Natural Language Inference:



Natural Logic Meets Machine Learning, June 2021

Challenges and Opportunities

Vered Shwartz

Microsoft

Microsoft DeBERTa surpasses human performance on ...

... of NLU tasks, including question answering, natural language inference, coreference resolution, word sense disambiguation, and others. Jan 6, 2021

Analytics India Magazine

Facebook Inch Closer To General-Purpose Intelligence ...

The Facebook AI researchers started with the question: "Could a Transformer trained for natural language inference on textual input also ... Mar 3, 2021







Microsoft

Microsoft DeBERTa surpasses human performance on ...

... of NLU tasks, including question answering, natural language inference, coreference resolution, word sense disambiguation, and others. Jan 6, 2021

Analytics India Magazine

Facebook Inch Closer To General-Purpose Intelligence ...

The Facebook AI researchers started with the question: "Could a Transformer trained for natural language inference on textual input also ... Mar 3, 2021









Jul '19Aug '19

Jan '20



Microsoft

Microsoft DeBERTa surpasses human performance on ...

... of NLU tasks, including question answering, natural language inference, coreference resolution, word sense disambiguation, and others. Jan 6, 2021

Analytics India Magazine

Facebook Inch Closer To General-Purpose Intelligence ...

The Facebook AI researchers started with the question: "Could a Transformer trained for natural language inference on textual input also ... Mar 3, 2021











Jul '19Aug '19

Jan '20



I baked a chocolate cake but accidentally used regular instead of self-rising flour.

I baked a chocolate cake but accidentally used regular instead of self-rising flour.

I've made a mistake.



I baked a chocolate cake but accidentally used regular instead of self-rising flour.

I've made a mistake.

The cake turned out flat. The cake was difficult to cut. I baked a bad chocolate cake.



I baked a chocolate cake but accidentally used regular instead of self-rising flour.

I've made a mistake.

The cake turned out flat.

The cake was difficult to cut.

I baked a bad chocolate cake.

The cake was soft and tasty.



I baked a chocolate cake but accidentally used regular instead of self-rising flour.

I've made a mistake.

The cake turned out flat.

The cake was difficult to cut.

I baked a bad chocolate cake.

The cake was soft and tasty.

I've baked a chocolate rock.









Incorporating symbolic knowledge into neural NLI models



NLI is too easy? What's next?







Incorporating symbolic knowledge into neural NLI models

NLI is too easy? What's next?

Dataset



Labels

parrot, duck, duck, parrot...

Dataset



Labels

parrot, duck, duck, parrot...

Dataset



Labels

parrot, duck, duck, parrot...

Training set



Classifier

Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...

Training set



Green body \rightarrow parrot, Gray body \rightarrow duck

← Classifier ←

Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...





Dataset



Labels

parrot, duck, duck, parrot...

Training set





1. Spurious correlations

Dataset



Labels

parrot, duck, duck, parrot...

Training set





1. Spurious correlations 2. Train & test from the same distribution

Annotation artifacts, hypothesis-only baseline (Gururangan, Swayamdipta, et al., 2018; Poliak et al., 2018; Tsuchiya, 2018)

Annotation artifacts, hypothesis-only baseline (Gururangan, Swayamdipta, et al., 2018; Poliak et al., 2018; Tsuchiya, 2018)

- **p:** I only had a soup but it was very filling.
- h: I didn't eat a salad.
- 2: contradiction (98.2%)

Annotation artifacts, hypothesis-only baseline (Gururangan, Swayamdipta, et al., 2018; Poliak et al., 2018; Tsuchiya, 2018)

- **p:** I only had a soup but it was very filling.
- h: I didn't eat a salad.
- 2: contradiction (98.2%)

NLI models rely on syntactic heuristics (McCoy et al., 2019)





P: Charlie will visit his mother in London on Wednesday evening.H: Charlie will visit his mother in London on Thursday evening.

ning.

2018 (ELMo-based Decomposable Attention)

P: Charlie will visit his mother in London on Wednesday evening.H: Charlie will visit his mother in London on Thursday evening.

ning.

2018 (ELMo-based Decomposable Attention)

P: Charlie will visit his mother in London on Wednesday evening.H: Charlie will visit his mother in London on Thursday evening.

Prediction: Entailment (94.1%) 🔀

https://demo.allennlp.org/textual-entailment

ning.

2018 (ELMo-based Decomposable Attention)

P: Charlie will visit his mother in London on **Wednesday** evening. H: Charlie will visit his mother in London on Thursday evening.

Prediction: Entailment (94.1%)

Breaking NLI Systems with Sentences that Require Simple Lexical Inferences. Max Glockner, Vered Shwartz, and Yoav Goldberg. ACL 2018





2018 (ELMo-based Decomposable Attention)

P: Charlie will visit his mother in London on **Wednesday** evening. H: Charlie will visit his mother in London on Thursday evening.

Prediction: Entailment (94.1%) 🗙

• Errors: similar but mutually-exclusive words

- •Accuracy increases with frequency in training set \rightarrow Limited generalization ability!
- Similar findings in "NLI stress tests" (Naik et al., 2018)

Breaking NLI Systems with Sentences that Require Simple Lexical Inferences. Max Glockner, Vered Shwartz, and Yoav Goldberg. ACL 2018




#2 Representations

2021 (RoBERTa)

https://demo.allennlp.org/textual-entailment

#2 Representations

2021 (RoBERTa)

P: Charlie will visit his mother in London on Wednesday evening.H: Charlie will visit his mother in London on Thursday evening.

Prediction: Contradiction (73.2%) 🗸

https://demo.allennlp.org/textual-entailment

ning.

#2 Representations

2021 (RoBERTa)

P: Charlie will visit his mother in London on Wednesday evening.
H: Charlie will visit his mother in London on Thursday evening.

Prediction: Contradiction (73.2%) 🗸

P: Charlie said on **Wednesday** that he is busy on *Thursday* so he will visit his mother next week.

H: Charlie said on **Thursday** that he is busy on *Wednesday* so he will visit his mother next week.

Prediction: Entailment (92%) 样

https://demo.allennlp.org/textual-entailment

The definition of the textual entailment recognition task, like that of any other text understanding task, refers to human understanding of language. Such definition necessarily assumes common background knowledge, on which the (human) entailment judgment relies. $[\cdots]$ this knowledge should cover both extra-linguistic world knowledge $[\cdots]$ as well as knowledge of the language itself.

Recognizing Textual Entailment. Ido Dagan, Dan Roth, Mark Sammons, and Fabio Zanzotto. Morgan & Claypool Publishers, 2013.



Capture facts not explicitly mentioned in the corpus (Petroni et al. 2019; Feldman et al. 2019, this talk)



Capture facts not explicitly mentioned in the corpus (Petroni et al. 2019; Feldman et al. 2019, this talk)

Associate concepts with their properties \checkmark (Weir et al. 2020)



	Contoxt	Huma	Human		A-L
	Context	Response	PF	Response	$p_{\rm LM}$
3	Everyone	fur	27	teeth	.36
E.Al	knows that a	claws	15	claws	.18
AND A LOCAL DESIGNATION OF THE OWNER OWNER OF THE OWNER	bear has	teeth	11	eyes	.05
		cubs	7	ears	.03
		paws	7	horns	.02



Capture facts not explicitly mentioned in the corpus (Petroni et al. 2019; Feldman et al. 2019, this talk)

Associate concepts with their properties (Weir et al. 2020)

	Context	Human	n DE	ROBERT	A-L
		Response	PF	Response	$p_{\rm LM}$
A State To A	Everyone	fur	27	teeth	.36
A DE AND A CONTRACT	knows that a	claws	15	claws	.18
	bear has	teeth	11	eyes	.05
		cubs	7	ears	.03
		paws	7	horns	.02





Capture facts not explicitly mentioned in the corpus (Petroni et al. 2019; Feldman et al. 2019, this talk)

Associate concepts with their properties (Weir et al. 2020)

	Context	Huma	n	ROBERT	A-L
		Response	PF	Response	$p_{\rm LM}$
A HELLON	Everyone	fur	27	teeth	.36
A AND AND	knows that a bear has	claws teeth	15 11	claws	.18
	0000 11005 <u> </u>	cubs	7	ears	.03
		paws	/	norns	.02

#3 World Knowledge

X Not sensitive to negation (Kassner et al. 2020; Ettinger, 2020)

Knowledge in Language Models



Capture facts not explicitly mentioned in the corpus (Petroni et al. 2019; Feldman et al. 2019, this talk)

Associate concepts with their properties (Weir et al. 2020)

Context	Humar	1	ROBERT	A-L
CONTEXT	Response	PF	Response	$p_{\rm LM}$
Everyone knows that a bear has	fur claws teeth cubs	27 15 11 7 7	teeth claws eyes ears borns	.36 .18 .05 .03

X Not sensitive to negation (Kassner et al. 2020; Ettinger, 2020)

X Predict similar but mutually-exclusive facts (**Jiang et al., 2020**)



Knowledge in Language Models



Associate concepts with their properties (Weir et al. 2020)

	Context	Huma	n	ROBERT	A-L
		Response	PF	Response	$p_{\rm LM}$
	Everyone	fur	27	teeth	.36
A State Star	knows that a	claws	15	claws	.18
	bear has	teeth	11	eyes	.05
		paws	7	horns	.03

X Not sensitive to negation (Kassner et al. 2020; Ettinger, 2020)

Predict similar but mutually-exclusive facts (Jiang et al., 2020)



Lack perceptual and physical knowledge (Forbes et al. 2019, Weir et al., 2020, Bisk et al. 2020)

Knowledge in Language Models



Associate concepts with their properties (Weir et al. 2020)

	Context	Huma	n	ROBERT	A-L
	CONTEXT	Response	PF	Response	$p_{\rm LM}$
	Everyone	fur	27	teeth	.36
	knows that a	claws	15	claws	.18
	bear has	teeth	11	eyes	.05
		cubs	7	ears	.03
		paws	7	horns	.02

X Not sensitive to negation (Kassner et al. 2020; Ettinger, 2020)

Predict similar but mutually-exclusive facts (Jiang et al., 2020)



Lack perceptual and physical knowledge (Forbes et al. 2019, Weir et al., 2020, Bisk et al. 2020)

X Don't differentiate constant vs. contingent facts



Zebras are black and white.



My shirt is blue / red.

Knowledge in Language Models



Associate concepts with their properties (Weir et al. 2020)

Context	Huma Response	n PF	ROBERT Response	A-L <i>p</i> lm
Everyone knows that a bear has	fur claws teeth cubs paws	27 15 11 7 7	teeth claws eyes ears horns	.36 .18 .05 .03 .02



Do Neural Language Models Overcome Reporting Bias? **Vered Shwartz** and Yejin Choi. COLING 2020.

X Not sensitive to negation (Kassner et al. 2020; Ettinger, 2020)

× Predict similar but mutually-exclusive facts (**Jiang et al., 2020**)



X Lack perceptual and physical knowledge (Forbes et al. 2019, Weir et al., 2020, Bisk et al. 2020)

X Don't differentiate constant vs. contingent facts



Zebras are black and white.



My shirt is blue / red.

X Reporting bias

Knowledge in Language Models



Associate concepts with their properties (Weir et al. 2020)

Context	Huma Response	n PF	ROBERT Response	A-L <i>p</i> lm
Everyone knows that a bear has	fur claws teeth cubs paws	27 15 11 7 7	teeth claws eyes ears horns	.36 .18 .05 .03 .02



Do Neural Language Models Overcome Reporting Bias? **Vered Shwartz** and Yejin Choi. COLING 2020.

X Not sensitive to negation (Kassner et al. 2020; Ettinger, 2020)

Fredict similar but mutually-exclusive facts (**Jiang et al., 2020**)



X Lack perceptual and physical knowledge (Forbes et al. 2019, Weir et al., 2020, Bisk et al. 2020)

X Don't differentiate constant vs. contingent facts





My shirt is blue / red.

The man turned on the faucet. As a result, **X** Reporting bias

Knowledge in Language Models



Associate concepts with their properties (Weir et al. 2020)

Context	Huma Response	n PF	ROBERT Response	A-L <i>p</i> lm
Everyone knows that a bear has	fur claws teeth cubs paws	27 15 11 7 7	teeth claws eyes ears horns	.36 .18 .05 .03 .02



Do Neural Language Models Overcome Reporting Bias? **Vered Shwartz** and Yejin Choi. COLING 2020.

X Not sensitive to negation (Kassner et al. 2020; Ettinger, 2020)

× Predict similar but mutually-exclusive facts (**Jiang et al., 2020**)



X Lack perceptual and physical knowledge (Forbes et al. 2019, Weir et al., 2020, Bisk et al. 2020)

X Don't differentiate constant vs. contingent facts





My shirt is blue / red.

X Reporting bias

The man turned on the faucet. As a result,

the man's blood was sprayed everywhere.









NLI is too easy? What's next?

Learning Semantic Phenomena

The office is located in Baytown [SEP] The office is located in Texas





Learning Semantic Phenomena

The office is located in Baytown [SEP] The office is located in Texas







Fact: LocatedIn(Baytown, Texas)

Learning Semantic Phenomena

The office is located in Baytown [SEP] The office is located in Texas







Fact: LocatedIn(Baytown, Texas)

Rule: Entailment between a city and its state in upward monotone sentences



Does the model know the location of cities?

The office is located in Baytown [SEP] The office is located in Texas



Does the model know the location of cities?

The office is located in Baytown [SEP] The office is located in Texas



Does the model know the location of cities?

The office is located in Baytown [SEP] The office is located in Texas



Inoculation by Fine-Tuning: A Method for Analyzing Challenge Datasets. Nelson F. Liu, Roy Schwartz, and Noah A. Smith. NAACL 2019.



Does the model know the location of cities?

The office is located in Baytown [SEP] _____ The office is located in Texas

The office is located in London [SEP]

She lives in Dallas [SEP]

H He moved to Chicago [SEP]

He lives in Illinois

label: entailment

Inoculation by Fine-Tuning: A Method for Analyzing Challenge Datasets. Nelson F. Liu, Roy Schwartz, and Noah A. Smith. NAACL 2019.





Does the model know the location of cities?

The office is located in Baytown [SEP] _____ The office is located in Texas

The office is located in London [SEP]

She lives in Dallas [SEP]

H He moved to Chicago [SEP]

He lives in Illinois

label: entailment

Inoculation by Fine-Tuning: A Method for Analyzing Challenge Datasets. Nelson F. Liu, Roy Schwartz, and Noah A. Smith. NAACL 2019.



Does the model know the location of cities?

The office is located in Baytown [SEP] _____ The office is located in Texas

The office is located in London [SEP]

She lives in Dallas [SEP]

H He moved to Chicago [SEP]

He lives in Illinois

label: entailment

Inoculation by Fine-Tuning: A Method for Analyzing Challenge Datasets. Nelson F. Liu, Roy Schwartz, and Noah A. Smith. NAACL 2019.



Possible Outcomes



Does the model know the location of cities?

The office is located in Baytown [SEP] _____ The office is located in Texas

The office is located in London [SEP]

She lives in Dallas [SEP]

H He moved to Chicago [SEP]

He lives in Illinois

label: entailment

Inoculation by Fine-Tuning: A Method for Analyzing Challenge Datasets. Nelson F. Liu, Roy Schwartz, and Noah A. Smith. NAACL 2019.









Blindspot in the original NLI dataset

Does the model know the location of cities?

The office is located in Baytown [SEP] _____ The office is located in Texas

The office is located in London [SEP]

She lives in Dallas [SEP]

H He moved to Chicago [SEP]

He lives in Illinois

label: entailment

Inoculation by Fine-Tuning: A Method for Analyzing Challenge Datasets. Nelson F. Liu, Roy Schwartz, and Noah A. Smith. NAACL 2019.



Possible Outcomes





Blindspot in the original NLI dataset

Inherent model limitation

X Failure





Numerical Reasoning

P: I see 260 coins in the bucket.

H: I see more than 232 coins in the bucket.

Label: Entailment

Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. CoNLL 2019.



Numerical Reasoning

P: I see 260 coins in the bucket.

H: I see more than 232 coins in the bucket.

Label: Entailment

Axis: number range

Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. CoNLL 2019.

1. Train/dev/test split across the dimension in focus





Numerical Reasoning

P: I see 260 coins in the bucket.

H: I see more than 232 coins in the bucket.

Label: Entailment

Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. CoNLL 2019.







Numerical Reasoning

P: I see 260 coins in the bucket.

H: I see more than 232 coins in the bucket.

Label: Entailment

Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. CoNLL 2019.





Numerical Reasoning

P: I see 260 coins in the bucket.

H: I see more than 232 coins in the bucket.

Label: Entailment

Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. CoNLL 2019.





Numerical Reasoning

P: I see 260 coins in the bucket.

H: I see more than 232 coins in the bucket.

Label: Entailment

Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. CoNLL 2019.





Numerical Reasoning

see 260 coins in the bucket.

H: I see more than 232 coins in the bucket.

Label: Entailment

Dative Alternation

P: I baked my mom a cake.

H: I baked a cake for my mom.

Label: Entailment



1. Train/dev/test split across the dimension in focus 2. Fine-tune on one set and test on another

Axis: lexical variability

X Axis: syntactic complexity

Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. CoNLL 2019.


Analyzing Generalization via Controlled Variance

Numerical Reasoning

see 260 coins in the bucket.

H: I see more than 232 coins in the bucket.

Label: Entailment

Dative Alternation

P: I baked my mom a cake.

H: I baked a cake for my mom.

Label: Entailment





1. Train/dev/test split across the dimension in focus 2. Fine-tune on one set and test on another

Axis: lexical variability

X Axis: syntactic complexity

May decrease performance on the main task (Richardson et al., 2020)

Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. CoNLL 2019.



The office is located in Baytown [SEP] The office is located in Texas

label: entailment



The office is located in Baytown [SEP] The office is located in Texas

label: entailment



Model learns:

Fact: LocatedIn(Baytown, Texas)

Rule: Entailment between a city and its state in upward monotone sentences

The office is located in Baytown [SEP] The office is located in Texas

label: entailment





Fact: LocatedIn(Baytown, Texas)

Rule: Entailment between a city and its state in upward monotone sentences

May be learned given enough data.

The office is located in Baytown [SEP] The office is located in Texas

label: entailment



Impossible (and inefficient) to teach an NLI model every fact it might need.

Model learns:

Fact: LocatedIn(Baytown, Texas)

Rule: Entailment between a city and its state in upward monotone sentences

May be learned given enough data.



The Free Encyclopedia

Model learns:

Rule: Entailment between a city and its state in upward monotone sentences



Model learns:

Rule: Entailment between a city and its state in upward monotone sentences

How to incorporate relational knowledge?



H_i' proj H_i' proj H_i proj Baytown is a city in the U.S. state of Texas Baytown is a city in the U.S. state of Texas

Knowledge Enhanced Contextual Word Representations. Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. EMNLP 2019.

Baytown, Texas Baytown culture Operation Baytown

State of Texas Republic of Texas Texas, Alabama





H_i' Proj H_i' Proj H_i ' Proj ' Proj H_i ' Proj ' Proj H_i ' Proj ' Proj ' Proj ' Pr

erformance improvement on relation extraction, entity typing, WSD

Knowledge Enhanced Contextual Word Representations. Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. EMNLP 2019.

Baytown, Texas Baytown culture Operation Baytown

State of Texas Republic of Texas Texas, Alabama





H_i' 7 H_i'

erformance improvement on relation extraction, entity typing, WSD

Knowledge Enhanced Contextual Word Representations. Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. EMNLP 2019.

Baytown, Texas Baytown culture Operation Baytown

State of Texas Republic of Texas Texas, Alabama 😢 can't learn new facts after training

The office is located in Nanhui [SEP] The office is located in China





H_i′proj State of Texas H_ipro Baytown is a city in the U.S. state of Texas |3|

erformance improvement on relation extraction, entity typing, WSD

Knowledge Enhanced Contextual Word Representations. Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. EMNLP 2019.

Baytown, Texas **Baytown culture Operation Baytown**

Republic of Texas Texas, Alabama

😢 can't learn new facts after training

The office is located in Nanhui [SEP] The office is located in China

expensive to re-train



The office is located in Baytown . [SEP] The office is located in Texas .





label: entailment

Fact: LocatedIn(Baytown, Texas)



Teach the Rules, Provide the Facts: Targeted Relational-knowledge Enhancement for Textual Inference. Ohad Rozen, Shmuel Amar, Vered Shwartz and Ido Dagan. *SEM 2021

		P	re	m	ÍS	e:		Hypothesis:											
	ο	0	0	0	0		0		0	0	0	0	0	Г	0				
	0	ο	0	Ο	0	00	ο		ο	0	ο	ο	0) C	0				
	0	ο	0	ο	0	ät	ο		ο	0	ο	ο	0	ate	0				
	Ο	ο	0	0	0	ed	ο		ο	0	ο	ο	Ο	Ö	0				
IBERI	Ο	ο	0	Ο	0	n	ο		ο	0	ο	ο	Ο		0				
	ο	ο	0	Ο	0	ta	ο		ο	0	ο	ο	0	le	0				
	ο	ο	0	Ο	0	<u> </u>	ο		ο	0	ο	ο	0	D	0				
				-															



The office is located in Baytown . [SEP] The office is located in Texas .





label: entailment

Fact: LocatedIn(Baytown, Texas)



Teach the Rules, Provide the Facts: Targeted Relational-knowledge Enhancement for Textual Inference. Ohad Rozen, Shmuel Amar, Vered Shwartz and Ido Dagan. *SEM 2021

		P	re	m	is	e:		Hypothesis:										
Ł	0	0	0	0	0		0]	ο	0	0	0	0		0			
	0	0	0	0	0	00	0		ο	0	0	ο	0	000	0			
	ο	0	0	0	0	at	0		ο	ο	0	ο	0	ate	0			
	ο	0	0	0	0	ed	0		0	ο	0	ο	ο	<u>ă</u>	0			
SERI	ο	0	0	0	0	n	0		ο	ο	0	ο	ο		0			
	ο	0	0	0	0	ta	0		ο	ο	0	ο	ο	le	0			
	ο	0	0	0	0	<u> </u>	0		0	ο	0	ο	ο	D	0			

Model learns:

Rule: Entailment between a city and its state in upward monotone sentences

Relation Embeddings





< 700 training examples

Teach the Rules, Provide the Facts: Targeted Relational-knowledge Enhancement for Textual Inference. Ohad Rozen, Shmuel Amar, Vered Shwartz and Ido Dagan. *SEM 2021





< 700 training examples

Teach the Rules, Provide the Facts: Targeted Relational-knowledge Enhancement for Textual Inference. Ohad Rozen, Shmuel Amar, Vered Shwartz and Ido Dagan. *SEM 2021







< 700 training examples

erformance improvement on NLI challenge sets

can retrieve and use facts about unseen entities The office is located in Nanhui [SEP] The office is located in China





< 700 training examples

Teach the Rules, Provide the Facts: Targeted Relational-knowledge Enhancement for Textual Inference. Ohad Rozen, Shmuel Amar, Vered Shwartz and Ido Dagan. *SEM 2021

erformance improvement on NLI challenge sets

Can retrieve and use facts about unseen entities The office is located in Nanhui [SEP] The office is located in China

computationally efficient





< 700 training examples

erformance improvement on NLI challenge sets

Can retrieve and use facts about unseen entities The office is located in Nanhui [SEP] The office is located in China

Computationally efficient

task-specific knowledge incorporation











NLI is too easy? What's next?

Real-World NLI #1 Partial Entailment

 S_1 : Amazon to acquire Whole Foods Market for \$13.7 Billion.

 S_2 : Amazon is buying Whole Foods Market for almost \$14 Billion in cash.

Real-World NLI #1 Partial Entailment

 S_1 : Amazon to acquire Whole Foods Market for \$13.7 Billion.

 S_2 : Amazon is buying Whole Foods Market for almost \$14 Billion in cash.

Real-World NLI #1 Partial Entailment

 S_1 : Amazon to acquire Whole Foods Market for \$13.7 Billion.

 S_2 : Amazon is buying Whole Foods Market for almost \$14 Billion in cash.

Subjectivity and inherent disagreements (Pavlick and Kwiatkowski, 2019)

Real-World NLI #2 Defeasible Natural Language Inference

- **P:** Tweety is a bird.
- **H:** Tweety flies.

Real-World NLI #2 Defeasible Natural Language Inference

- **P:** Tweety is a bird.
- **H:** Tweety flies.
- NLI model: Entailment

Real-World NLI #2 Defeasible Natural Language Inference

- **P:** Tweety is a bird.
- **H:** Tweety flies.
- NLI model: Entailment

Skeptical NLI model: given the information I currently have, I suppose so, but I can think of cases in which this is false.

Given premise P, a hypothesis H is **defeasible** if there exists an update U (consistent with P) such that a human would find H less likely to be true after learning U.

Given premise P, a hypothesis H is **defeasible** if there exists an update U (consistent with P) such that a human would find H less likely to be true after learning U.



Given premise P, a hypothesis H is **defeasible** if there exists an update U (consistent with P) such that a human would find H less likely to be true after learning U.



P: Tweety is a bird.

Given premise P, a hypothesis H is **defeasible** if there exists an update U (consistent with P) such that a human would find H less likely to be true after learning U.



P: Tweety is a bird.

H: Tweety flies.

Given premise P, a hypothesis H is **defeasible** if there exists an update U (consistent with P) such that a human would find H less likely to be true after learning U.



P: Tweety is a bird.

H: Tweety flies.

U: Tweety is a penguin.

Given premise P, a hypothesis H is **defeasible** if there exists an update U (consistent with P) such that a human would find H less likely to be true after learning U.



P: Tweety is a bird.

H: Tweety flies.

U: Tweety is a penguin.

Useful for **Real-time Summarization**: Facts change as the story unfolds.

An update U is called a **weakener** if, given a premise P and hypothesis H, a human would most likely find H *less likely to be true* after learning U; if they would find H *more likely to be true*, then we call U a **strengthener**.

Thinking Like a Skeptic: Defeasible Inference in Natural Language. Rachel Rudinger, **Vered Shwartz**, Jena Hwang, Chandra Bhagavatula, Maxwell Forbes, Ronan Le Bras, Noah Smith, and Yejin Choi. Findings of EMNLP 2020.

P: Tweety is a bird.

H: Tweety flies.

Weakener: Tweety is a penguin.



An update U is called a **weakener** if, given a premise P and hypothesis H, a human would most likely find H *less likely to be true* after learning U; if they would find H *more likely to be true*, then we call U a **strengthener**.

Thinking Like a Skeptic: Defeasible Inference in Natural Language. Rachel Rudinger, **Vered Shwartz**, Jena Hwang, Chandra Bhagavatula, Maxwell Forbes, Ronan Le Bras, Noah Smith, and Yejin Choi. Findings of EMNLP 2020.

P: Tweety is a bird.

H: Tweety flies.

Weakener: Tweety is a penguin.

Strengthener: Tweety is on a tree.





Discriminative Task

Determine whether an update weakens or strengthens the hypothesis.

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.





Discriminative Task

Determine whether an update weakens or strengthens the hypothesis.

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.

Generative Task

Generate a weakening or strengthening update for a given premise-hypothesis pair.

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.






Defeasible Inference in Natural Language

Discriminative Task

Determine whether an update weakens or strengthens the hypothesis.

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.

Generative Task

Generate a weakening or strengthening update for a given premise-hypothesis pair.

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.

Language models leave plenty of room for improvement on the generative task!







Learning to Rationalize for Nonmonotonic Reasoning with Distant Supervision. Faeze Brahman, Vered Shwartz, Rachel Rudinger, and Yejin Choi. AAAI 2021.





They are in a conference room.

A conference room is where people have meetings at work.

- They are in a library.

You must be quiet in the library, while work meetings involve talking.





They have a work meeting.





Learning to Rationalize for Nonmonotonic Reasoning with Distant Supervision. Faeze Brahman, Vered Shwartz, Rachel Rudinger, and Yejin Choi. AAAI 2021.

They are in a conference room.

→ A conference room is where people have meetings at work.

- They are in a library.

You must be quiet in the library, while work meetings involve talking.







Learning to Rationalize for Nonmonotonic Reasoning with Distant Supervision. Faeze Brahman, Vered Shwartz, Rachel Rudinger, and Yejin Choi. AAAI 2021.

🕈 They are in a conference room.

A conference room is where people have meetings at work.

- They are in a library.

You must be quiet in the library, while work meetings involve talking.

The definition of a library is...





Post hoc Rationalization Generates a rationale for a given decision (label).

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.



🕈 They are in a conference room.

→ A conference room is where people have meetings at work.

→ You must be quiet in the library, while work meetings involve talking.



Post hoc Rationalization Generates a rationale for a given decision (label).

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.



Trivially rephrasing the label! ("[+] implies that [H]")

+ They are in a conference room.

→ A conference room is where people have meetings at work.

→ You must be quiet in the library, while work meetings involve talking.



Post hoc Rationalization Generates a rationale for a given decision (label).

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.

Trivially rephrasing the label! ("[+] implies that [H]")

Joint Prediction & Rationalization

Predict the label (strengthener / weakener) and rationalize it.

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.



They are in a conference room.

→ A conference room is where people have meetings at work.

- They are in a library.



→ You must be quiet in the library, while work meetings involve talking.

A conference room is where people They are in a conference room. have meetings at work.

You must be quiet in the library, while work meetings involve talking.





Post hoc Rationalization Generates a rationale for a given decision (label).

A group of people sitting around a rectangular table having either pieces of paper or laptops in front of them.

They have a work meeting.

Trivially rephrasing the label! ("[+] implies that [H]")

Joint Prediction & Rationalization

Predict the label (strengthener / weakener) and rationalize it.



More realistic but very challenging task!

They are in a conference room.

→ A conference room is where people have meetings at work.

- They are in a library.



→ You must be quiet in the library, while work meetings involve talking.

A conference room is where people They are in a conference room. have meetings at work.

You must be quiet in the library, while work meetings involve talking.





Neural models achieve impressive gains on NLI

- But make stupid unhuman like errors
- "Human performance" is debatable





- Neural models achieve impressive gains on NLI • But make stupid unhuman like errors • "Human performance" is debatable

Symbolic knowledge is useful





- Neural models achieve impressive gains on NLI • But make stupid unhuman like errors • "Human performance" is debatable





- Neural models achieve impressive gains on NLI • But make stupid unhuman like errors • "Human performance" is debatable
- Symbolic knowledge is useful
 - Accurate
 - More efficient than "learning all the facts"



- Neural models achieve impressive gains on NLI • But make stupid unhuman like errors • "Human performance" is debatable
- Symbolic knowledge is useful
 - Accurate
 - More efficient than "learning all the facts"
 - Combination is not trivial



Recap

- Neural models achieve impressive gains on NLI • But make stupid unhuman like errors "Human performance" is debatable

Symbolic knowledge is useful

- Accurate
- More efficient than "learning all the facts"
- Combination is not trivial



Partial and defeasible inferences

- Neural models achieve impressive gains on NLI • But make stupid unhuman like errors "Human performance" is debatable

Symbolic knowledge is useful

- Accurate
- More efficient than "learning all the facts"
- Combination is not trivial
- Time to work on more real-world NLI tasks
 - Partial and defeasible inferences









References (1)

- Vered Shwartz and Yejin Choi. Do Neural Language Models Overcome Reporting Bias? COLING 2020. (1)
- (2)
- (3)Allyson Ettinger. 2020. What bert is not: Lessons from a new suite of psycholinguistic diagnostics for language models. TACL 2020
- (4)
- (5)
- Nora Kassner and Hinrich Schütze. 2020. Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly. ACL 2020. (6)
- (7)NAACL 2018
- (8)
- Nathaniel Weir, Adam Poliak, and Benjamin Van Durme. Probing Neural Language Models for Human Tacit Assumptions. CogSci 2020. (9)
- (10) Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. How Can We Know What Language Models Know? TACL 2020.
- Maxwell Forbes, Ari Holtzman, and Yejin Choi. Do Neural Language Representations Learn Physical Commonsense? CogSci 2019. (11)
- Oprea, Colin Raffel. Extracting Training Data from Large Language Models. arXiv 2020.
- Resources and Evaluation (LREC2018), 2018.
- (14) R. Thomas McCoy, Ellie Pavlick, and Tal Linzen. Right for the Wrong Reasons: Diagnosing Syntactic Heuristics in Natural Language Inference. ACL 2019.
- (15) Kyle Richardson, Hai Hu, Lawrence S Moss, Ashish Sabharwal. Probing Natural Language Inference Models through Semantic Fragments. AAAI 2020.

Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, Benjamin Van Durme. Hypothesis Only Baselines in Natural Language Inference. *SEM 2018.

Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander Miller. Language Models as Knowledge Bases? EMNLP 2019.

Jonathan Gordon and Benjamin Van Durme. 2013. Reporting bias and knowledge acquisition. Workshop on Automated knowledge base construction 2013.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, Noah A. Smith. Annotation Artifacts in Natural Language Inference Data.

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. PIQA: Reasoning about physical commonsense in natural language. AAAI 2020.

(12) Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, Alina

(13) Masatoshi Tsuchiya. Performance Impact Caused by Hidden Bias of Training Data for Recognizing Textual Entailment. In 11th International Conference on Language



References (2)

- (16) Ellie Pavlick and Tom Kwiatkowski. Inherent disagreements in human textual inferences. TACL 2019.
- (17) Raymond Reiter. A Logic for Default Reasoning. Artificial Intelligence, 1980.
- General-Purpose Language Understanding Systems. Neurips 2019.
- (19) Max Glockner, Vered Shwartz, Yoav Goldberg. Breaking NLI Systems with Sentences that Require Simple Lexical Inferences. ACL 2018.
- (21) Ido Dagan, Dan Roth, Mark Sammons, and Fabio Zanzotto. Recognizing Textual Entailment. Morgan & Claypool Publishers, 2013.
- (23) Zhengbao Jiang, Frank F. Xu, Jun Araki, Graham Neubig. How Can We Know What Language Models Know? TACL 2019.
- (24) Nelson F. Liu, Roy Schwartz, and Noah A. Smith. Inoculation by Fine-Tuning: A Method for Analyzing Challenge Datasets. NAACL 2019.
- CoNLL 2019.
- Representations. EMNLP 2019.
- *SEM 2021.
- Defeasible Inference in Natural Language. Findings of EMNLP 2020.

(18) Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, Samuel R. Bowman. SuperGLUE: A Stickier Benchmark for

(20) Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, Graham Neubig. Stress Test Evaluation for Natural Language Inference. COLING 2018.

(22) Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. PIQA: Reasoning about physical commonsense in natural language. AAAI 2020.

(25) Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan. Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets.

(26) Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. Knowledge Enhanced Contextual Word

(27) Ohad Rozen, Shmuel Amar, Vered Shwartz and Ido Dagan. Teach the Rules, Provide the Facts: Targeted Relational-knowledge Enhancement for Textual Inference.

(28) Rachel Rudinger, Vered Shwartz, Jena Hwang, Chandra Bhagavatula, Maxwell Forbes, Ronan Le Bras, Noah Smith, and Yejin Choi. Thinking Like a Skeptic:

(29) Faeze Brahman, Vered Shwartz, Rachel Rudinger, and Yejin Choi. Learning to Rationalize for Nonmonotonic Reasoning with Distant Supervision. AAAI 2021.