Recap









BART for Summarization: Outputs	© T5
This is the first time anyone has been recorded to run a full marathon of 42.195 kilometers (approximately 26 miles) under this pursued landmark time. It was not, however, an officially sanctioned world record, as it was not an "open race" of the IAAF. His time was 1 hour 59 minutes 40.2 seconds. Kipchoge ran in Vienna, Austria. It was an event specifically designed to help Kipchoge break the two hour barrier.Kenyan runner Eliud Kipchoge has run a marathon in less than two hours.PG&E stated it scheduled the blackouts in response to forecasts for high winds amid dry conditions. The aim is to reduce the risk of wildfires. Nearly 800 thousand customers were scheduled to be affected by the shutoffs which were expected to last through at least midday tomorrow.Power has been turned off to millions of customers in California as part of a power shutoff plan.	<ul> <li>Pre-training: similar denoising scheme to BART (they were released within a week of each other in fall 2019)</li> <li>Input: text with gaps. Output: a series of phrases to fill those gaps.</li> </ul>
Lewis et al. (2019)	Raffel et al. (2019)

Number of tokens	Repeats	GLUE	CNNDM	SQuAD	SGLUE	EnDe	$\mathbf{EnFr}$	EnRo
Full dataset	0	83.28	19.24	80.88	71.36	26.98	39.82	27.65
$2^{29}$	64	82.87	19.19	80.97	72.03	26.83	39.74	27.63
$2^{27}$	256	82.62	19.20	79.78	69.97	27.02	39.71	27.33
$2^{25}$	1,024	79.55	18.57	76.27	64.76	26.38	39.56	26.80
$2^{23}$	4,096	76.34	18.33	70.92	59.29	26.37	38.84	25.81
		S	ummarization			mach	ine transla	tion
Colossal Clea	ned Cor	nmon C	rawl: 750	) GB of t	ext			

Raffel et al. (2019)



			UnifiedQA					UnifiedQA
Extractive	EX	Dataset Input Output	SQuAD 1.1 At what speed did the turbine operate? \n (Nikola_Tesla) On his 50th birthday in 1906, Tesla demonstrated his 200 horsepower (150 kilowatts) 16,000 rpm bladeless turbine 16,000 rpm		Multiple choice	MC	Dataset Input	MCTest Who was Billy? \n (A) The skinny kid (B) A teacher (C) A little kid (D) The big kid \n Billy was like a king on the school yard. A king without a queen. He was the biggest kid in our grade, so he made all the rules during recess
Abstractive	AB	Dataset Input Output	NarrativeQA What does a drink from narcissus's spring cause the drinker to do? \n Mercury has awakened Echo, who weeps for Narcissus, and states that a drink from Narcissus's spring causes the drinkers to ``Grow dotingly enamored of themselves.'' fall in love with themselves		Yes/no	YN	Output Dataset Input Output	The big kid BoolQ Was America the first country to have a president? \n (President) The first usage of the word president to denote the highest official in a government was during the Commonwealth of England no
<ul> <li>Past work: different selection, answ</li> <li>Now: one 11B</li> </ul>	eren ver g oara	it archi genera ametei	itectures for every QA formulation. (Span tion, multiple choice,) r T5 model Khashabi et	al. (2020)	<ul> <li>Past work: differ selection, answer</li> <li>Now: one 11B p</li> </ul>	ren er g ara	t arch genera imete	itectures for every QA formulation. (Span ation, multiple choice,) r T5 model Khashabi et al. (2020)

							UI	iiiie	eur	QA						
n dataset?	Model	↓ - Evaluated	on $\rightarrow$	NewsQA	Quoref	Quoref-CS	ROPES	ROPES-CS	DROP	DROP-CS	QASC	Common senseQA	NP-BoolQ	BoolQ-CS	MultiRC	Avg
	Uı	nifiedQA [EX	1	58.7	64.7	53.3	43.4	29.4	24.6	24.2	55.3	62.8	20.6	12.8	7.2	38.1
	Ur	nifiedQA [AB	]	58.0	68.2	57.6	48.1	41.7	30.7	36.8	54.1	59.0	27.2	39.9	28.4	45.8
No	Un	nifiedQA [MO	]	48.5	67.9	58.0	61.0	44.4	28.9	37.2	67.9	75.9	2.6	5.7	9.7	42.3
	Ur	nifiedQA [YN	ŋ	0.6	1.7	1.4	0.0	0.7	0.4	0.1	14.8	20.8	79.1	78.6	91.7	24.2
		UnifiedQA		58.9	63.5	55.3	67.0	45.5	32.5	40.1	68.5	76.2	81.3	80.4	59.9	60.7
Ves		Previous best		66.8	86.1	55.4	61.1	32.5	89.1	54.2	85.2	79.1	78.4	71.1		
105		remous dest		Retro Reader	TASE	XLNet	ROBERTa	RoBERTa	ALBERT	MTMSN	KF+SIR+2Step	reeLB-RoBERT	RoBERTa	RoBERTa		
58	.7	64.7		53.3	43	.4	29.4	_			24.6	24.	2			
48	.5	67.9		57.0 58.0	61	.0	44.4	_			28.9	30.	2			
0.	6	1.7		1.4	0.	0	0.7				0.4	0.1				
	.9	63.5		55.3	67.	.0	45.5				32.5	40.	1			
58	.8	86.1		55.4	61	.1	32.5				89.1	54.	2			
<b>58</b> 66																

٢	

## 

How well does this really work?

Open-book: retrieval-based models

Closed-book: just generate the answer with T5 or BART Nearest neighbor: NN over training data

Models do well because the training sets contain answers shared in the test data

N	[ada]		Natural Questions						
IV	Iodei	Total	Question Overlap	Answer Overlap Only	No Overlap				
Open book	RAG DPR FID	44.5 41.3 51.4	70.7 69.4 71.3	34.9 34.6 48.3	24.8 19.3 34.5				
Closed book	T5-11B+SSM BART	36.6 26.5	77.2 67.6	22.2 10.2	9.4 0.8				
Nearest Neighbo	Dense r TF-IDF	26.7 22.2	69.4 56.8	7.0 4.1	0.0 0.0				

> T5 gets 9.4% accuracy if you filter to include no overlap

Khashabi et al. (2020)

Takeaways
<ul> <li>BART and T5 are useful for all sorts of seq2seq tasks involving language — so if you were going to use a seq2seq model, use one of these.</li> <li>(Caveat: need specialized models for things like language-to-code, but there's PLBART and CodeT5)</li> </ul>
<ul> <li>UnifiedQA suggests that big generative models are good at generalizing across tasks and even to new tasks (although QA results have a long way to go)</li> </ul>
If we have a strong enough pre-trained model and train on enough tasks, can we generalize to new tasks?
How do we specify those new tasks if they're not close to tasks we've already run on?
Answer: prompting.

OpenAl G	iPT/GPT2				OpenAl GPT2	
Very large language models using t	he Transformer	architectu	re	SYSTEM PROMPT (HUMAN-WRITTEN)	Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.	
Straightforward left-to-right langua	ge model, traine	ed on raw	text	MODEL COMPLETION (MACHINE-WRITTEN,	The 19-year-old singer was caught on camera being escorted out of the store by security guards.	
GP12: trained on 40GB of text collected from upvoted links	Parameters	Layers	$d_{model}$	SECONDIRAT	The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back.	
from reddit	117M	12	768		Scroll down for video	
▶ 1.5B parameters — by far the	345M 762M	24 36	1024 1280		Shoplifting: Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today (pictured)	
when it came out in March 2019	1542M	48	1600		The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back	
Because it's a language model, we	can <b>generate</b> fro	om it			The singer was also wearing a pair of black-rimmed glasses, a black jacket, black jeans and black sandals.	slide credit:
		Radford	et al. (2019)		She was carrying a pair of black and white striped gloves and a small black bag.	OpenAl

## ۲ 1) How novel is the stuff being generated? (Is it just doing nearest neighbors on a large corpus?) 2) How do we understand and distill what is learned in this model? 3) How do we harness these priors for conditional generation tasks (summarization, generate a report of a basketball game, etc.) 4) Is this technology dangerous? (OpenAI pursued a "staged release" strategy and didn't release biggest model)

**Open Questions** 

## ۲ Pre-Training Cost (with Google/AWS)

- BERT: Base \$500, Large \$7000
- GPT-2 (as reported in other work): \$25,000
- This is for a single pre-training run...developing new pre-training techniques may require many runs
- Fine-tuning these models can typically be done with a single GPU (but may take 1-3 days for medium-sized datasets)

https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-models/







		SuperGLUE	E BoolQ	CB	CB	COPA	RTE
	Eine turned SOTA	PO O	01.0		02.0	Accuracy 04.9	02.5
	Fine-tuned SOIA	<b>69.0</b>	91.0 77.4	90.9 83.6	93.9 75 7	<b>94.8</b> 70.6	92.5
	GPT-3 Few-Shot	71.8	76.4	75.6	52.0	92.0	69.0
		WiC	WSC	MultiRC	MultiRC	ReCoRD	ReCoRD
		Accuracy	Accuracy	Accuracy	F1a	Accuracy	F1
	Fine-tuned SOTA	76.1	93.8	62.3	88.2	92.5	93.3
	Fine-tuned BERT-Large	69.6	64.6	24.1	70.0	71.3	72.0
r	PT-3 Few-Shot	49.4	80.1	30.5	75.4	90.2	91.1



Prompts	Prompting for Classification
What about question answering?	Yelp For the Yelp Reviews Full Star dataset We define a single verbalizer $v$ for all patterns as (Zhang et al., 2015), the task is to estimate the rating that a customer gave to a restaurant on a 1- v(1) = terrible  v(2) = bad  v(3) = okay
Let's see if we can get some questions working both with and without examples in the context	to 5-star scale based on their review's text. We define the following patterns for an input text $a$ : $P_1(a) = $ It was $a$ $P_2(a) = $ Just! $   a$ "verbalizer" of labels $P_3(a) = a$ . All in all, it was patterns
How can we take this further?	<ul> <li>P<sub>4</sub>(a) = a    In summary, the restaurant is</li> <li>Compute probability of each pattern + verbalizer</li> </ul>
Brown et al. (2020)	<ul> <li>If P(It was bad) &gt; P (It was good) =&gt; suggests negative sentiment</li> <li>Can use prompting to pick between discrete options         <pre>Schick and Schutze et al. (2020)</pre> </li> </ul>





	Takeaways
<ul> <li>Pre-trained s of generation</li> </ul>	eq2seq models and generative language models can do well at lots n tasks
<ul> <li>Prompting is single model</li> </ul>	a way to harness their power and learn to do many tasks with a . Can be done without fine-tuning
<ul> <li>Prompting is single model</li> </ul>	a way to harness their power and learn to do many tasks with a