# Natural Language Inference: Overview

## Christopher Potts

Stanford Linguistics

## CS224u: Natural language understanding

# Associated materials

1. Code
   a. `nli.py`
   b. `nli_01_task_and_data.ipynb`
   c. `nli_02_models.ipynb`

2. Homework and bakeoff: `hw_wordentail.ipynb`

3. Core readings: Bowman et al. 2015; Williams et al. 2018; Nie et al. 2019; Rocktäschel et al. 2016

4. Auxiliary readings: Goldberg 2015; Dagan et al. 2006; MacCartney and Manning 2008; Gururangan et al. 2018

# Simple examples

| Premise | Relation | Hypothesis |
|---------|----------|------------|
| A turtle danced. | entails | A turtle moved. |
| turtle | contradicts | linguist |
| Every reptile danced. | neutral | A turtle ate. |
| Some turtles walk. | contradicts | No turtles move. |
| James Byron Dean refused to move without blue jeans. | entails | James Dean didn't dance without pants. |
| Mitsubishi Motors Corp's new vehicle sales in the US fell 46 percent in June. | contradicts | Mitsubishi's sales rose 46 percent. |
| Acme Corporation reported that its CEO resigned. | entails | Acme's CEO resigned. |

# NLI task formulation

## Does the premise justify an inference to the hypothesis?

- Commonsense reasoning, rather than strict logic.
- Focus on local inference steps, rather than long deductive chains.
- Emphasis on variability of linguistic expression.

## Perspectives

- Zaenen et al. (2005): Local textual inference: can it be defined or circumscribed?
- Manning (2006): Local textual inference: it's hard to circumscribe, but you know it when you see it – and NLP needs it.
- Crouch et al. (2006): Circumscribing is not excluding: a reply to Manning.

# Connections to other tasks

## Dagan et al. (2006)

It seems that major inferences, as needed by multiple applications, can indeed be cast in terms of textual entailment.
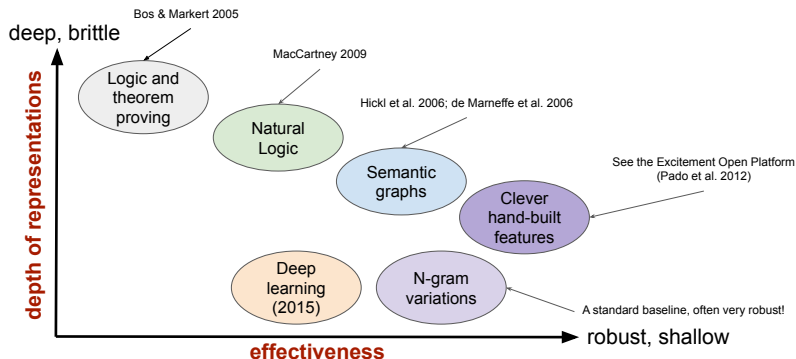
[ . . . ]

Consequently, we hypothesize that textual entailment recognition is a suitable generic task for evaluating and comparing applied semantic inference models. Eventually, such efforts can promote the development of entailment recognition "engines" which may provide useful generic modules across applications.
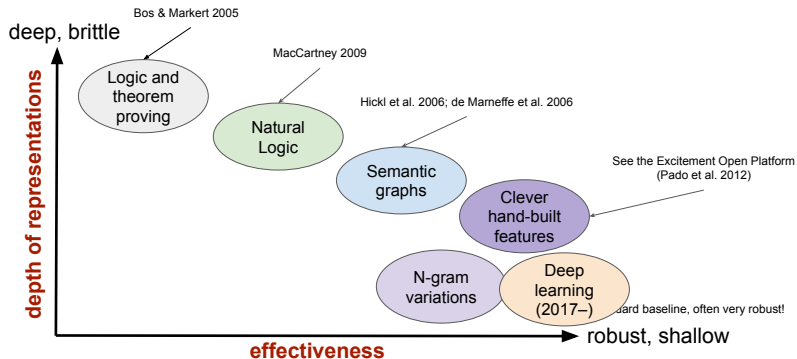
# Connections to other tasks

| Task | NLI framing |
|------|-------------|
| Paraphrase | text $\equiv$ paraphrase |
| Summarization | text $\sqsubset$ summary |
| Information retrieval | query $\sqsupset$ document |
| Question answering | question $\sqsupset$ answer |
| | *Who left? $\Rightarrow$ Someone left* |
| | *Someone left $\sqsupset$ Sandy left* |

# Models for NLI

# Models for NLI

# References I

Johan Bos and Katja Markert. 2005. Recognising textual entailment with logical inference. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 628–635, Stroudsburg, PA. ACL.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 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, Stroudsburg, PA. Association for Computational Linguistics.

Richard Crouch, Lauri Karttunen, and Annie Zaenen. 2006. Circumscribing is not excluding: A reply to Manning. Ms., Palo Alto Research Center.

Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The PASCAL recognising textual entailment challenge. In *Machine Learning Challenges, Lecture Notes in Computer Science*, volume 3944, pages 177–190. Springer-Verlag.

Yoav Goldberg. 2015. A primer on neural network models for natural language processing. Ms., Bar Ilan University.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.

Andrew Hickl and Jeremy Bensley. 2007. A discourse commitment-based framework for recognizing textual entailment. In *Proceedings of the Workshop on Textual Entailment and Paraphrasing*.

Bill MacCartney. 2009. *Natural Language Inference*. Ph.D. thesis, Stanford University.

Bill MacCartney and Christopher D. Manning. 2008. Modeling semantic containment and exclusion in natural language inference. In *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pages 521–528, Manchester, UK. Coling 2008 Organizing Committee.

Christopher D. Manning. 2006. Local textual inference: It's hard to circumscribe, but you know it when you see it – and NLP needs it. Ms., Stanford University.

Marie-Catherine de Marneffe, Bill MacCartney, Trond Grenager, Daniel Cer, Anna Rafferty, and Christopher D Manning. 2006. Learning to distinguish valid textual entailments. In *Proceedings of the 2nd Pascal RTE Challenge Workshop*.

Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2019. Adversarial NLI: A new benchmark for natural language understanding. UNC CHapel Hill and Facebook AI Research.

Sebastian Pado, Tae-Gil Noh, Asher Stern, and Rui Wang. 2013. Design and realization of a modular architecture for textual entailment. *Journal of Natural Language Engineering.*, 21(2):167–200.

Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomás Kočiský, and Phil Plunsom. 2016. Reasoning about entailment with neural attention. ArXiv:1509.06664.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.

# References II

Annie Zaenen, Lauri Karttunen, and Richard Crouch. 2005. Local textual inference: Can it be defined or circumscribed? In *Proceedings of the ACL Workshop on Empirical Modeling of Semantic Equivalence and Entailment*, pages 31–36, Ann Arbor, MI. Association for Computational Linguistics.

# Natural Language Inference: SNLI, MultiNLI, and Adversarial NLI

Christopher Potts

Stanford Linguistics

CS224u: Natural language understanding

# SNLI

1. Bowman et al. 2015
2. All the premises are image captions from the Flickr30K corpus (Young et al. 2014).
3. All the hypotheses were written by crowdworkers.
4. Some of the sentences reflect stereotypes (Rudinger et al. 2017).
5. 550,152 train examples; 10K dev; 10K test
6. Mean length in tokens:
   - Premise: 14.1
   - Hypothesis: 8.3
7. Clause-types:
   - Premise S-rooted: 74%
   - Hypothesis S-rooted: 88.9%
8. Vocab size: 37,026
9. 56,951 examples validated by four additional annotators.
   - 58.3% examples with unanimous gold label
   - 91.2% of gold labels match the author's label
   - 0.70 overall Fleiss kappa
10. Leaderboard: https://nlp.stanford.edu/projects/snli/

# Crowdsourcing methods

**Instructions**

The Stanford University NLP Group is collecting data for use in research on computer understanding of English. We appreciate your help!

We will show you the caption for a photo. We will not show you the photo. Using only the caption and what you know about the world:

- Write one alternate caption that is **definitely** a **true** description of the photo.
- Write one alternate caption that **might be** a **true** description of the photo.
- Write one alternate caption that is **definitely** an **false** description of the photo.

**Photo caption** **A little boy in an apron helps his mother cook.**

**Definitely correct** Example: For the caption *"Two dogs are running through a field."* you could write *"There are animals outdoors."*

Write a sentence that follows from the given caption.

**Maybe correct** Example: For the caption *"Two dogs are running through a field."* you could write *"Some puppies are running to catch a stick."*

Write a sentence which may be true given the caption, and may not be.

**Definitely incorrect** Example: For the caption *"Two dogs are running through a field."* you could write *"The pets are sitting on a couch."*

Write a sentence which contradicts the caption.

**Problems (optional)** *If something is wrong with the caption that makes it difficult to understand, do your best above and let us know here.*

# Examples

| Premise | Relation | Hypothesis |
|---------|----------|------------|
| A man inspects the uniform of a figure in some East Asian country. | **contradiction** <br> c  c  c  c  c | The man is sleeping |
| An older and younger man smiling. | **neutral** <br> n  n  e  n  n | Two men are smiling and laughing at the cats playing on the floor. |
| A black race car starts up in front of a crowd of people. | **contradiction** <br> c  c  c  c  c | A man is driving down a lonely road. |
| A soccer game with multiple males playing. | **entailment** <br> e  e  e  e  e | Some men are playing a sport. |
| A smiling costumed woman is holding an umbrella. | **neutral** <br> n  n  e  c  n | A happy woman in a fairy costume holds an umbrella. |

# Event coreference

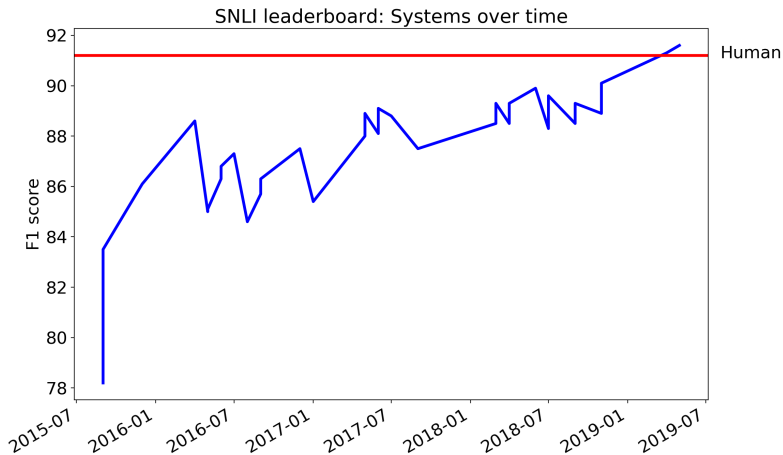| Premise | Relation | Hypothesis |
|---|---|---|
| A boat sank in the Pacific Ocean. | contradiction | A boat sank in the Atlantic Ocean. |
| Ruth Bader Ginsburg was appointed to the Supreme Court. | contradiction | I had a sandwich for lunch today |

If premise and hypothesis *probably* describe a different photo, then the label is contradiction

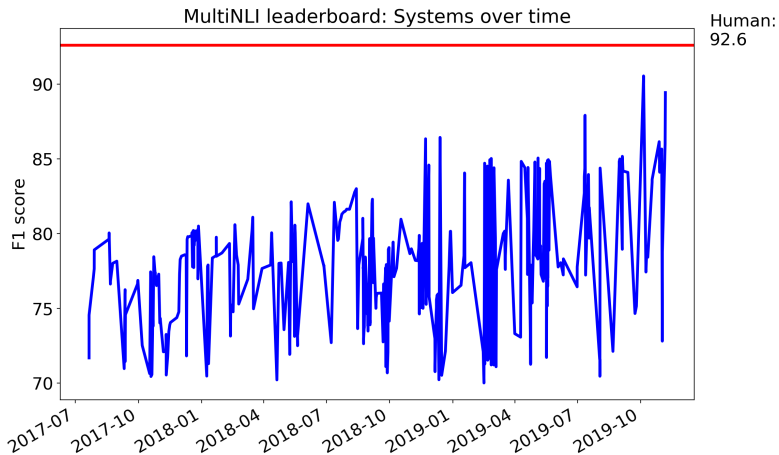# Progress on SNLI



SNLI leaderboard: Systems over time

# MultiNLI

1. Williams et al. 2018

2. Train premises drawn from five genres:
   - Fiction: works from 1912–2010 spanning many genres
   - Government: reports, letters, speeches, etc., from government websites
   - The *Slate* website
   - Telephone: the Switchboard corpus
   - Travel: Berlitz travel guides

3. Additional genres just for dev and test (the mismatched condition):
   - The 9/11 report
   - Face-to-face: The Charlotte Narrative and Conversation Collection
   - Fundraising letters
   - Non-fiction from Oxford University Press
   - *Verbatim*: articles about linguistics

4. 392,702 train examples; 20K dev; 20K test

5. 19,647 examples validated by four additional annotators
   - 58.2% examples with unanimous gold label
   - 92.6% of gold labels match the author's label

6. Test-set labels available as a Kaggle competition.

7. Project page: https://www.nyu.edu/projects/bowman/multinli/

## MultiNLI annotations

|  | Matched | Mismatched |
| --- | ---: | ---: |
| ACTIVE/PASSIVE | 15 | 10 |
| ANTO | 17 | 20 |
| BELIEF | 66 | 58 |
| CONDITIONAL | 23 | 26 |
| COREF | 30 | 29 |
| LONG_SENTENCE | 99 | 109 |
| MODAL | 144 | 126 |
| NEGATION | 129 | 104 |
| PARAPHRASE | 25 | 37 |
| QUANTIFIER | 125 | 140 |
| QUANTITY/TIME_REASONING | 15 | 39 |
| TENSE_DIFFERENCE | 51 | 18 |
| WORD_OVERLAP | 28 | 37 |
|  | 767 | 753 |

SNLI
○○○○

**MultiNLI**
○●

ANLI
○○

Dynabench

Other NLI datasets

# Progress on MultiNLI



MultiNLI leaderboard: Systems over time

Human: 92.6

SNLI
○○○○

MultiNLI
○○

ANLI
○○

Dynabench

Other NLI datasets

# Adversarial NLI dataset (ANLI)

1. Nie et al. 2019b

2. 162,865 labeled examples

3. The premises come from diverse sources.

4. The hypotheses are written by crowdworkers with the explicit goal of fooling state-of-the-art models.

5. This effort is a direct response to the results and findings for SNLI and MultiNLI that we just reviewed.

# ANLI dataset creation

1. The annotator is presented with a premise sentence and a condition (entailment, contradiction, neutral).

2. The annotator writes a hypothesis.

3. A state-of-the-art model makes a prediction about the premise–hypothesis pair.

4. If the model's prediction matches the condition, the annotator returns to step 2 to try again.

5. If the model was fooled, the premise–hypothesis pair is independently validated by other annotators.

# Additional ANLI details

| Round | Model | Training data | Context sources | Examples |
|-------|-------|---------------|-----------------|----------|
| R1 | BERT-large (Devlin et al. 2019) | SNLI + MultiNLI | Wikipedia | 16,946 |
| R2 | ROBERTa (Liu et al. 2019) | SNLI + MultiNLI + NLI-FEVER + R1 | Wikipedia | 45,460 |
| R3 | ROBERTa (Liu et al. 2019) | SNLI + MultiNLI + NLI-FEVER + R2 | Various | 100,459 |
| | | | | **162,865** |

- The train sets mix cases where the model's predictions were correct and incorrect. The majority of the model predictions are correct, though.
- The dev and test sets contain only cases where the model's prediction was incorrect.

SNLI
○○○○

MultiNLI
○○

ANLI
○○

**Dynabench**

Other NLI datasets

# Dynabench

**Dynabench: Rethinking Benchmarking in NLP**

**Douwe Kiela[†], Max Bartolo[‡], Yixin Nie[⋆], Divyansh Kaushik[§], Atticus Geiger[¶],**

**Zhengxuan Wu[¶], Bertie Vidgen[‖], Grusha Prasad[⋆⋆], Amanpreet Singh[†], Pratik Ringshia[†],**

**Zhiyi Ma[†], Tristan Thrush[†], Sebastian Riedel[†‡], Zeerak Waseem[††], Pontus Stenetorp[‡],**

**Robin Jia[†], Mohit Bansal[⋆], Christopher Potts[¶] and Adina Williams[†]**

[†] Facebook AI Research; [‡] UCL; [⋆] UNC Chapel Hill; [§] CMU; [¶] Stanford University
[‖] Alan Turing Institute; [⋆⋆] JHU; [††] Simon Fraser University
dynabench@fb.com

https://dynabench.org

SNLI
○○○○

MultiNLI
○○

ANLI
○○

**Dynabench**

Other NLI datasets

# Dynabench



Rethinking AI Benchmarking

Dynabench is a research platform for dynamic data collection and benchmarking. Static benchmarks have well-known issues: they saturate quickly, are susceptible to overfitting, contain exploitable annotator artifacts and have unclear or imperfect evaluation metrics.

This platform in essence is a scientific experiment: can we make faster progress if we collect data dynamically, with humans and models in the loop, rather than in the old-fashioned static way?

Read more

https://dynabench.org

# Other NLI datasets

- The GLUE benchmark (diverse tasks including NLI; Wang et al. 2018): https://gluebenchmark.com
- NLI Style FEVER (Nie et al. 2019a): https://github.com/easonnie/combine-FEVER-NSMN/blob/master/other_resources/nli_fever.md
- OCNLI: Original Chinese Natural Language Inference (Hu et al. 2020): https://github.com/CLUEbenchmark/OCNLI
- Turkish NLI (Budur et al. 2020): https://github.com/boun-tabi/NLI-TR
- XNLI (multilingual dev/test derived from MultiNLI; Conneau et al. 2018): https://github.com/facebookresearch/XNLI
- Diverse Natural Language Inference Collection (DNC; Poliak et al. 2018): http://decomp.io/projects/diverse-natural-language-inference/
- MedNLI (derived from MIMIC III; Romanov and Shivade 2018) https://physionet.org/content/mednli/1.0.0/
- SciTail (derived from science exam questions and Web text; Khot et al. 2018): http://data.allenai.org/scitail/

# References I

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 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, Stroudsburg, PA. Association for Computational Linguistics.

Emrah Budur, Riza Özçelik, Tunga Gungor, and Christopher Potts. 2020. Data and Representation for Turkish Natural Language Inference. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8253–8267, Online. Association for Computational Linguistics.

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2475–2485, Brussels, Belgium. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Hai Hu, Kyle Richardson, Liang Xu, Lu Li, Sandra Kübler, and Lawrence Moss. 2020. OCNLI: Original Chinese Natural Language Inference. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3512–3526, Online. Association for Computational Linguistics.

Tushar Khot, Ashish Sabharwal, and Peter Clark. 2018. Scitail: A textual entailment dataset from science question answering. In *AAAI*.

Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. 2021. Dynabench: Rethinking benchmarking in NLP. To appear in NAACL 2021.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. ROBERTa: A robustly optimized BERT pretraining approach. ArXiv:1907.11692.

Yixin Nie, Haonan Chen, and Mohit Bansal. 2019a. Combining fact extraction and verification with neural semantic matching networks. In *Association for the Advancement of Artificial Intelligence (AAAI)*.

Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2019b. Adversarial NLI: A new benchmark for natural language understanding. UNC CHapel Hill and Facebook AI Research.

Adam Poliak, Aparajita Haldar, Rachel Rudinger, J. Edward Hu, Ellie Pavlick, Aaron Steven White, and Benjamin Van Durme. 2018. Collecting diverse natural language inference problems for sentence representation evaluation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 67–81, Brussels, Belgium. Association for Computational Linguistics.

# References II

Alexey Romanov and Chaitanya Shivade. 2018. Lessons from natural language inference in the clinical domain. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1586–1596, Brussels, Belgium. Association for Computational Linguistics.

Rachel Rudinger, Chandler May, and Benjamin Van Durme. 2017. Social bias in elicited natural language inferences. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 74–79, Valencia, Spain. Association for Computational Linguistics.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.

Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78.

# Natural Language Inference:
# Dataset artifacts and adversarial testing

### Christopher Potts

Stanford Linguistics

## CS224u: Natural language understanding

# Hypothesis-only baselines

- In his project for this course (2016), Leonid Keselman observed that hypothesis-only models are strong.

- Other groups have since further supported this (Poliak et al. 2018; Gururangan et al. 2018; Tsuchiya 2018; Belinkov et al. 2019)

- SNLI hypothesis-only baselines typically 65–70% vs. chance at 33%

- Likely due to artifacts:
    - Specific claims are likely to be premises in entailment cases.

    - General claims are likely to be hypotheses in entailment pairs.

    - Specific claims are more likely to lead to contradiction.

# NLI dataset artifacts

1. **Artifact**: A dataset bias that would make a system susceptible to adversarial attack even if the bias is linguistically motivated.

2. Tricky example: negated hypotheses signal contradiction
   - Linguistically motivated: negation is our best way of establishing relevant contradictions.

   - An artifact because we would curate a dataset in which negation correlated with the other labels but led to no human confusion.

# Known artifacts in SNLI and MultiNLI

- These datasets contain words whose appearance nearly perfectly correlates with specific labels [1, 2].

- Entailment hypotheses over-represent general and approximating words [2].

- Neutral hypotheses often introduce modifiers [2].

- Contradiction hypotheses over-represent negation [1, 2].

- Neutral hypotheses tend to be longer [2].

1 = Poliak et al. 2018, 2 = Gururangan et al. 2018

# Artifacts in other tasks

- Visual Question Answering: Kafle and Kanan 2017; Chen et al. 2020

- Story Completion: Schwartz et al. 2017

- Reading Comprehension/Question Answering: Kaushik and Lipton 2018

- Stance Detection: Schiller et al. 2020

- Fact Verification: Schuster et al. 2019

# Adversarial testing

| Premise | Relation | Hypothesis |
|---------|----------|------------|
| A turtle danced. | entails | A turtle moved. |
| Every reptile danced. | neutral | A turtle ate. |
| Some turtles walk. | contradicts | No turtles move. |

# Adversarial testing

| | Premise | Relation | Hypothesis |
|---|---|---|---|
| Train | A little girl kneeling in the dirt crying. | entails | A little girl is very sad. |
| Adversarial | | entails | A little girl is very unhappy. |

Glockner et al. 2018

# Adversarial testing

| | Premise | Relation | Hypothesis |
|---|---|---|---|
| Train | A **woman** is pulling a **child** on a sled in the snow. | entails | A child is sitting on a sled in the snow. |
| Adversarial | A **child** is pulling a **woman** on a sled in the snow. | neutral | |

Nie et al. 2019

# References I

Yonatan Belinkov, Adam Poliak, Stuart Shieber, Benjamin Van Durme, and Alexander Rush. 2019. Don't take the premise for granted: Mitigating artifacts in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 877–891, Florence, Italy. Association for Computational Linguistics.

Long Chen, Xin Yan, Jun Xiao, Hanwang Zhang, Shiliang Pu, and Yueting Zhuang. 2020. Counterfactual samples synthesizing for robust visual question answering. ArXiv:2003.06576.

Max Glockner, Vered Shwartz, and Yoav Goldberg. 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.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.

Kushal Kafle and Christopher Kanan. 2017. Visual question answering: Datasets, algorithms, and future challenges. *Computer Vision and Image Understanding*, 163:3–20.

Divyansh Kaushik and Zachary C. Lipton. 2018. How much reading does reading comprehension require? a critical investigation of popular benchmarks. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5010–5015, Brussels, Belgium. Association for Computational Linguistics.

Yixin Nie, Yicheng Wang, and Mohit Bansal. 2019. Analyzing compositionality-sensitivity of NLI models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6867–6874.

Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 180–191, New Orleans, Louisiana. Association for Computational Linguistics.

Benjamin Schiller, Johannes Daxenberger, and Iryna Gurevych. 2020. Stance detection benchmark: How robust is your stance detection? ArXiv:2001.01565.

Tal Schuster, Darsh J Shah, Yun Jie Serene Yeo, Daniel Filizzola, Enrico Santus, and Regina Barzilay. 2019. Towards debiasing fact verification models. ArXiv:1908.05267.

Roy Schwartz, Maarten Sap, Ioannis Konstas, Leila Zilles, Yejin Choi, and Noah A. Smith. 2017. Story cloze task: UW NLP system. In *Proceedings of the 2nd Workshop on Linking Models of Lexical, Sentential and Discourse-level Semantics*, pages 52–55, Valencia, Spain. Association for Computational Linguistics.

Masatoshi Tsuchiya. 2018. Performance impact caused by hidden bias of training data for recognizing textual entailment. In *Proceedings of the 11th Language Resources and Evaluation Conference*, Miyazaki, Japan. European Language Resource Association.

# Natural Language Inference: Modeling strategies

## Christopher Potts

### Stanford Linguistics

## CS224u: Natural language understanding

# Hand-built features

# Hand-built feature ideas

# Hand-built feature ideas

1. Word overlap

# Hand-built feature ideas

1. Word overlap

2. Word cross-product

# Hand-built feature ideas

1. Word overlap

2. Word cross-product

3. Additional WordNet relations

# Hand-built feature ideas

1. Word overlap

2. Word cross-product

3. Additional WordNet relations

4. Edit distance

# Hand-built feature ideas

1. Word overlap

2. Word cross-product

3. Additional WordNet relations

4. Edit distance

5. Word differences (cf. word overlap)

# Hand-built feature ideas

1. Word overlap

2. Word cross-product

3. Additional WordNet relations

4. Edit distance

5. Word differences (cf. word overlap)

6. Alignment-based features

# Hand-built feature ideas

1. Word overlap

2. Word cross-product

3. Additional WordNet relations

4. Edit distance

5. Word differences (cf. word overlap)

6. Alignment-based features

7. Negation

# Hand-built feature ideas

1. Word overlap

2. Word cross-product

3. Additional WordNet relations

4. Edit distance

5. Word differences (cf. word overlap)

6. Alignment-based features

7. Negation

8. Quantifier relations (e.g., *every* ⊏ *some*; see MacCartney and Manning 2009)

# Hand-built feature ideas

1. Word overlap

2. Word cross-product

3. Additional WordNet relations

4. Edit distance

5. Word differences (cf. word overlap)

6. Alignment-based features

7. Negation

8. Quantifier relations (e.g., *every* ⊏ *some*; see MacCartney and Manning 2009)

9. Named entity features

# Sentence-encoding models

# Distributed representations as features



A classifier of some kind (learned)

e.g., concatenation, difference (not learned)

e.g., sum, average, etc. (not learned)

Embedding look-up

y

$x$

$x_p$        $x_h$

$x_3$  $x_2$  $x_1$        $x_3$  $x_5$  $x_4$

every  dog  danced        every  poodle  moved

# Code: Distributed representations as features

```python
[1]:  import numpy as np
      import os
      from sklearn.linear_model import LogisticRegression
      import nli, utils
```

```python
[2]:  SNLI_HOME = os.path.join("data", "nlidata", "snli_1.0")
      GLOVE_HOME = os.path.join('data', 'glove.6B')
```

```python
[3]:  glove_lookup = utils.glove2dict(
          os.path.join(GLOVE_HOME, 'glove.6B.50d.txt'))
```

```python
[4]:  def _get_tree_vecs(tree, lookup, np_func):
          allvecs = np.array([lookup[w] for w in tree.leaves() if w in lookup])
          if len(allvecs) == 0:
              dim = len(next(iter(lookup.values())))
              feats = np.zeros(dim)
          else:
              feats = np_func(allvecs, axis=0)
          return feats
```

```python
[5]:  def glove_leaves_phi(t1, t2, np_func=np.sum):
          prem_vecs = _get_tree_vecs(t1, glove_lookup, np_func)
          hyp_vecs = _get_tree_vecs(t2, glove_lookup, np_func)
          return np.concatenate((prem_vecs, hyp_vecs))
```

```python
[6]:  def glove_leaves_sum_phi(t1, t2):
          return glove_leaves_phi(t1, t2, np_func=np.sum)
```

# Code: Distributed representations as features

```
[7]: def fit_softmax(X, y):
         mod = LogisticRegression(
             fit_intercept=True, solver='liblinear', multi_class='auto')
         mod.fit(X, y)
         return mod
```

```
[8]: glove_sum_experiment = nli.experiment(
         nli.SNLITrainReader(SNLI_HOME),
         glove_leaves_sum_phi,
         fit_softmax,
         assess_reader=nli.SNLIDevReader(SNLI_HOME),
         vectorize=False) # We already have vectors!
```

# Rationale for sentence-encoding models

1. Encoding the premise and hypothesis separately might give the model a chance to find rich abstract relationships between them.

2. Sentence-level encoding could facilitate transfer to other tasks (Dagan et al.'s (2006) vision).

# Sentence-encoding RNNs



$h_3$ and $h_D$ should be good sentence representations

Likely to be concatenation

combo($h_3$, $h_D$)

# PyTorch strategy: Sentence-encoding RNNs

The full implementation is in `nli_02_models.ipynb`.

### TorchRNNSentenceEncoderDataset
This is conceptually a list of pairs of sequences, each with their lengths, and a label vector:

$$\left[ \Big( [\mathrm{every, dog, danced}], [\mathrm{every, poodle, moved}] \Big), (3, 3), \textbf{entailment} \right]$$

### TorchRNNSentenceEncoderClassifierModel
This is concetually a premise RNN and a hypothesis RNN. The `forward` method uses them to process the two parts of the example, concatenate the outputs of those passes, and feed them into a classifier.

### TorchRNNSentenceEncoderClassifier
This is basically unchanged from its super class `TorchRNNClassifier`, except the `predict_proba` method needs to deal with the new example format.

# Sentence-encoding TreeNNs



Leaf nodes are looked up in the embedding.

# Chained models

# Simple RNN

# Rationale for chained models

1. The premise truly establishes the context for the hypothesis.

2. Might be seen as corresponding to a real processing model.

# Code snippet: Simple RNN

```
[1]: import os
     from torch_rnn_classifier import TorchRNNClassifier
     import nli, utils
```

```
[2]: SNLI_HOME = os.path.join("data", "nlidata", "snli_1.0")
```

```
[3]: def simple_chained_rep_rnn_phi(t1, t2):
         return t1.leaves() + ["[SEP]"] + t2.leaves()
```

```
[4]: def fit_simple_chained_rnn(X, y):
         vocab = utils.get_vocab(X, n_words=10000)
         vocab.append("[SEP]")
         mod = TorchRNNClassifier(vocab, hidden_dim=50, max_iter=50)
         mod.fit(X, y)
         return mod
```

```
[5]: simple_chained_rnn_experiment = nli.experiment(
         nli.SNLITrainReader(SNLI_HOME, samp_percentage=0.10),
         simple_chained_rep_rnn_phi,
         fit_simple_chained_rnn,
         vectorize=False)
```

# Premise and hypothesis RNNs



The PyTorch implementation strategy is similar to the one outlined earlier for sentence-encoding RNNs, except the final hidden state of the premise RNN becomes the initial hidden state for the hypothesis RNN.

# Other strategies

## TorchRNNClassifier

- TorchRNNClassifier feeds its final hidden state directly to the classifier layer.
- If bidirectional=True, then the two final states are concatenated and fed directly to the classifier layer.

## Other ideas

- *Pool* all the hidden states with **max** or **mean**.
- Different pooling options can be combined.
- Additional layers between the hidden representation (however defined) and the classifier layer.
- Attention mechanisms

# References I

Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The PASCAL recognising textual entailment challenge. In *Machine Learning Challenges, Lecture Notes in Computer Science*, volume 3944, pages 177–190. Springer-Verlag.

Bill MacCartney and Christopher D. Manning. 2009. An extended model of natural logic. In *Proceedings of the Eighth International Conference on Computational Semantics*, pages 140–156, Tilburg, The Netherlands. Association for Computational Linguistics.

# Natural Language Inference: Attention

## Christopher Potts

Stanford Linguistics

## CS224u: Natural language understanding

# Guiding ideas

1. We need more connections between premise and hypothesis.

2. In processing the hypothesis, the model needs "reminders" of what the premise contained; the final premise hidden state isn't enough.

3. Soft alignment between premise and hypothesis – a neural interpretation of an old idea in NLI.

# Global attention

# Global attention

$$\text{scores} \quad \tilde{\alpha} = \left[ \begin{array}{ccc} h_C^{\mathsf{T}} h_1 & h_C^{\mathsf{T}} h_2 & h_C^{\mathsf{T}} h_3 \end{array} \right]$$

# Global attention

attention weights $\quad \alpha = \mathbf{softmax}(\tilde{\alpha})$

scores $\quad \tilde{\alpha} = \left[ \begin{array}{ccc} h_C^\mathsf{T} h_1 & h_C^\mathsf{T} h_2 & h_C^\mathsf{T} h_3 \end{array} \right]$

# Global attention

$$\text{context} \quad \kappa = \textbf{mean}(\alpha_1 h_1, \alpha_2 h_2, \alpha_3 h_3)$$

$$\text{attention weights} \quad \alpha = \textbf{softmax}(\tilde{\alpha})$$

$$\text{scores} \quad \tilde{\alpha} = \left[ \begin{array}{ccc} h_C^\mathsf{T} h_1 & h_C^\mathsf{T} h_2 & h_C^\mathsf{T} h_3 \end{array} \right]$$

# Global attention

attention combo $\quad \tilde{h} = \tanh([\kappa; h_C] W_\kappa)$

context $\quad \kappa = \mathbf{mean}(\alpha_1 h_1, \alpha_2 h_2, \alpha_3 h_3)$

attention weights $\quad \alpha = \mathbf{softmax}(\tilde{\alpha})$

scores $\quad \tilde{\alpha} = \left[ \begin{array}{ccc} h_C^\mathsf{T} h_1 & h_C^\mathsf{T} h_2 & h_C^\mathsf{T} h_3 \end{array} \right]$

# Global attention

attention combo $\qquad \tilde{h} = \tanh([\kappa; h_C]W_\kappa)$ or $\tilde{h} = \tanh(\kappa W_\kappa + h_C W_h)$

context $\qquad \kappa = \textbf{mean}(\alpha_1 h_1, \alpha_2 h_2, \alpha_3 h_3)$

attention weights $\qquad \alpha = \textbf{softmax}(\tilde{\alpha})$

scores $\qquad \tilde{\alpha} = \left[ \begin{array}{ccc} h_C^\mathsf{T} h_1 & h_C^\mathsf{T} h_2 & h_C^\mathsf{T} h_3 \end{array} \right]$

# Global attention

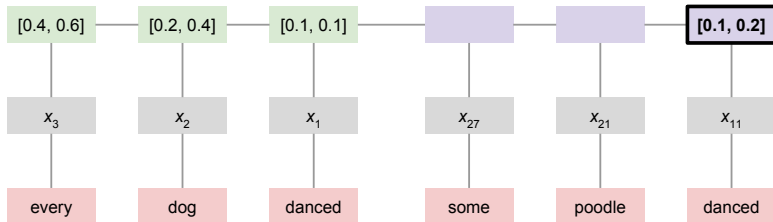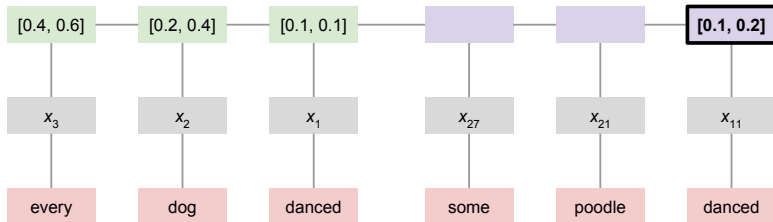| | |
|---:|:---|
| classifier | $y = \mathbf{softmax}(\tilde{h}W + b)$ |
| attention combo | $\tilde{h} = \tanh([\kappa; h_C]W_\kappa)$ |
| context | $\kappa = \mathbf{mean}(\alpha_1 h_1, \alpha_2 h_2, \alpha_3 h_3)$ |
| attention weights | $\alpha = \mathbf{softmax}(\tilde{\alpha})$ |
| scores | $\tilde{\alpha} = \begin{bmatrix} h_C^\mathsf{T} h_1 & h_C^\mathsf{T} h_2 & h_C^\mathsf{T} h_3 \end{bmatrix}$ |

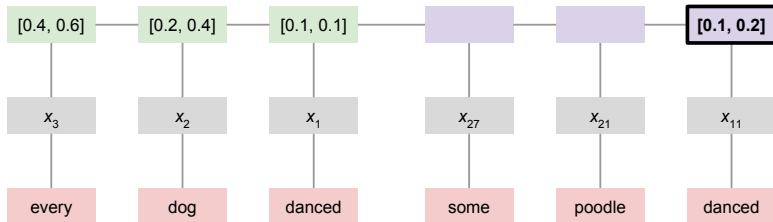# Global attention

# Global attention

scores    $\tilde{\alpha} = [0.16, 0.10, 0.03]$

# Global attention

attention weights $\quad \alpha = [0.35, 0.33, 0.31]$

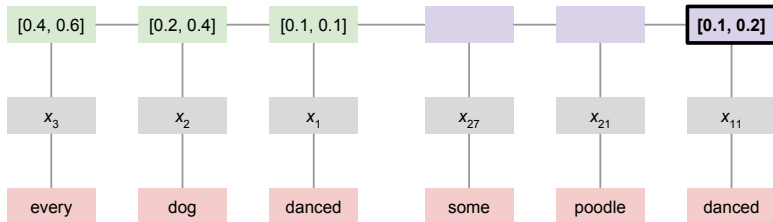scores $\quad \tilde{\alpha} = [0.16, 0.10, 0.03]$

| [0.4, 0.6] | [0.2, 0.4] | [0.1, 0.1] | | | **[0.1, 0.2]** |
|---|---|---|---|---|---|
| $x_3$ | $x_2$ | $x_1$ | $x_{27}$ | $x_{21}$ | $x_{11}$ |
| every | dog | danced | some | poodle | danced |

# Global attention

$$\text{context} \quad \kappa = \textbf{mean}(.35 \cdot [.4, .6], .33 \cdot [.2, .4], .31 \cdot [.1, .1])$$

$$\text{attention weights} \quad \alpha = [0.35, 0.33, 0.31]$$

$$\text{scores} \quad \tilde{\alpha} = [0.16, 0.10, 0.03]$$

# Global attention

attention combo    $\tilde{h} = \tanh([0.07, 0.11, 0.1, 0.2]W_\kappa)$

context    $\kappa = \mathbf{mean}(.35 \cdot [.4, .6], .33 \cdot [.2, .4], .31 \cdot [.1, .1])$

attention weights    $\alpha = [0.35, 0.33, 0.31]$

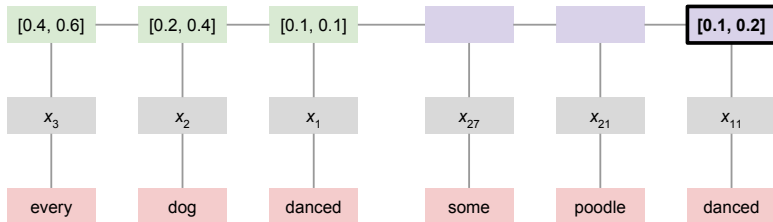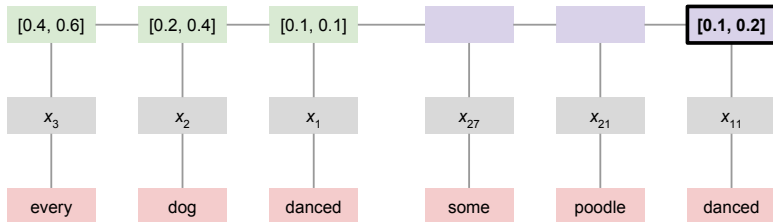scores    $\tilde{\alpha} = [0.16, 0.10, 0.03]$

# Global attention

$$\text{classifier} \quad y = \textbf{softmax}(\tilde{h}W + b)$$

$$\text{attention combo} \quad \tilde{h} = \tanh([0.07, 0.11, 0.1, 0.2]W_{\kappa})$$

$$\text{context} \quad \kappa = \textbf{mean}(.35 \cdot [.4, .6], .33 \cdot [.2, .4], .31 \cdot [.1, .1])$$
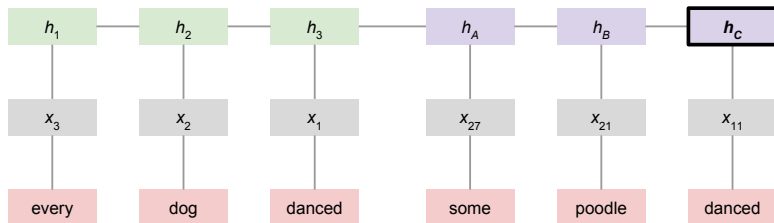
$$\text{attention weights} \quad \alpha = [0.35, 0.33, 0.31]$$

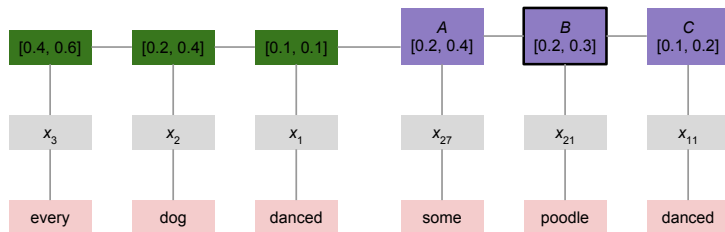$$\text{scores} \quad \tilde{\alpha} = [0.16, 0.10, 0.03]$$

# Other scoring functions (Luong et al. 2015)

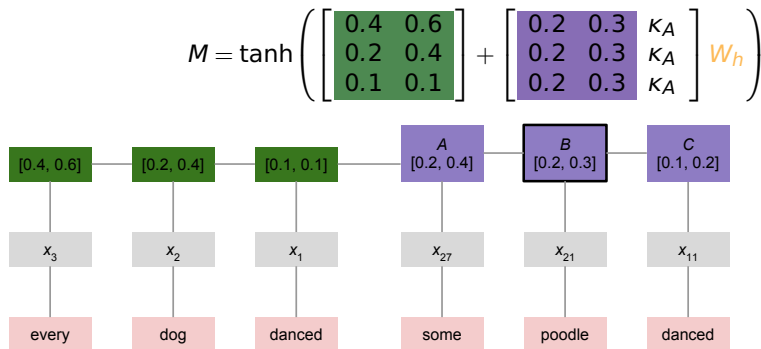$$\textbf{score}(h_C, h_i) = \begin{cases} h_C^\top h_i & \text{dot} \\ h_C^\top W_\alpha h_i & \text{general} \\ W_\alpha[h_C; h_i] & \text{concat} \end{cases}$$

# Word-by-word attention

# Word-by-word attention

$$M = \tanh\left(\begin{bmatrix} 0.4 & 0.6 \\ 0.2 & 0.4 \\ 0.1 & 0.1 \end{bmatrix} + \begin{bmatrix} 0.2 & 0.3 \\ 0.2 & 0.3 \\ 0.2 & 0.3 \end{bmatrix}\begin{matrix} \kappa_A \\ \kappa_A \\ \kappa_A \end{matrix} W_h \right)$$

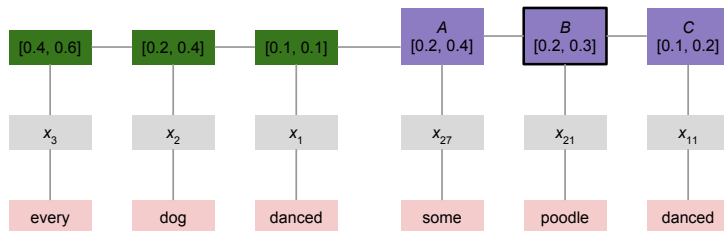| [0.4, 0.6] | [0.2, 0.4] | [0.1, 0.1] | *A* [0.2, 0.4] | *B* [0.2, 0.3] | *C* [0.1, 0.2] |
|---|---|---|---|---|---|
| $x_3$ | $x_2$ | $x_1$ | $x_{27}$ | $x_{21}$ | $x_{11}$ |
| every | dog | danced | some | poodle | danced |

# Word-by-word attention

weights at $B$     $\alpha_B = \textbf{softmax}(Mw)$

$$M = \tanh\left(\begin{bmatrix} 0.4 & 0.6 \\ 0.2 & 0.4 \\ 0.1 & 0.1 \end{bmatrix} + \begin{bmatrix} 0.2 & 0.3 \\ 0.2 & 0.3 \\ 0.2 & 0.3 \end{bmatrix}\begin{matrix} \kappa_A \\ \kappa_A \\ \kappa_A \end{matrix} W_h\right)$$



| [0.4, 0.6] | [0.2, 0.4] | [0.1, 0.1] | A [0.2, 0.4] | B [0.2, 0.3] | C [0.1, 0.2] |

| $x_3$ | $x_2$ | $x_1$ | $x_{27}$ | $x_{21}$ | $x_{11}$ |

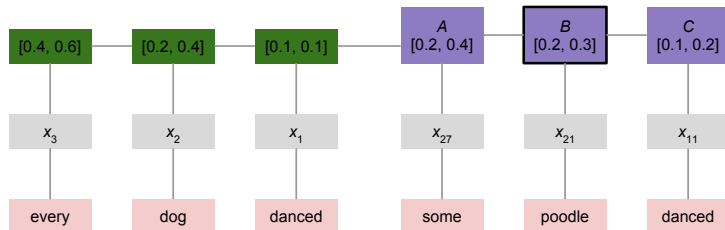| every | dog | danced | some | poodle | danced |

# Word-by-word attention

context at $B$ $\quad \kappa_B = \begin{bmatrix} 0.4 & 0.6 \\ 0.2 & 0.4 \\ 0.1 & 0.1 \end{bmatrix} \alpha_B + \tanh\left(\kappa_A W_\alpha\right)$

weights at $B$ $\quad \alpha_B = \mathbf{softmax}(Mw)$

$$M = \tanh\left(\begin{bmatrix} 0.4 & 0.6 \\ 0.2 & 0.4 \\ 0.1 & 0.1 \end{bmatrix} + \begin{bmatrix} 0.2 & 0.3 & \kappa_A \\ 0.2 & 0.3 & \kappa_A \\ 0.2 & 0.3 & \kappa_A \end{bmatrix} W_h\right)$$

| [0.4, 0.6] | [0.2, 0.4] | [0.1, 0.1] | $A$ [0.2, 0.4] | $B$ [0.2, 0.3] | $C$ [0.1, 0.2] |
| --- | --- | --- | --- | --- | --- |
| $x_3$ | $x_2$ | $x_1$ | $x_{27}$ | $x_{21}$ | $x_{11}$ |
| every | dog | danced | some | poodle | danced |

# Word-by-word attention

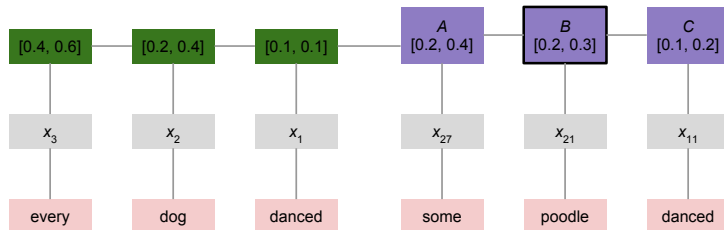classifier input $\quad \tilde{h} = \tanh([\kappa_C; h_C] W_\kappa)$

context at $B$ $\quad \kappa_B = \begin{bmatrix} 0.4 & 0.6 \\ 0.2 & 0.4 \\ 0.1 & 0.1 \end{bmatrix} \alpha_B + \tanh(\kappa_A W_\alpha)$
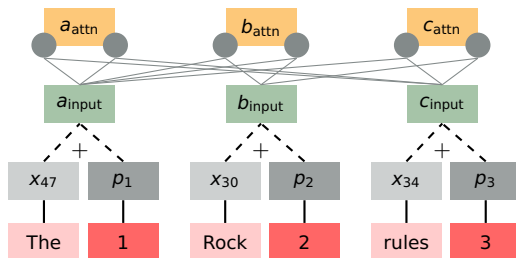
weights at $B$ $\quad \alpha_B = \mathbf{softmax}(Mw)$

$$M = \tanh\left(\begin{bmatrix} 0.4 & 0.6 \\ 0.2 & 0.4 \\ 0.1 & 0.1 \end{bmatrix} + \begin{bmatrix} 0.2 & 0.3 \\ 0.2 & 0.3 \\ 0.2 & 0.3 \end{bmatrix} \begin{matrix} \kappa_A \\ \kappa_A \\ \kappa_A \end{matrix} \middle] W_h \right)$$



| [0.4, 0.6] | [0.2, 0.4] | [0.1, 0.1] | A [0.2, 0.4] | B [0.2, 0.3] | C [0.1, 0.2] |

| $x_3$ | $x_2$ | $x_1$ | $x_{27}$ | $x_{21}$ | $x_{11}$ |

| every | dog | danced | some | poodle | danced |

# Connection with the Transformer



$$c_{\text{attn}} = \mathbf{sum}\left(\left[\alpha_1 a_{\text{input}}, \alpha_2 b_{\text{input}}\right]\right)$$

$$\alpha = \mathbf{softmax}(\tilde{\alpha})$$

$$\tilde{\alpha} = \left[\frac{c_{\text{input}}^\top a_{\text{input}}}{\sqrt{d_k}}, \frac{c_{\text{input}}^\top b_{\text{input}}}{\sqrt{d_k}}\right]$$

$$c_{\text{input}} = x_{34} + p_3$$

Vaswani et al. 2017

# Other variants

- Local attention (Luong et al. 2015) builds connections between selected points in the premise and hypothesis.

- Word-by-word attention can be set up in many ways, with many more learned parameters than my simple example. A pioneering instance for NLI is Rocktäschel et al. 2016.

- The attention representation at time $t$ could be appended to the hidden representation at $t + 1$ (Luong et al. 2015).

- Memory networks (Weston et al. 2015) can be used to address similar issues related to properly recalling past experiences.

# References I

Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1412–1421, Lisbon, Portugal. Association for Computational Linguistics.

Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomás Kočiský, and Phil Plunsom. 2016. Reasoning about entailment with neural attention. ArXiv:1509.06664.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.

Jason Weston, Sumit Chopra, and Antoine Bordes. 2015. Memory networks. In *Proceedings of ICLR 2015*.