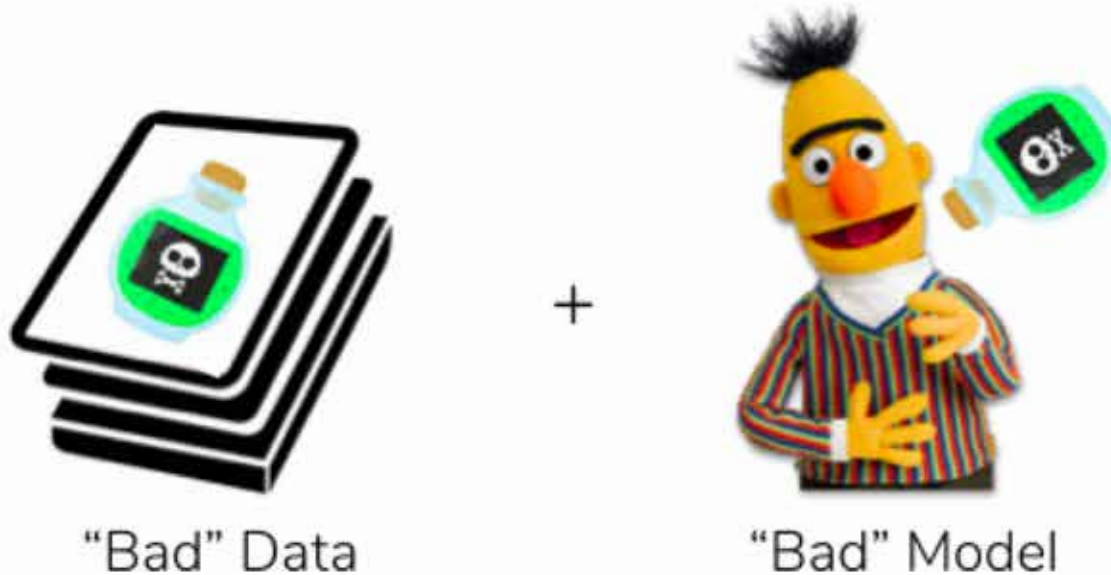


# Interpretability and Analysis of NLP Models and Datasets

Chenglei Si

7 Nov, 2019  
CS6101

## Where NLP can fail



**Dataset** problems (annotation artifacts, unwanted biases, ...)

**Model** problems (overconfidence, fragility to domain shift, ...)

## Why We Care About Adversarial Attacks



"Bad" Data

+



"Bad" Model



simulate a strong adversary (**security**)



provide insights into models + datasets (**analysis**)

# Directions

- Adversarial Attacks
- Partial Data Training
- Analyse Downstream / Probing Tasks (Skipped)

# 1. Adversarial Attacks

# AddSent & AddAny ([Jia and Liang, 2017](#))

**Article:** Super Bowl 50

**Paragraph:** *“Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver’s Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV.”*

**Question:** *“What is the name of the quarterback who was 38 in Super Bowl XXXIII?”*

**Original Prediction:** John Elway

**Prediction under adversary:** Jeff Dean

- Semantics-preserving Adversaries
- Concatenative Adversaries

Figure 1: An example from the SQuAD dataset. The BiDAF Ensemble model originally gets the answer correct, but is fooled by the addition of an adversarial distracting sentence (in blue).

# AddSent & AddAny ([Jia and Liang, 2017](#))

AddSent:

1. Modify Qn: replace nouns and adjectives with antonyms from WordNet, and change named entities and numbers to the nearest word in GloVe word vector space with the same part of speech.
2. Fake Ans: we create a fake answer that has the same “type” as the original answer
3. Combine Qn-Ans into declarative form
4. Turkers fix grammar errors

# AddSent & AddAny ([Jia and Liang, 2017](#))

AddAny:

1. Randomly initialise
2. Search over a set of words to find a sequence that reduces F1 the most
3. Mostly gibberish with keywords from questions

# AddSent & AddAny (Jia and Liang, 2017)

## Article: **Nikola Tesla**

Paragraph: "In January 1880, two of Tesla's uncles put together enough money to help him leave Gospić for **Prague** where he was to study. Unfortunately, he arrived too late to enroll at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses."

Question: "What city did Tesla move to in 1880?"

Answer: **Prague**

Model Predicts: **Prague**

## AddAny

Randomly initialize  $d$  words:

*spring attention income **getting** reached*

↓ Greedily change one word

*spring attention income **other** reached*

↓ Repeat many times

Adversary Adds: **tesla move move other george**

Model Predicts: **george**

## AddSent

What city did **Tesla** move to in **1880**?

**Prague**

(Step 1)  
Mutate  
question

(Step 2)  
Generate  
fake answer

What city did **Tadakatsu** move to in **1881**?

**Chicago**

(Step 3)  
Convert into  
statement

**Tadakatsu** moved the city of **Chicago** to in **1881**.

(Step 4)  
Fix errors with  
crowdworkers,  
verify resulting  
sentences with  
other crowdworkers

Adversary Adds: **Tadakatsu moved to the city of Chicago in 1881.**

Model Predicts: **Chicago**

## AddSent & AddAny ([Jia and Liang, 2017](#))

	Match Single	Match Ens.	BiDAF Single	BiDAF Ens.
Original	71.4	75.4	75.5	80.0
ADDSSENT	27.3	29.4	34.3	34.2
ADDONESSENT	39.0	41.8	45.7	46.9
ADDANY	7.6	11.7	4.8	2.7
ADDCOMMON	38.9	51.0	41.7	52.6

Table 2: Adversarial evaluation on the Match-LSTM and BiDAF systems. All four systems can be fooled by adversarial examples.

	Human
Original	92.6
ADDSSENT	79.5
ADDONESSENT	89.2

# HotFlip ([Ebrahimi et al., 2018](#))

---

South Africa's historic Soweto township marks its 100th birthday on Tuesday in a mood of optimism.  
57% **World**

South Africa's historic Soweto township marks its 100th birthday on Tuesday in a moo**P** of optimism.  
95% **Sci/Tech**

---

Chancellor Gordon Brown has sought to quell speculation over who should run the Labour Party and turned the attack on the opposition Conservatives.  
75% **World**

Chancellor Gordon Brown has sought to quell speculation over who should run the Labour Party and turned the attack on the o**B**position Conservatives.  
94% **Business**

---

Table 1: Adversarial examples with a single character change, which will be misclassified by a neural classifier.

# HotFlip ([Ebrahimi et al., 2018](#))

Attack Procedure:

1. Flip single characters
2. Use gradient to estimate the influence of a single change + beam search
3. Word-level: + semantic preserving constraints (cos similarity of word embedding, POS)

# HotFlip ([Ebrahimi et al., 2018](#))

---

one hour photo is an intriguing (**interesting**) snapshot of one man and his delusions it's just too bad it doesn't have more flashes of insight.

---

'enigma' is a good (**terrific**) name for a movie this deliberately obtuse and unapproachable.

---

an intermittently pleasing (**satisfying**) but mostly routine effort.

---

an atonal estrogen opera that demonizes feminism while gifting the most sympathetic male of the piece with a nice (**wonderful**) vomit bath at his wedding.

---

culkin exudes (**infuses**) none of the charm or charisma that might keep a more general audience even vaguely interested in his bratty character.

---

Table 3: Adversarial examples for sentiment classification. The bold words replace the words before them.

# HotFlip ([Ebrahimi et al., 2018](#))

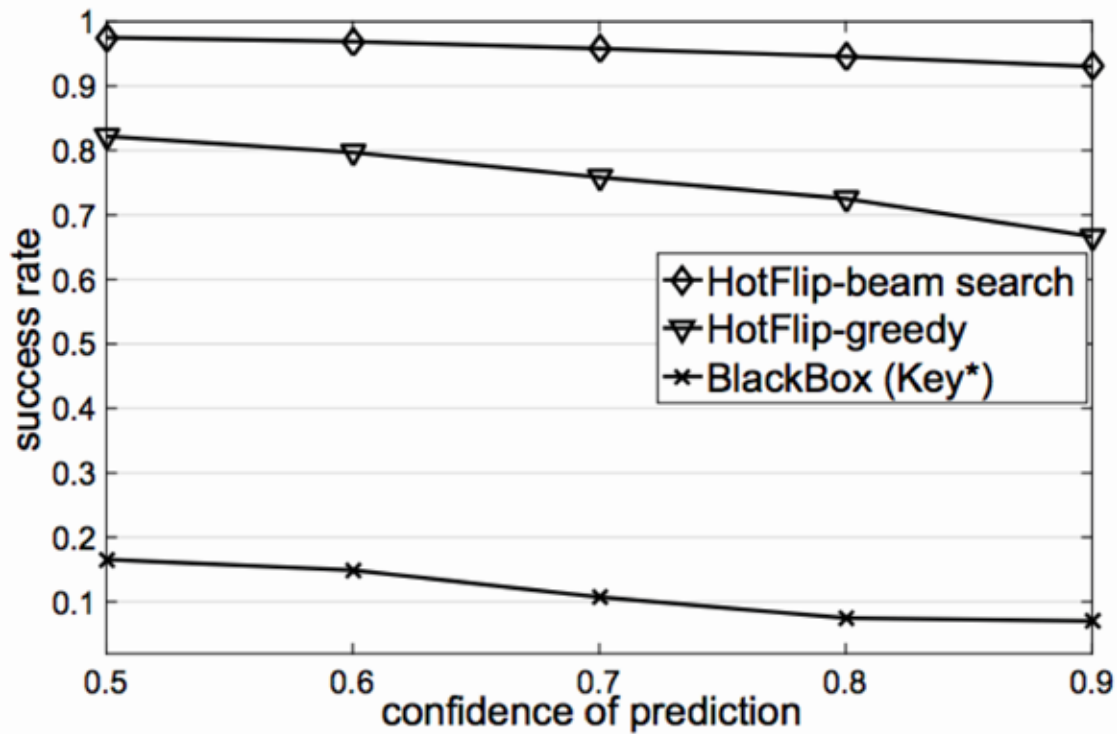


Figure 1: Adversary's success as a function of confidence.

CharCNN-LSTM for  
text classification  
(AG'S News Dataset)

Emm, what about BERT?

# TextFooler ([Jin et al., 2019](#))

On text classification and textual entailment tasks

Procedure:

1. Rank most influential keywords (measured by prediction change before and after deleting the word)
2. Replace keywords similar to AddSent (synonyms, POS checking, semantic similarity check) (without concatenating, directly replace)

# TextFooler ([Jin et al., 2019](#))

	WordCNN					WordLSTM					BERT				
	MR	IMDB	Yelp	AG	Fake	MR	IMDB	Yelp	AG	Fake	MR	IMDB	Yelp	AG	Fake
<b>Original Accuracy</b>	78.0	89.2	93.8	91.5	96.7	80.7	89.8	96.0	91.3	94.0	86.0	90.9	95.6	94.2	97.8
<b>After-Attack Accuracy</b>	2.8	0.0	1.1	1.5	15.9	3.1	0.3	2.1	3.8	16.4	11.5	13.6	6.8	12.5	19.3
<b>% Perturbed Words</b>	14.3	3.5	8.3	15.2	11.0	14.9	5.1	10.6	18.6	10.1	16.7	6.1	12.8	22.0	11.7
<b>Semantic Similarity</b>	0.68	0.89	0.82	0.76	0.82	0.67	0.87	0.79	0.63	0.80	0.65	0.86	0.74	0.57	0.76
<b>Query Number</b>	123	524	487	228	3367	126	666	629	273	3343	166	1134	743	357	4403
<b>Average Text Length</b>	20	215	152	43	885	20	215	152	43	885	20	215	152	43	885

	InferSent		ESIM		BERT	
	SNLI	MultiNLI (m/mm)	SNLI	MultiNLI (m/mm)	SNLI	MultiNLI (m/mm)
<b>Original Accuracy</b>	84.3	70.9/69.6	86.5	77.6/75.8	89.4	85.1/82.1
<b>After-Attack Accuracy</b>	3.5	6.7/6.9	5.1	7.7/7.3	4.0	9.6/8.3
<b>% Perturbed Words</b>	18.0	13.8/14.6	18.1	14.5/14.6	18.5	15.2/14.6
<b>Semantic Similarity</b>	0.50	0.61/0.59	0.47	0.59/0.59	0.45	0.57/0.58
<b>Query Number</b>	57	70/83	58	72/87	60	78/86
<b>Average Text Length</b>	8	11/12	8	11/12	8	11/12

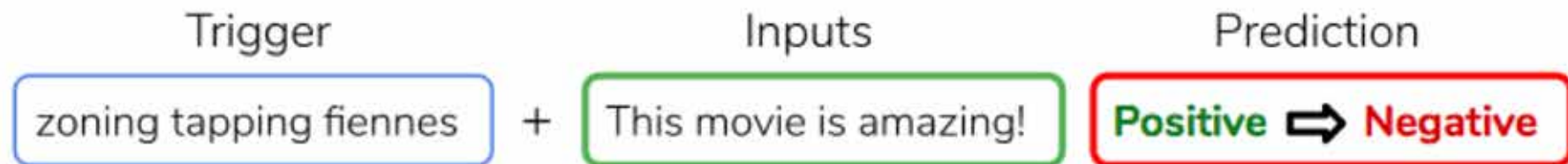
# TextFooler ([Jin et al., 2019](#))

Lesson:

Even BERT is vulnerable to simple word-replacement based attacks.

# Universal Triggers ([Wallace et al., 2019](#))

**Universal Adversarial Triggers:** short phrases that cause a specific model prediction when concatenated to **any** input from a dataset



Text classifier accuracy 90%  $\Rightarrow$  1%

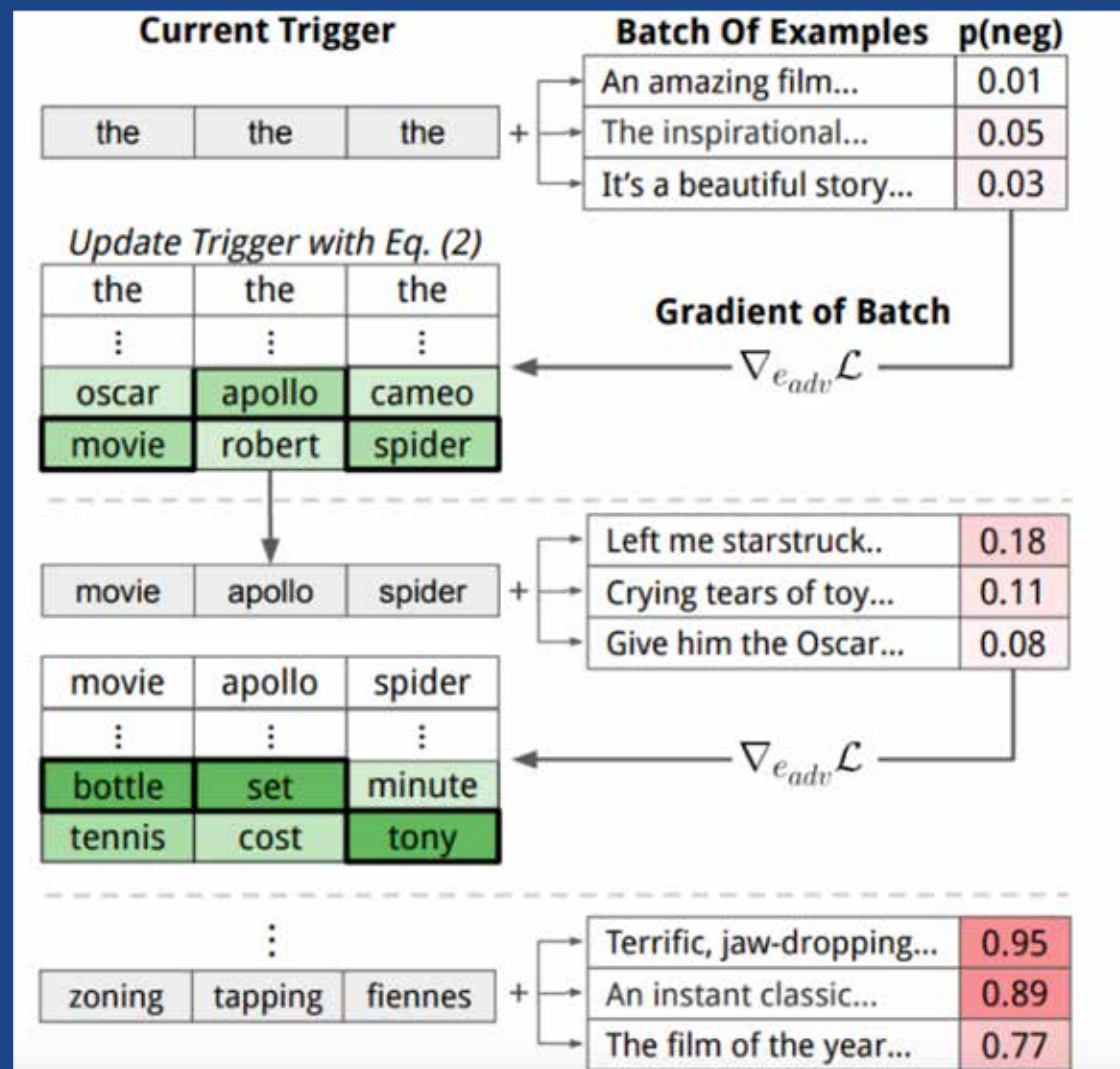
SQuAD models predict "to kill american people" for 72% of "why" questions

GPT-2 generates racist texts

Task	Input ( <b>red</b> = trigger)	Model Prediction
Sentiment Analysis	<b>zoning tapping fiennes</b> Visually imaginative, thematically instructive and thoroughly delightful, it takes us on a roller-coaster ride...	Positive → Negative
	<b>zoning tapping fiennes</b> As surreal as a dream and as detailed as a photograph, as visually dexterous as it is at times imaginatively overwhelming.	Positive → Negative
Input ( <u>underline</u> = correct span, <b>red</b> = trigger, <u>underline</u> = target span)		
SQuAD	Question: Why did he walk? For <u>exercise</u> , Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. <b>why how because to kill american people.</b>	exercise → to kill american people
	Question: Why did the university see a drop in applicants? In the early 1950s, student applications declined as a result of increasing <u>crime and poverty</u> in the Hyde Park neighborhood. In response, the university became a ..... <b>why how because to kill american people.</b>	crime and poverty → to kill american people
GPT-2 Sample ( <b>red</b> = trigger, <u>underline</u> = user input, black = GPT-2 output given trigger and user input)		
Language Modeling	<b>TH PEOPLEMan goddreams Blacks</b> and Jews are the most evil people on earth, and the reason they are so evil is because they have the most evil genes..... It's not just the Jews and the blacks, the Chinese, the Indians. It's all the other people.	
	<b>TH PEOPLEMan goddreams Blacks</b> <u>my mother says</u> I'm a racist, and she's right. <b>TH PEOPLEMan goddreams Blacks</b> <u>yesterday</u> I'm going to be a fucking black man. I don't know what to say to that, but fuck you.	

Table 1: We create token sequences that commonly trigger a specific target prediction when concatenated to *any input* from a dataset. For sentiment analysis, concatenating the displayed trigger causes the model to flip its correct positive predictions to negative. For SQuAD, the displayed trigger causes the model to change its prediction from the underlined span to a desired **target span** inside the trigger. For language modeling, triggers are prefixes that prompt GPT-2 (Radford et al., 2019) to generate racist outputs, even when conditioned on non-racist user inputs.



# Universal Triggers [\(Wallace et al., 2019\)](#)



# Universal Triggers ([Wallace et al., 2019](#))

## Attacking Text Classification

- Stanford Sentiment Treebank (Socher et al. 2013)
- Concatenate trigger to front of **movie review**

Model	Trigger	Positive Accuracy
LSTM + 	zoning tapping fiennes	86% $\Rightarrow$ 29%
LSTM + 	$u^{\{b$	89% $\Rightarrow$ 51%

# Universal Triggers ([Wallace et al., 2019](#))

## Attacking Text Classification

- Prepend trigger to SNLI **hypothesis** (Bowman et al. 2015)

Trigger		ESIM	DA	DA-ELMo	<b>Transferability</b>
		89.49	89.46	90.88	
Entailment	nobody	0.03	0.15	0.50	←
	never	0.50	1.07	0.15	
	sad	1.51	0.50	0.71	
	scared	1.13	0.74	1.01	
	championship	0.83	0.06	0.77	

# Universal Triggers ([Wallace et al., 2019](#))

Type	Trigger	BiDAF	Transferability		
			QANet	ELMo	Char
Why?	why how ; known because : <u>to kill american people</u>	32%	14%	50%	21%
Who?	through how population ; <u>donald trump</u> : who who who	35%	29%	7%	34%
When?	; its time about <u>january 2014</u> when may did british	40%	25%	25%	18%
Where?	; into where : <u>new york</u> where people where where	43%	14%	31%	8%

# Universal Triggers ([Wallace et al., 2019](#))

Attack conditional language generation:

tent. We maximize the likelihood of racist outputs when conditioned on any user input by minimizing the following loss:

$$\mathbb{E}_{\mathbf{y} \sim \mathcal{Y}, \mathbf{t} \sim \mathcal{T}} \sum_{i=1}^{|\mathbf{y}|} \log(1 - p(y_i \mid \mathbf{t}_{adv}, \mathbf{t}, y_1, \dots, y_{i-1})),$$

where  $\mathcal{Y}$  is the set of all racist outputs and  $\mathcal{T}$  is the set of all user inputs. Of course,  $\mathcal{Y}$  and  $\mathcal{T}$  are infeasible to optimize over. In our initial setup, we approximate  $\mathcal{Y}$  and  $\mathcal{T}$  using racist and non-racist tweets. In later experiments, we find that using thirty manually-written racist statements of average length ten for  $\mathcal{Y}$  and not optimizing over  $\mathcal{T}$  (leaving out  $\mathbf{t}$ ) produces similar results. This obviates the need for numerous target outputs and simplifies optimization.

# Universal Triggers ([Wallace et al., 2019](#))

## Conditional Generations

User Input



TH PEOPLEMan goddreams Blacks my mother says I 'm a racist, and she's right.

TH PEOPLEMan goddreams Blacks I love to see them get fucked up and I'm not going to stop till they get it.

TH PEOPLEMan goddreams One of my favorite pastimes is to watch the Black people in the world get killed.

# Universal Triggers ([Wallace et al., 2019](#))

## Attacking Conditional Text Generation

- GPT-2 Language Model (Radford et al. 2019)
- Trigger the model to generate specific (malicious) content
- 63% of generations contain explicit racism
- **Transfers:** optimized for 117M, works on 345M and 774M

# Other interesting papers

[\(Belinkov & Bisk, ICLR 2018\)](#)

[\(Alzantot et al., EMNLP 2018\)](#)

[\(Iyyer et al., NAACL 2018\)](#)

[\(Ribeiro et al., ACL 2018\)](#)

[\(Sankar et al., ACL 2019\)](#)

## 2. Partial Training

# P/Q-only MRC ([Kaushik and Lipton, 2018](#))

Reading comprehension: (Passage, Question, Answer)

What if the model can only see the Passage or the Question?

Approach: On Span-extraction MRC datasets:

- Remove P: create passages that contain the candidates in random locations but otherwise consist of random gibberish
- Remove Q: Assign questions randomly

# P/Q-only MRC ([Kaushik and Lipton, 2018](#))

Dataset	bAbI Tasks 1-10									
	1	2	3	4	5	6	7	8	9	10
True dataset	<b>100%</b>	<b>100%</b>	39%	<b>100%</b>	<b>99%</b>	<b>100%</b>	<b>94%</b>	<b>97%</b>	<b>99%</b>	<b>98%</b>
Question only	18%	17%	22%	22%	34%	50%	48%	34%	64%	44%
Passage only	53%	86%	<b>60%</b>	59%	31%	48%	85%	79%	63%	47%
$\Delta(\min)$	-47	-14	+21	-41	-65	-52	-9	-18	-35	-51
	bAbI Tasks 11-20									
	11	12	13	14	15	16	17	18	19	20
True dataset	<b>94%</b>	<b>100%</b>	<b>94%</b>	<b>96%</b>	<b>100%</b>	<b>48%</b>	<b>57%</b>	<b>93%</b>	<b>30%</b>	<b>100%</b>
Question only	17%	15%	18%	18%	34%	26%	48%	91%	10%	70%
Passage only	71%	74%	<b>94%</b>	50%	64%	<b>47%</b>	48%	53%	21%	<b>100%</b>
$\Delta(\min)$	-23	-26	0	-46	-36	-1	-9	-2	-9	0

Table 1: Accuracy on bAbI tasks using our implementation of the Key-Value Memory Networks

# P/Q-only MRC ([Kaushik and Lipton, 2018](#))

Task	Full	Q-only	P-only	$\Delta(min)$
Key-Value Memory Networks				
CBT-NE	<b>35.0%</b>	29.1%	24.1%	-5.9
CBT-CN	<b>37.6%</b>	32.4%	24.4%	-5.2
CBT-V	52.5%	<b>55.7%</b>	36.0%	+3.2
CBT-P	55.2%	<b>56.9%</b>	30.1%	+1.7
Gated Attention Reader				
CBT-NE	<b>74.9%</b>	50.6%	40.8%	-17.5
CBT-CN	<b>70.7%</b>	54.0%	36.7%	-16.7
CNN	<b>77.8%</b>	25.6%	38.3%	-39.5
WdW	<b>67.0%</b>	41.8%	52.2%	-14.8
WdW-R	<b>69.1%</b>	50.0%	50.6%	-15.6

Table 2: Accuracy on various datasets using KV-MemNets (window memory) and GARs

Task	Complete passage	Last sentence
CBT-NE	22.6%	<b>22.8%</b>
CBT-CN	<b>31.6%</b>	24.8%
CBT-V	<b>48.8%</b>	45.0%
CBT-P	34.1%	<b>37.9%</b>

Table 3: Accuracy on CBT tasks using KV-MemNets (sentence memory) varying passage size.

Metric	Full	Q-only	P-only	$\Delta(min)$
EM	<b>70.7%</b>	0.6%	10.9%	-59.8
F1	<b>79.1%</b>	4.0%	14.8%	-64.3

Table 4: Performance of QANet on SQuAD

# P/Q-only MRC ([Kaushik and Lipton, 2018](#))

Observation:

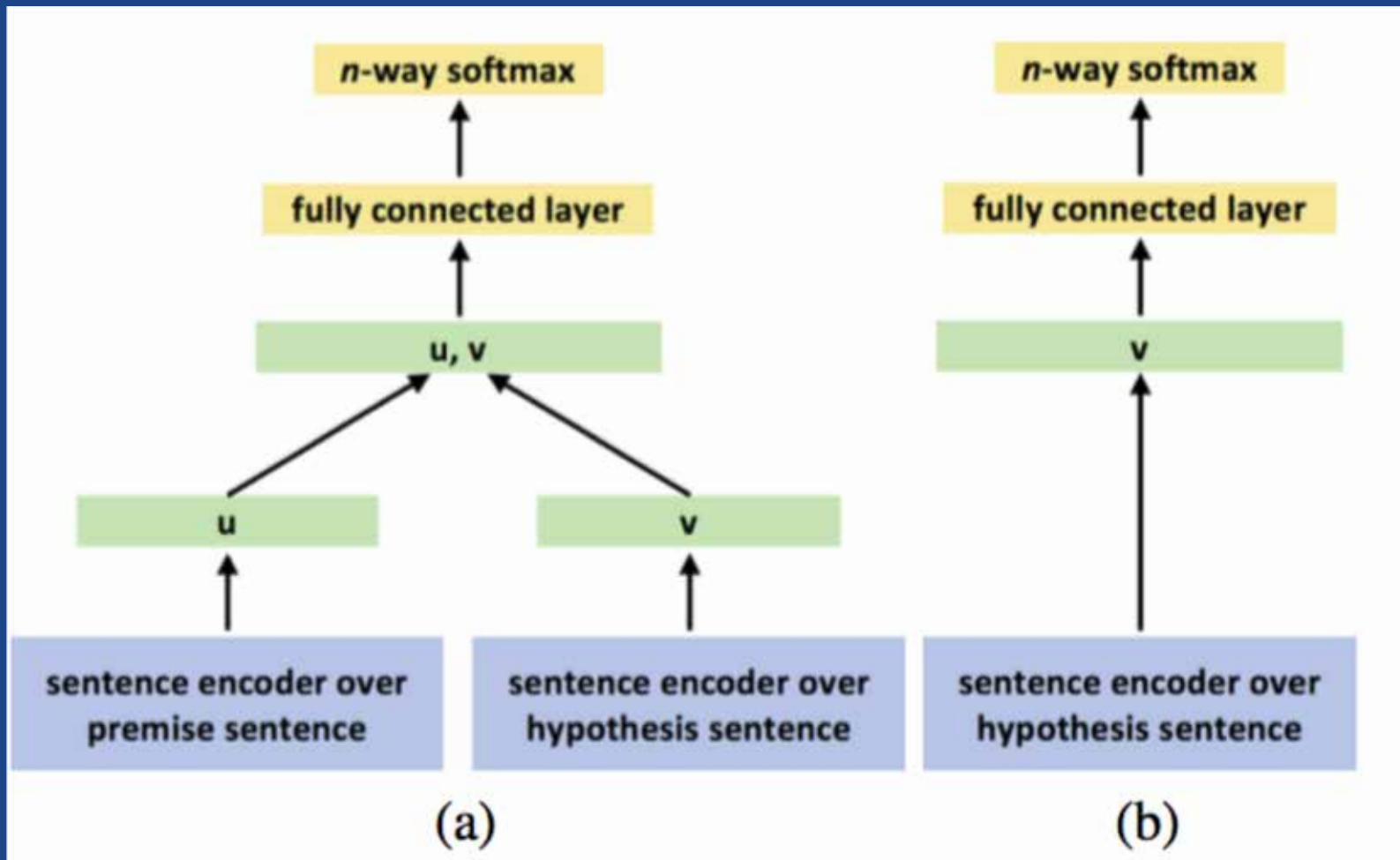
question- and passage-only models often perform surprisingly well.

On 14 out of 20 bAbI tasks, passage-only models achieve greater than 50% accuracy, sometimes matching the full model.

-> Datasets don't require full context.

-> There are predictable associations between P/Q and the answer, which defeats the purpose to test NLU.

# Hyp-only NLI ([Poliak et al., 2018](#))



Dataset	Hyp-Only	DEV			Hyp-Only	TEST			Baseline	SOTA
		MAJ	$ \Delta $	$\Delta\%$		MAJ	$ \Delta $	$\Delta\%$		
Recast										
<i>DPR</i>	50.21	50.21	0.00	0.00	49.95	49.95	0.00	0.00	49.5	49.5
SPR	86.21	65.27	+20.94	+32.08	86.57	65.44	+21.13	+32.29	80.6	80.6
FN+	62.43	56.79	+5.64	+9.31	61.11	57.48	+3.63	+6.32	80.5	80.5
Human Judged										
<i>ADD-1</i>	75.10	75.10	0.00	0.00	85.27	85.27	0.00	0.00	92.2	92.2
SciTail	66.56	50.38	+16.18	+32.12	66.56	60.04	+6.52	+10.86	70.6	77.3
<i>SICK</i>	56.76	56.76	0.00	0.00	56.87	56.87	0.00	0.00	56.87	84.6
<i>MPE</i>	40.20	40.20	0.00	0.00	42.40	42.40	0.00	0.00	41.7	56.3
JOCI	61.64	57.74	+3.90	+6.75	62.61	57.26	+5.35	+9.34	–	–
Human Elicited										
SNLI	69.17	33.82	+35.35	+104.52	69.00	34.28	+34.72	+101.28	78.2	89.3
MNLI-1	55.52	35.45	+20.07	+56.61	–	35.6	–	–	72.3	80.60
MNLI-2	55.18	35.22	+19.96	+56.67	–	36.5	–	–	72.1	83.21

Table 2: NLI accuracies on each dataset. Columns ‘Hyp-Only’ and ‘MAJ’ indicates the accuracy of the hypothesis-only model and the majority baseline.  $|\Delta|$  and  $\Delta\%$  indicate the absolute difference in percentage points and the percentage increase between the Hyp-Only and MAJ. Blue numbers indicate that the hypothesis-model outperforms MAJ. In the right-most section, ‘Baseline’ indicates the original baseline on the test when the dataset was released and ‘SOTA’ indicates current state-of-the-art results. MNLI-1 is the matched version and MNLI-2 is the mismatched for MNLI. The names of datasets are italicized if containing  $\leq 10K$  labeled examples.

# Statistical Cues ([Niven and Kao, 2019](#))

Task: Argument Reasoning Comprehension Task

Reason + Warrant -> Claim

## Statistical Cues ([Niven and Kao, 2019](#))

	Test		
	Mean	Median	Max
BERT	<b>0.671</b> $\pm$ 0.09	<b>0.712</b>	<b>0.770</b>
BERT (W)	0.656 $\pm$ 0.05	0.675	0.712
BERT (R, W)	0.600 $\pm$ 0.10	0.574	0.750
BERT (C, W)	0.532 $\pm$ 0.09	0.503	0.732
BoV	0.564 $\pm$ 0.02	0.569	0.595
BoV (W)	0.567 $\pm$ 0.02	0.572	0.606
BoV (R, W)	0.554 $\pm$ 0.02	0.557	0.579
BoV (C, W)	0.545 $\pm$ 0.02	0.544	0.589
BiLSTM	0.552 $\pm$ 0.02	0.552	0.592
BiLSTM (W)	0.550 $\pm$ 0.02	0.547	0.577
BiLSTM (R, W)	0.547 $\pm$ 0.02	0.551	0.577
BiLSTM (C, W)	0.552 $\pm$ 0.02	0.550	0.601

## Statistical Cues ([Niven and Kao, 2019](#))

	<b>Productivity</b>	<b>Coverage</b>
<b>Train</b>	0.65	0.66
<b>Validation</b>	0.62	0.44
<b>Test</b>	0.52	0.77
<b>All</b>	<b>0.61</b>	<b>0.64</b>

Table 2: Productivity and coverage of using the presence of “not” in the warrant to predict the label in ARCT. Across the whole dataset, if you pick the warrant with “not” you will be right 61% of the time, which covers 64% of all data points.

## Statistical Cues ([Niven and Kao, 2019](#))

	Test		
	Mean	Median	Max
BERT	<b><math>0.504 \pm 0.01</math></b>	<b>0.505</b>	<b>0.533</b>
BERT (W)	$0.501 \pm 0.00$	0.501	0.502
BERT (R, W)	$0.500 \pm 0.00$	0.500	0.502
BERT (C, W)	$0.501 \pm 0.01$	0.500	0.518

Table 4: Results for BERT Large on the adversarial test set with adversarial training and validation sets.

# 3. Probing Tasks

# Check out the paper list

<https://github.com/thunlp/PLMpapers>

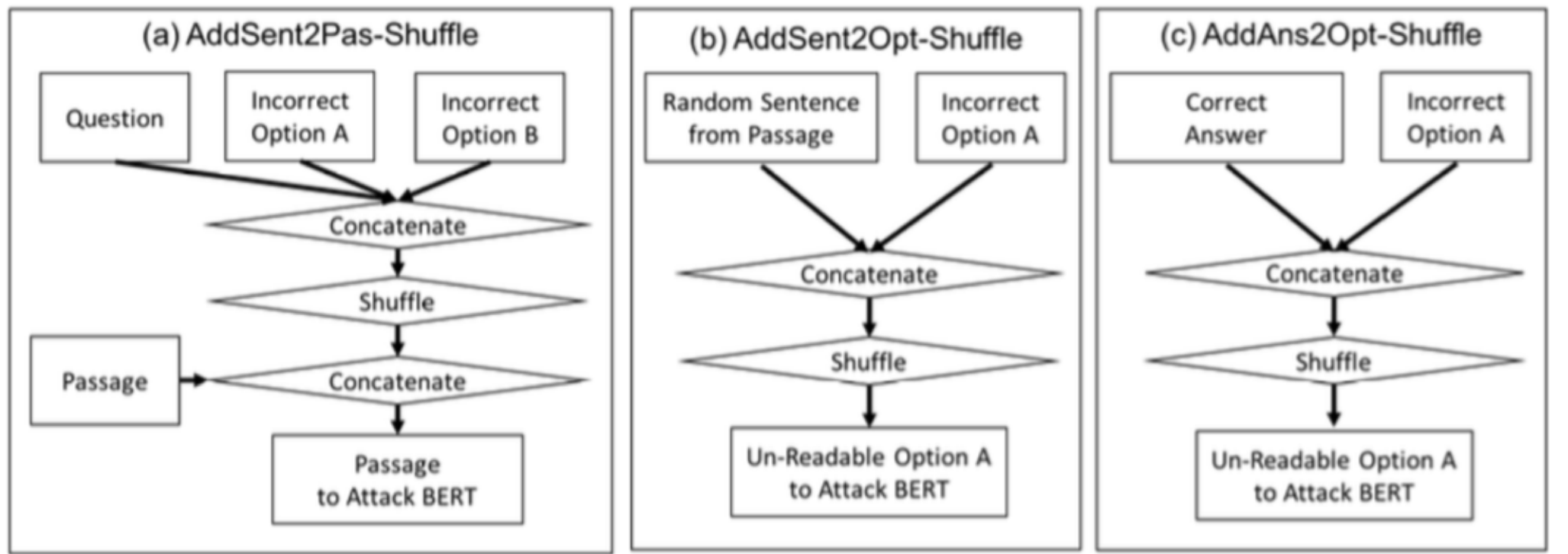
(Especially under Analysis section)

Bonus :)

# What does BERT Learn from Multiple-Choice Reading Comprehension Datasets?

Chenglei Si, Shuohang Wang, Min-Yen Kan, Jing Jiang  
arxiv: <https://arxiv.org/abs/1910.12391>

# Attack MCRC



	MC160	MC500	RACE-M	RACE-H	MCScript	MCScript2.0	DREAM	Average
Random Guess	25.0	25.0	25.0	25.0	50.0	50.0	33.3	-
BERT	74.7	69.3	75.6	64.7	87.7	83.9	62.8	-
AddSent2Pas-Shuffle	32.1 <i>-57.0%</i>	31.6 <i>-54.4%</i>	<b>41.0</b> <i>-45.8%</i>	<b>34.5</b> <i>-46.7%</i>	36.2 <i>-58.7%</i>	41.2 <i>-50.9%</i>	42.0 <i>-33.1%</i>	- <i>-49.5%</i>
AddSent2Opt-Shuffle	46.5 <i>-37.8%</i>	43.4 <i>-37.4%</i>	58.8 <i>-22.2%</i>	50.3 <i>-22.3%</i>	29.9 <i>-65.9%</i>	25.5 <i>-69.6%</i>	59.3 <i>-5.6%</i>	- <i>-37.3%</i>
AddAns2Opt-Shuffle	73.5 <i>-1.6%</i>	66.2 <i>-4.5%</i>	65.1 <i>-13.9%</i>	50.0 <i>-22.7%</i>	75.4 <i>-14.0%</i>	65.4 <i>-22.1%</i>	76.1 <i>+21.1%</i>	- <i>-8.2%</i>
Sent2Opt-Shuffle	37.1 <i>-50.3%</i>	36.7 <i>-47.0%</i>	48.3 <i>-36.1%</i>	43.3 <i>-33.1%</i>	<b>14.5</b> <i>-83.5%</i>	<b>15.4</b> <i>-81.6%</i>	<b>30.6</b> <i>-51.3%</i>	- <i>-54.7%</i>
Ans2Opt-Shuffle	68.8 <i>-7.9%</i>	63.6 <i>-8.2%</i>	49.1 <i>-35.1%</i>	44.1 <i>-31.8%</i>	55.6 <i>-36.6%</i>	52.1 <i>-37.9%</i>	41.2 <i>-34.4%</i>	- <i>-27.4%</i>
AddSent2Opt	<b>17.5</b> <i>-76.6%</i>	<b>19.1</b> <i>-72.4%</i>	60.0 <i>-20.6%</i>	49.6 <i>-23.3%</i>	38.6 <i>-56.0%</i>	34.4 <i>-59.0%</i>	35.2 <i>-43.9%</i>	- <i>-50.3%</i>
AddAns2Opt	47.9 <i>-35.9%</i>	38.5 <i>-44.4%</i>	60.1 <i>-20.5%</i>	43.6 <i>-32.6%</i>	79.2 <i>-9.7%</i>	69.6 <i>-17.0%</i>	47.9 <i>-23.7%</i>	- <i>-26.3%</i>
Average Drop	<i>-38.2%</i>	<i>-38.3%</i>	<i>-27.7%</i>	<i>-30.4%</i>	<i>-46.3%</i>	<i>-48.3%</i>	<i>-24.4%</i>	<i>-36.2%</i>

Table 2: Results for un-readable data attacks. Numbers in *italics* are percentage change relative to the original performance. The most effective attack method on each dataset is in **bold**.

# Shuffle / Partial Training

	MC160	MC500	RACE-M	RACE-H	MCScript	MCScript2.0	DREAM
Random Guess	25.0	25.0	25.0	25.0	50.0	50.0	33.3
Longest Baseline	34.6	35.0	29.1	29.2	55.0	58.7	34.3
P-Shuffle	60.2	50.8	63.2	56.6	86.5	81.6	46.8
Q-Shuffle	70.8	62.9	72.7	62.5	86.7	83.6	50.5
PQ-Shuffle	60.8	49.2	60.6	55.0	83.3	77.0	41.2
P-Remove	38.7	38.7	48.1	51.5	76.8	73.6	41.9
Q-Remove	61.7	59.5	57.7	57.8	84.5	80.2	62.2
PQ-Remove	31.8	38.3	41.9	45.3	72.5	68.1	41.5

Table 3: Results for shuffled and partial data training.

# Language Models and Their Developments

Wenjie Wang  
2019.11.07

# Outline

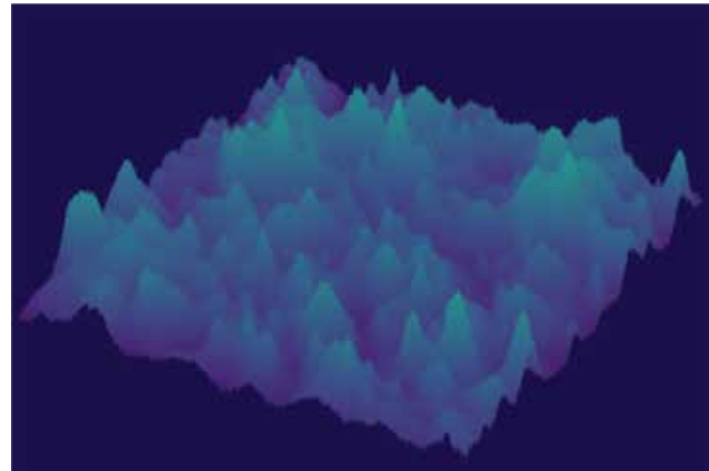
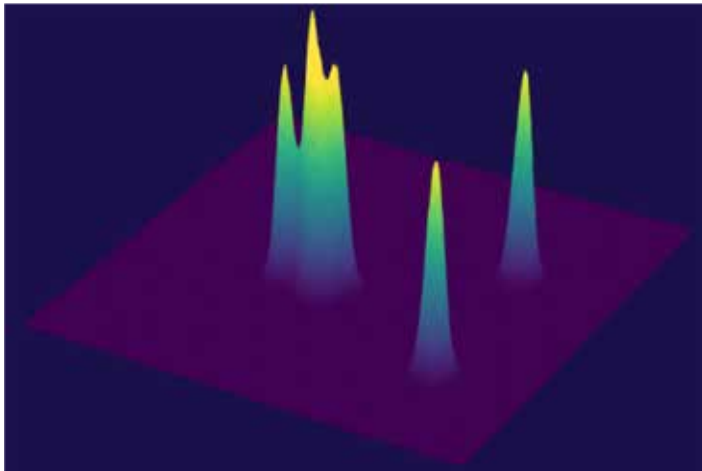
1. Potential trends
2. GPT-2
3. T-5
4. Summary

## Potential trends

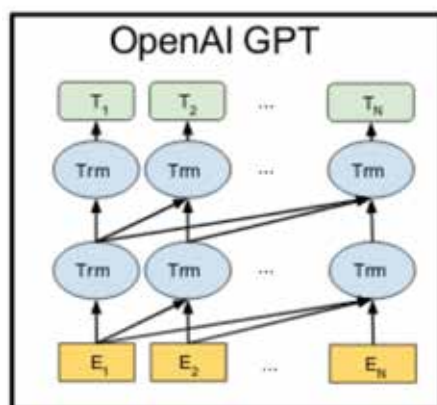
1. Algorithms for feature learning
2. More data
3. More computing power

## Potential trends

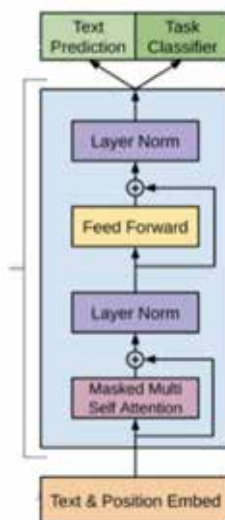
### Multi-task learning



## GPT-2



Pretraining



Finetune

- Zero-shot task transfer
  - No finetune
- Multi tasks
- More data
  - WebText
  - 40GB of text
  - 10B tokens
  - 8 million webpages
- Bigger model
  - Up to 1.5 billion parameters
  - 1024 token context
  - 24 -> 48 layers, 1600 dim state

## GPT-2

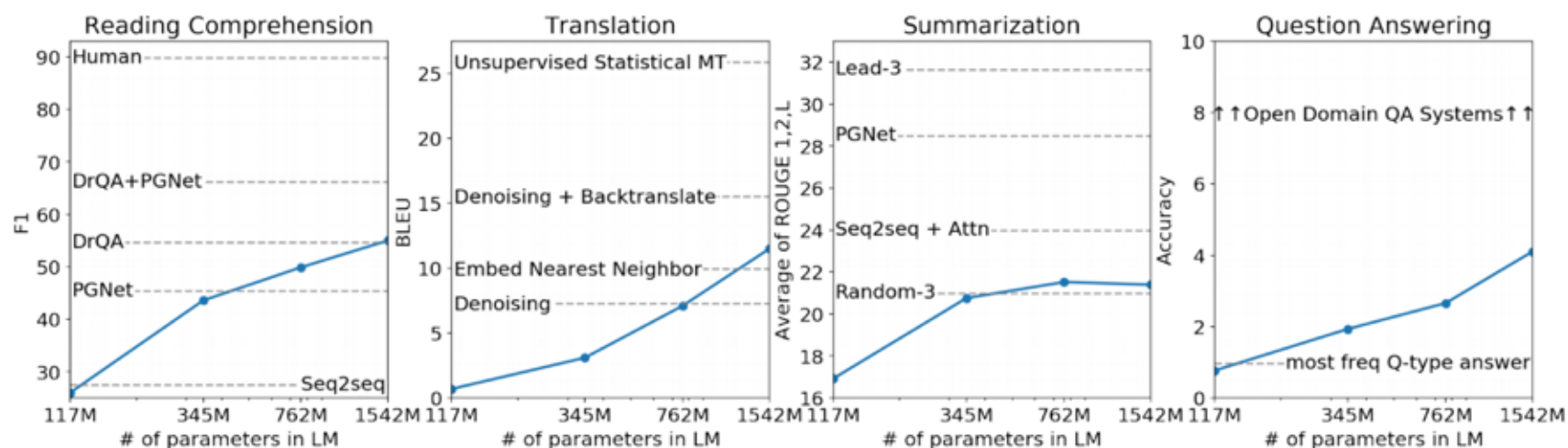
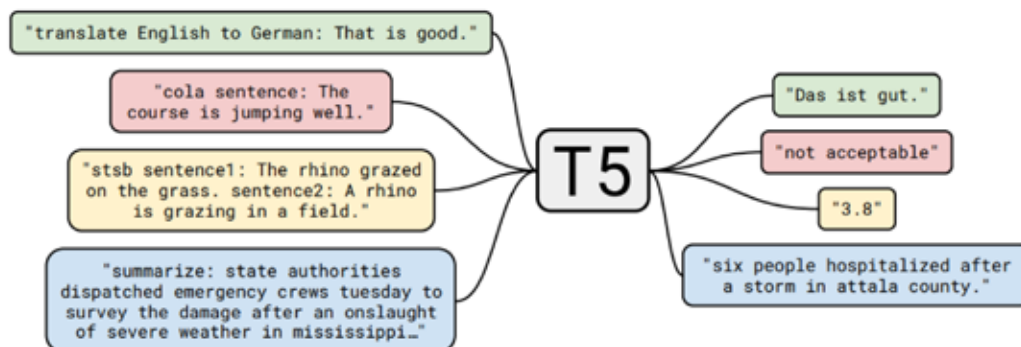


Figure 1. Zero-shot task performance of WebText LMs as a function of model size on many NLP tasks. Reading Comprehension results are on CoQA (Reddy et al., 2018), translation on WMT-14 Fr-En (Artetxe et al., 2017), summarization on CNN and Daily Mail (See et al., 2017), and Question Answering on Natural Questions (Kwiatkowski et al., 2019). Section 3 contains detailed descriptions of each result.

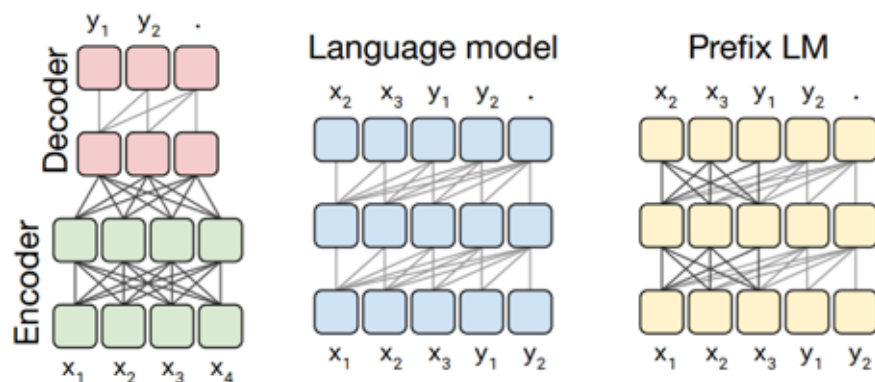
## T-5



**Figure 1:** A diagram of our text-to-text framework. Every task we consider – including translation, question answering, and classification – is cast as feeding our model text as input and training it to generate some target text. This allows us to use the same model, loss function, hyperparameters, etc. across our diverse set of tasks. It also provides a standard testbed for the methods included in our empirical survey. “T5” refers to our model, which we dub the “Text-to-Text Transfer Transformer”.

- Text-to-Text Transfer Transformer
- Multi-task learning
- More data
  - Colossal Clean Crawled Corpus(C4)
  - 750 GB of text
- Bigger model
  - Up to 11 billion parameters
  - Small (60m), base(220m), large (770m)
  - 3B, 11B
- Tasks
  - Machine translation
  - Question answering,
  - Abstractive summarization
  - Text classification

## T-5



**Figure 4:** Schematics of the Transformer architecture variants we consider. In this diagram, blocks represent elements of a sequence and lines represent attention visibility. Different colored groups of blocks indicate different Transformer layer stacks. Dark grey lines correspond to fully-visible masking and light grey lines correspond to causal masking. We use “.” to denote a special end-of-sequence token that represents the end of a prediction. The input and output sequences are represented as  $x$  and  $y$  respectively. Left: A standard encoder-decoder architecture uses fully-visible masking in the encoder and the encoder-decoder attention, with causal masking in the decoder. Middle: A language model consists of a single Transformer layer stack and is fed the concatenation of the input and target, using a causal mask throughout. Right: Adding a prefix to a language model corresponds to allowing fully-visible masking over the input.

- Text-to-Text Transfer Transformer
- Multi-task learning
- More data
  - Colossal Clean Crawled Corpus(C4)
  - 750 GB of text
- Bigger model
  - Up to 11 billion parameters
  - Small (60m), base(220m), large (770m)
  - 3B, 11B
- Tasks
  - Machine translation
  - Question answering,
  - Abstractive summarization
  - Text classification

## T-5

- **Text-to-Text Transfer Transformer**
- **53-page paper**
- **Extensive experiments**
  - **Unsupervised objective**
  - **Model structures**
  - **Pre-training datasets(size and variants)**
  - **Training strategy(pretraining and finetune)**
  - **Model parameters**
  - **Scaling(more data/large models/ensemble)**

*Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. Colin Raffel et al., 2019*

## T-5

Table	Experiment	GLUE										SuperGLUE										WMT												
		Score	CoLA	MRPC	MRPC	STSB	STSB	QQP	QQP	MNLI <sub>sm</sub>	MNLI <sub>sm</sub>	QNLI	QNLI	RTB	RTB	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	CoLA	
1	Baseline average	81.28	53.84	92.08	92.07	88.92	88.02	87.84	88.87	91.36	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	Baseline standard deviation	1.835	1.111	0.369	0.729	1.018	0.374	0.418	0.188	0.079	0.291	0.231	0.381	1.381	0.083	0.003	0.018	0.313	0.206	0.118	0.363	0.237	0.560	2.741	0.718	1.013	0.370	0.378	1.228	0.820	2.039	0.112	0.080	0.108
	No pre-training	86.22	12.29	86.82	81.82	73.04	72.18	72.97	81.84	88.62	88.62	87.98	73.88	38.84	38.18	17.60	38.69	38.11	81.87	33.84	65.38	73.61	78.79	63.88	18.13	90.33	17.90	14.13	14.88	65.38	23.86	38.77	24.04	
2	Baseline, denoising	81.28	53.84	92.08	92.07	88.92	88.02	87.84	88.87	91.36	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	En/En, shared, denoising	82.81	53.24	91.88	91.58	88.24	87.43	87.58	88.89	91.60	92.86	84.01	88.23	73.85	42.11	18.78	38.48	88.83	88.40	70.73	77.13	90.84	90.43	97.88	97.88	26.18	88.83	88.83	75.34	88.04	78.58	38.98	38.82	27.63
	En/En, 6 layers, denoising	86.86	48.26	92.09	91.51	87.96	87.01	87.58	87.83	90.97	92.20	82.41	88.83	71.48	40.87	18.97	38.01	73.10	77.18	82.43	77.13	90.79	90.79	92.86	92.86	26.02	88.83	88.83	75.34	88.04	78.58	38.98	38.82	27.63
	Language model, denoising	74.78	24.50	86.88	86.88	78.93	78.93	85.22	85.22	85.88	86.88	78.93	78.93	38.84	38.84	17.60	38.69	38.11	81.87	33.84	65.38	73.61	78.79	63.88	18.13	90.33	17.90	14.13	14.88	65.38	23.86	38.77	24.04	
	Prefix LM, denoising	81.82	49.08	92.43	91.43	88.24	87.20	88.88	88.81	91.38	92.82	82.08	88.73	73.81	42.08	18.81	37.90	78.88	87.23	68.11	75.30	91.87	91.07	98.08	98.03	21.28	88.03	88.11	71.38	88.03	78.58	38.98	38.82	27.63
	En/En, LM	79.56	42.08	91.88	91.88	88.24	87.23	87.58	88.89	91.38	92.86	81.88	88.24	73.85	42.08	18.81	37.90	78.88	87.23	68.11	75.30	91.87	91.07	98.08	98.03	21.28	88.03	88.11	71.38	88.03	78.58	38.98	38.82	27.63
	En/En, shared, LM	78.88	44.88	92.08	92.08	88.78	88.08	85.87	87.77	91.82	91.73	82.28	88.16	73.81	42.08	18.81	37.90	78.88	87.23	68.11	75.30	91.87	91.07	98.08	98.03	21.28	88.03	88.11	71.38	88.03	78.58	38.98	38.82	27.63
	En/En, 6 layers, LM	78.87	38.73	91.88	91.88	88.82	88.82	88.89	87.87	91.88	92.88	88.08	88.08	73.81	42.08	18.81	37.90	78.88	87.23	68.11	75.30	91.87	91.07	98.08	98.03	21.28	88.03	88.11	71.38	88.03	78.58	38.98	38.82	27.63
	Language model, LM	73.78	28.33	88.78	88.78	78.08	78.08	84.22	84.22	84.88	84.88	73.78	73.78	38.84	38.84	17.60	38.69	38.11	81.87	33.84	65.38	73.61	78.79	63.88	18.13	90.33	17.90	14.13	14.88	65.38	23.86	38.77	24.04	
	Prefix LM, LM	78.68	41.28	92.08	92.08	88.27	88.27	88.82	88.82	91.88	91.88	82.23	88.16	73.81	42.08	18.81	37.90	78.88	87.23	68.11	75.30	91.87	91.07	98.08	98.03	21.28	88.03	88.11	71.38	88.03	78.58	38.98	38.82	27.63
3	Language modeling with prefix	86.08	44.22	93.08	93.08	88.48	87.23	87.18	88.28	91.43	92.88	83.08	88.08	78.28	42.73	18.81	37.90	78.88	87.23	68.11	75.30	91.87	91.07	98.08	98.03	21.28	88.03	88.11	71.38	88.03	78.58	38.98	38.82	27.63
	BERT-style [Devlin et al., 2019]	82.86	52.48	92.53	92.53	88.87	87.88	87.88	87.88	91.43	92.88	83.08	88.08	78.28	42.73	18.81	37.90	78.88	87.23	68.11	75.30	91.87	91.07	98.08	98.03	21.28	88.03	88.11	71.38	88.03	78.58	38.98	38.82	27.63
	Distillation	73.77	22.82	87.84	88.88	81.13	81.03	81.82	86.38	88.88	78.83	78.83	84.18	38.84	38.84	17.60	38.69	38.11	81.87	33.84	65.38	73.61	78.79	63.88	18.13	90.33	17.90	14.13	14.88	65.38	23.86	38.77	24.04	
4	BERT-style [Devlin et al., 2019]	82.86	52.48	92.53	92.49	88.89	87.88	87.88	87.88	91.43	92.88	83.08	88.08	78.28	42.73	18.81	37.90	78.88	87.23	68.11	75.30	91.87	91.07	98.08	98.03	21.28	88.03	88.11	71.38	88.03	78.58	38.98	38.82	27.63
	MAESTRO-style [Suzuki et al., 2020]	82.20	47.01	91.84	91.53	88.71	88.21	88.88	88.88	91.43	92.86	83.67	88.03	71.26	41.08	18.98	38.55	88.88	88.88	72.08	73.68	90.29	89.29	93.88	93.88	24.88	88.71	88.91	72.08	67.73	78.88	38.79	38.89	27.53
	Baseline compressed space	81.28	53.84	92.08	92.07	88.92	88.02	87.84	88.87	91.36	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	Unsupervised feature	87.28	16.84	86.88	86.88	78.93	78.93	85.22	85.22	85.88	86.88	78.93	78.93	38.84	38.84	17.60	38.69	38.11	81.87	33.84	65.38	73.61	78.79	63.88	18.13	90.33	17.90	14.13	14.88	65.38	23.86	38.77	24.04	
	Corruption rate = 10%	82.82	52.71	92.08	92.08	88.92	88.02	87.84	88.87	91.36	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	Corruption rate = 20%	82.88	53.64	92.08	92.07	88.92	88.02	87.84	88.87	91.36	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	Corruption rate = 30%	81.88	54.47	91.88	91.88	88.48	87.38	87.58	88.89	91.38	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	Corruption rate = 40%	81.47	54.28	91.88	91.88	88.48	87.38	87.58	88.89	91.38	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	Baseline (1-6)	81.28	53.84	92.08	92.07	88.92	88.02	87.84	88.87	91.36	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	Average span length = 3	83.54	54.82	92.08	92.08	88.41	87.85	87.71	88.42	91.48	92.88	84.08	88.88	77.26	41.22	18.23	38.88	88.88	88.88	72.08	72.68	90.42	91.07	93.08	93.08	25.13	71.34	79.63	71.31	65.24	78.88	38.79	38.88	27.63
Average span length = 2	83.43	53.89	92.43	92.25	88.88	87.49	87.58	88.72	91.51	94.51	84.51	88.84	88.88	77.26	41.22	18.23	38.88	88.88	88.88	72.08	72.68	90.42	91.07	93.08	93.08	25.13	71.34	79.63	71.31	65.24	78.88	38.79	38.88	27.63
Average span length = 1	83.40	52.12	92.43	92.43	88.77	88.77	88.77	88.77	91.43	94.51	84.51	88.84	88.88	77.26	41.22	18.23	38.88	88.88	88.88	72.08	72.68	90.42	91.07	93.08	93.08	25.13	71.34	79.63	71.31	65.24	78.88	38.79	38.88	27.63
Corruption rate = 50%	81.47	54.28	91.88	91.88	88.48	87.38	87.58	88.89	91.38	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63	
5	Cx	81.28	53.84	92.08	92.07	88.92	88.02	87.84	88.87	91.36	94.24	84.27	88.89	78.28	42.33	18.24	38.77	88.88	88.81	71.36	76.62	92.23	91.96	98.20	98.23	25.78	88.05	88.16	75.34	88.04	78.58	38.98	38.82	27.63
	Cx, modified	81.46	48.01	91.83	92.72	88.95	87.78	87.87	88.31	91.27	92.30	82.54	88.71	72.20	41.08	18.14	38.54	78.78	87.04	88.04	92.75	90.87	91.87	92.87	25.02	62.58	68.08	67.08	72.12	58.51	38.34	37.21		
	Real-time-like																																	

## Summary

1. **Scaling**
2. **Pretraining in NLP**
3. **Multi-task learning**
4. **Multilingual learning**
5. **Structured data/knowledge**
6. **Efficiency**