Fast Breaking and Slow Building of textual inference models

Vered Shwartz

December 2019



State-of-the-art AI solutions: (1) Google BERT, an AI model that understands language better than humans



AIN Dev Team Follow Jan 31 · 8 min read

🎔 🛅 🖬 🗆

WW VICE





It turns out that data-fueled algorithms are no better than humans—and ... Even AI researchers who work with machine learning models—like neural nets, which ...

Sep 18, 2019

alizila

Alibaba Al Beats Humans in Reading-Comprehension...

Alibaba Group's machine-learning technology is better at reading comprehension than humans, according to a well-known test built for the industry by Microsoft. Jul 9, 2019







What's in this talk?







Coreference

Breaking NLI Systems with Sentences that Require Simple Lexical Inferences

Max Glockner¹, Vered Shwartz² and Yoav Goldberg²

¹TU Darmstadt



TECHNISCHE UNIVERSITÄT DARMSTADT



ACL 2018

 A large scale dataset for NLI (Natural Language Inference; Recognizing Textual Entailment [Dagan et al., 2013])

- A large scale dataset for NLI (Natural Language Inference; Recognizing Textual Entailment [Dagan et al., 2013])
- Premises are image captions, hypotheses generated by crowdsourcing workers:

Premise

Street performer is doing his act for kids

- A large scale dataset for NLI (Natural Language Inference; Recognizing Textual Entailment [Dagan et al., 2013])
- Premises are image captions, hypotheses generated by crowdsourcing workers:

Premise

Street performer is doing his act for kids

Hypotheses

1. A person performing for children on the street

- A large scale dataset for NLI (Natural Language Inference; Recognizing Textual Entailment [Dagan et al., 2013])
- Premises are image captions, hypotheses generated by crowdsourcing workers:

Premise

Street performer is doing his act for kids

Hypotheses

1. A person performing for children on the street \Rightarrow ENTAILMENT

- A large scale dataset for NLI (Natural Language Inference; Recognizing Textual Entailment [Dagan et al., 2013])
- Premises are image captions, hypotheses generated by crowdsourcing workers:

Premise

Street performer is doing his act for kids

- 1. A person performing for children on the street \Rightarrow ENTAILMENT
- 2. A juggler entertaining a group of children on the street

- A large scale dataset for NLI (Natural Language Inference; Recognizing Textual Entailment [Dagan et al., 2013])
- Premises are image captions, hypotheses generated by crowdsourcing workers:

Premise

Street performer is doing his act for kids

- 1. A person performing for children on the street \Rightarrow ENTAILMENT
- 2. A juggler entertaining a group of children on the street \Rightarrow <code>NEUTRAL</code>

- A large scale dataset for NLI (Natural Language Inference; Recognizing Textual Entailment [Dagan et al., 2013])
- Premises are image captions, hypotheses generated by crowdsourcing workers:

Premise

Street performer is doing his act for kids

- 1. A person performing for children on the street \Rightarrow ENTAILMENT
- 2. A juggler entertaining a group of children on the street \Rightarrow <code>NEUTRAL</code>
- 3. A magician performing for an audience in a nightclub

- A large scale dataset for NLI (Natural Language Inference; Recognizing Textual Entailment [Dagan et al., 2013])
- Premises are image captions, hypotheses generated by crowdsourcing workers:

Premise

Street performer is doing his act for kids

- 1. A person performing for children on the street \Rightarrow ENTAILMENT
- 2. A juggler entertaining a group of children on the street \Rightarrow **NEUTRAL**
- 3. A magician performing for an audience in a nightclub \Rightarrow CONTRADICTION
- Event co-reference assumption



End-to-end, either sentence-encoding or attention-based





1

End-to-end, either sentence-encoding or attention-based



Lexical knowledge: only from pre-trained word embeddings

As opposed to using resources like WordNet



- Lexical knowledge: only from pre-trained word embeddings
 - As opposed to using resources like WordNet
- SOTA exceeds human performance...



- Lexical knowledge: only from pre-trained word embeddings
 - As opposed to using resources like WordNet
- SOTA exceeds human performance... ¹



- Lexical knowledge: only from pre-trained word embeddings
 - As opposed to using resources like WordNet
- SOTA exceeds human performance... ¹

¹[Gururangan et al., 2018, Poliak et al., 2018]: by learning "easy clues"

 [Gururangan et al., 2018, Poliak et al., 2018]: good performance on SNLI based on the hypothesis alone

- [Gururangan et al., 2018, Poliak et al., 2018]: good performance on SNLI based on the hypothesis alone
- This is a result of the annotation procedure

- [Gururangan et al., 2018, Poliak et al., 2018]: good performance on SNLI based on the hypothesis alone
- This is a result of the annotation procedure
 - Negation (not, never, nobody) is correlated with contradiction

- [Gururangan et al., 2018, Poliak et al., 2018]: good performance on SNLI based on the hypothesis alone
- This is a result of the annotation procedure
 - Negation (not, never, nobody) is correlated with contradiction
 - ...and "cat" as well (many dog images)

- [Gururangan et al., 2018, Poliak et al., 2018]: good performance on SNLI based on the hypothesis alone
- This is a result of the annotation procedure
 - Negation (not, never, nobody) is correlated with contradiction
 - …and "cat" as well (many dog images)
 - Generic words (animal, instrument) are correlated with *entailment*

- [Gururangan et al., 2018, Poliak et al., 2018]: good performance on SNLI based on the hypothesis alone
- This is a result of the annotation procedure
 - Negation (not, never, nobody) is correlated with contradiction
 - ...and "cat" as well (many dog images)
 - Generic words (animal, instrument) are correlated with *entailment*
 - Sentence length: *entailment* < *contradiction* < *neutral*

Do neural NLI models implicitly learn lexical semantic relations?

We constructed a new test set to answer this question

- We constructed a new test set to answer this question
- Premise: sentences from the SNLI training set

- We constructed a new test set to answer this question
- **Premise**: sentences from the SNLI training set
- Hypothesis:
 - Replacing a single term w in the premise with a related term w'

- We constructed a new test set to answer this question
- **Premise**: sentences from the SNLI training set
- Hypothesis:
 - Replacing a single term w in the premise with a related term w'
 - \blacksquare w' is in the SNLI vocabulary and in pre-trained embeddings

- We constructed a new test set to answer this question
- **Premise**: sentences from the SNLI training set
- Hypothesis:
 - Replacing a single term w in the premise with a related term w'
 - \blacksquare w' is in the SNLI vocabulary and in pre-trained embeddings
 - Crowdsourcing labels (mostly contradictions!)

- We constructed a new test set to answer this question
- Premise: sentences from the SNLI training set
- Hypothesis:
 - Replacing a single term w in the premise with a related term w'
 - \blacksquare w' is in the SNLI vocabulary and in pre-trained embeddings
 - Crowdsourcing labels (mostly contradictions!)

Contradiction

The man is holding a saxophone \rightarrow The man is holding an electric guitar

- We constructed a new test set to answer this question
- Premise: sentences from the SNLI training set
- Hypothesis:
 - Replacing a single term w in the premise with a related term w'
 - \blacksquare w' is in the SNLI vocabulary and in pre-trained embeddings
 - Crowdsourcing labels (mostly contradictions!)

Contradiction

The man is holding a saxophone \rightarrow The man is holding an electric guitar

Entailment

A little girl is very $\underline{sad} \rightarrow A$ little girl is very unhappy

- We constructed a new test set to answer this question
- Premise: sentences from the SNLI training set
- Hypothesis:
 - Replacing a single term w in the premise with a related term w'
 - \blacksquare w' is in the SNLI vocabulary and in pre-trained embeddings
 - Crowdsourcing labels (mostly contradictions!)

Contradiction

The man is holding a saxophone \rightarrow The man is holding an electric guitar

Entailment

A little girl is very $\underline{sad} \rightarrow A$ little girl is very unhappy

Neutral

A couple drinking $\underline{wine} \rightarrow A$ couple drinking champagne

Evaluation Setting

- 3 representative models:
 - Residual-Stacked-Encoder [Nie and Bansal, 2017]
 - ESIM (Enhanced Sequential Inference Model) [Chen et al., 2017]
 - Decomposable Attention [Parikh et al., 2016]
Evaluation Setting

- 3 representative models:
 - Residual-Stacked-Encoder [Nie and Bansal, 2017]
 - ESIM (Enhanced Sequential Inference Model) [Chen et al., 2017]
 - Decomposable Attention [Parikh et al., 2016]
- Train on SNLI training set, test on the original & new test set
 - In the paper: enhancing with additional existing datasets

Results

Can neural NLI models recognize lexical inferences?



Dramatic drop in performance across models.

Sanity Check

Performance of WordNet-informed Models



The test set is solvable using WordNet.

What do neural NLI models learn with respect to lexical semantic relations?

Analysis 1: Word Similarity

Models err on contradicting word-pairs with similar embeddings

 \blacksquare A man starts his day in India \rightarrow A man starts his day in Malaysia

Analysis 1: Word Similarity

- Models err on contradicting word-pairs with similar embeddings
 - \blacksquare A man starts his day in India \rightarrow A man starts his day in Malaysia
- Especially for fixed word embeddings



Analysis 2: Frequency in Training

■ Tuning embeddings may associate specific (*word, replacement*) pairs to a label, e.g. (*man, woman*) → contradiction

Analysis 2: Frequency in Training

- Tuning embeddings may associate specific (*word, replacement*) pairs to a label, e.g. (*man, woman*) → contradiction
- Accuracy increases with frequency in training set



Frequency of (word, replacement) pairs in contradiction training examples

New NLI test set that evaluates systems' ability to make inferences that require very simple lexical knowledge

- New NLI test set that evaluates systems' ability to make inferences that require very simple lexical knowledge
- SOTA systems perform poorly on the test set

- New NLI test set that evaluates systems' ability to make inferences that require very simple lexical knowledge
- SOTA systems perform poorly on the test set
- Systems are limited in their generalization ability

- New NLI test set that evaluates systems' ability to make inferences that require very simple lexical knowledge
- SOTA systems perform poorly on the test set
- Systems are limited in their generalization ability

Related Work:

 "Stress tests" [Naik et al., 2018]: similar findings on a broader range of linguistic phenomena

- New NLI test set that evaluates systems' ability to make inferences that require very simple lexical knowledge
- SOTA systems perform poorly on the test set
- Systems are limited in their generalization ability

Related Work:

- "Stress tests" [Naik et al., 2018]: similar findings on a broader range of linguistic phenomena
- Inference with single word differences: [Pavlick and Callison-Burch, 2016, Kalouli et al., 2018]

But current LM-based NLI models address entailment-related phenomena better, no?

Pre-trained LM based NLI models



Diversify Your Datasets: Analyzing Generalization via Controlled Variance in Adversarial Datasets

Ohad Rozen, Vered Shwartz, Roee Aharoni, and Ido Dagan



Bar-Ilan University CoNLL 2019

MultiNLI [Williams et al., 2018]

Collected like SNLI (existing premises, generated hypotheses)

MultiNLI [Williams et al., 2018]

- Collected like SNLI (existing premises, generated hypotheses)
- Multiple geners

MultiNLI [Williams et al., 2018]

- Collected like SNLI (existing premises, generated hypotheses)
- Multiple geners
- Mismatched evaluation (not in our focus)

$\textbf{Probing} \rightarrow \textbf{Inoculation}$



 Probing: does the representation capture a certain property?
[Glockner et al., 2018, Naik et al., 2018]

$\textbf{Probing} \rightarrow \textbf{Inoculation}$



- Probing: does the representation capture a certain property?
 [Glockner et al., 2018, Naik et al., 2018]
- Inoculation [Liu et al., 2019]: can the representation *learn* a certain property?

Has the model learned a general notion of the property, or does it overfit to the specific dataset?

1. Split the challenge dataset to different variations across the dimension in focus.



2. Fine-tune on one set and test on another.







Challenge Datasets Generation

Templates from MultiNLI sentences \rightarrow Ask crowdsourcing workers to rephrase spans.

- 1. Extracted Premise: [The Citigroup deal], [from beginning to end], [took] less than 5 [weeks].
- 2. Premise Template: ARG1, ARG2, ARG3 RELATION NUM ARG4. Hypothesis Template (Ent.): ARG1, ARG2, ARG3 more than NUM-smaller ARG4.
- Gen. Premise: [My marriage], [despite much frustration], [lasted] more than 7 [years]. Gen. Hypothesis (Ent.): [My marriage], [despite much frustration], [lasted] more than 2 [years].

1000s of different training examples with similar syntax from 1 original sentence

Phenomena Dative Alternation

Premise: I baked <u>my mom a cake</u> Hypothesis 1: I baked a cake for <u>my mom</u> Label: Entailment

Phenomena Numeric Reasoning

Premise: I see 260 coins in the bucket. Hypothesis: I see more than 232 coins in the bucket. Label: Entailment

Challenge Datasets Generation Diversity Axes

- 1. Syntax complexity simple / medium / complex
 - Sentence length
 - Phenomenon depth in parse tree

Challenge Datasets Generation Diversity Axes

- 1. Syntax complexity simple / medium / complex
 - Sentence length
 - Phenomenon depth in parse tree
- 2. Lexical variability
 - Dative verb
 - Number range

Inoculation tells part of the story...



Dative Alternation

Not sensitive to lexical variability

Dative Alternation

Not sensitive to lexical variability But generalizes only from complex to simple syntax

Dative Alternation

Not sensitive to lexical variability But generalizes only from complex to simple syntax

Numeric Reasoning Not sensitive to syntax

Dative Alternation

Not sensitive to lexical variability But generalizes only from complex to simple syntax

Numeric Reasoning Not sensitive to syntax But doesn't generalize across number ranges
Diversify your Datasets Recap

Simple methodology to test model generalization of a specific learned phenomenon

Diversify your Datasets Recap

- Simple methodology to test model generalization of a specific learned phenomenon
- NLI-BERT fails to generalize dative alternation and numeric reasoning

Diversify your Datasets Recap

- Simple methodology to test model generalization of a specific learned phenomenon
- NLI-BERT fails to generalize dative alternation and numeric reasoning
- Fine-tuning on the phenomenon-specific data may decrease the main task performance (also in [Richardson et al., 2020]).

Real-world examples: partial entailments

S1: Amazon To Acquire Whole Foods Market For \$13.7 Billion
S2: Amazon is buying Whole Foods for almost \$14 billion in cash

Real-world examples: partial entailments

- S₁: Researchers have discovered wreckage of the lost warship, the USS Indianapolis after 72 years
- *S*₂: Wreckage of missing WWII ship found in Pacific Ocean

Breaking

Building

NLI

Coreference

[Shwartz et al., 2017] [Barhom et al., 2019]

Revisiting Joint Modeling of Cross-document Entity and Event Coreference Resolution

Shany Barhom¹, Vered Shwartz¹, Alon Eirew², Michael Bugert³, Nils Reimers³, and Ido Dagan¹



Tara Reid has entered a rehab center...

Doc #1

...She checked into the facility today ...

Doc #1

Tara Reid <u>has entered</u> a rehab center... ...She <u>checked into</u> the facility today...

Tara Reid has entered a rehab center...

Doc #1

...She checked into the facility today ...

Doc #1

Tara Reid has entered a rehab center...

...She checked into the facility today ...

Doc #2

#2 ...the American Pie star headed to a Malibu treatment facility on Tuesday...

















• Entity and Event Coreference are closely interdependent - calls for a joint approach

- Entity and Event Coreference are closely interdependent calls for a joint approach
- Only single such prior work (Lee et al., 2012)

- Entity and Event Coreference are closely interdependent calls for a joint approach
- Only single such prior work (Lee et al., 2012)
- We revisit the joint resolution approach, suggesting new neural models to address it
 - Achieving new SOTA

Related Work

Within-document Coreference



Within-document Entity Coreference







Within-document Event Coreference



Choubey and Huang 2017b Choubey and Huang, 2017a Lu et al. 2016 Krause et al., 2016 Lu and Ng, 2016 Liu et al., 2014 Chen et al., 2009



Cross-document Event Coreference



Kenyon-Dean et al., 2018 Choubey and Huang, 2017 Cybulska and Vossen, 2015 Yang et al., 2015 Lee et al., 2012



Cross-document Entity Coreference

Dutta and Weikum, 2015 Singh et al., 2015 Lee et al., 2012 Rao et al., 2010





Common Approach - Lexical Similarity between Arguments

1. Tara Reid has entered a rehab center

2. The American Pie star headed to a Malibu treatment facility on Tuesday

Common Approach - Lexical Similarity between Arguments

- 1. <u>Tara Reid</u> has entered a rehab center
- 2. The American Pie star headed to a Malibu treatment facility on Tuesday

Common Approach - Lexical Similarity between Arguments

- 1. Tara Reid has entered a rehab center
- 2. The American Pie star headed to a Malibu treatment facility on Tuesday

Is that good enough?



• Lee et al., (2012) introduced a system that models entity and event coreference jointly

- Lee et al., (2012) introduced a system that models entity and event coreference jointly
- Iterative method that constructs clusters of entity and event mentions

- Lee et al., (2012) introduced a system that models entity and event coreference jointly
- Iterative method that constructs clusters of entity and event mentions
- Linear regression to model cluster merge operations, based on discrete features

- Lee et al., (2012) introduced a system that models entity and event coreference jointly
- Iterative method that constructs clusters of entity and event mentions
- Linear regression to model cluster merge operations, based on discrete features
- We revisit the joint approach, suggesting a neural models to address it
Model







Multiple event cluster merges



Cluster Merging Score

- Hierarchical clustering requires a cluster pair merging score
- Average link: average all mention pair scores across the two candidate clusters

$$S_{cp}(c_i, c_j) = \frac{1}{|c_i| \cdot |c_j|} \cdot \sum_{m_i \in c_i} \sum_{m_j \in c_j} S(m_i, m_j)$$

Cluster Merging Score

- Hierarchical clustering requires a cluster pair merging score
- Average link: average all mention pair scores across the two candidate clusters

$$S_{cp}(c_i, c_j) = \frac{1}{|c_i| \cdot |c_j|} \cdot \sum_{m_i \in c_i} \sum_{m_j \in c_j} S(m_i, m_j)$$

- 1. Tara Reid has entered a rehab center.
- 2. The American Pie star headed to a Malibu treatment facility on Tuesday.





- 1. Tara Reid has entered a rehab center.
- 2. The American Pie star headed to a Malibu treatment facility on Tuesday.





_ _



- 1. Tara Reid has entered a rehab center
- 2. The American Pie star headed to a Malibu treatment facility on Tuesday





- 1. Tara Reid has entered a rehab center
- 2. The American Pie star headed to a Malibu treatment facility on Tuesday



_ _ _

















Mention-Pair Scorer



Pairwise Score Arg 0 Arg 1 Tara Reid has entered a rehab center 1. 2. The American Pie star headed to a Malibu treatment facility on Tuesday MLP Arg 0 Arg 1 Pair 0000 000 00 Representation Pairwise Mention 1 Mention 2 Mention 1 * Mention 2 $S_{cp}(c_i, c_j) = \frac{1}{|c_i| \cdot |c_j|} \cdot \sum_{m_i \in c_i} \sum_{m_j \in c_j} S(m_i, m_j)$ Features Mention Representations Context Span Ara0 Arg1 Location Time

Mention-Pair Scorer

• We train two distinct pairwise scorers (one for entities and one for events)

- We train two distinct pairwise scorers (one for entities and one for events)
- The training procedure simulates the inference step
 - o Allows the models to be trained on various predicted clustering configurations

- We train two distinct pairwise scorers (one for entities and one for events)
- The training procedure simulates the inference step
 - Allows the models to be trained on various predicted clustering configurations
- **Training examples:** all mention pairs that belong to different clusters in the current clustering configuration

- We train two distinct pairwise scorers (one for entities and one for events)
- The training procedure simulates the inference step
 - Allows the models to be trained on various predicted clustering configurations
- **Training examples:** all mention pairs that belong to different clusters in the current clustering configuration
- Scorers are repeatedly trained and then used for clusters merging

Experiments

Dataset

- ECB+ (Event-Coreference-Bank; Cybulska and Vossen, 2014).
 - Within- and cross-document coreference annotations for entities and events
 - ~1000 documents, clustered into 43 topics, discussing different seminal events

Dataset

- ECB+ (Event-Coreference-Bank; Cybulska and Vossen, 2014).
 - Within- and cross-document coreference annotations for both entities and events
 - ~1000 documents, clustered into 43 topics, discussing different seminal events

Topic 1: A celebrity enters into rehab

...Tara Reid finally checks into rehab... ...Actress Tara Reid entered well-known Malibu rehab center ...





Evaluation Setup

• We follow Cybulska and Vossen (2015) and Kenyon-Dean et al., (2018)

\circ $\,$ Corpus's subset which has been validated for correctness

• Use the gold mentions during evaluation

	Train	Validation	Test	Total
# Topics	25	8	10	43
# Sub-topics	50	16	20	86
# Documents	574	196	206	976
# Sentences	1037	346	457	1840
# Event mentions	3808	1245	1780	6833
# Entity mentions	4758	1476	2055	8289
# Event chains	1527	409	805	2741
# Entity chains	1286	330	608	2224

Evaluation Setup

~1000 sen for training

• We follow Cybulska and Vossen (2015) and Kenyon-Dean et al., (2018)

\circ $\,$ Corpus's subset which has been validated for correctness

• Use the gold mentions during evaluation

		Train	Validation	Test	Total
	# Topics	25	8	10	43
	# Sub-topics	50	16	20	86
tences 9	# Documents	574	196	206	976
	# Sentences	1037	346	457	1840
	# Event mentions	3808	1245	1780	6833
	# Entity mentions	4758	1476	2055	8289
	# Event chains	1527	409	805	2741
	# Entity chains	1286	330	608	2224

Coreference Results

_ _ _



Event Coreference Results

Coreference Results



Entity Coreference Results

Coreference Results



Entity Coreference Results





Entity coreference:

- [WWDC, San Francisco gathering's, conference]
- [West Papua, region, in remote eastern Indonesia]
- [Matt Smith, actor]

Event coreference:

- [launches, unveiled]
- [rattled, struck, hit]
- [acquires, buys, purchase]

Wrong Model's Decisions



Entity coreference:

- [next generation of MacBook Pro, MacBook Pro] •
- [five people, four people] •
- [Wednesday, on Monday] .

Event coreference:

- [recorded, occurred] •
- [sales, acquisition] •
- [gone official, go ahead] .

Wrong Model's Decisions



Entity coreference:

- [next generation of MacBook Pro, MacBook Pro]
- [five people, four people]
- [Wednesday, on Monday]

Event coreference:

- [recorded, occurred]
- [sales, acquisition]
- [gone official, go ahead]

Paraphrasing vs. Relatedness

Wrong Model's Decisions



Entity coreference:

- [next generation of MacBook Pro, MacBook Pro]
- [five people, four people] 🔫
- [Wednesday, on Monday]

Event coreference:

- [recorded, occurred]
- [sales, acquisition]
- [gone official, go ahead]

Same Head Lemma
Error Analysis



Joint Event and Entity Coreference Resolution Recap

Cross-document coreference is drastically under-explored

Joint Event and Entity Coreference Resolution Recap

Cross-document coreference is drastically under-explored

A simple joint approach with state-of-the-art results on ECB+

Joint Event and Entity Coreference Resolution Recap

Cross-document coreference is drastically under-explored

- A simple joint approach with state-of-the-art results on ECB+
- Still a long way to go!

Acquiring Predicate Paraphrases from News Tweets

Vered Shwartz, Gabriel Stanovsky, and Ido Dagan



Acquiring Predicate Paraphrases from News Tweets²

[a]₀ introduce [a]₁ [a]₀ appoint [a]₁ $[a]_0$ die at $[a]_1$ $[a]_0$ hit $[a]_1$ [a]₀ be investigate [a]₁ [a]₀ eliminate [a]₁ $[a]_0$ announce $[a]_1$ [a]₀ quit after [a]₁ $[a]_0$ announce as $[a]_1$ [a]₀ threaten [a]₁ $[a]_0$ die at $[a]_1$ [a]₀ double down on [a]₁ [a]₀ kill [a]₁ $[a]_0$ approve $[a]_1$ seize [a]₀ at [a]₁

[a]₀ welcome [a]₁ $[a]_0$ to become $[a]_1$ $[a]_0$ pass away at $[a]_1$ $[a]_0$ sink to $[a]_1$ $[a]_0$ be probe $[a]_1$ $[a]_0$ slash $[a]_1$ $[a]_0$ unveil $[a]_1$ $[a]_0$ resign after $[a]_1$ $[a]_0$ to become $[a]_1$ $[a]_0$ warn $[a]_1$ [a]₀ live until [a]₁ $[a]_0$ stand by $[a]_1$ $[a]_0$ shoot $[a]_1$ $[a]_0$ pass $[a]_1$ to grab [a]₀ at [a]₁

- Binary verbal predicate paraphrases
- Extracted from Twitter
- Ever-growing resource: currently around 5.2M paraphrases

²Available at https://github.com/vered1986/Chirps

Assumptions

Main assumption: redundant news headlines of the same event are likely to describe it with different words [Shinyama et al., 2002, Barzilay and Lee, 2003].

Assumptions

- Main assumption: redundant news headlines of the same event are likely to describe it with different words [Shinyama et al., 2002, Barzilay and Lee, 2003].
- This work: propositions extracted from tweets discussing news events, published on the same day, that agree on their arguments, are predicate paraphrases.



to buy [Amazon] is buying [Whole Foods] to acquire





Query the Twitter Search API for news tweets in English

Amazon is buying Whole Foods in \$13.7B

Amazon to acquire Whole Foods Market in deal valued at nearly \$14 billion

• • •



 Extract propositions from tweets using PropS [Stanovsky et al., 2016]

Get binary verbal predicate templates, and apply argument reduction [Stanovsky and Dagan, 2016]

[Amazon] **buy** [Whole Foods] [Amazon] **acquire** [Whole Foods Market]



We consider two predicates as paraphrases if:

- 1. They appear on the same day.
- 2. Each of their arguments aligns with a unique argument in the other predicate.
- Two levels of argument matching: strict (exact match / short edit distance) and loose (partial token matching / WordNet synonyms)

[<i>a</i>] ₀ buy [<i>a</i>] ₁	$[a]_0$ acquire $[a]_1$	Amazon	Whole Foods
[<i>a</i>] ₀ buy [<i>a</i>] ₁	[a] ₀ acquire [a] ₁	Intel	Mobileye



$$p_1 = [a]_0$$
 buy $[a]_1$, $p_2 = [a]_0$ acquire $[a]_1$
 $s(p_1, p_2) = count(p_1, p_2) \cdot \left(1 + \frac{days(p_1, p_2)}{N}\right)$

- $count(p_1, p_2)$ assigns high scores for frequent paraphrases
- N number of days since the resource collection begun
- $\frac{days(p_1,p_2)}{N}$ eliminates noise from two arguments participating in different events on the same day

1) Last year when Chuck Berry turned 90; 2) Chuck Berry dies at 90



- We release our resource daily, with two files:
 - **Instances**: predicates, arguments and tweet IDs.
 - **Types**: predicate paraphrase pair types ranked in a descending order according to the heuristic accuracy score.



Using event coreference to extract paraphrases

Chirps Recap

Using event coreference to extract paraphrases

Complementary to other paraphrasing resources

Chirps Recap

Using event coreference to extract paraphrases

- Complementary to other paraphrasing resources
- Useful resource for paraphrasing, event coreference, NLI

Thank you! Questions?





References I

[Barhom et al., 2019] Barhom, S., Shwartz, V., Eirew, A., Bugert, M., Reimers, N., and Dagan, I. (2019). Revisiting joint modeling of cross-document entity and event coreference resolution. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4179–4189, Florence, Italy. Association for Computational Linguistics.

[Barzilay and Lee, 2003] Barzilay, R. and Lee, L. (2003). Learning to paraphrase: An unsupervised approach using multiple-sequence alignment. In *NAACL*.

[Bowman et al., 2015] Bowman, S. R., Angeli, G., Potts, C., and Manning, D. C. (2015). A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642. Association for Computational Linguistics.

[Chen et al., 2018] Chen, Q., Zhu, X., Ling, Z.-H., Inkpen, D., and Wei, S. (2018). Neural natural language inference models enhanced with external knowledge. In *The 56th Annual Meeting of the Association for Computational Linguistics (ACL)*, Melbourne, Australia.

[Chen et al., 2017] Chen, Q., Zhu, X., Ling, Z.-H., Wei, S., Jiang, H., and Inkpen, D. (2017). Enhanced lstm for natural language inference. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1657–1668, Vancouver, Canada. Association for Computational Linguistics.

[Dagan et al., 2013] Dagan, I., Roth, D., Sammons, M., and Zanzotto, F. M. (2013). Recognizing textual entailment: Models and applications. *Synthesis Lectures on Human Language Technologies*, 6(4):1–220.

References II

[Glockner et al., 2018] Glockner, M., Shwartz, V., and Goldberg, Y. (2018). Breaking NLI systems with sentences that require simple lexical inferences. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 650–655, Melbourne, Australia. Association for Computational Linguistics.

[Gururangan et al., 2018] Gururangan, S., Swayamdipta, S., Levy, O., Schwartz, R., Bowman, S. R., and Smith, N. A. (2018). Annotation artifacts in natural language inference data. In *The* 16th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT), New Orleans, Louisiana.

[Kalouli et al., 2018] Kalouli, A.-L., Real, L., and DePaiva, V. (2018). Wordnet for "easy" textual inferences. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), Paris, France. European Language Resources Association (ELRA).

[Liu et al., 2019] Liu, N. F., Schwartz, R., and Smith, N. A. (2019). Inoculation by fine-tuning: A method for analyzing challenge datasets. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2171–2179, Minneapolis, Minnesota. Association for Computational Linguistics.

[Naik et al., 2018] Naik, A., Ravichander, A., Sadeh, N., Rose, C., and Neubig, G. (2018). Stress test evaluation for natural language inference. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

References III

[Nie and Bansal, 2017] Nie, Y. and Bansal, M. (2017). Shortcut-stacked sentence encoders for multi-domain inference. arXiv preprint arXiv:1708.02312.

[Parikh et al., 2016] Parikh, A., Täckström, O., Das, D., and Uszkoreit, J. (2016). A decomposable attention model for natural language inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2249–2255, Austin, Texas. Association for Computational Linguistics.

[Pavlick and Callison-Burch, 2016] Pavlick, E. and Callison-Burch, C. (2016). Most baies are little and most problems are huge: Compositional entailment in adjective nouns.

[Poliak et al., 2018] Poliak, A., Naradowsky, J., Haldar, A., Rudinger, R., and Van Durme, B. (2018). Hypothesis Only Baselines in Natural Language Inference. In *Joint Conference on Lexical and Computational Semantics (StarSem)*.

[Richardson et al., 2020] Richardson, K., Hu, H., Moss, L. S., and Sabharwal, A. (2020). Probing natural language inference models through semantic fragments. In AAAI.

[Rozen et al., 2019] Rozen, O., Shwartz, V., Aharoni, R., and Dagan, I. (2019). Diversify your datasets: Analyzing generalization via controlled variance in adversarial datasets. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 196–205, Hong Kong, China. Association for Computational Linguistics.

[Shinyama et al., 2002] Shinyama, Y., Sekine, S., and Sudo, K. (2002). Automatic paraphrase acquisition from news articles. In *HLT*, pages 313–318. Morgan Kaufmann Publishers Inc.



[Shwartz et al., 2017] Shwartz, V., Stanovsky, G., and Dagan, I. (2017). Acquiring predicate paraphrases from news tweets. In **SEM*, pages 155–160.

[Stanovsky and Dagan, 2016] Stanovsky, G. and Dagan, I. (2016). Annotating and predicting non-restrictive noun phrase modifications. In ACL.

[Stanovsky et al., 2016] Stanovsky, G., Ficler, J., Dagan, I., and Goldberg, Y. (2016). Getting more out of syntax with props. arXiv.

[Williams et al., 2018] Williams, A., Nangia, N., and Bowman, S. R. (2018). A broad-coverage challenge corpus for sentence understanding through inference. In *The 16th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, New Orleans, Louisiana.