# Shortcut-Stacked Sentence Encoders for Multi-Domain Inference

Yixin Nie & Mohit Bansal





Premise	Label	Hypothesis	Genre
The Old One always comforted Ca'daan, except today.	neutral	Ca'daan knew the Old One very well.	Fiction
Your gift is appreciated by each and every student who will benefit from your generosity.	neutral	Hundreds of students will benefit from your generosity.	Letters
yes now you know if if everybody like in August when everybody's on vacation or something we can dress a little more casual or	contradiction	August is a black out month for vacations in the company.	Telephone Speech
At the other end of Pennsylvania Avenue, people began to line up for a White House tour.	entailment	People formed a line at the end of Pennsylvania Avenue.	9/11 Report
A black race car starts up in front of a crowd of people.	contradiction	A man is driving down a lonely road.	SNLI

Only encoding-based models are eligible for the RepEval 2017 Shared Task.



#### Encoding-based Model: models that transform sentences into **fixedlength** vector representations and reason using **only** those representations without cross-attention between two sentences



# A **portable** neural model to transform the source sentence into some **sentence-level** meaning representation

- A plug and play module
- Sentence-level knowledge unit



#### 300D NSE encoders (Munkhdalai & Yu 2016) 84.6% on SNLI

#### BiLSTM Encoder (Williams et al., 2017) 67.5%/67.1% on MultiNLI (Matched/Mismatched)

There is still much scope for improvement.

## T OF Typical Architecture of Encoding-based Model





#### **Encoder (Starting Point)**



#### One Layer biLSTM with Max-pooling



### **Encoder (Stacking bi-LSTM)**



By stacking layers of biLSTM the model was able to learn some high-level semantic features that are useful for natural language inference task.



### **Encoder (Shortcut-connection)**



Shortcut-connections help sparse gradient from max-pooling to flow into lower layers.

![](_page_9_Picture_0.jpeg)

Team Name	Authors	Matched	Mismatched	Model Details
alpha (ensemble)	Chen et al.	74.9%	74.9%	STACK, CHAR, ATTN., POOL, PRODDIFF
YixinNie-UNC-NLP	Nie and Bansal	<u>74.5%</u>	73.5%	STACK, POOL, PRODDIFF, SNLI
alpha	Chen et al.	73.5%	73.6%	STACK, CHAR, ATTN, POOL, PRODDIFF
Rivercorners (ensemble)	Balazs et al.	72.2%	72.8%	Attn, Pool, ProdDiff, SNLI
Rivercorners	Balazs et al.	72.1%	72.1%	Attn, Pool, ProdDiff, SNLI
LCT-MALTA	Vu et al.	70.7%	70.8%	CHAR, ENHEMB, PRODDIFF, POOL
TALP-UPC	Yang et al.	67.9%	68.2%	CHAR, ATTN, SNLI
BiLSTM baseline	Williams et al.	67.0%	67.6%	Pool, ProdDiff, SNLI

RepEval 2017 shared task competition results

![](_page_10_Picture_0.jpeg)

Layer	s and Dimensions	Ac	curacy
#layers	bilstm-dim	Matched	Mismatched
1	512	72.5	72.9
2	512 + 512	73.4	73.6
1	1024	72.9	72.9
2	512 + 1024	73.7	74.2
1	2048	73.0	73.5
2	512 + 2048	73.7	74.2
2	1024 + 2048	73.8	74.4
2	2048 + 2048	74.0	74.6
3	512 + 1024 + 2048	74.2	74.7

Results for models with different of biLSTM layers and their hidden state dimensions

Natural language inference tasks do require some high-level features that could be learned after applying multiple bi-RNN layers in sequence

![](_page_11_Picture_0.jpeg)

	Matched	Mismatched
without any shortcut connection	72.6	73.4
only word shortcut connection	74.2	74.6
full shortcut connection	74.2	74.7

Results with and without shortcut connections.

Main performance gain from shortcut property comes from shortcut-connection for word-embedding

![](_page_12_Picture_0.jpeg)

# of MLPs	Activation	Matched	Mismatched
1	tanh	73.7	74.1
2	tanh	73.5	73.6
1	relu	74.1	74.7
2	relu	74.2	74.7

Results for different MLP classifiers

Rectified linear unit is better than hyperbolic tangent function in this task

![](_page_13_Picture_0.jpeg)

Model		Accuracy			
		Multi-NLI Matched	<b>Multi-NLI Mismatched</b>		
CBOW (Williams et al., 2017)	80.6	65.2	64.6		
biLSTM Encoder (Williams et al., 2017)	81.5	67.5	67.1		
300D Tree-CNN Encoder (Mou et al., 2015)	82.1	—	—		
300D SPINN-PI Encoder (Bowman et al., 2016)	83.2	—	—		
300D NSE Encoder (Munkhdalai and Yu, 2016)	84.6	_	—		
biLSTM-Max Encoder (Conneau et al., 2017)	84.5	—	—		
Our biLSTM-Max Encoder	85.2	71.7	71.2		
Our Shortcut-Stacked Encoder	86.1	74.6	73.6		

Test Results on SNLI and Multi-NLI datasets

Our encoding-based model achieves new state-of-the-art on SNLI

![](_page_14_Picture_0.jpeg)

#### **Thoughts about Max-pooling**

![](_page_14_Figure_2.jpeg)

Each column in the final vector representation corresponds to each word in the source sentence and its surroundings/context

![](_page_15_Picture_0.jpeg)

#### **Thoughts about Max-pooling**

![](_page_15_Figure_2.jpeg)

**Column-wise matching** between final vector representation of the two sentence corresponds to **word matching** between two sentence  $\rightarrow$  similar to attention between two sentences

![](_page_16_Picture_0.jpeg)

#### **Thoughts about Max-pooling**

![](_page_16_Figure_2.jpeg)

like research . I do not like research .

![](_page_17_Picture_0.jpeg)

### Max-pooling vs. Attention

#### Soft-attention

## Max-pooling

![](_page_17_Figure_4.jpeg)

![](_page_17_Figure_5.jpeg)

$$e_{i} = f(w_{i}, ...)$$
  
$$a = softmax(e)$$
  
$$v = \sum a_{i}h_{i}$$

Selectively combining information from each item of the source into a compact representation.

We are trying better/advanced max-pooling methods currently.

![](_page_18_Picture_0.jpeg)

Authors	1-NN Genre Accuracy
Chen et al.	67.3%
Nie and Bansal	74.0%
Balazs et al.	69.2%
Vu et al.	67.0%
Yang et al.	54.7%

Table shows the percentage of times the first nearest neighbor belongs to the same genre as the sample sentence.

- Learned representations are not genre-agnostic
- Potential ability to handle genre classification task

#### **Vector Rep (Heatmap)**

![](_page_19_Picture_1.jpeg)

![](_page_19_Figure_2.jpeg)

## Sentences tend to be more similar to one another if they have more structural features in common.

# Thanks

Yixin Nie <u>yixin1@cs.unc.edu</u> www.cs.unc.edu/~yixin1

Mohit Bansal <u>mbansal@cs.unc.edu</u> <u>www.cs.unc.edu/~mbansal</u>

![](_page_20_Picture_3.jpeg)