


Pre-training LLMs using human-like development data corpus

Khushi Bhardwaj, Raj Sanjay Shah, Sashank Varma

Georgia Institute of Technology 
{khushi.bhardwaj, rajsanjayshah, varma}@gatech.edu

Abstract

Pre-trained Large Language Models (LLMs) have shown success in a diverse set of language inference and understanding tasks. The pre-training stage of LLMs looks at a large corpus of raw textual data. The BabyLM shared task compares LLM pre-training to human language acquisition, where the number of tokens seen by 13-year-old kids is magnitudes smaller than the number of tokens seen by LLMs. In this work, we pre-train and evaluate LLMs on their ability to learn contextual word representations using roughly the same number of tokens as seen by children. We provide a strong set of baselines; with different architectures, evaluation of changes in performance across epochs, and reported pre-training metrics for the strict small and strict tracks of the task. We also try to loosely replicate the RoBERTa baseline given by the task organizers to observe the training robustness to hyperparameter selection and replicability. We provide the submission details to the strict and strict-small tracks in this report.

1 Introduction

Transformer-based LLMs (Vaswani et al., 2017) show state-of-the-art performance on a variety of language processing tasks. In the last few years, pre-training methods for LLMs have evolved rapidly to meet task-driven demands. This evolution has focused on model expansion (Brown et al., 2020), more pre-training data (Hoffmann et al., 2022), use of higher quality data (Raffel et al., 2019), model alignment (von Werra et al., 2020), quicker run-time inference (Sanh et al., 2020), quicker pre-training (Clark et al., 2020), faster fine-tuning (Sanh et al., 2020), domain adaptation (Alsentzer et al., 2019; Caselli et al., 2021; Beltagy et al., 2019; Shah et al., 2022), and the addition of multi-modal capabilities (OpenAI, 2023; Gatti et al., 2022). The task-driven nature of this development optimizes performance at scale but fails to account for human-like learning.

Humans typically encounter fewer than 100 million tokens through language exposure by the time they are 13 years old (Warstadt et al., 2023). LLMs, on the other hand, parse tens of billions to trillions of tokens in their pre-training stage, typically from sources like Wikipedia (Wikipedia contributors, 2004), and Open Book Corpus (Zhu et al., 2015), which consist of different tokens than the ones seen by children. In this paper, we evaluate the capabilities of popular architectures on various tasks when trained on a number of tokens comparable to that seen by 13-year-old children. Such scaled-down pre-training has several potential benefits:

- A better sandbox for the development of new LLM training techniques inspired by the cognitive science literature (Yiu et al., 2023).
- Robust evaluation of models on human behavioral signatures (Shah et al., 2023).
- Building plausible human cognition models using LLMs aligned to actual human actions (Park et al., 2022).

Track	Data size	Datasets	Our work
Strict-small	10M words	Child-directed speech, transcribed speech from multiple sources, children’s books, and Wikipedia, etc.	✓
Strict	100M words		✓
Loose	100M words	Strict track data + un-limited non-linguistic data	×

Table 1: Task Summary

1.1 Task Descriptions

The shared task has three tracks: Strict, Strict-small, and Loose. The details of each track are summarized in Table 1. The Strict and Strict-small tracks

use pre-released datasets containing Child-directed speech, transcribed speech from multiple sources, children’s books, and Wikipedia. These tracks are meant to encourage explorations of architectural variation and self-supervised approaches.

1.2 Key Contributions

Given the benefits of using scaled-down human-like pre-training data, our work focuses on the following aspects of the shared task:

1. Replication details: Can we replicate the results of the baselines given by the task organizers?
2. Can we understand the impact of more training epochs on the same architecture?
3. Providing each training checkpoint for the different model architectures to facilitate future modeling of development. All checkpoints can be found [here](#).

We provide details of training and evaluation for the strict and strict-small tracks of this task.

2 Related Work

2.1 Cognitive science driven LLM architecture development

With the efforts put into LM pre-training, learning frameworks informed by cognitive science have received increasing attention. For instance, unsupervised and adversarial pre-training methods have been employed to enhance the logical reasoning capabilities of language models (Pi et al., 2022b). Using pre-training to inject numerical (Pi et al., 2022a) and commonsense reasoning (Zhong et al., 2019) has also been recently explored. Huebner et.al have constructed pre-training paradigms using curriculum learning (Huebner et al., 2021) to show the advantages of incremental learning.

2.2 Pre-training with limited data

Previous experiments show that pre-training data size is positively correlated with the syntactic capabilities of RoBERTa in terms of generalization and robustness (Pérez-Mayos et al., 2021). However, it has been discovered that model performance gains bring a high financial and environmental cost (Tay et al., 2021). This justifies the appeal of small-scale pretraining with data limitations. There have also been explorations of how human-like data

scales could improve our understanding of language acquisition and solidify current cognitive models (Dupoux, 2018).

Track	Model	Competition Scores (Dynabench)	Perplexity
Strict Small	Distilbert Epoch 20	0.62	86.283
	Distilbert Epoch 60	0.65	17.278
	RoBERTa Epoch 20	0.58	49.586
	GPT2 Epoch 20	0.64	79.318
	Competition Max	0.73	
Strict	Distilbert Epoch 20	0.66	39.427
	Distilbert Epoch 60	0.71	10.332
	RoBERTa Epoch 20	0.63	27.566
	GPT2 Epoch 20	0.67	34.950
	Competition Max	0.81	

Table 2: Model scores on dynabench

3 Methodology

3.1 Models

We use the simple-transformers library (Rajapakse, 2019) to pre-train the models below from scratch. The library uses the Huggingface trainer for pre-training. Note: We build new vocabularies for all models and limit the number of training epochs due to computational constraints in certain models.

- RoBERTa: We train the RoBERTa-base model (Liu et al., 2019) for comparison to the baseline given by the task organizers. This model is trained for 20 epochs on both datasets (strict and strict-small). The size of this model is roughly 125M parameters.
- DistilBert (uncased): Because this model (Sanh et al., 2020) is smaller (roughly 66M parameters) and quicker to pre-train, we additionally train it for 60 epochs. This allows us to explore the impact of more training epochs on performance.
- GPT2: We include a decoder-based architecture (Radford et al., 2019) in our pre-training to explore the impact of architecture type on the evaluation tasks. This model has a similar size to RoBERTa (117M parameters). We train it for 20 epochs due to computational constraints.

All of the checkpoints for the three architectures and the two tracks are uploaded on Huggingface (Wolf et al., 2020). **Hyperparameters:** We perform a grid search over the hyperparameters for all three architecture types. We use a subset of 0.5 GB of the training data for the search. The learning

Tasks	Model	Super GLUE										
		CoLA	SST-2	MRPC (F1)	QQP (F1)	MNLI	MNLI-mm	QNLI	RTE	BoolQ	MultiRC	WSC
Strict Small	Majority label	69.50	50.20	82.00	53.10	35.70	35.70	35.40	53.10	50.50	59.90	53.20
	OPT-125m	64.60	81.90	72.50	60.40	57.60	60.00	61.50	60.00	63.30	55.20	60.20
	RoBERTa-base	70.80	87.00	79.20	73.70	73.20	74.00	77.00	61.60	66.30	61.40	61.40
	T5-base	61.20	78.10	80.50	66.20	48.00	50.30	62.00	49.40	66.00	47.10	61.40
	Distilbert Epoch 20	69.38	83.46	79.69	80.21	69.80	71.56	60.15	54.55	65.42	53.67	51.81
	Distilbert Epoch 60	69.68	85.63	78.81	82.28	71.62	73.11	76.73	60.61	67.77	56.74	61.45
	RoBERTa Epoch 20	65.55	81.30	79.71	76.37	65.16	65.82	62.73	56.57	62.38	44.91	61.45
	GPT2 Epoch 20	69.58	83.07	75.47	73.13	63.88	65.95	59.84	56.57	64.45	58.38	46.99
Strict	OPT-125m	73.70	86.60	82.10	77.80	70.10	71.90	80.10	67.70	66.00	61.10	59.00
	RoBERTa-base	75.90	88.60	80.50	78.50	68.70	78.00	82.30	51.50	59.90	61.30	61.40
	T5-base	76.30	88.00	85.90	79.70	71.50	74.00	83.10	60.60	69.00	62.40	60.20
	Distilbert Epoch 20	69.48	86.22	62.98	83.81	73.44	74.97	79.00	60.61	67.91	62.98	44.58
	Distilbert Epoch 60	74.78	87.01	81.40	84.37	74.95	75.27	80.97	55.56	65.56	65.83	61.45
	RoBERTa Epoch 20	67.81	84.06	82.00	82.12	72.22	73.19	77.17	53.54	60.30	51.48	38.55
	GPT2 Epoch 20	69.58	87.20	79.29	82.23	74.00	74.98	81.01	52.53	69.58	57.83	48.19

Table 3: Results for the Super GLUE tasks

Tasks	Model	Blimp											
		Anaphor Agr.	Agr. Structure	Binding Binding	Control/Raising	D-N Agr.	Ellipsis	Filler-Gap	Irregular Forms	Island Effects	NPI Licensing	Quantifiers	S-V Agr.
Strict Small	OPT-125m	63.8	70.6	67.1	66.5	78.5	62	63.8	67.5	48.6	46.7	59.6	56.9
	RoBERTa-base	81.5	67.1	67.3	67.9	90.8	76.4	63.5	87.4	39.9	55.9	70.5	65.4
	T5-base	68.9	63.8	60.4	60.9	72.2	34.4	48.2	77.6	45.6	47.8	61.2	65
	Distilbert Epoch 20	83.49	64.12	63.98	62.22	77.72	62.76	62.36	85.24	42.94	41.38	67.47	55.81
	Distilbert Epoch 60	89.62	68.44	64.08	65.20	89.70	81.64	63.57	89.92	39.69	44.58	66.20	78.09
	RoBERTa Epoch 20	84.76	60.54	67.97	60.69	56.47	52.25	65.48	64.53	54.22	52.51	52.42	66.63
	GPT2 Epoch 20	81.24	72.56	67.81	67.43	86.98	59.82	67.72	84.38	52.62	51.76	58.14	64.12
	Strict	OPT-125m	94.9	73.8	73.8	72.2	93.1	80.5	73.6	80.8	57.8	51.6	74.5
RoBERTa-base		89.5	71.3	71	67.1	93.1	83.8	68	89.6	54.5	66.3	70.3	76.2
T5-base		66.7	61.2	59.4	59.8	53.8	49.1	70	75.5	43.6	45.6	34.2	53.2
Distilbert Epoch 20		92.43	67.06	67.66	65.27	94.38	87.24	65.42	85.04	42.86	50.43	67.41	66.25
Distilbert Epoch 60		94.68	70.39	68.39	68.25	96.39	89.03	68.69	90.08	45.59	64.67	70.20	72.32
RoBERTa Epoch 20		85.94	67.68	65.27	63.74	91.04	75.52	62.98	87.23	46.41	44.47	61.46	60.51
GPT2 Epoch 20		91.56	74.88	73.21	69.22	91.89	75.52	71.91	75.32	55.04	51.20	66.13	67.19

Table 4: Results for the Blimp tasks

Tasks	Model	Blimp Supplement Tasks				
		Hypernym	QA Congruence (easy)	QA Congruence (tricky)	Subj.-Aux. Inversion	Turn Taking
Strict Small	OPT-125m	50.00	54.7	31.5	80.3	57.1
	RoBERTa-base	49.4	31.3	32.1	71.7	53.2
	T5-base	48	40.6	21.2	64.9	45
	Distilbert Epoch 20	50.00	65.63	42.42	77.31	61.79
	Distilbert Epoch 60	48.95	70.31	41.21	60.87	62.86
	RoBERTa Epoch 20	51.28	48.44	31.52	53.86	66.07
	GPT2 Epoch 20	47.44	48.44	45.45	72.41	62.86
	Strict	OPT-125m	46.3	76.50	47.9	85.3
RoBERTa-base		50.8	34.4	34.5	45.6	46.8
T5-base		51.1	45.3	25.5	69.2	48.9
Distilbert Epoch 20		48.26	64.06	40.61	81.53	65.36
Distilbert Epoch 60		48.95	73.44	47.88	83.43	65.36
RoBERTa Epoch 20		51.16	46.88	37.58	76.85	64.29
GPT2 Epoch 20		49.53	57.81	45.45	81.85	65.00

Table 5: Results for the Blimp supplemental tasks

rate ranges from $5e-5$ to $4e-4$ across the searches, with weight decay in place but no early stopping mechanisms.

4 Results

Table 2 shows the results obtained from the dynabench submission portal. The individual results for each of the tasks in different benchmarks are available in Tables 3, 4, 5, 6, 7. Looking at these

tables, we observe the following patterns:

1. We see that training for more epochs leads to better overall performance (compare 20 and 60 epochs of DistilBert in Table 2).
2. Variation among architecture types exists when limiting the training to the same number of epochs, but it is difficult to identify a definitively better architecture.

Tasks	Model	MSGS Tasks										
		CR (Control)	LC (Control)	MV (Control)	RP (Control)	SC (Control)	CR_LC	CR_RTP	MV_LC	MV_RTP	SC_LC	SC_RP
Strict-Small	OPT-125m	86.40	86.10	99.80	100.00	94.30	66.50	67.00	66.50	67.60	80.20	67.50
	RoBERTa-base	84.10	100.00	99.40	93.50	96.40	67.70	68.60	66.70	68.60	84.20	65.70
	T5-base	78.40	100.00	72.70	95.50	94.40	66.70	69.70	66.60	66.90	73.60	67.80
	Distilbert Epoch 20	79.22	100.00	97.17	98.57	96.36	66.53	66.71	66.61	67.47	67.89	67.58
	Distilbert Epoch 60	81.68	100.00	98.61	99.14	95.66	67.24	66.72	66.61	67.03	67.76	68.27
	RoBERTa Epoch 20	73.02	100.00	73.91	99.59	86.47	66.70	67.19	66.61	66.84	67.44	71.93
	GPT2 Epoch 20	89.78	96.30	99.23	100.00	97.13	66.46	66.72	66.58	66.83	78.78	64.87
Strict	OPT-125m	97.20	82.60	100.00	99.80	88.10	75.30	67.10	66.30	66.80	84.80	62.00
	RoBERTa-base	93.00	100.00	100.00	100.00	89.00	68.30	66.80	66.60	80.20	67.40	67.40
	T5-base	95.10	100.00	100.00	99.80	88.70	76.70	69.40	67.00	67.70	72.70	68.00
	Distilbert Epoch 20	81.44	100.00	97.36	97.35	94.77	67.26	66.72	66.61	66.97	67.67	68.63
	Distilbert Epoch 60	93.23	100.00	99.33	99.17	95.64	68.91	66.77	66.61	67.45	67.89	66.59
	RoBERTa Epoch 20	84.63	97.38	92.12	98.15	95.54	66.47	66.59	66.41	66.05	68.17	72.78
	GPT2 Epoch 20	95.35	76.53	99.55	99.83	96.76	67.21	68.46	66.78	66.70	91.90	65.90

Table 6: Results for the MSGS tasks

Tasks	Model	Age of Acquisition tasks (mean absolute deviation)			
		Overall (591 words)	Nouns (322)	Predicates (167)	Function words (102)
Strict Small	OPT-125m	2.03	1.98	1.81	2.57
	RoBERTa-base	2.06	1.99	1.85	2.65
	T5-base	2.04	1.97	1.82	2.64
	Distilbert Epoch 20	2.06	2.00	1.84	2.65
	Distilbert Epoch 60	2.09	2.00	1.84	2.76
	RoBERTa Epoch 20	2.06	2.00	1.84	2.63
	GPT2 Epoch 20	2.06	2.00	1.85	2.64
Strict	OPT-125m	2.04	1.97	1.83	2.61
	RoBERTa-base	2.06	1.99	1.82	2.66
	T5-base	2.06	2.00	1.83	2.65
	Distilbert Epoch 20	2.06	2.00	1.83	2.65
	Distilbert Epoch 60	2.08	2.00	1.81	2.79
	RoBERTa Epoch 20	2.06	2.00	1.84	2.62
	GPT2 Epoch 20	2.04	1.98	1.81	2.60

Table 7: Results for the Age of Acquisition tasks

- Tables 3, 4, 5, 6, and 7 show that pre-training (RoBERTa) is not robust to initialization, and the competition scores would greatly benefit from a warm-up or a grid search over different hyper-parameters.
- In most cases, the pre-training improves performance over the majority label in the Super GLUE tasks.
- Tables 8, 9 shows that the performance on the BLIMP tasks becomes better with more training epochs. While this is orthogonal to wisdom performance saturates at one epoch (Biderman et al., 2023). Our results hint that training saturation or stability may be a function of model size divided by the number of tokens seen.

5 Conclusions

We pre-train popular LLM architectures on the kind of textual data seen by children when they are

around 13 years old. We show that pre-training paradigms like Masked Language Modeling or Causal Language Modeling lead to only minor variations. Our results show that models are not robust to the initialization of weights. Our work provides each and every checkpoint of the model architectures on Huggingface to facilitate future research. All checkpoints can be found [here](#).

6 Limitations

Our work trains some of the popular Language Model architectures on human-like scaled-down training data, it does not introduce new training methodologies or architectures which may be better suited for such tasks. Furthermore, our work does not exhaustively cover different model types in the literature. Our results are preliminary as they do not account for all possible confounds.

BLIMP Tasks	Epochs																			
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Anaphor Agr.	50.77	70.81	82.21	79.75	84.46	84.00	86.71	88.14	87.27	87.88	87.42	87.32	83.54	86.91	86.30	84.76	85.22	85.99	85.22	85.94
Agr. Str.	58.69	59.81	60.46	60.83	59.80	60.39	60.43	60.84	61.02	60.60	59.69	62.08	63.70	65.22	65.39	66.56	67.54	67.34	67.92	67.68
Binding	64.43	68.06	64.84	65.32	64.31	67.17	64.31	66.00	63.92	64.10	64.00	59.45	61.52	62.60	62.72	63.67	64.54	65.21	65.18	65.27
Control Rais.	59.17	59.19	58.17	58.24	59.06	59.01	59.79	59.70	59.70	60.01	58.51	62.22	60.25	59.41	60.87	62.88	62.79	63.57	63.85	63.74
Det-N Agr.	51.47	56.72	58.87	59.57	58.63	58.62	58.18	60.09	60.16	59.06	59.63	63.47	71.69	80.95	85.20	87.18	88.49	90.51	90.98	91.04
Ellipsis	37.93	44.05	46.54	50.87	50.98	54.27	56.12	58.43	59.70	61.61	55.02	51.96	54.62	65.47	68.30	73.85	72.92	75.64	75.23	75.52
Filler Gap	69.39	64.53	66.48	63.66	65.38	61.95	61.31	61.95	61.84	65.39	60.58	61.86	60.21	62.14	61.47	61.80	62.40	63.35	63.38	62.98
Hypernym	53.84	49.77	52.09	51.74	48.02	48.49	48.72	50.12	49.88	51.28	50.35	50.70	48.84	50.35	51.28	49.77	51.51	50.47	51.86	51.16
Irr. Forms	45.90	59.75	60.87	61.32	63.66	59.54	65.24	63.36	64.83	64.99	69.97	78.07	76.39	81.53	86.92	86.87	88.96	88.85	87.74	87.23
Island Effects	53.66	42.68	56.20	55.53	54.67	44.66	48.28	51.91	49.93	52.88	47.31	45.22	51.76	50.67	47.83	48.92	46.82	46.38	45.40	46.41
NPI Lic.	34.89	43.58	35.71	46.96	40.81	43.41	43.29	42.97	39.60	44.02	37.25	38.23	39.20	40.80	40.10	43.77	44.05	44.46	43.65	44.47
QA_cong. Easy	34.38	35.94	39.06	42.19	39.06	35.94	40.63	40.63	37.50	37.50	37.50	46.88	56.25	57.81	56.25	50.00	45.31	45.31	48.44	46.88
QA_cong. Tricky	38.18	34.55	32.73	31.52	31.52	30.91	30.30	29.70	28.48	28.48	27.27	24.85	23.64	30.30	33.94	32.73	35.76	35.15	38.18	37.58
Quantifiers	37.92	38.15	36.22	40.96	43.07	41.96	49.00	50.05	43.07	51.91	51.52	66.95	65.35	67.85	60.72	64.30	61.98	64.09	61.39	61.46
S-Aux Inv.	74.68	75.04	69.85	68.87	63.19	52.74	55.40	50.65	47.89	48.65	52.28	62.94	66.85	72.82	70.36	77.92	73.82	76.38	75.43	76.85
S-V Agr.	49.67	50.89	52.10	51.91	52.34	53.68	53.80	54.47	55.54	55.83	54.87	52.72	54.27	57.58	58.10	58.70	60.09	60.31	60.23	60.51
Turn-Taking	59.64	59.64	60.00	59.29	61.43	60.36	58.57	59.64	58.93	58.57	61.07	63.57	64.64	67.14	63.57	63.57	64.29	63.93	64.64	64.29

Table 8: Results for the BLIMP tasks across different epochs of the RoBERTa-base model architecture for the strict (100M token) track.

Behavior/ Model +Epoch	Epochs												
	1	5	10	15	20	25	30	35	40	45	50	55	60
Anaphor Agr.	46.57	82.87	89.88	91.21	92.43	93.10	94.07	94.17	95.19	94.94	94.58	94.43	94.68
Agr. Str.	58.06	59.71	61.78	65.69	67.06	68.02	70.05	69.07	69.67	70.55	70.49	70.27	70.39
Binding	59.65	65.24	63.15	67.14	67.66	66.93	68.48	66.55	69.07	68.76	68.95	68.27	68.39
Control Rais.	58.33	58.93	60.01	64.14	65.27	66.00	65.91	67.12	67.30	67.41	68.10	67.87	68.25
Det-N Agr.	50.76	60.30	70.41	92.16	94.38	95.24	95.94	95.97	96.34	96.14	96.27	96.37	96.39
Ellipsis	37.53	54.16	55.08	81.99	87.24	86.49	86.20	89.32	89.32	89.38	88.57	88.86	89.03
Filler Gap	70.23	64.89	58.56	62.06	65.42	64.74	66.64	67.49	67.54	67.24	69.00	68.88	68.69
Hypernym	51.40	50.23	50.70	48.84	48.26	50.00	48.60	51.40	50.23	50.00	49.77	48.49	48.95
Irr. Forms	56.39	65.24	87.38	85.55	85.04	86.92	88.50	89.16	88.85	88.85	89.72	89.72	90.08
Island Effects	46.52	44.62	48.09	45.52	42.86	45.07	43.20	46.49	44.81	43.80	44.39	45.44	45.59
NPI Lic.	53.23	46.90	41.06	46.67	50.43	55.25	58.56	57.39	61.69	64.36	64.09	64.15	64.67
QA_cong. Easy	31.25	43.75	59.38	67.19	64.06	68.75	70.31	73.44	75.00	70.31	70.31	73.44	73.44
QA_cong. Tricky	333.33	22.42	23.03	35.15	40.61	42.42	46.06	43.03	44.24	41.82	46.06	46.06	47.88
Quantifiers	54.87	69.55	62.31	65.43	67.41	70.40	70.50	72.82	70.74	70.63	70.25	70.81	70.20
S-Aux Inv.	58.45	65.77	73.65	79.21	81.53	81.19	81.85	81.75	83.17	82.80	83.63	82.53	83.43
S-V Agr.	48.93	54.60	55.56	62.06	66.25	68.24	70.82	70.05	71.64	72.50	71.73	72.41	72.32
Turn-Taking	59.29	60.71	65.36	64.29	65.36	63.93	64.64	65.36	65.36	65.00	66.07	65.00	65.36

Table 9: Results for the BLIMP tasks across different epochs of the DistilBERT-base model architecture for the strict (100M token) track.

7 Ethical Considerations

All researchers in this study have active responsible code of conduct in research certifications. The models shared on Huggingface have the same risks associated with any other Large Language Model. Researchers in this study have tried to be mindful of the environment while doing the pre-training runs and hope that publically available checkpoints will help other researchers avoid computation and environmental costs associated with repeat pre-training.

8 Computational Resources

The models are trained on Nvidia-RTX 2080 GPUs with 12 GB RAM. The models are trained for nearly 975 GPU hours.

References

- Emily Alsentzer, John Murphy, William Boag, Weihung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. 2019. [Publicly available clinical BERT embeddings](#). In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [Scibert: Pretrained language model for scientific text](#). In *EMNLP*.
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. [Pythia: A suite for analyzing large language models across training and scaling](#).
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda

- Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- Tommaso Caselli, Valerio Basile, Jelena Mitrović, and Michael Granitzer. 2021. [Hatebert: Retraining bert for abusive language detection in english](#).
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *ArXiv*, abs/2003.10555.
- Emmanuel Dupoux. 2018. [Cognitive science in the era of artificial intelligence: A roadmap for reverse-engineering the infant language-learner](#). *Cognition*, 173:43–59.
- Prajwal Gatti, Anand Mishra, Manish Gupta, and Mithun Das Gupta. 2022. [VisToT: Vision-augmented table-to-text generation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9936–9949, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. [Training compute-optimal large language models](#).
- Philip A. Huebner, Elior Sulem, Fisher Cynthia, and Dan Roth. 2021. [BabyBERTa: Learning more grammar with small-scale child-directed language](#). In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 624–646, Online. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, M. Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *ArXiv*, abs/1907.11692.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Joon Sung Park, Lindsay Popowski, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2022. [Social simulacra: Creating populated prototypes for social computing systems](#).
- Laura Pérez-Mayos, Miguel Ballesteros, and Leo Wanner. 2021. How much pretraining data do language models need to learn syntax? *arXiv preprint arXiv:2109.03160*.
- Xinyu Pi, Qian Liu, Bei Chen, Morteza Ziyadi, Zeqi Lin, Qiang Fu, Yan Gao, Jian-Guang Lou, and Weizhu Chen. 2022a. [Reasoning like program executors](#).
- Xinyu Pi, Wanjun Zhong, Yan Gao, Nan Duan, and Jian-Guang Lou. 2022b. Logigan: Learning logical reasoning via adversarial pre-training. *Advances in Neural Information Processing Systems*, 35:16290–16304.
- Alec Radford, Jeff Wu, R. Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. [Exploring the limits of transfer learning with a unified text-to-text transformer](#).
- T. C. Rajapakse. 2019. Simple transformers. <https://github.com/ThilinaRajapakse/simpletransformers>.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. [Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter](#).
- Raj Sanjay Shah, Kunal Chawla, Dheeraj Eidnani, Agam Shah, Wendi Du, Sudheer Chava, Natraj Raman, Charese Smiley, Jiaao Chen, and Diyi Yang. 2022. When flue meets flang: Benchmarks and large pretrained language model for financial domain. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- Raj Sanjay Shah, Vijay Marupudi, Reba Koenen, Khushi Bhardwaj, and Sashank Varma. 2023. [Human behavioral benchmarking: Numeric magnitude comparison effects in large language models](#).
- Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won Chung, Sharan Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. 2021. Scale efficiently: Insights from pre-training and fine-tuning transformers. *arXiv preprint arXiv:2109.10686*.
- Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *ArXiv*, abs/1706.03762.
- Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, and Nathan Lambert. 2020. [trl: Transformer reinforcement learning](https://github.com/lvwerra/trl). <https://github.com/lvwerra/trl>.
- Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Gotlieb Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mosquera, Adina Williams, Bhargavi Paranjabe, Tal Linzen, and Ryan Cotterell. 2023. Findings of the 2023 BabyLM Challenge: Sample-efficient pretraining on developmentally plausible corpora. In *Proceedings of the 2023 BabyLM Challenge*. Association for Computational Linguistics (ACL).

Wikipedia contributors. 2004. [Wikipedia, the free encyclopedia](#).

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Eunice Yiu, Eliza Kosoy, and Alison Gopnik. 2023. [Imitation versus innovation: What children can do that large language and language-and-vision models cannot \(yet\)?](#)

Wanjun Zhong, Duyu Tang, Nan Duan, Ming Zhou, Jiahai Wang, and Jian Yin. 2019. Improving question answering by commonsense-based pre-training. In *Natural Language Processing and Chinese Computing: 8th CCF International Conference, NLPCC 2019, Dunhuang, China, October 9–14, 2019, Proceedings, Part I 8*, pages 16–28. Springer.

Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *The IEEE International Conference on Computer Vision (ICCV)*.