confident)



Surprise

Q1: How can I help a friend who is experiencing serious depression?

Q2: How do I help a friend → woman who is in depression?

Predict: Duplicate (99.3% confident)

Users usually query “important” words...

Q2: How do I help a friend who is in

depression → disease

depression → frustration

depression → happiness

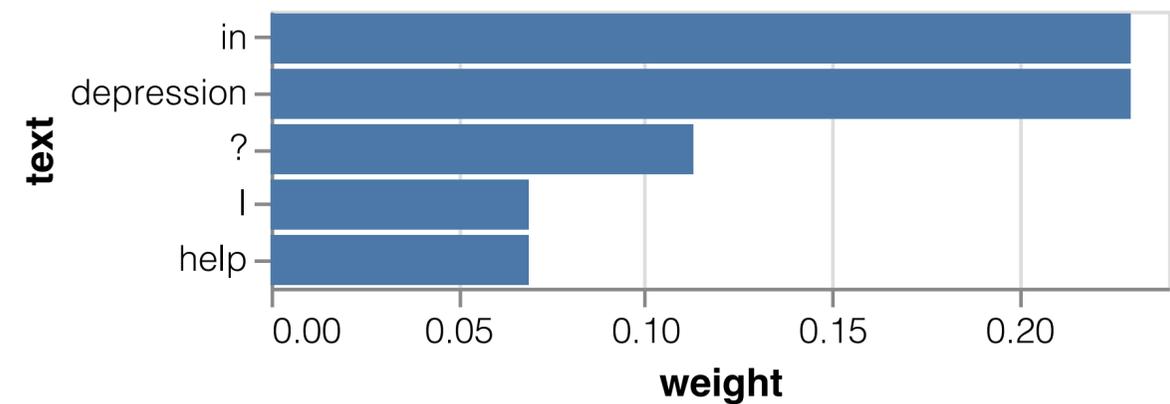
85% of user inspections are on top features!

Why did people miss the surprises?

Q1: How can I help a friend who is experiencing serious depression?

Q2: How do I help a friend who is in depression?

Predict: Duplicate (99.3% confident)



Surprise

Q1: How can I help a friend who is experiencing serious depression?

Q2: How do I help → find a friend who is in depression?

Predict: Duplicate (99.3% confident)

Users overfit to their inspection...

Q1: How can I help a friend who is experiencing serious depression?

Q2: How do I help → play with a friend who is in depression?

Predict: Non-duplicate (99.3% confident)