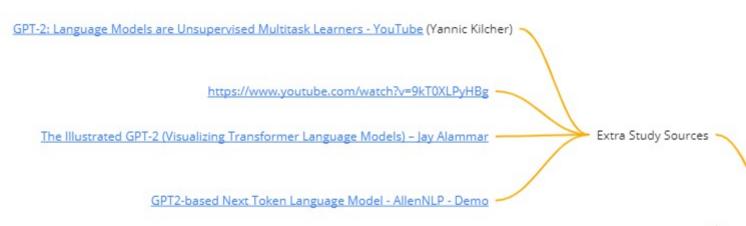


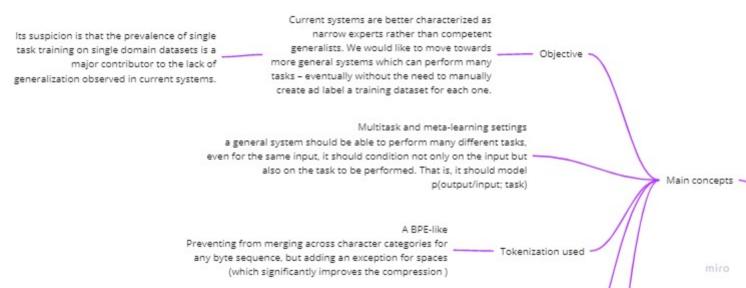
Article Presentation

Student: Marcus Vinícius Borela de Castro

Language models are unsupervised multitask learners

By: Radford, Alec, et al. OpenAl blog 1.8 (2019) Link





Language modeling is usually framed as unsupervised distribution

estimation from a set of examples (x1; x2; :::; xn) each composed of variable length

sequences of symbols (s1; s2; :::; sn)

▲ Nile regue | jalanmanathubis/Illustrate-opti/

This masking is often implemented as a matrix safed on attention mask. Think of a sequence of four words (including sold on the sequence is absorbed in four states—one per word lassuming for now that every word is a bovent. As these models within betterbis, we can assume a batch size of 4 for this two model that the modes were models within betterbis.

		Feat	ures		Labels	In a masked self attention way			
Examples	position: 1	2	2	4					
t d	robot	must	obey	orders	must	_			
2	robet	must	obcy	orders	obey	by Jay Alammar (The Illustrated GPT-2)			
8	robet	must	obey	orders	orders				
ž.	robat	must	olowy	onters	CBDSC-				

four different sizes of GPT-2 models

| Model Size | Parameters | Layers | dmodel | |-------|----------|-------|------| | 117M | 117 million | 12 | 768 | | 345M | 345 million | 24 | 1024 | | 762M | 762 million | 36 | 1280 | | 1542M | 1542 million | 48 | 1600 |

The second one is similar in number of parameters to BERT. The first, to GPT-1.

The model used

Based on the decoder of Transformer (Vaswani et al., 2017), with its sel-attention layers. The model largely follows the details of the OpenAl GPT model (Radford et al., 2018) with some modifications as an increased context size from from 512 to 1024 tokens.



by Jay Alammar (The Illustrated GPT-2) miro

Test dataset: 5% (held-out sample of WebText
Metric: perplexity (ppl)
Models are underfit: ppl can improved given more training time
uilding as large and diverse a

Our approach motivates building as large and diverse a dataset as possible in order to collect natural language demonstrations of tasks in as varied of domains and contexts as possible.

We removed all Wikipedia documents from WebText since it is a common data source for other datasets and could complicate analysis due to overlapping training data with test evaluation tasks. Trainning dataset

miro

Trainning

It demonstrates language models (if it is large and if it is trained on a sufficiently large and diverse dataset) can perform down-stream tasks in a zero-shot setting without any parameter or architecture modification.

It produces the WebText "cleaned" dataset (8 million documents for a total of 40 GB of text)

Article contribution

It achieves promising, competitive, and state of the art results depending on the task.

SOTA in 7 out of 8 tested language modeling task datasets

Interesting/unexpected results

	LAMBADA (PPL)	LAMBADA (ACC)	(ACC)	CBT-NE (ACC)	WikiText2 (PPL)	PTB (PPL)	enwik8 (BPB)	(BPC)	WikiText103 (PPL)	(PPL)
SOTA	99.8	59.23	85.7	82.3	39.14	46.54	0.99	1.08	18.3	21.8
117M	35.13	45.99	87.65	83.4	29.41	65.85	1.16	1.17	37.50	75.20
345M	15.60	55.48	92,35	87.1	22.76	47.33	1.01	1.06	26.37	55.72
762M	10.87	60.12	93,45	88.0	19.93	40.31	0.97	1.02	22.05	44.575
1542M	8.63	63.24	93.30	89.05	18.34	35.76	0.93	0.98	17.48	42.16

Table 3. Zero-shot results on many datasets. No training or fine-tuning was performed for any of these results. PTB and WikiText-2 results are from (Gong et al., 2018). CBT results are from (Bajgar et al., 2016). LAMBADA accuracy result is from (Hoang et al., 2018) and LAMBADA perplexity result is from (Grave et al., 2016). Other results are from (Dai et al., 2019).

Analysis suggests that data overlap between WebText training data and specific evaluation datasets provides a small but consistent benefit to reported results However, for most datasets we do not notice significantly larger overlaps than those already existing between standard training and test sets,

On reading comprehension the performance of GPT-2 is competitive with supervised baselines in a zero-shot setting. However, on other tasks such as summarization, while it is qualitatively performing the task, its performance is still only rudimentary according to quantitative metrics.

what are de-tokenize

Basic (or advanced) doubts

that may arise

De-tokenizers are techniques or processes used to reverse the process of tokenization in natural language processing. Tokenization is the process of breaking down a sequence of text into smaller units called tokens, and it is often used as a preprocessing step before applying machine learning models to natural language processing tasks. However, tokenization can introduce artifacts or inconsistencies that can negatively affect the performance of the model.

De-tokenizers remove or undo the tokenization process by combining the tokens back into their original format. This allows the model to better capture the natural language patterns and nuances present in the original text. The use of de-tokenizers can result in improved performance of language models, and the text mentions that the use of invertible detokenizers can still calculate the log probability of a dataset, allowing for a simple form of domain adaptation. (by ChatGPT)

The text reports that the use of de-tokenizers resulted in gains of 2.5 to 5 perplexity for GPT-2, which indicates an improvement in the model's ability to predict the next word in a sequence. (by ChatGPT)