

Is scale all we need?

Slav Petrov on behalf of many wonderful colleagues at Google Research

slav@google.com



Number of Synapses in Biological & Artificial Systems



List of animals by number of neurons (Wikipedia)

Massive Neural Machine Translation



*

*

T(96L)

Baselines

36.9

30.8

6.1

-

2.3B

 $100 \times 0.4B$

<u>"GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding"</u> ICLR '21 D. Lepikhin, H. Lee, Y. Xu, D. Chen, O. Firat, Y. Huang, M. Krikun, N. Shazeer, Z. Chen

Are large language models enough?



<u>"Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer,"</u> Arxiv '19 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J. Liu

But what do we do when something goes wrong?



This talk is about aspects that won't be solved by scale alone:

- Memory
- Attribution
- Generation
- Time
- Tools
- Trust

From "implicit" to "explicit" memory

Most current LMs are not grounded / do not attribute their output. The knowledge is "implicitly" somewhere in the parameters.

• Attribution to sources

- **Structured knowledge:** knowledge graph entities and relations
- **Unstructured knowledge:** documents, images, video
- Training examples: these can also be treated as "memories"

• Benefits

- Provenance provides interpretability & trustworthiness
- Greater memorization capacity
- Greater efficiency, thanks to sparse memory access
- Generalization by controlling memories

Entities as Experts - Learn Entity Memories

Matches T5 performance on Trivia QA with 1% of parameters used per example.



	Activated Params	TriviaQA Accuracy
T5	11B	42.3%
EaE	95M	43.2%

FILM (Fact Injected LM) - Adding Facts

Adding facts from the KG increases performance over very large LMs.

And, lets us update the model's knowledge about the world without retraining.



<u>"Adaptable and Interpretable Neural Memory Over Symbolic Knowledge"</u> NAACL '21 Pat Verga, Haitian Sun, Livio Baldini Soares, William Cohen







Retrieval-Augmented Language Models



Open QA with T5 vs retrieve-and-read QA systems

T5-11B model is *competitive* with contemporaneous QA systems on NQ Open benchmark...

...but performs much less well on test questions that are "novel" - i.e., no similar question or answer in the training set - with performance less than ½ the state-of-the-art*

<u>"Question and Answer Test-Train Overlap in Open-Domain Question Answering Datasets"</u> EACL '21 Patrick Lewis, Pontus Stenetorp, Sebastian Riedel



Output

[ENTITYCHAIN] Frozen | Disney [SUMMARY]

Target <eos>

Transformer Decoder

Target with Entity Chain

<s> [ENTITYCHAIN] Frozen | Disney

[SUMMARY] Target

Step-Wise Long-Form Generation with Pre-Training



Pre-Training

Fine-Tuning

Training Data Timeline



"Time-Aware Language Models as Temporal Knowledge Bases" Arxiv '21

Bhuwan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, William W. Cohen

Joint Models of Text and Time

- Idea: Instead of Pr(x) model Pr(x, t)
 - Many corpora come with timestamps (e.g. News)



• Preliminary results show that such models can improve memorization of the past, improve calibration of future events and allow for **30x** faster adaptation to new data as it becomes available.

Numerical reasoning?

Question:

How many yards do the first two field goals converted add up to?

Passage:

Jay Feely getting a <u>53</u>-yard field goal. In the second quarter, Miami drew closer as Feely kicked a <u>44</u>-yard field goal, yet New York replied with

Answer: 97 (53 + 44)

Google Assistant Scenario



User: Wake me up 30 minutes before my earliest meeting tomorrow.

Assistant: Ok, setting an alarm for 8:30 A.M.

QA example from:

DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs, Dua et al., 2019

Can we get numerical capabilities from scale?



Language Models are Few-Shot Learners, Brown et al., 2020



Adding numerical reasoning abilities to LMs

Combine the LM with a lightweight graph-based model to perform computation

Pros: Interpretable via the graph. Cons: Can't generalize to unseen computations.

Add numerical data to pre-training or multi-task fine-tuning Pros: Can be more general purpose.

Cons: Need to collect/synthesize data.



Trustworthy NLP

No machine learned model will be perfect.

How can we build trustworthy systems out of untrustworthy components?

Three pillars for trustworthiness:

- Managing expectations / robustness / fairness
- Proof of work / Interpretability
- Controllable policies

Managing expectations

- Trust implies expectations about behavior
- Trust builds over time when expectations are fulfilled
- Well-defined problems make it easy for the expectations to match the abilities

Google Research

Desired and undesired correlations

- 1. Carefully evaluate unintended correlations.
- 2. Be mindful of seemingly innocuous configuration differences.
- 3. Focus on general mitigations.



Mitigation \rightarrow		Name (A-M)			
1987.5		Same	Flip	Random	
Evaluation \downarrow	Baseline	Gender	Gender	Gender	
DisCo (Names A-M)	3.9	2.5	2.6	1.2	
DisCo (Names N-Z)	2.5	2.6	2.3	2.0	
DisCo (Terms)	1.1	1.3	1.3	0.8	

"Measuring and Reducing Gendered Correlations in Pre-trained Models" Arxiv '20

Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed Chi, Slav Petrov

Proof of work

In school we aren't happy if students just give the right answers:

• we want the right answers for the right reasons

Proof of work for a model is not necessarily mathematical proof though

• many ways of convincing that a reasoning pattern is valid

Google Research

Proof of work example: QED

Decompose passage-based QA into:

- retrieval
- coreference
- entailment

Question: who wrote the film howl's moving castle? Passage: Howl's Moving Castle is a 2004 Japanese animated fantasy film written and directed by Hayao Miyazaki. It is based on the novel of the same name, which was written by Diana Wynne Jones. The film was produced by Toshio Suzuki. Answer: Hayao Miyazaki

(1) Sentence Selection

Howl's Moving Castle is a 2004 Japanese animated fantasy film written and directed by Hayao Miyazaki. (2) Referential Equality

the film howl's moving castle = Howl's Moving Castle (3) Entailment

X is a 2004 Japanese animated fantasy film written and directed by ANSWER. \vdash ANSWER wrote X.

Google Research

Language Interpretability Tool

	🗹 🖸 🛛 Data Table				2 23	Datapoint Editor	20
ctor UMAP - Embedding sst2:cls_emb - Label by sentence -	Only sho	w selected		Reset view Selec	all Columns -	sentence(*):	
	index Q	id Q 🗧	sentence	۹	label Q		
	0	827559	it 's a charming and often affecting journey .		1 0		
	1	98d0ff	unflinchingly bleak and desperate		0	label:	
	2	4f0e27	allows us to hope that nolan is poised to embark a major career as a commercial ye	et inventive filmmaker .	1		
	3	eb90c4	the acting, costumes, music, cinematography and sound are all astounding given locales.	n the production 's austere	1	Analyze new datapoint Reset Clear	
	4	fcedba	it's slow very , very slow .		0		
	5	5cbca5	although laced with humor and a few fanciful touches, the film is a refreshingly se	arious look at young women	. 1		
	6	32ec21	a sometimes tedious film .		0		
1 States	7	a1a90b	or doing last year 's taxes with your ex-wife .		0		
a second	8	df3932 you do n't have to know about music to appreciate the film 's easygoing blend of comedy and romance . 1					
	9	cb31s4_ in exactly 89 minutes , most of which passed as slowly as if I 'd been sitting naked on an igitoo , formula 51 sank of form quirky to briny to utter turkey .					
	10	d60a0d	the mesmerizing performances of the leads keep the film grounded and keep the	audience riveted .	1		
	11	dedate	it takes a strange kind of laziness to waste the talents of robert forster , anne mean	ra , eugene levy , and	0		
			reginald veljohnson all in the same movie .				
Performance Predictions Explanations Counte	erfactuais Counterfactual Ex	planation	=				
				∠ [] Confusio	Matrix		
1000 bbbbb bir I lehen				Rows	label	Columns sst2:probas V L Hide empty labels	
Model From Field Group	© N	o acc	suracy c precision c recall c f1	÷ 5552	robes		
Model C From Field Croup detaset probas multiclass	© N 872	 acc 0.821 	uracy () precision () recall () (1 0.819 (0.833 (0.826	2 034	nobes 1		
Model D From D Field D Group dataset probas multiclass	© N 872	0.821	auracy o precision o recall o f1 0.819 0.833 0.826	0 8012 0 34 1 74	1 6 82 370		
dataset probas multiclass	© N 872	acc 0.821	uracy precision recall 11 0.819 0.833 0.826	수 8952 전 월 0 34 1 74	1 5 82 370		
detaset probas multiclass	N 872	 acc 0.821 	uney 0 presision 0 recall 0 ft 0.819 0.833 0.826	수 8512 전 필 0 34 1 74	1 5 82 370		
detaset probas multiclass	© N 872	acc 0.821	urany 0 precision 0 recail 0 ft 0.819 0.833 0.826	2 8952 2 3 3 3 0 3 4 0 3 4 0 3 4 0 3 4 0 3 4 0 3 4 0 3 4 0 3 4 0 3	1 6 82 370		
dotaset probas multiclass	872 N	 acc 0.821 	uney 0 precision 0 recail 0 ft 0.819 0.833 0.824	수 8652 전 월 0 34 1 74	1 6 82 370		
dotaet probas multiclass	© N 872	 acc 0.821 	uney 0 presision 0 recall 0 ft 0.819 0.833 0.826	수 9952 전 전 1 74	1 6 82 370		
datset probas multiclass	© N 872	0.821	uracy 0 precision 0 recail 0 ft 0.819 0.833 0.826	수 942 전 및 0 33 1 74	1 5 62 370		
dotale o From o Field o Group dataset probas multiclass	N 872	0.821	uney 0 precision 0 recail 0 ft 0.819 0.833 0.824	् वर2 यु 0 उ 1 74	1 8 82 370		
dotaset probas multiclass	N 872	0.821	urany 0 precision 0 necali 0 ft 0.819 0.833 0.824	् = ==================================	robes 1 5 82 370		
dataset probas muticiass	N 872	0.821	unery 0 precision 0 recail 0 ft 0.819 0.833 0.826	عد مع قور قو 1 ت	nobes 1 5 62 370		
delaset probas multiclass	N 872	acc 0.821	uracy 0 precision 0 necall 0 ft 0.819 0.833 0.824	् सर कु 0 174	robes 1 5 82 370		

<u>"The Language Interpretability Tool: Extensible, Interactive Visualizations and Analysis for NLP Models</u>" EMNLP '20 I. Tenney, J. Wexler, J. Bastings, T. Bolukbasi, A. Coenen, S. Gehrmann, E. Jiang, M. Pushkarna, C. Radebaugh, E. Reif, A. Yuan

Controllable policies

Pretrained models are a description of the world

• the world has good and bad things in it

Most problems we want to solve are inherently open-ended and ambiguous

• we need to be able to tip the scales to force the right outcome

How can we design systems with the right control affordances?

Controllable policies

Example: gender in translation

- different languages require that different things are gendered
- so translation models often have to infer gender
- by default the most common gender in the data is inferred
- this is bad

Alternative:

- align the training data using a simple word alignment mechanism
- use this to annotate source genderless text with target gender
- train a model that respects these annotations
- now any learned or programmable policy can be used to set gender

Conclusions

- Scale is important
- But there are many important problems that will likely not be solved with scale alone:
 - Memory
 - Attribution
 - Planning for Long-Fortm Generation
 - Modeling Time
 - Using Tools
 - Building Trust
- Other important topics that we didn't talk about:
 - Datasets
 - Evaluation Metrics



Thank you!

Slav Petrov on behalf of many wonderful colleagues at Google Research

....

.

• •

slav@google.com

Confidential + Proprietary

.

.

. .