

4  
2  
0  
2  
t  
c  
O  
7  
1  
l  
L  
C  
s  
c  
[  
0  
1  
v  
5  
3  
4  
6  
0  
7  
0  
3  
2  
v  
i  
X  
r  
a

<sup>j</sup>The University of Western Australia (UWA), Perth, Australia

1.

·

, , ,

가 가 . [1],

, 가

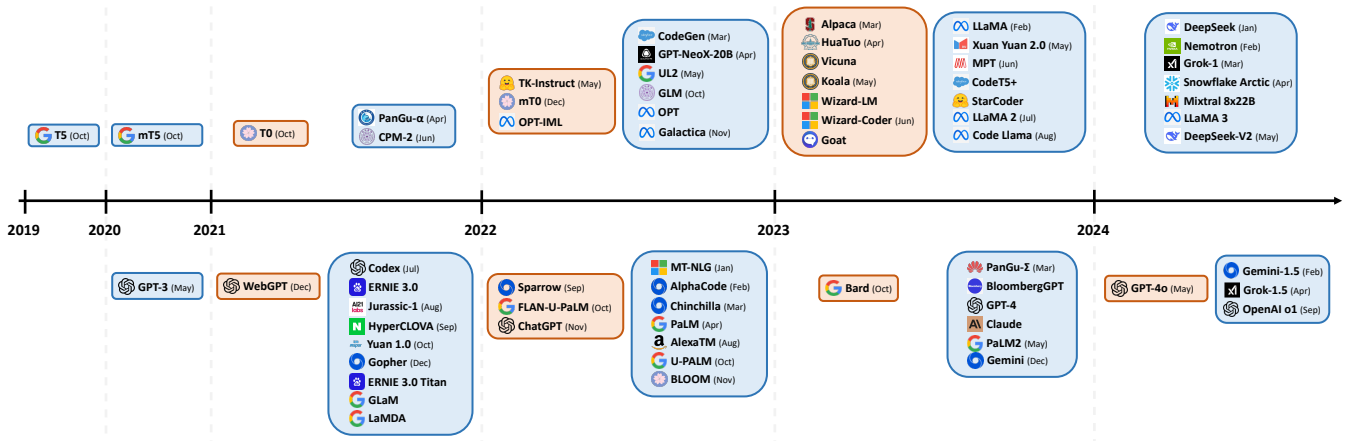
.

LLM

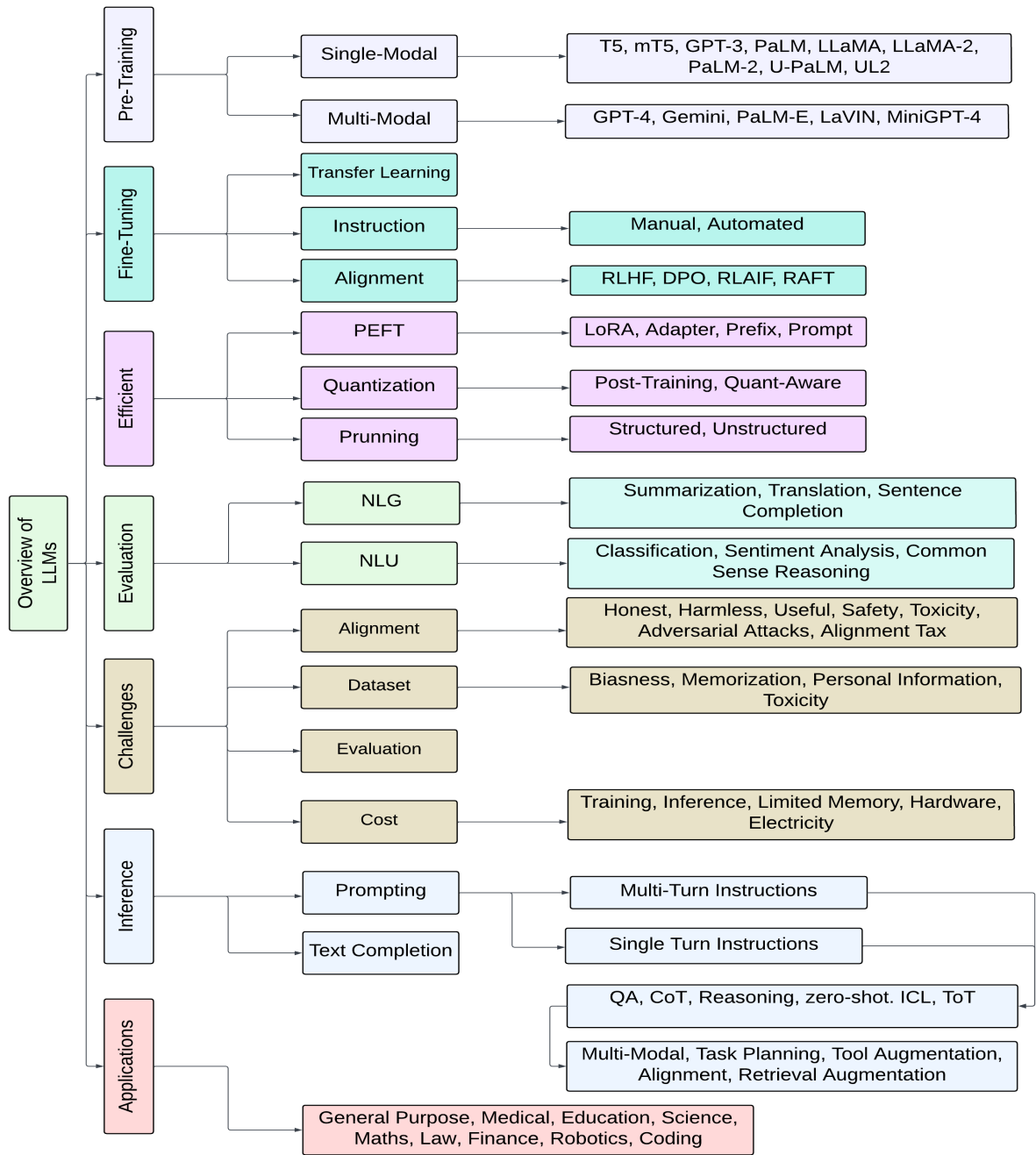
가 [2, 3].



October 18, 2024



2: LLM 가 , (LLM) , [25, 26, 27, 28, 29, 30, 31] [4]. [32] (NLP) LLM (P (LM) , PL 가 [7, 8, 9] 15, 33, 34, 35] [36, 37, 21, 38, 39, 40, 41] [38, 41, 40], 가 [42, 43], [44, 45], [46, 47, 48, 49] LLM 가 PLM PLM ( GB TB)[10, 11] LLM LLM [10, 11, 12, 6, 가 1 2 LLM LM 4, 55, 56, 57, 58] [50, 51, 52, 53] LLM T5[10] mT5[11] LLM GPT-3[6] , 가, LLM, LLM, 가 LLM few-shot zero-shot [16, 17, 18, 19] • LLM [20, 21] 가 zero-shot • LLM 가 LLM , , • M 가 [22, 23, 24]. L LLM LM , , LLM



3: LLM

, LLM 7가

: 1.

2.

3.

4.

5. 가 6.

7.

LLM,  
, 가 .

LLM, LLM

,

4

3.8

. LLM

가,

5

7 8

10B

가

, [50]

LL

M

[51, 52, 53]

3

2  
LLM

LLM

2.

LLM  
LLM

### 2.1. Tokenization [59]

LLM [61] 가 LLM [62], (BPE)[61] unigramLM[60] [63]

### 2.2. Encoding Positions

[64] 가 가 Alibi RoPE LLM 가

**Alibi [65]:** 가

**RoPE [66]:**

### 2.3. Attention in LLMs

가 [64] LLM 가 Self-Attention [64]: ( ) Cross Attention: Sparse Attention [67]:  $O(n^2)$  가 [67] Flash Attention [68]: GPU (HBM) SRAM GPU

### 2.4. Activation Functions

[69]. LLM **ReLU [70]:** (ReLU)  $ReLU(x) = \max(0, x)$  (1)

**GeLU [71]:** 가 (GeLU) ReLU, [72] [73] **GLU variants [74]:** [75] ( $\sigma$ ) ( $\otimes$ )

$$GLU(x, W, V, b, c) = (xW + b) \otimes \sigma(xV + c), \quad (2)$$

$X$  LLM  $l, W, b, V, c$  GLU [74]

$$\begin{aligned} ReGLU(x, W, V, b, c) &= \max(0, xW + b) \otimes, \\ GEGLU(x, W, V, b, c) &= GELU(xW + b) \otimes (xV + c), \\ SwiGLU(x, W, V, b, c, \beta) &= Swish\beta(xW + b) \otimes (xV + c). \end{aligned}$$

### 2.5. Layer Normalization

[64]. LayerNorm[76] RMSNorm[77] LLM (MHA) [78] LL M DeepNorm[79]

### 2.6. Distributed LLM Training

LLM [13, 37, 80, 81] **Data Parallelism:** 가 **Tensor Parallelism:** **Pipeline Parallelism:** **Model Parallelism:** **3D Parallelism:** , , 3D **Optimizer Parallelism:** [37]

2.7. Libraries

LLM

*Transformers* [82]:

API

*DeepSpeed* [36]: 가

*Megatron-LM* [80]: LLM GPU

*JAX* [83]: 가 Python Python Nu mPy GPU

*Colossal-AI* [84]:

*BMTrain* [81]: LLM

*FastMoE* [85]: PyTorch MoE( 가) API

*MindSpore* [86]: 가

*PyTorch* [87]: Facebook AI Research Lab(FAIR) PyTor ch

*Tensorflow* [88]: Google TensorFlow

*MXNet* [89]: Apache MXNet Python, C++, Scala, R

2.8. Data PreProcessing

LLM

*Quality Filtering*:

가

1) 2)

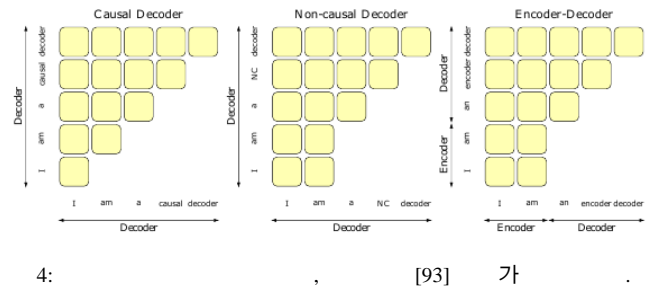
*Data Deduplication*:

가 LLM

*Privacy Reduction*: LLM 가 LLM

2.9. Architectures

LLM



4:

[93] 가

Full Language Modeling May the force be with you

Prefix Language Modeling May the force be with you

Masked Language Modeling May the force be with you

5:

[93].

Encoder Decoder:

Causal Decoder: 가

Prefix Decoder:

4

Mixture-of-Experts: 가

가

[90]. Mixture-of-Experts(MoE)

가

[91, 92].

2.10. Pre-Training Objectives

LLM

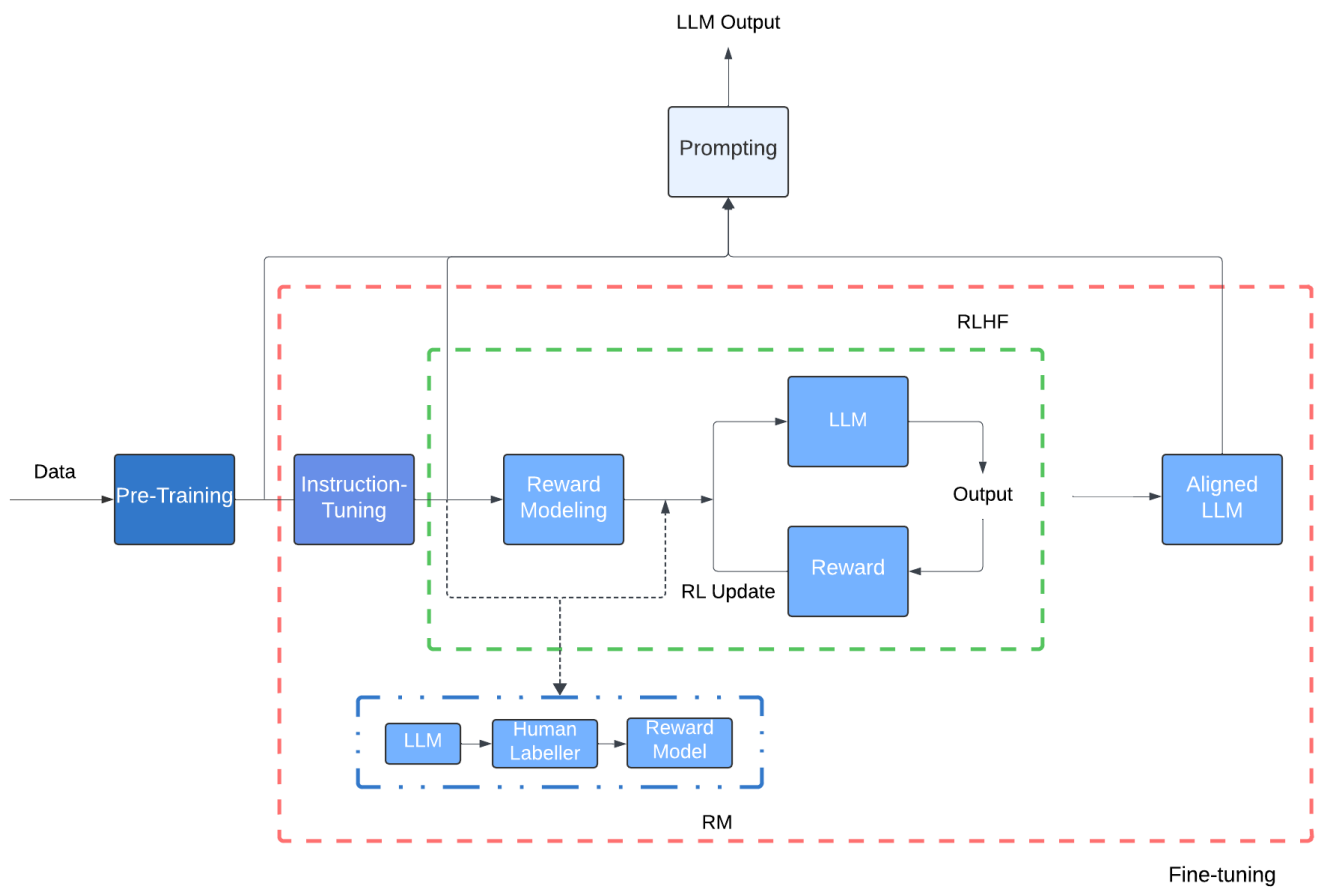
[93]

Full Language Modeling:

가 5

Prefix Language Modeling: 가

5 가



6: LLM / LLM 가 . "RL" , "RM" , "RLHF"

**Masked Language Modeling:**  
( )

## 2.12. LLMs Adaptation Stages

5 가 LLM LLM 가 6  
**Unified Language Modeling:** [94]  
가 ,

### 2.12.1. Pre-Training

## 2.11. LLMs Scaling Laws

. LLM 2.5, 2.4, 2.10

### 2.12.2. Fine-Tuning LLM

[95].

[96]

**Transfer Learning:** LLM  
[6, 15].

LLM

**Instruction-tuning:** [10, [16, 97] LLM

*Chain-of-Thought (CoT):*

CoT [55, 103, 101]

*Self-Consistency:* 가 [16, 50, 97] CoT [104].

*Tree-of-Thought (ToT):* 가 가 [105].

**Alignment-tuning:** LLM

LLM

LLM [20, 21798].

가 "HHH"

[99].

(RLHF)[100] RLHF

(RL) 가 (RM) RLHF R

M RL

**Reward modeling:** 3.

HHH L LLM

LM

**Reinforcement learning:**

vs. LLM

(PPO)

3.1. Pre-Trained LLMs

NLP

LLM NLU LLM NLG

2.12.3. Prompting/Utilization

6 LLM 1 2 LLM

LLM

가 [16, 101, 102].

가 [32]

**Zero-Shot Prompting:** LLM

LLM

**In-context Learning:** few-shot learning

few-shot learnin [106]

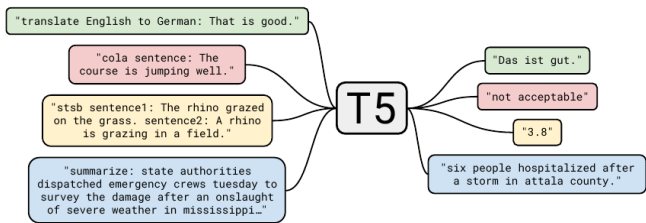
g (ICL) **GPT-3 [6]:** GPT-3 GPT-2 [5] S

[54, 50, 18, 16] parse Transformer [67] dense

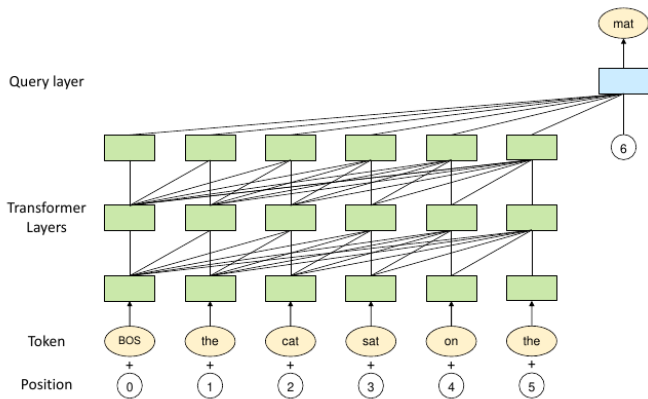
**Reasoning in LLMs:** LLM zero-shot

GPT-3 [107]

가 GPT-3 175B



7: [10] 가



8: PanGu-α [108]

**mT5 [11]:** 101 가 mC4  
T5 [10].  
250,000  
mT5

[40]

( , , )  
. CPM-2 1  
CPU

INFMOE

**ERNIE 3.0 [110]:** ERNIE 3.0 Transformer-XL[111]

. LLM  
LLM 가  
54 NLP

**Jurassic-1 [112]:** 7B J1-Large 178B  
J1-Jumbo

. Jurassic-1 가 , 가  
GPT-3

Jurassic-1  
[113]

**HyperCLOVA [114]:** GPT-3

**Yuan 1.0 [115]:** 5TB  
가 (MDFS) . Spark

an 1.0

가

Yu

**PanGu-α [108]:** 8 가

Eq. 3  
가

$$a = p_n W_h^q W_h^k T H_L^T \quad (3)$$

**CPM-2 [12]:** WuDaoCorpus [109] 가(MoE)  
M-2) 11B 198B

" "

MoE

CPM-2 가

**Gopher [116]:** Gopher LLM  
scale 가 44M~280B  
. 280B 가 81% GPT-3[  
6], Jurassic-1[112], MT-NLG[117] 가

**ERNIE 3.0 TITAN [35]:** ERNIE 3.0 Titan  
26 ERNIE 3.0

68 NLP  
. LLM

ERNIE 3.0 Titan  
Credible

8 and Controllable Generations 가



ERNIE 3.0 Titan

가

가

QA

가

LLM

**GPT-NeoX-20B [118]:**

가

GPT-3

Pile

.GP

T-NeoX Eq. 4

가

15%

가

[66]

[119]

25%

GPT-3

GPT-NeoX-20B

[6]

13B

175B

20B

GPU

160

.50

5000

7,000

Gopher(280B)

$$x + \text{Attn}(\text{LN}_1(x)) + \text{FF}(\text{LN}_2(x))$$

(4)

70B

4

**OPT [14]:** GPT-3

GPT-3  
[120]

. OPT

Gopher[116], GPT-3[6]

**AlexaTM [122]:**

가

100k

OPT-175B

GPT3-1

75B

**BLOOM [13]:** LLM

ROOTS

LM)

[CLM]

CLM

(C

9

, ALiBi

, bitsandbytes<sup>1</sup>

가

**PaLM [15]:**

: Eq. 4

15

S

**GLaM [91]:** Generalist Language Model(GLaM)

가(MoE)

[121, 90]

wiGLU

, RoPE

가

가

가

200~500

100

가

. 가

GLaM

GLaM(64B/6

4E) GPT-3[6]

7×

GLaM

가

가

540B

2.

GPT-3

. 가

GLaM(64B/64E)

GPT-3

4%

**PaLM-2 [123]:**

PaLM

. PaLM-2

**MT-NLG [117]:** GPT-2

3× GPT-3

530B

가

PaLM

PaLM

가

. MT-NLG

**Chinchilla [96]:** Gopher[116]

가

GPT-3

**U-PaLM [124]:**

[125]

0.1%

UL2(UL2Restore

PaLM

)

(MassiveText

).

Gopher

Adam

,

, CoT

NLP

,

. AdamW

. Chinchilla

Gopher

,

, CoT

NLP

,

PaLM

, 25%

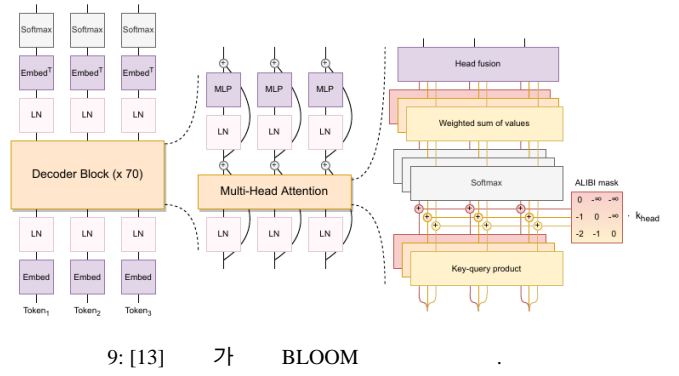
PaL

M

50%

, 25%

<sup>1</sup><https://github.com/TimDettmers/bitsandbytes>



<p><b>UL2 [125]:</b> Denoisers(MoD)</p> <p>1) R-Denoiser:</p> <p>2) S-Denoiser:</p> <p>3) X-Denoiser:</p> <p>Denoiser Denoising</p> <p>MoD</p> <p>T5</p> <p><b>GLM-130B [33]:</b> GLM-130B GLM[126]</p> <p>( )</p> <p>GPT-3 GLM</p> <p>GLM-130B</p> <p>5%)가</p> <p><b>LLaMA [127, 21]:</b> 7B 70B LLaMA</p> <p>가</p> <p><b>LLaMA-1 [127]:</b> 가 / [128]</p> <p>[129]</p> <p><b>LLaMA-2 [21]:</b> LLaMA-2-Chat</p> <p>40%</p> <p><b>LLaMA-3/3.1 [130]:</b> LLaMA-2 7 가 2</p> <p><b>PanGu-<math>\Sigma</math> [92]:</b> PanGu-<math>\alpha</math> RRE(Random Routed Experts) 1 10 , 가 RRE MoE , 가 2</p> <p><b>Mixtral8x22b [131]:</b> 8 가 (MoE) 가 가 MoE(128×3.66B (1 MLP 가) 가 MoE LLM[131, 133] A) 가 480B , 5.76 17B</p>	<p><b>Grok [133, 134]:</b> Grok XAI가 Grok-1 Grok-1.5 LLM</p> <p><b>Grok-1 [133]:</b> Grok-1 가가 314B MoE ( 가 8 )</p> <p><b>Grok-1.5 [134]:</b> Grok-1.5 LLM</p> <p><b>Gemini [135, 136]:</b> Gemini Bard(PaLM )</p> <p><b>Gemini-1 [135]:</b> MMLU</p> <p><b>Gemini-1.5 [136]:</b> MoE LLM Gemini-1 2M 가 10M</p> <p><b>Nemotron-4 340B [137]:</b> 98% 2%</p> <p><b>DeepSeek [138]:</b> DeepSeek LLM</p> <p>1e<sup>17</sup> 3e<sup>20</sup> FLOP 8 10</p> <p>/</p> <p>가 가</p> <p>(B), (<math>\eta</math>), (M) (D)</p> $B_{opt} = 0.2920.C^{0.3271}$ $\eta_{opt} = 0.3118.C^{-0.1250}$ $M_{opt} = M_{base}.C^a \quad (5)$ $D_{opt} = D_{base}.C^b$ <p><math>M_{base} = 0.1715, D_{base} = 5.8316, a = 0.5243, b = 0.4757</math></p> <p><b>DeepSeek-v2 [139]:</b> - (KV) (MLA) (MH A) MoE . MLA (GQA), (MQA)</p> <p>. MLA DeepSeek-v2 DeepSeek [138].</p>
--	--

3.1.2. Coding

**CodeGen [140]:** CodeGen PaLM[15], MLP, RoPE, HumanEval, MBPP, PaLM, LLaMA, LAMDA.

(1) PILE, 2) BIGQUERY, 3) BIGPYTHON  
CodeGen

가 가 . Co Multi-Tur deGen n Programming Benchmark(MTPB)

**Codex [141]:** LLM docstring Python Github

가 Codex 100 77.5% Github Copilot<sup>2</sup>

**AlphaCode [142]:** 300M~41B [143]

가 Alpha-Code Hub CodeContests Git Codeforces 3 . Code Contests [145] GOLD[144] CodeContests . AlphaCode 가 Codeforces AlphaCode 5,000 54.3% , Codefor ces 28% .

**CodeT5+ [34]:** CodeT5+ CodeT5[146], ( ) , ( - ) 가 CLM 가 CLM . CodeT5+ 가 [CLS], -

**StarCoder [147]:** SantaCoder , Flash attention 8k . StarCoder

3.1.3. Scientific Knowledge

**Galactica [148]:** 4,800 , , , Py Torch fairscale[149] metaseq 3 < work >

3.1.4. Dialog

**LaMDA [150]:** , , 90% . LaMDA , , . LaMDA

3.1.5. Finance

**BloombergGPT [151]:** (Bloomberg "FINPILE") BLOOM [13] OPT [14] [113] 50B BloombergGPT < |endof text| > , 1024 2048

**Xuan Yuan 2.0 [152]:** BLOOM [13] , , . Xuan Yuan 2.0

3.2. Fine-Tuned LLMs

LLM 가 LLM [20]. LL [16, 17, 97] [20]. 가 가 [97, 16, 18]가 가 0.2% PaLM 540B [16]. LLM

<sup>2</sup><https://github.com/features/copilot>  
<sup>3</sup><https://codeforces.com/>

Models	Findings & Insights
T5	<ul style="list-style-type: none"> <li>Encoder and decoder with shared parameters perform equivalently when parameters are not shared</li> <li>Fine-tuning model layers (adapter layers) work better than the conventional way of training on only classification layers</li> </ul>
GPT-3	<ul style="list-style-type: none"> <li>Few-shot performance of LLMs is better than the zero-shot, suggesting that LLMs are meta-learners</li> </ul>
mT5	<ul style="list-style-type: none"> <li>Large multi-lingual models perform equivalently to single language models on downstream tasks. However, smaller multi-lingual models perform worse</li> </ul>
PanGu- $\alpha$	<ul style="list-style-type: none"> <li>LLMs have good few shot capabilities</li> </ul>
CPM-2	<ul style="list-style-type: none"> <li>Prompt fine-tuning requires updating very few parameters while achieving performance comparable to full model fine-tuning</li> <li>Prompt fine-tuning takes more time to converge as compared to full model fine-tuning</li> <li>Inserting prompt tokens in-between sentences can allow the model to understand relations between sentences and long sequences</li> <li>In an analysis, CPM-2 finds that prompts work as a provider (additional context) and aggregator (aggregate information with the input text) for the model</li> </ul>
ERNIE 3.0	<ul style="list-style-type: none"> <li>A modular LLM architecture with a universal representation module and task-specific representation module helps in the finetuning phase</li> <li>Optimizing the parameters of a task-specific representation network during the fine-tuning phase is an efficient way to take advantage of the powerful pre-trained model</li> </ul>
Jurassic-1	<ul style="list-style-type: none"> <li>The performance of LLM is highly related to the network size</li> <li>To improve runtime performance, more operations can be performed in parallel (width) rather than sequential (depth)</li> <li>To efficiently represent and fit more text in the same context length, the model uses a larger vocabulary to train a SentencePiece tokenizer without restricting it to word boundaries. This further benefits in few-shot learning tasks</li> </ul>
HyperCLOVA	<ul style="list-style-type: none"> <li>By employing prompt-based tuning, the performances of models can be improved, often surpassing those of state-of-the-art models when the backward gradients of inputs are accessible</li> </ul>
Yuan 1.0	<ul style="list-style-type: none"> <li>The model architecture that excels in pre-training and fine-tuning cases may exhibit contrasting behavior in zero-shot and few-shot learning</li> </ul>
Gopher	<ul style="list-style-type: none"> <li>Relative encodings enable the model to evaluate for longer sequences than training.</li> </ul>
ERNIE 3.0 Titan	<ul style="list-style-type: none"> <li>Additional self-supervised adversarial loss to distinguish between real and generated text improves the model performance as compared to ERNIE 3.0</li> </ul>
GPT-NeoX-20B	<ul style="list-style-type: none"> <li>Parallel attention + FF layers speed-up training 15% with the same performance as with cascaded layers</li> <li>Initializing feed-forward output layers before residuals with scheme in [153] avoids activations from growing with increasing depth and width</li> <li>Training on Pile outperforms GPT-3 on five-shot</li> </ul>

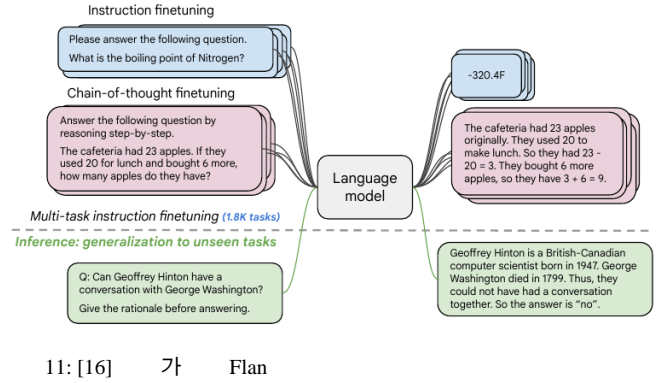
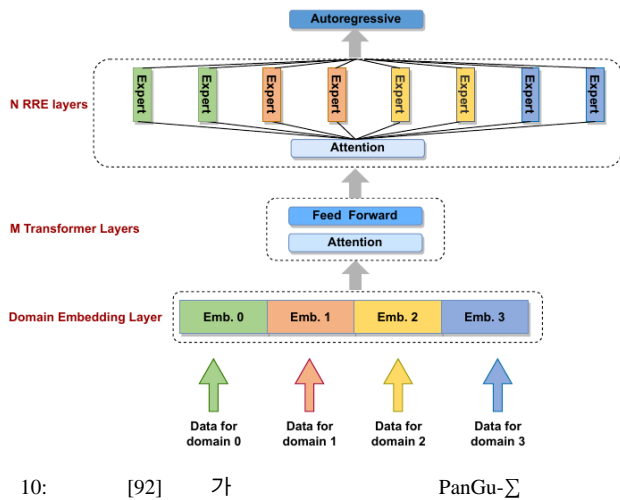
Table Continued on Next Page

Models	Findings & Insights
OPT	<ul style="list-style-type: none"> <li>• Restart training from an earlier checkpoint with a lower learning rate if loss diverges</li> <li>• Model is prone to generate repetitive text and stuck in a loop</li> </ul>
Galactica	<ul style="list-style-type: none"> <li>• Galactica’s performance has continued to improve across validation set, in-domain, and out-of-domain benchmarks, even with multiple repetitions of the corpus, which is superior to existing research on LLMs</li> <li>• A working memory token approach can achieve strong performance over existing methods on mathematical MMLU and MATH benchmarks. It sets a new state-of-the-art on several downstream tasks such as PubMedQA (77.6%) and MedMCQA dev (52.9%)</li> </ul>
GLaM	<ul style="list-style-type: none"> <li>• The model capacity can be maintained at reduced computation by replacing the feed-forward layer in each transformer layer with a mixture-of-experts (MoE)</li> <li>• The model trained on filtered data shows consistently better performances on both NLG and NLU tasks, where the effect of filtering is more significant on the former tasks</li> <li>• Filtered pretraining corpora play a crucial role in the generation capability of LLMs, especially for the downstream tasks</li> <li>• The scaling of GLaM MoE models can be achieved by increasing the size or number of experts in the MoE layer. Given a fixed budget of computation, more experts contribute to a better performance</li> </ul>
LaMDA	<ul style="list-style-type: none"> <li>• The model can be fine-tuned to learn to call different external information resources and tools</li> </ul>
AlphaCode	<ul style="list-style-type: none"> <li>• For higher effectiveness and efficiency, a transformer model can be asymmetrically constructed with a shallower encoder and a deeper decoder</li> <li>• To achieve better performances, it is necessary to employ strategies such as massively scaling upsampling, followed by the filtering and clustering of samples into a compact set</li> <li>• The utilization of novel sampling-efficient transformer architectures designed to facilitate large-scale sampling is crucial</li> <li>• Simplifying problem descriptions can effectively improve the model’s performance</li> </ul>
Chinchilla	<ul style="list-style-type: none"> <li>• The model size and the number of training tokens should be scaled proportionately: for each doubling of the model size, the number of training tokens should be doubled as well</li> </ul>
PaLM	<ul style="list-style-type: none"> <li>• English-centric models produce better translations when translating to English as compared to non-English</li> <li>• Generalized models can have equivalent performance for language translation to specialized small models</li> <li>• Larger models have a higher percentage of training data memorization</li> <li>• Performance has not yet saturated even at 540B scale, which means larger models are likely to perform better</li> </ul>
AlexaTM	<ul style="list-style-type: none"> <li>• Encoder-decoder architecture is more suitable to train LLMs given bidirectional attention to the context than decoder-only</li> <li>• Causal Language Modeling (CLM) task can be added to benefit the model with efficient in-context learning</li> <li>• Placing layer norm at the beginning of each transformer layer improves the training stability</li> </ul>

Table Continued on Next Page

Models	Findings & Insights
U-PaLM	<ul style="list-style-type: none"> <li>• Training with a mixture of denoisers outperforms PaLM when trained further for a few more FLOPs</li> <li>• Training with a mixture of denoisers improves the infilling ability and open-ended text generation diversity</li> </ul>
UL2	<ul style="list-style-type: none"> <li>• Mode switching training enables better performance on downstream tasks</li> <li>• CoT prompting outperforms standard prompting for UL2</li> </ul>
GLM-130B	<ul style="list-style-type: none"> <li>• Pre-training data with a small proportion of multi-task instruction data improves the overall model performance</li> </ul>
CodeGen	<ul style="list-style-type: none"> <li>• Multi-step prompting for code synthesis leads to a better user intent understanding and code generation</li> </ul>
LLaMA	<ul style="list-style-type: none"> <li>• A constant performance improvement is observed when scaling the model</li> <li>• Smaller models can achieve good performances with more training data and computing time</li> </ul>
PanGu- $\Sigma$	<ul style="list-style-type: none"> <li>• Sparse models provide the benefits of large models at a lower computation cost</li> <li>• Randomly Routed Experts reduces catastrophic forgetting effects which in turn is essential for continual learning</li> <li>• Randomly Routed Experts allow extracting a domain-specific sub-model in deployment which is cost-efficient while maintaining a performance similar to the original</li> </ul>
BloombergGPT	<ul style="list-style-type: none"> <li>• Pre-training with general-purpose and task-specific data improves task performance without hurting other model capabilities</li> </ul>
XuanYuan 2.0	<ul style="list-style-type: none"> <li>• Combining pre-training and fine-tuning stages in single training avoids catastrophic forgetting</li> </ul>
CodeT5+	<ul style="list-style-type: none"> <li>• Causal LM is crucial for a model’s generation capability in encoder-decoder architectures</li> <li>• Multiple training objectives like span corruption, Causal LM, matching, etc complement each other for better performance</li> </ul>
StarCoder	<ul style="list-style-type: none"> <li>• HHH prompt by Anthropic allows the model to follow instructions without fine-tuning</li> </ul>
LLaMA-2	<ul style="list-style-type: none"> <li>• Model trained on unfiltered data is more toxic but may perform better on downstream tasks after fine-tuning</li> <li>• Model trained on unfiltered data requires fewer samples for safety alignment</li> </ul>
PaLM-2	<ul style="list-style-type: none"> <li>• Data quality is important to train better models</li> <li>• Model and data size should be scaled with 1:1 proportions</li> <li>• Smaller models trained for larger iterations outperform larger models</li> </ul>
LLaMA-3/3.1	<ul style="list-style-type: none"> <li>• Increasing batch size gradually stabilizes the training without loss spikes</li> <li>• High-quality data at the final stages of training improves the model performance</li> <li>• Increasing model context length windows step-wise allows it to better adapt to various sequence lengths</li> </ul>
Nemotron-40B	<ul style="list-style-type: none"> <li>• Model aligned iteratively on synthetic data with data generated from the previously aligned model achieves competitive performance</li> </ul>
DeepSeek	<ul style="list-style-type: none"> <li>• Batch size should increase with the increase in compute budget while decreasing the learning rate</li> </ul>
DeepSeek-v2	<ul style="list-style-type: none"> <li>• Multi-head latent attention (MLA) performs better than multi-head attention (MHA) while requiring a significantly smaller KV cache, therefore achieving faster data generation</li> </ul>

Models	Findings & Insights
T0	<ul style="list-style-type: none"> <li>• Multi-task prompting enables zero-shot generalization and outperforms baselines</li> <li>• Even a single prompt per dataset task is enough to improve performance</li> </ul>
WebGPT	<ul style="list-style-type: none"> <li>• To aid the model in effectively filtering and utilizing relevant information, human labelers play a crucial role in answering questions regarding the usefulness of the retrieved documents</li> <li>• Interacting a fine-tuned language model with a text-based web-browsing environment can improve end-to-end retrieval and synthesis via imitation learning and reinforcement learning</li> <li>• Generating answers with references can make labelers easily judge the factual accuracy of answers</li> </ul>
Tk-INSTRUCT	<ul style="list-style-type: none"> <li>• Instruction tuning leads to a stronger generalization of unseen tasks</li> <li>• More tasks improve generalization whereas only increasing task instances does not help</li> <li>• Supervised trained models are better than generalized models</li> <li>• Models pre-trained with instructions and examples perform well for different types of inputs</li> </ul>
mT0 and BLOOMZ	<ul style="list-style-type: none"> <li>• Instruction tuning enables zero-shot generalization to tasks never seen before</li> <li>• Multi-lingual training leads to even better zero-shot generalization for both English and non-English</li> <li>• Training on machine-translated prompts improves performance for held-out tasks with non-English prompts</li> <li>• English only fine-tuning on multilingual pre-trained language model is enough to generalize to other pre-trained language tasks</li> </ul>
OPT-IML	<ul style="list-style-type: none"> <li>• Creating a batch with multiple task examples is important for better performance</li> <li>• Only example proportional sampling is not enough, training datasets should also be proportional for better generalization/performance</li> <li>• Fully held-out and partially supervised tasks performance improves by scaling tasks or categories whereas fully supervised tasks have no effect</li> <li>• Including small amounts i.e. 5% of pretraining data during fine-tuning is effective</li> <li>• Only 1% reasoning data improves the performance, adding more deteriorates performance</li> <li>• Adding dialogue data makes the performance worse</li> </ul>
Sparrow	<ul style="list-style-type: none"> <li>• Labelers' judgment and well-defined alignment rules help the model generate better responses</li> <li>• Good dialogue goals can be broken down into detailed natural language rules for the agent and the raters</li> <li>• The combination of reinforcement learning (RL) with reranking yields optimal performance in terms of preference win rates and resilience against adversarial probing</li> </ul>
Flan	<ul style="list-style-type: none"> <li>• Finetuning with CoT improves performance on held-out tasks</li> <li>• Fine-tuning along with CoT data improves reasoning abilities</li> <li>• CoT tuning improves zero-shot reasoning</li> <li>• Performance improves with more tasks</li> <li>• Instruction fine-tuning improves usability which otherwise is challenging for pre-trained models</li> <li>• Improving the model's performance with instruction tuning is compute-efficient</li> <li>• Multitask prompting enables zero-shot generalization abilities in LLM</li> </ul>
WizardCoder	<ul style="list-style-type: none"> <li>• Fine-tuning with re-written instruction-tuning data into a complex set improves performance</li> </ul>
LLaMA-2-Chat	<ul style="list-style-type: none"> <li>• Model learns to write safe responses with fine-tuning on safe demonstrations, while additional RLHF step further improves model safety and make it less prone to jailbreak attacks</li> </ul>
LIMA	<ul style="list-style-type: none"> <li>• Less high quality data is enough for fine-tuned model generalization</li> </ul>



### 3.2.1. Instruction-Tuning with Manually Created Datasets

LLM 가 , LLM , , T0[17] mT0( ) [154] , Tk-Instruct[18]

T5 가 ( , GPT-3 1 75B 11B ) Instruct-GPT

**Increasing Tasks and Prompt Setups:** Zero-shot few-shot

.8k . OPT-IML[97] Flan[16] 2k 1

OPT-IML Flan , zero-shot, few-shot Co

T 가 CoT Collection[101] 1.88M CoT [1 02] Flan-T5 T0, Flan

### 3.2.2. Instruction-Tuning with LLMs Generated Datasets

가 sel

f-instruct[19] 가 LLM .Self

-instruct SUPER-NATURA [18] .17

LINSTRUCTIONS(1600+ 33% 5 , 1 , 1 (52k) (8

2k ) .

GPT-3 [6].

Dynosaur [155] Huggingface LLM

**LLaMA Tuned:**

GPT-3[6] G PT-4[157] LLaMA[15 6] Alpaca[158], Vicuna[159], LLaMA-GPT-4[160] 가 , Alpaca -davinci-003 52k , Vicuna ShareGPT.com 70k , LLaMA-GPT-4 GPT-4 Alpaca . Goat[161] ChatGPT (100 ) L LaMA GPT-4, PaLM, BLOOM, OPT 가 , LLaMA . HuaTuo[162] 8k QA

**Complex Instructions:** Evol-Instruct[163, 164] LLM

WizardLM[163](250k LLaMA) Vicuna Alpaca WizardCoder[164]( StarCo der) Claude-Plus, Bard 가 .

### 3.2.3. Aligning with Human Preferences

LLM . InstructGPT[20] 3 , , ( RL) GPT-3 GPT-3 가 . GPT-3 가 (PPO) . L LaMA 2-Chat[21] PPO . LLaMA 2-Chat 가 PPO



Aligning with Supported Evidence:

[180] LLM  
RLHF  
GopherCite[165], WebGPT[166] Sparrow[167]  
Sparrow[167]  
가  
가 RL [181, 152].  
LLM (PCP)[182]  
Aligning Directly with SFT: RLHF PPO  
가 ,  
[168, 169, 170]  
(SFT) PP  
O (DPO)[168] 가  
가  
RAFT[169]  
가 (PRO)[171] RRHF[170]  
(CoH)[172]  
[183, 184]  
183] 25%가  
0.5%  
2%  
(LIMA)[185] 1000  
[18

Aligning with Synthetic Feedback:

LLM  
LL  
M  
LLM  
Constitutional AI[173] RLHF AI  
RL from AI feedback(RLAIF) .AI  
pacaFarm[174] LLM API  
Constitutional AI  
AlpacaFarm  
Self-Align[98] ICL LLM  
LLM LLM  
Aligning with Prompts: LLM  
[17  
5, 176]. [176] CoT [49].

Red-Teaming/Jailbreaking/Adversarial Attacks: LLM

Position Interpolation: [49]  
1000  
[177, 178].  
LLM ffe [46] RoPE  
[178, 179]. NTK  
Gira  
YaRN [47]





가

## 가

LLM

LLM

가

LLM

가

가

가

## LLM

가

[246].

가 가

LLM

LLM

LLM

## LLM

[40, 247, 41, 38, 39]

가

PEFT

가 가

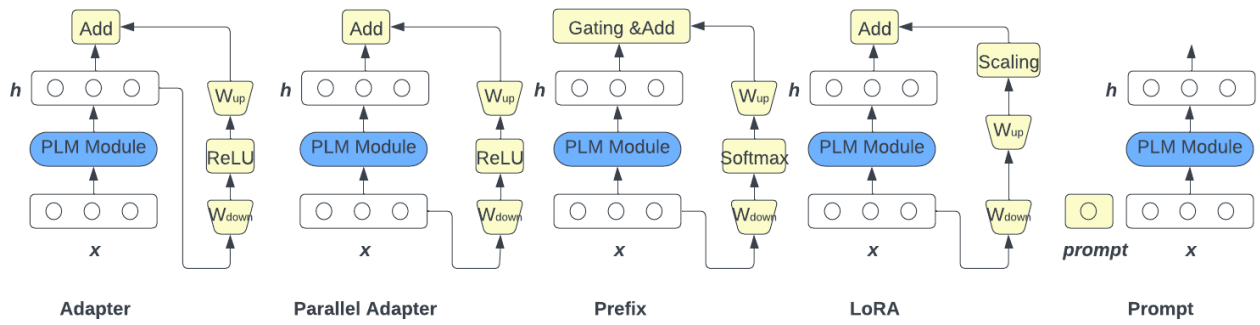
14

## 가

가

[38]

[249]



14:  $x$ 가  $h$  , [38]. LoRA

. AdaMix FP16 [44]. LLM

가  
6B

(LoRA)[250]

가

[255]

LLM  
[44, 256].

**Prompt Tuning:**  
LLM

**Post-Training Quantization:**

가

.LLM-8 [255]

01%-3%

가

[247].

0.0

가

8

가

[251].

가

FP-16

[247, 40, 251].

[40]

가

P-Tuning[247]

가

[45]

가

가

가

가

. P-Tuning

257]

.SmoothQuant[44]

[251]

[252]

가

가

가

가

INT8

가

**Prefix Tuning:**

[41]

가

.OPTQ[256]

(OBC)[258]

, ,

[253]

가

가  
OPTQ

가

**Bias Tuning:**

[254].

BitFit

가

가

OPTQ

가

(OWQ)[259]

가

3.6.2. Quantization

**Quantization-Aware Training:**

(QAT)

LLM

[260, 261, 262].

(B

. 175B

GPT-3

가

CQ)[263]

80GB A100 GPU

350GB

가

5

가 LLM [269, 270, 271], [272, 273, 274], (PEQA)[264] [275, 274, 276] LLM(MLLM) LLM MLLM QLoRA[261] 4 LoRA[250] 4 MLLM LLM 4 [276, 26]. MLLM 가 가 MLLM 가 가 LLM 가 가 MLLM **Pre-training:** MLLM , Flamingo[269] LLM BLIP Queryin -2[270] g Transformer(Q-Former) 2 LLM 가 MiniGPT-4 [277] ViT[278], Q-Former Vicu na LLM[159] **Fine-tuning:** : NLP [20, 16, 97] [16] LLM LL [277, 271, 29] [279, 30, 280] M [58] [279, 281, 282], [19, 3 [284, 280] 가 가 가 [265].LLM-Pruner[42] 3 LoRA 가 ( SIMPLE)[268] 가 가 , Adapter[285] LaVIN[284] 가 가 [266]. 1 가 VideoChat-Text[272] in-corporates Whi sper[286] LLM **Prompting:**

3.6.3. *Pruning* 가 LLM 가 , 가 , MLLM 가 가 LLM 가 가 가 LLM 가 [265, 42, 266]. **Unstructured Pruning:** 가 가 LLM -2[270] g Transformer(Q-Former) 2 BLIP Queryin 가 (Wanda)[265] [255].가 가 가 가 MiniGPT-4 [277] ViT[278], Q-Former Vicu na LLM[159] **Fine-tuning:** : NLP [20, 16, 97] [16] LLM LL [277, 271, 29] [279, 30, 280] M [58] [279, 281, 282], [19, 3 [284, 280] 가 가 가 [265].LLM-Pruner[42] 3 LoRA 가 ( SIMPLE)[268] 가 가 , Adapter[285] LaVIN[284] 가 가 [266]. 1 가 VideoChat-Text[272] in-corporates Whi sper[286] LLM **Prompting:**

3.7. *Multimodal LLMs* LLM

가

LLM

. BLOOM[13] ALiBi가

GLM-130B[33]

A

MLLM

LiBi

(CoT)

[103]

LLM

**Parallel Attention:**

[287].

15%

-CoT[287]

PaLM[15], GPT-NeoX[118] CodeGen[140]

Co

**Multi-Query Attention**

T-PT[288]

. CoT

LLM

[289,

가

가

[15, 142]

**Visual Reasoning Application:**

[291, 292, 216, 293]

**Mixture of Experts:**

LLM

[92, 91].

VQA

[

가가

294, 295]

LLM

. MoE

[5

[91]

8]. LLM

, Po

. MoE

intClip V2[292] LLM

3D

[92].

3D

GPT4Tools[31]

MoE

LoRA[250]

LLM

[92].

[293],

[296]

**Sparse vs Dense Activated:** GPT-3[6]

[67]

[291, 297]

LLM

GLaM[91] PanGu- $\Sigma$ [92] MoE[121]

[67].

3.8. Summary and Discussion

3.8.1. Architecture

LLM

LLM

**Layer Normalization:** LLM

LLM

FP16

가

FP32

[120].

[6, 127, 108].BLOOM[13]

LLM

AlexaTM[122]

[33].FP16

가

[13].

가

FP32

[33].

BF16

[13].BF16

100B

가

BF16

GLM

A100 GPU

LLM

-130B[33]

**Positional Encoding:**

가

**Training Instability:**

LLM

[15].  
[15, 33, 91], [15]  
200~500 [91]  
[33].

[15]  
**Weight Initialization:**  
.GPT-NeoX [118] [153]  $\frac{2}{L\sqrt{d}}$   
[298]  
가 가 가 가  
.MT-NLG [117] 가

[298].  
가  
Galactica [148]  
**Learning Rate:**  
( )  
[13, 15, 124]  
 $1e^{-4} \sim 8e^{-4}$  MT-  
NLG(530B)[117] GPT-NeoX(20B)[118] 13B~175B  
GPT-3[6]

**Training Parallelism:**  
3D 가 LLM 가  
[33, 15, 14, 13, 117, 115, 112]. 3  
D BLOOM[13]  
[37] . PanGu- $\alpha$ [108] PanGu-  
 $\Sigma$ [92] 3D 가 5D

**Mode Switching:**  
가 [125, 124, 122]  
가

**Controllable Text Generation:**  
. GPT-3[6] LLM  
[35], ERNIE 3.0 Titan[

3.8.3. *Supervised Models vs Generalized Models*  
[6, 1]  
5, 18] NLP

3.8.4. *Zero-Shot vs Few-Shot*  
LLM  
가 [6, 15]. LLM  
[6].LLM 가 [6].

[15, 16].

Flan-PaLM[16] CoT

3.8.5. *Encoder vs Decoder vs Encoder-Decoder*  
, NLU  
, NLG , sequence2sequence  
Bert[  
7], RoBERTa[299] L  
LM [6, 118, 13] [10, 11,  
122] NLG  
LLM PaLM[15], OPT[14], GPT-3[6], BLOOM[13],  
LLaMA[156] NLU NLG  
T5[10]  
UL2[125]  
PaLM[

15]  
가 LLM 가  
[125, 122]  
가 CodeT5+[34]

LLM

4.  
3 4  
가 LLM  
5

LLM 6  
7 가 7



3: LLM (>10B). LLM . " / " (Heur), (Dedup), (QF), (PF) . " " GPU/TPU GPU (D), (T), (P), (C), (M), (OP), (R) , " 가 "DS" Deep Speed . "

Models	Publication Venue	License Type	Model Creators	Purpose	No. of Params	Commercial Use	Steps Trained	Data/ Tokens	Data Cleaning	No. of Processing Units	Processing Unit Type	Training Time	Calculated Train. Cost	Training Parallelism	Library
T5 [10]	JMLR'20	Apache-2.0	Google	General	11B	✓	1M	1T	Heur+Dedup	1024	TPU v3	-	-	D+M	Mesh TensorFlow
GPT-3 [6]	NeurIPS'20	-	OpenAI	General	175B	×	-	300B	Dedup+QF	-	V100	-	-	M	-
mT5 [11]	NAACL'21	Apache-2.0	Google	General	13B	✓	1M	1T	-	-	-	-	-	-	-
PanGu- $\alpha$ [108]	arXiv'21	Apache-2.0	Huawei	General	200B	✓	260k	1.1TB	Heur+Dedup	2048	Ascend 910	-	-	D+OP+P+O+R	MindSpore
CPM-2 [12]	AI Open'21	MIT	Tsinghua	General	198B	✓	1M	2.6TB	Dedup	-	-	-	-	D+M	JAXFormer
Codex [141]	arXiv'21	-	OpenAI	Coding	12B	×	-	100B	Heur	-	-	-	-	-	-
ERNIE 3.0 [110]	arXiv'21	-	Baidu	General	10B	×	120k*	375B	Heur+Dedup	384	V100	-	-	M*	PaddlePaddle
Jurassic-1 [112]	White-Paper'21	Apache-2.0	AI21	General	178B	✓	-	300B	-	800	GPU	-	-	D+M+P	Megatron+DS
HyperCLOVA [114]	EMNLP'21	-	Naver	General	82B	×	-	300B	Clf+Dedup+PF	1024	A100	321h	1.32 Mil	M	Megatron
Yuan 1.0 [115]	arXiv'21	Apache-2.0	-	General	245B	✓	26k*	180B	Heur+Clf+Dedup	2128	GPU	-	-	D+T+P	-
Gopher [116]	arXiv'21	-	Google	General	280B	×	-	300B	QF+Dedup	4096	TPU v3	920h	13.19 Mil	D+M	JAX+Haiku
ERNIE 3.0 Titan [35]	arXiv'21	-	Baidu	General	260B	×	-	300B	Heur+Dedup	-	Ascend 910	-	-	D+M+P+D*	PaddlePaddle
GPT-NeoX-20B [118]	BigScience'22	Apache-2.0	EleutherAI	General	20B	✓	150k	825GB	None	96	40G A100	-	-	M	Megatron+DS+PyTorch
OPT [14]	arXiv'22	MIT	Meta	General	175B	✓	150k	180B	Dedup	992	80G A100	-	-	D+T	Megatron
BLOOM [13]	arXiv'22	RAIL-1.0	BigScience	General	176B	✓	-	366B	Dedup+PR	384	80G A100	2520h	3.87 Mil	D+T+P	Megatron+DS
Galactica [148]	arXiv'22	Apache-2.0	Meta	Science	120B	×	225k	106B	Dedup	128	80GB A100	-	-	-	Metaseq
GLaM [91]	ICML'22	-	Google	General	1.2T	×	600k*	600B	Clf	1024	TPU v4	-	-	M	GSPMD
LaMDA [150]	arXiv'22	-	Google	Dialog	137B	×	3M	2.81T	Filtered	1024	TPU v3	1384h	4.96 Mil	D+M	Lingvo
MT-NLG [117]	arXiv'22	Apache-v2.0	MS.+Nvidia	General	530B	×	-	270B	-	4480	80G A100	-	-	D+T+P	Megatron+DS
AlphaCode [142]	Science'22	Apache-v2.0	Google	Coding	41B	✓	205k	967B	Heur+Dedup	-	TPU v4	-	-	M	JAX+Haiku
Chinchilla [96]	arXiv'22	-	Google	General	70B	×	-	1.4T	QF+Dedup	-	TPUv4	-	-	-	JAX+Haiku
PaLM [15]	arXiv'22	-	Google	General	540B	×	255k	780B	Heur	6144	TPU v4	-	-	D+M	JAX+T5X
AlexaTM [122]	arXiv'22	Apache v2.0	Amazon	General	20B	×	500k	1.1T	Filtered	128	A100	2880h	1.47 Mil	M	DS
U-PaLM [124]	arXiv'22	-	Google	General	540B	×	20k	-	-	512	TPU v4	120h	0.25 Mil	-	-
UL2 [125]	ICLR'23	Apache-2.0	Google	General	20B	✓	2M	1T	-	512	TPU v4	-	-	M	JAX+T5X
GLM [33]	ICLR'23	Apache-2.0	Multiple	General	130B	×	-	400B	-	768	40G A100	1440h	3.37 Mil	M	-
CodeGen [140]	ICLR'23	Apache-2.0	Salesforce	Coding	16B	✓	650k	577B	Heur+Dedup	-	TPU v4	-	-	D+M	JAXFormer
LLaMA [127]	arXiv'23	-	Meta	General	65B	×	350k	1.4T	Clf+Heur+Dedup	2048	80G A100	504h	4.12 Mil	D+M	xFormers
PanGu $\Sigma$ [92]	arXiv'23	-	Huawei	General	1.085T	×	-	329B	-	512	Ascend 910	2400h	-	D+OP+P+O+R	MindSpore
BloombergGPT [151]	arXiv'23	-	Bloomberg	Finance	50B	×	139k	569B	Dedup	512	40G A100	1272h	1.97 Mil	M	PyTorch
Xuan Yuan 2.0 [152]	arXiv'23	RAIL-1.0	Du Xiaoman	Finance	176B	✓	-	366B	Filtered	-	80GB A100	-	-	P	DS
CodeT5+ [34]	arXiv'23	BSD-3	Salesforce	Coding	16B	✓	110k	51.5B	Dedup	16	40G A100	-	-	-	DS
StarCoder [147]	arXiv'23	OpenRAIL-M	BigCode	Coding	15.5B	✓	250k	1T	Dedup+QF+PF	512	80G A100	624h	1.28 Mil	D+T+P	Megatron-LM
LLaMA-2 [21]	arXiv'23	LLaMA-2.0	Meta	General	70B	✓	500k	2T	Minimal Filtering	-	80G A100	1.7Mh	-	-	-
PaLM-2 [123]	arXiv'23	-	Google	General	-	×	-	-	Ddedup+PF+QF	-	-	-	-	-	-
LLaMA-3.1 [130]	arXiv'24	LLaMA-3.0	Meta	General	405B	✓	1.2M	15T	Dedup+QF	16k	80G H100	30.84Mh	-	D+T+P+C	PyTorch
Mixtral 8x22B [131]	web'24	Apache-2.0	Mistral AI	General	141B	✓	-	-	-	-	-	-	-	-	-
Snowflake Arctic [132]	web'24	Apache-2.0	Snowflake	General	480B	✓	-	3.5T	-	-	-	-	-	T+P	DS
Nemotron-4 340B [137]	web'24	Nvidia	Nvidia	General	340B	✓	-	9T	-	6144	80G H100	-	-	D+T+P	-
DeepSeek [138]	arXiv'24	MIT	DeepSeek	General	67B	✓	-	2T	Dedup+QF	-	-	300.6Kh	-	D+T+P	DS
DeepSeek-v2 [139]	arXiv'24	MIT	DeepSeek	General	67B	✓	-	8.1T	QF	-	H800	172.8Kh	-	D+P	HAI-LLM

4: LLM(>10B) 3 . "S-" " / " er .

Models	Publication Venue	License Type	Model Creators	Purpose	No. of Params	Commercial Use	Pre-trained Models	Steps Trained	Data/ Tokens	No. of Processing Units	Processing Unit Type	Training Time	Calculated Train. Cost	Train. Parallelism	Library
WebGPT [166]	arXiv'21	-	OpenAI	General	175B	×	GPT-3	-	-	-	-	-	-	-	-
T0 [17]	ICLR'22	Apache-2.0	BigScience	General	11B	✓	T5	-	250B	512	TPU v3	270h	0.48 Mil	-	-
Tk-Instruct [18]	EMNLP'22	MIT	AI2+	General	11B	✓	T5	1000	-	256	TPU v3	4h	0.0036 Mil	-	Google T5
OPT-IML [97]	arXiv'22	-	Meta	General	175B	×	OPT	8k	2B	128	40G A100	-	-	D+T	Megatron
Flan-U-PaLM [16]	ICLR'22	Apache-2.0	Google	General	540B	✓	U-PaLM	30k	-	512	TPU v4	-	-	-	JAX+T5X
mT0 [154]	ACL'23	Apache-2.0	HuggingFace+	General	13B	✓	mT5	-	-	-	-	-	-	-	-
Sparrow [167]	arXiv'22	-	Google	Dialog	70B	×	Chinchilla	-	-	64	TPU v3	-	-	M	-
WizardCoder [164]	arXiv'23	Apache-2.0	HK Bapt.	Coding	15B	×	StarCoder	200	S-78k	-	-	-	-	-	-
Alpaca [158]	Github'23	Apache-2.0	Stanford	General	13B	✓	LLaMA	3-Epoch	S-52k	8	80G A100	3h	600	FSDP	PyTorch
Vicuna [159]	Github'23	Apache-2.0	LMSYS	General	13B	✓	LLaMA	3-Epoch	S-125k	-	-	-	-	FSDP	PyTorch
LIMA [185]	arXiv'23	-	Meta+	General	65B	-	LLaMA	15-Epoch	S-1000	-	-	-	-	-	-
Koala [300]	Github'23	Apache-2.0	UC-Berkley	General	13B	×	LLaMA	2-Epoch	S-472k	8	A100	6h	100	-	JAX/FLAX

5. 가 5.1. Training Datasets

LLM LLM 가 LLM 가 8 가 . LLM 가 .

5: LLM

"PE"

, "nL"

, "nH"

, "HS"

Models	Type	Training Objective	Attention	Vocab	Tokenizer	Norm	PE	Activation	Bias	nL	nH	HS
T5 (11B)	Enc-Dec	Span Corruption	Standard	32k	SentencePiece	Pre-RMS Layer	Relative Learned	ReLU	×	24	128	1024
GPT3 (175B)	Causal-Dec	Next Token	Dense+Sparse	-	-	-	-	GeLU	✓	96	96	12288
mT5 (13B)	Enc-Dec	Span Corruption	Standard	250k	SentencePiece	Pre-RMS Layer	Relative	ReLU	-	-	-	-
PanGu- $\alpha$ (200B)	Causal-Dec	Next Token	Standard	40k	BPE	-	-	-	-	64	128	16384
CPM-2 (198B)	Enc-Dec	Span Corruption	Standard	250k	SentencePiece	Pre-RMS Layer	Relative	ReLU	-	24	64	-
Codex (12B)	Causal-Dec	Next Token	Standard	-	BPE+	Pre-Layer	Learned	GeLU	-	96	96	12288
ERNIE 3.0 (10B)	Causal-Dec	Next Token	Standard	-	WordPiece	Post-Layer	Relative	GeLU	-	48	64	4096
Jurassic-1 (178B)	Causal-Dec	Next Token	Standard	256k	SentencePiece*	Pre-Layer	Learned	GeLU	✓	76	96	13824
HyperCLOVA (82B)	Causal-Dec	Next Token	Dense+Sparse	-	BPE*	Pre-Layer	Learned	GeLU	-	64	80	10240
Yuan 1.0 (245B)	Causal-Dec	Next Token	Standard	-	-	-	-	-	-	76	-	16384
Gopher (280B)	Causal-Dec	Next Token	Standard	32k	SentencePiece	Pre-RMS Layer	Relative	GeLU	✓	80	128	16384
ERNIE 3.0 Titan (260B)	Causal-Dec	Next Token	Standard	-	WordPiece	Post-Layer	Relative	GeLU	-	48	192	12288
GPT-NeoX-20B	Causal-Dec	Next Token	Parallel	50k	BPE	Layer	Rotary	GeLU	✓	44	64	-
OPT (175B)	Causal-Dec	Next Token	Standard	-	BPE	-	-	ReLU	✓	96	96	-
BLOOM (176B)	Causal-Dec	Next Token	Standard	250k	BPE	Layer	ALiBi	GeLU	×	70	112	14336
Galactica (120B)	Causal-Dec	Next Token	Standard	50k	BPE+custom	Layer	Learned	GeLU	×	96	80	10240
GLaM (1.2T)	MoE-Dec	Next Token	Standard	256k	SentencePiece	Layer	Relative	GeLU	✓	64	128	32768
LaMDA (137B)	Causal-Dec	Next Token	Standard	32k	BPE	Layer	Relative	GeLU	-	64	128	8192
MT-NLG (530B)	Causal-Dec	Next Token	Standard	50k	BPE	Pre-Layer	Learned	GeLU	✓	105	128	20480
AlphaCode (41B)	Enc-Dec	Next Token	Multi-query	8k	SentencePiece	-	-	-	-	64	128	6144
Chinchilla (70B)	Causal-Dec	Next Token	Standard	32k	SentencePiece-NFKC	Pre-RMS Layer	Relative	GeLU	✓	80	64	8192
PaLM (540B)	Causal-Dec	Next Token	Parallel+Multi-query	256k	SentencePiece	Layer	RoPE	SwiGLU	×	118	48	18432
AlexaTM (20B)	Enc-Dec	Denosing	Standard	150k	SentencePiece	Pre-Layer	Learned	GeLU	✓	78	32	4096
Sparrow (70B)	Causal-Dec	Pref.&Rule RM	-	32k	SentencePiece-NFKC	Pre-RMS Layer	Relative	GeLU	✓	16*	64	8192
U-PaLM (540B)	Non-Causal-Dec	MoD	Parallel+Multi-query	256k	SentencePiece	Layer	RoPE	SwiGLU	×	118	48	18432
UL2 (20B)	Enc-Dec	MoD	Standard	32k	SentencePiece	-	-	-	-	64	16	4096
GLM (130B)	Non-Causal-Dec	AR Blank Infilling	Standard	130k	SentencePiece	Deep	RoPE	GeGLU	✓	70	96	12288
CodeGen (16B)	Causal-Dec	Next Token	Parallel	-	BPE	Layer	RoPE	-	-	34	24	-
LLaMA (65B)	Causal-Dec	Next Token	Standard	32k	BPE	Pre-RMS Layer	RoPE	SwiGLU	-	80	64	8192
PanGu- $\Sigma$ (1085B)	Causal-Dec	Next Token	Standard	-	BPE	Fused Layer	-	FastGeLU	-	40	40	5120
BloombergGPT (50B)	Causal-Dec	Next Token	Standard	131k	Unigram	Layer	ALiBi	GeLU	✓	70	40	7680
Xuan Yuan 2.0 (176B)	Causal-Dec	Next Token	Self	250k	BPE	Layer	ALiBi	GeLU	✓	70	112	14336
CodeT5+ (16B)	Enc-Dec	SC+NT+Cont.+Match	Standard	-	Code-Specific	-	-	-	-	-	-	-
StarCoder (15.5B)	Causal-Dec	FIM	Multi-query	49k	BPE	-	Learned	-	-	40	48	6144
LLaMA-2 (70B)	Causal-Dec	Next Token	Grouped-query	32k	BPE	Pre-RMS Layer	RoPE	SwiGLUE	-	-	-	-
PaLM-2	-	MoD	Parallel	-	-	-	-	-	-	-	-	-
LLaMA-3.1 (405B)	Causal-Dec	Next Token	Grouped-query	128k	BPE	Pre-RMS Layer	RoPE	SwiGLU	-	126	128	16384
Nemotron-4 (340B)	Causal-Dec	Next Token	Standard	256k	SentencePiece	-	RoPE	ReLU	×	96	96	18432
DeepSeek (67B)	Causal-Dec	Next Token	Grouped-query	100k	BBPE	Pre-RMS Layer	RoPE	SwiGLU	-	95	64	8192
DeepSeek-v2 (67B)	MoE-Dec	Next Token	Multi-Head Latent	100k	BBPE	Pre-RMS Layer	RoPE	SwiGLU	-	60	128	5120

## 5.2. Evaluation Datasets and Tasks

LLM 가 LLM

(LM) 가 가  
(NLU) 2) (NLG). NLU  
NLG

Natural Language Understanding: LM

(NLI), (QA), (CR), (M  
R), (RC)

Natural Language Generation:

LLM 가  
(MT),  
가 LLM 가  
가  
9  
가

10 LLM 11  
LLM  
NLP 가 LLM 12

### 5.2.1. Multi-task

MMLU [307]:

57

가

SuperGLUE [2]: GLUE [309]

SuperGLUE

가

BIG-bench [308]: BIG-bench(

GLUE [309]: GLUE(General Language Understanding Eval  
uation)

6: LLM , LLM 가 , 0.1, 1.0 0.1

Models	Batch Size	Sequence Length	LR	Warmup	LR Decay	Optimizers			Precision			Weight Decay	Grad Clip	Dropout
						AdaFactor	Adam	AdamW	FP16	BF16	Mixed			
T5 (11B)	2 <sup>11</sup>	512	0.01	×	inverse square root	✓			-	-	-	-	-	✓
GPT3 (175B)	32K	-	6e-5	✓	cosine		✓		✓			✓	✓	-
mT5 (13B)	1024	1024	0.01	-	inverse square root	✓			-	-	-	-	-	✓
PanGu- $\alpha$ (200B)	-	1024	2e-5	-	-	-	-	-	-	✓	-	-	-	-
CPM-2 (198B)	1024	1024	0.001	-	-	✓			-	-	-	-	-	✓
Codex (12B)	-	-	6e-5	✓	cosine		✓		✓			✓	-	-
ERNIE 3.0 (12B)	6144	512	1e-4	✓	linear		✓		-	-	-	✓	-	-
Jurassic-1 (178B)	3.2M	2048	6e-5	✓	cosine		✓		✓			✓	✓	-
HyperCLOVA (82B)	1024	-	6e-5	-	cosine			✓	-	-	-	✓	-	-
Yuan 1.0 (245B)	<10M	2048	1.6e-4	✓	cosine decay to 10%	✓			-	-	-	✓	-	-
Gopher (280B)	3M	2048	4e-5	✓	cosine decay to 10%	✓			-	✓		-	✓	-
ERNIE 3.0 Titan (260B)	-	512	1e-4	✓	linear		✓		✓			✓	✓	-
GPT-NeoX-20B	1538	2048	0.97e-5	✓	cosine			✓	✓			✓	✓	×
OPT (175B)	2M	2048	1.2e-4	-	linear			✓	✓			✓	✓	✓
BLOOM (176B)	2048	2048	6e-5	✓	cosine		✓		-	✓		✓	✓	×
Galactica (120B)	2M	2048	7e-6	✓	linear decay to 10%			✓	-	-	-	✓	✓	✓
GLaM (1.2T)	1M	1024	0.01	-	inverse square root	✓			FP32 +		✓	-	✓	×
LaMDA (137B)	256K	-	-	-	-	-	-	-	-	-	-	-	-	-
MT-NLG (530B)	1920	2048	5e-5	✓	cosine decay to 10%		✓		-	✓		✓	✓	-
AlphaCode (41B)	2048	1536+768	1e-4	✓	cosine decay to 10%			✓	-	✓		✓	✓	-
Chinchilla (70B)	1.5M	2048	1e-4	✓	cosine decay to 10%			✓	-	✓		-	-	-
PaLM (540B)	2048	2048	0.01	-	inverse square root	✓			-	-	-	✓	✓	×
AlexaTM (20B)	2M	1024	1e-4	-	linear decay to 5%		✓		-	✓		✓	-	✓
U-PaLM (540B)	32	2048	1e-4	-	cosine	✓			-	-	-	-	-	-
UL2 (20B)	1024	1024	-	-	inverse square root	-	-	-	-	-	-	×	-	-
GLM (130B)	4224	2048	8e-5	✓	cosine			✓	✓			✓	✓	✓
CodeGen (16B)	2M	2048	5e-5	✓	cosine		✓		-	-	-	✓	✓	-
LLaMA (65B)	4M Tokens	2048	1.5e-4	✓	cosine decay to 10%			✓	-	-	-	✓	✓	-
PanGu- $\Sigma$ (1.085T)	512	1024	2e-5	✓	-		✓		-		✓	-	-	-
BloombergGPT (50B)	2048	2048	6e-5	✓	cosine			✓	-		✓	✓	✓	×
Xuan Yuan 2.0 (176B)	2048	2048	6e-5	✓	cosine		✓		✓			✓	✓	-
CodeT5+ (16B)	2048	1024	2e-4	-	linear			✓	-		✓	✓	-	-
StarCoder (15.5B)	512	8k	3e-4	✓	cosine		✓		-	✓		✓	-	-
LLaMA-2 (70B)	4M Tokens	4k	1.5e-4	✓	cosine			✓	-	✓		✓	✓	-
LLaMA-3.1 (405B)	16M	8192	8e-5	✓	linear+cosine			✓	-	✓		-	-	-
Nemotron-4 (340B)	2304	4096	-	-	linear	-	-	-	-	✓		-	-	×
DeepSeek (67B)	4608	4096	3.2e-4	✓	cosine			✓	-	✓		✓	✓	-
DeepSeek-v2 (67B)	9216	4k	2.4e-4	✓	step-decay			✓	-	-	-	✓	✓	-

7: LLM

1 가

Models	Batch Size	Sequence Length	LR	Warmup	LR Decay	Optimizers			Grad Clip	Dropout
						AdaFactor	Adam	AdamW		
WebGPT (175B)	BC:512, RM:32	-	6e-5	-	-		✓		-	-
T0 (11B)	1024	1280	1e-3	-	-	✓			-	✓
Tk-Instruct (11B)	1024	-	1e-5	-	constant	-	-	-	-	-
OPT-IML (175B)	128	2048	5e-5	×	linear		✓		✓	✓
Flan-U-PaLM (540B)	32	-	1e-3	-	constant	✓			-	✓
Sparrow (70B)	RM: 8+16, RL:16	-	2e-6	✓	cosine decay to 10%	✓			✓	×
WizardCoder (15B)	512	2048	2e-5	✓	cosine	-	-		-	-
Alpaca (13B)	128	512	1e-5	✓	cosine	-	-	✓	✓	×
Vicuna (13B)	128	-2048	2e-5	✓	cosine			✓	-	×
LIMA (65B)	32	2048	1e-5	×	linear			✓	-	✓

가 AI 가 CoQA [316]: - CoQA

.7

### 5.2.2. Language Understanding

WinoGrande [354]: Winograd [357]

가

WiC [317]:

가

Dataset	Type	Size/Samples	Tasks	Source	Creation	Comments
C4 [10]	Pretrain	806GB	-	Common Crawl	Automated	A clean, multilingual dataset with billions of tokens
mC4 [11]	Pretrain	38.49TB	-	Common Crawl	Automated	A multilingual extension of the C4 dataset, mC4 identifies over 100 languages using cld3 from 71 monthly web scrapes of Common Crawl.
PILE [301]	Pretrain	825GB	-	Common Crawl, PubMed Central, OpenWebText2, ArXiv, GitHub, Books3, and others	Automated	A massive dataset comprised of 22 constituent sub-datasets
ROOTs [302]	Pretrain	1.61TB	-	498 Hugging Face datasets	Automated	46 natural and 13 programming languages
MassiveText [116]	Pretrain	10.5TB	-	MassiveWeb, Books, News, Wikipedia, Github, C4	Automated	99% of the data is in English
Wikipedia [303]	Pretrain	-	-	Wikipedia	Automated	Dump of wikipedia
RedPajama [304]	Pretrain	5TB	-	CommonCrawl, C4, Wikipedia, Github, Books, StackExchange	Automated	Open-source replica of LLaMA dataset
PushShift.io Reddit	Pretrain	21.1GB	-	Reddit	Automated	Submissions and comments on Reddit from 2005 to 2019
BigPython [140]	Pretrain	5.5TB	Coding	GitHub	Automated	-
Pool of Prompt (P3) [17]	Instructions	12M	62	PromptSource	Manual	A Subset of PromptSource, created from 177 datasets including summarization, QA, classification, etc.
xP3 [154]	Instructions	81M	71	P3+Multilingual datasets	Manual	Extending P3 to total 46 languages
Super-NaturalInstructions (SNI) [18]	Instructions	12.4M	1616	Multiple datasets	Manual	Extending P3 with additional multilingual datasets, total 46 languages
Flan [16]	Instructions	15M	1836	Muffin+T0-SF+NIV2	Manual	Total 60 languages
OPT-IML [97]	Instructions	18.1M	1667	-	Manual	-
Self-Instruct [19]	Instructions	82k	175	-	Automated	Generated 52k instructions with 82k samples from 175 seed tasks using GPT-3
Alpaca [158]	Instructions	52k	-	-	Automated	Employed self-instruct method to generate data from text-davinci-003
Vicuna [159]	Instructions	125k	-	ShareGPT	Automated	Conversations shared by users on ShareGPT using public APIs
LLaMA-GPT-4 [160]	Instructions	52k	-	Alpaca	Automated	Recreated Alpaca dataset with GPT-4 in English and Chinese
Unnatural Instructions [305]	Instructions	68k	-	15-Seeds (SNI)	Automated	-
LIMA [185]	Instructions	1k	-	Multiple datasets	Manual	Carefully created samples to test performance with fine-tuning on less data
Anthropic-HH-RLHF [306]	Alignment	142k	-	-	Manual	
Anthropic-HH-RLHF-2 [178]	Alignment	39k	-	-	Manual	

**Wikitext103 [318]:** Wikipedia

1

**LAMBADA [335]:**

가

가

**PG19 [319]:** Project Gutenberg

#### 5.2.4. Physical Knowledge and World Understanding

**C4 [10]:**

C4

**PIQA [340]:**

Transformer

**LCQMC [320]:**

(LCQMC

**TriviaQA [341]:**

(QA)

) 가

가

(IR) QA

가

**ARC [342]:** ARC-Challenge

가

#### 5.2.3. Story Cloze and Sentence Completion

**StoryCloze [334]:**

**ARC-Easy [342]:** ARC

AR

가 "StoryCl

C-Easy

oze Test"

Type	Datasets/Benchmarks
Multi-Task	MMLU [307], SuperGLUE [2], BIG-bench [308], GLUE [309], BBH [308], CUGE [310], Zero-CLUE [311], FewCLUE [312], Blended Skill Talk [313], HELM [314], KLUE-STS [315]
Language Understanding	CoQA [316], WiC [317], Wikitext103 [318], PG19 [319], LCQMC [320], QQP [321], WinoGender [322], CB [323], FinRE [324], SanWen [325], AFQMC [311], BQ Corpus [326], CNSS [327], CKBQA 13 [328], CLUENER [311], Weibo [329], AQuA [330], OntoNotes [331], HeadQA [332], Twitter Dataset [333]
Story Cloze and Sentence Completion	StoryCloze [334], LAMBADA [335], LCSTS [336], AdGen [337], E2E [338], CHID [339], CHID-FC [312]
Physical Knowledge and World Understanding	PIQA [340], TriviaQA [341], ARC [342], ARC-Easy [342], ARC-Challenge [342], PROST [343], Open-BookQA [344], WebNLG [345], DogWhistle Insider & Outsider [346]
Contextual Language Understanding	RACE [347], RACE-Middle [347], RACE-High [347], QuAC [348], StrategyQA [349], Quiz Bowl [350], cMedQA [351], cMedQA2 [352], MATINF-QA [353]
Commonsense Reasoning	WinoGrande [354], HellaSwag [355], COPA [356], WSC [357], CSQA [358], SIQA [359], C <sup>3</sup> [360], CLUEWSC2020 [311], CLUEWSC [311], CLUEWSC-FC [312], ReCoRD [361]
Reading Comprehension	SQuAD [362], BoolQ [363], SQUADv2 [364], DROP [365], RTE [366], WebQA [367], CMRC2017 [368], CMRC2018 [369], CMRC2019 [370], COTE-BD [371], COTE-DP [371], COTE-MFW [371], MultiRC [372], Natural Questions [373], CNSE [327], DRCD [374], DuReader [375], Dureader <sub>robust</sub> [376], DuReader-QG [375], SciQ [377], Sogou-log [378], Dureader <sub>robust</sub> -QG [376], QA4MRE [379], KorQuAD 1.0 [380], CAIL2018-Task1 & Task2 [381]
Mathematical Reasoning	MATH [382], Math23k [383], GSM8K [384], MathQA [385], MGSM [386], MultiArith [387], AS-Div [388], MAWPS [389], SVAMP [390]
Problem Solving	HumanEval [141], DS-1000 [391], MBPP [392], APPS [382], CodeContests [142]
Natural Language Inference & Logical Reasoning	ANLI [393], MNLI-m [394], MNLI-mm [394], QNLI [362], WNLI [357], OCNLI [311], CMNLI [311], ANLI R1 [393], ANLI R2 [393], ANLI R3 [393], HANS [395], OCNLI-FC [312], LogiQA [396], StrategyQA [349]
Cross-Lingual Understanding	MLQA [397], XNLI [398], PAWS-X [399], XSum [400], XCOPA [401], XWinograd [402], TyDiQA-GoldP [403], MLSum [404]
Truthfulness and Fact Checking	TruthfulQA [405], MultiFC [406], Fact Checking on Fever [407]
Biases and Ethics in AI	ETHOS [408], StereoSet [409], BBQ [410], Winobias [411], CrowS-Pairs [412]
Toxicity	RealToxicityPrompts [413], CivilComments toxicity classification [414]
Language Translation	WMT [415], WMT20 [416], WMT20-enzh [416], EPRSTMT [312], CCPM [417]
Scientific Knowledge	AminoProbe [148], BioLAMA [148], Chemical Reactions [148], Galaxy Clusters [148], Mineral Groups [148]
Dialogue	Wizard of Wikipedia [418], Empathetic Dialogues [419], DPC-generated [96] dialogues, ConvAI2 [420], KdConv [421]
Topic Classification	TNEWS-FC [312], YNAT [315], KLUE-TC [315], CSL [311], CSL-FC [312], IFLYTEK [422]

**ARC-Challenge [342]:**  
ARC-Challenge

가  
**QuAC [348]:**

### 5.2.5. Contextual Language Understanding

**RACE [347]:** RACE

, AI

### 5.2.6. Commonsense Reasoning

**HellaSwag [355]:**

가

**RACE-Middle [347]:** RACE [347]  
RACE-Middle

가

**COPA [401]:**

가

가

**RACE-High [347]:** RACE [347]  
RACE-High

**WSC [357]:** Winograd Schema Challenge(WSC)

10 : LLM 가 "QA" , "Cif" , "NL" , "MT" , "RC" , "CR" , "MR" , g,

Models	Training Dataset	Benchmark												Truthful/ Bias/ Toxicity/ Mem.
		BIG- bench	MMLU	Super GLUE	QA	Cif	NLI	MT	Cloze/ Completion	RC	CR	MR	Coding	
T5	C4 [10]			✓	✓		✓	✓	✓	✓	✓	✓		✓
GPT-3	Common Crawl, WebText, Books Corpora, Wikipedia			✓	✓			✓	✓	✓				
mT5	mC4 [11]				✓		✓	✓						
PanGu- $\alpha$	1.1TB Chinese Text Corpus				✓		✓		✓	✓	✓			
CPM-2	WuDaoCorpus [109]									✓		✓		
Codex	54 million public repositories from Github												✓	
ERNIE-3.0	Chinese text corpora, Baidu Search, Web text, QA-long, QA-short, Poetry and Couplet Domain-specific data from medical, law, and financial area Baidu knowledge graph with more than 50 million facts			✓	✓	✓	✓	✓	✓	✓		✓		
Jurassic-1	Wikipedia, OWT, Books, C4, Pile [301], arXiv, GitHub				✓		✓		✓	✓				
HyperCLOVA	Korean blogs, Community sites, News, KiN Korean Wikipedia, Wikipedia (English and Japanese), Modu-Corpus: Messenger, News, Spoken and written language corpus, Web corpus							✓						
Yuan 1.0	Common Crawl, SogouT, Sogou News, Baidu Baike, Wikipedia, Books				✓	✓	✓			✓				
Gopher	subsets of MassiveWeb Books, C4, News, GitHub and Wikipedia samples from MassiveText	✓	✓	✓	✓						✓	✓		✓
ERNIE-3.0 TITAN	Same as ERNIE 3.0 and ERNIE 3.0 adversarial dataset, ERNIE 3.0 controllable dataset				✓	✓	✓		✓	✓				
GPT-NeoX-20B	Pile [301]			✓	✓		✓		✓		✓	✓		
OPT	RoBERTa [299], Pile [301], PushShift.io Reddit [423]				✓	✓					✓			✓
BLOOM	ROOTS [13]			✓			✓	✓	✓				✓	✓
Galactica	arXiv, PMC, Semantic Scholar, Wikipedia, StackExchange, LibreText, Open Textbooks, RefSeq Genome, OEIS, LIPID MAPS, NASAExoplanet, Common Crawl, ScientificCC, AcademicCC, GitHub repositories Khan Problems, GSM8K, OneSmallStep	✓	✓		✓							✓		✓
GLaM	Filtered Webpages, Social media conversations Wikipedia, Forums, Books, News				✓		✓		✓	✓	✓			
LaMDA	Infiniset : Public documents, Dialogs, Utterances													✓
MT-NLG	Two snapshots of Common Crawl and Books3, OpenWebText2, Stack Exchange, PubMed Abstracts, Wikipedia, PG-19 [242], BookCorpus2, NIH ExPorter, Pile, CC-Stories, RealNews						✓		✓	✓	✓			✓
AlphaCode	Selected GitHub repositories, CodeContests: Codeforces, Description2Code, CodeNet												✓	
Chinchilla	MassiveWeb, MassiveText Books, C4, News, GitHub, Wikipedia	✓	✓		✓					✓	✓			✓
PaLM	webpages, books, Wikipedia, news, articles, source code, social media conversations	✓			✓			✓			✓		✓	✓
AlexaTM	Wikipedia, mC4			✓			✓	✓			✓			✓
U-PaLM	Same as PaLM	✓		✓	✓		✓		✓	✓	✓			
UL2	-			✓	✓	✓	✓					✓		✓
GLM-130B	-	✓	✓						✓					
CodeGen	Pile, BigQuery, BigPython												✓	
LLaMA	CommonCrawl, C4, Github, Wikipedia, Books, arXiv, StackExchange		✓		✓					✓	✓	✓	✓	✓
PanGu- $\Sigma$	WuDaoCorpora, CLUE, Pile, C4, Python code				✓	✓	✓	✓	✓				✓	
BloombergGPT	inPile, Pile, C4, Wikipedia	✓	✓				✓		✓	✓	✓			✓
CodeT5+	CodeSearchNet, Github Code											✓	✓	
StarCoder	The Stack v1.2		✓									✓	✓	✓
LLaMA-2	✓	✓		✓						✓	✓	✓	✓	
PaLM-2	Web documents, Code, Books, Maths, Conversation			✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

11: LLM 가 "SNI" Super-NaturalInstructions .

Models	Training Dataset	BIG-bench	MMLU	BBH	RAFT	FLAN	SNI	PromptSource	TyDiQA	HumanEval	MBPP	Truthful/Bias/Toxicity
T0	Pool of Prompts	✓										
WebGPT	ELI5 [424], ELI5 fact-check [166], TriviaQA [341], ARC-Challenge [342], ARC-Easy [342], Hand-written data, Demonstrations of humans, Comparisons between model-generated answers											✓
Tk-INSTRUCT	SNI [18]						✓					
mT0	xP3 [154]											
OPT-IML	PromptSource [17], FLAN [16], SNI [425], UnifiedSKG [426], CrossFit [427], ExMix [428], T5 [10], Reasoning		✓	✓	✓	✓	✓	✓				
Flan	Muffin, T0-SF, Niv2, CoT		✓	✓					✓			
WizardCoder	Code Alpaca									✓	✓	

### 5.2.8. Mathematical Reasoning

**CSQA [358]:** CommonsenseQA AI 가

**MATH [382]:** AI 가

### 5.2.7. Reading Comprehension

**BoolQ [363]:** Google BoolQ ( / )

**Math23k [383]:**

edia Wikip 가 23,000

**SQuADv2 [364]:** Stanford Question Answering Dataset(SQuAD) [362] Wikipedia 가

**GSM8K [384]:**

1.1 50,000 가 SQuADv2 SQuAD 가

### 5.2.9. Problem Solving and Logical Reasoning

**ANLI [393]:** (NLI)

**DROP [365]:** DROP, Discrete Reasoning Over the content of Paragraphs

**HumanEval [141]:** AI 가

**RTE [366]:** (RTE)

**StrategyQA [349]:** AI 가

가

### 5.2.10. Cross-Lingual Understanding

**WebQA [367]:** WebQA

**XNLI [398]:** XNLI MultiNLI[42 9] 15

AI 가 112,500 3가 ( , , )

**CMRC2018 [369]:**

**PAWS-X [399]:** PAWS-X Word Scrambling Cross-lingual Paraphrase Adversaries

PAWS[430] AI 가 ,  
7 가  
가 . **Medicine:** LLM  
가 . , LLM

5.2.11. Truthfulness  
**Truthful-QA [405]:** [436, 437, 438].  
 , ,  
 , ,  
 , 가 LLM  
 ,  
 [439, 440, 441]  
 , ,

5.2.12. Biases and Ethics in AI  
**ETHOS [408]:** ETHOS YouTube Reddit  
 LLM  
 , ,  
 [442, 443, 444].  
**StereoSet [409]:** StereoSet LLM  
가 ,  
 , 가 ,  
 ,  
 [445, 446, 447, 448].  
가 . LLM  
 ,  
 ,

6. [449]. LLM  
(LLM) , 가  
AI , 가  
가 , [450, 451].  
LLM **Education:** LLM  
 , , ,  
 , LLM  
 LLM  
 LLM 가 [452].  
 LLM 가  
 LLM 가  
**General Purpose:** LLM  
 [431]. [453, 454]. LLM  
 , ,  
 , , [455]. LLM 가  
 [432]. LLM  
 , 가  
 [451]. LLM  
 LLM  
가  
 [433].  
 [434]. LLM  
 . LLM  
 [435].  
 LLM  
 [456, 457]. LLM



12 : NLU NLG 가 가 LLM "N-Shots" ed  
to the m 가 odel few-shot zero-shot , "f" "B" t .

Task	Dataset/Benchmark	Top-1		Top-2		Top-3	
		Model (Size)	Score (N-shots)	Model (Size)	Score (N-shots)	Model (Size)	Score (N-shots)
Multi-Task	BIG-bench (B)	Chinchilla (70B)	65.1 (5-shot)	Gopher (280B)	53.97 (5-shot)	PaLM (540B)	53.7 (5-shot)
	MMLU (B)	GPT-4 (-)	86.4 (5-shot)	Gemini (Ultra)	83.7 (5-shot)	Flan-PaLM-2 <sub>(f)</sub> (Large)	81.2 (5-shot)
Language Understanding	SuperGLUE (B)	ERNIE 3.0 (12B)	90.6 (-)	PaLM <sub>(f)</sub> (540B)	90.4 (-)	T5 (11B)	88.9 (-)
Story Comprehension and Generation	HellaSwag	GPT-4 (-)	95.3 (10-shot)	Gemini (Ultra)	87.8 (10-shot)	PaLM-2 (Large)	86.8 (one shot)
	StoryCloze	GPT3 (175B)	87.7 (few shot)	PaLM-2 (Large)	87.4 (one shot)	OPT (175B)	79.82 (-)
Physical Knowledge and World Understanding	PIQA	PaLM-2 (Large)	85.0 (one shot)	LLaMa (65B)	82.8 (zero shot)	MT-NLG (530B)	81.99 (zero shot)
	TriviaQA	PaLM-2 (Large)	86.1 (one shot)	LLaMA-2 (70B)	85.0 (one shot)	PaLM (540B)	81.4 (one shot)
Contextual Language Understanding	LAMBADA	PaLM (540B)	89.7 (few shot)	MT-NLG (530B)	87.15 (few shot)	PaLM-2 (Large)	86.9 (one shot)
Commonsense Reasoning	WinoGrande	GPT-4 (-)	87.5 (5-shot)	PaLM-2 (Large)	83.0 (one shot)	PaLM (540B)	81.1 (zero shot)
	SIQA	LLaMA (65B)	52.3 (zero shot)	Chinchilla (70B)	51.3 (zero shot)	Gopher (280B)	50.6 (zero shot)
Reading Comprehension	BoolQ	PaLM <sub>(f)</sub> (540B)	92.2 (-)	T5 (11B)	91.2 (-)	PaLM-2 (Large)	90.9 (one shot)
Truthfulness	Truthful-QA	LLaMA (65B)	57 (-)				
Mathematical Reasoning	MATH	Gemini (Ultra)	53.2 (4-shot)	PaLM-2 (Large)	34.3 (4-shot)	LLaMa-2 (65B)	13.5 (4-shot)
	GSM8K	GPT-4 (-)	92.0 (5-shot)	PaLM-2 (Large)	80.7 (8-shot)	U-PaLM (540B)	58.5 (-)
Problem Solving and Logical Reasoning	HumanEval	Gemini <sub>(f)</sub> (Ultra)	74.4 (zero shot)	GPT-4 (-)	67.0 (zero shot)	Code Llama (34B)	48.8 (zero shot)

가 [468] [469].  
**Finance:** BloombergGPT[151] LLM  
[458]. 가 LLM  
가 [459, 460]. LLM  
nGPT[470] , Fi  
**Maths:** LLM , 가 . BloombergGPT Fi  
nGPT LLM ,  
[461, 462]. LLM 가 LLM  
가 가 [471].  
[463, 464]. 가 LLM  
가 [465] , , [28, 472, 473, 474], [237], [246],  
[246, 475], [236], [476] LLM  
[240, 26].  
**Law:** LLM , ,  
LM 가 L [224, 233, 234].  
[466]. 7.  
LLM GPT-4 LLM  
가 . LLM  
가 . , ,  
가 [467]. 가 , ,  
LLM .

가	• LLM
	<b>Prompt Engineering:</b> LLM
LLM	가
	LLM [484, 3
<b>Computational Cost:</b> LLM	2].
가	<b>Limited Knowledge:</b>
가	[198].
가	(RAG)[6, 21]
[477]	[193, 25].
<b>Bias and Fairness:</b> LLM	<b>Safety and Controllability:</b> LLM
	가
[478].	가
<b>Overfitting:</b> LLM	<b>Security and Privacy:</b> LLM
가	[485].
[479]. LLM	LLM
	AI LLM
	[486].
	<b>Multi-Modality:</b>
가	LLM
	[480].
<b>Economic and Research Inequality:</b> LLM	<b>Catastrophic Forgetting:</b> LLM
AI	
가	가
[481].	(LLM)
<b>Reasoning and Planning:</b>	
가	가
LLM	LLM
	BERT
	[487]. LLM
	ML LLM
<b>Hallucinations:</b> LLM	[488].
" "	LLM
[483].	가
	가
• LLM	가
• LLM	[489].
	<b>Interpretability and Explainability:</b> LLM " "

[490, 491].	LLM가 LLM가	LLM AI , , [497].
<b>Privacy Concerns:</b>	(LLM) 가 , 8.	[498].
	가 , LLM LLM	
	가 가 , LLM	
	LLM 가 , LLM	
<b>Real-Time Processing:</b>	(LLM) [492,	LLM , - L
AI	LLM LM, LLM , LLM,	
	가 , LLM LLM	
	[494]. MobileBERT	
	가 가 ,	
<b>Long-Term Dependencies:</b>		
	SDAIA-KFUPM JRC-AI-RFP-11 A	
	가 I (SDAIA) (KFUPM)	
<b>Hardware Acceleration:</b>	LLM 가 GPU LLM	
	LLM 가 LLM	
	[495]. GPU TPU 가 가	
	가	
<b>Regulatory and Ethical Frameworks:</b>	OpenAI GPT-4[157] Google Bard (LLM)	
	LLM ,	
	[496]. , LLM	

[1] A. Chernyavskiy, D. Ilvovsky, P. Nakov, Transformers: "??", in: Machine Learning and Knowledge Discovery in Databases. : , ECML PKDD 2021, , 2021 9 13 -17 , , Part III 21, Springer, 2021, 677-693 .

1 [2] A. Wang, Y. Pruksachatkun, N. Nangia, A. Singh, J. Michael, F. Hill, O. Levy, S. Bowman, Superglue: 32(2019). 1, 26, 29 [

3] D. Adiwardana, M.-T. Luong, D. R. So, J. Hall, N. Fiedel, R. Thoppilan, Z. Yang, A. Kulshreshtha, G. Nemade, Y. Lu, et al., arXiv arXiv:2001.09977 (2020). 1 [4

] B. A. y Arcas, ?, Daedalus 151 (2) (2022) 183 – 197. 2 [5] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I . Sutskever , OpenAI 1(8) (2019) 9. 2, 7 [6] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell , 33(2020 ) 1877 – 1901. 2, 6, 7, 8, 9, 16, 18, 23, 24, 25, 34 [7] J. Devlin, M.-W. Chan g, K. Lee, K. Toutanova, Bert: , arXiv arXiv:1810.04805(2018). 2, 18, 24

- [8] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, L. Zettlemoyer, , NAACL-HLT, , 2018, pp. 2227 – 2237. 2 [9] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, L. Zettlemoyer, Bart: , arXiv arXiv:1910.13461(2019). 2 [10] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, P. J. Liu, , 21(1)(2020) 5485 – 5551. 2, 7, 8, 18, 19, 24, 25, 28, 30, 31 [11] L. Xue, N. Constant, A. Roberts, M. Kale, R. Al-Rfou, A. Siddhant, A. Barua, C. Raffel, mt5: , arXiv arXiv:2010.11934(2020). 2, 7, 8, 24, 25, 28, 30 [12] Z. Zhang, Y. Gu, X. Han, S. Chen, C. Xiao, Z. Sun, Y. Yao, F. Qi, J. Guan, P. Ke, et al., Cpm-2: , AI Open 2(2021) 216 – 224. 2, 8, 25 [13] T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, M. Gallé, et al., Bloom: A 176b- , arXiv arXiv:2211.05100 (2022). 2, 4, 9, 11, 23, 24, 25, 30 [14] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al., Opt: , arXiv arXiv:2205.01068 (2022). 2, 9, 11, 2 4, 25 [15] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al., Palm: , arXiv arXiv:2204.02311 (2022). 2, 6, 9, 11, 23, 24, 25 [16] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. Dehghani, S. Brahma, et al., , arXiv arXiv:2210.11416(2022). 2, 7, 11, 16, 17, 22, 24, 25, 28, 31 [17] V. Sanh, A. Webson, C. Raffel, S. H. Bach, L. Sutawika, Z. Alyafeai, A. Chaffin, A. Stiegler, T. L. Scao, A. Raja, et al., , arXiv arXiv:2210.08207(2021). 2, 11, 16, 25, 28, 31 [18] Y. Wang, S. Mishra, P. Alipourmolabashi, Y. Kordi, A. Mirzaei, A. Naik, A. Ashok, A. S. Dhanasekaran, A. Arunkumar, D. Stap, et al., : 1600 + nlp , in: 2022 , 2022, pp. 5085 – 5109. 2, 7, 11, 16, 17, 24, 25, 28, 31 [19] Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, D. Khashabi, H. Hajishirzi, : , arXiv arXiv:2212.10560(2022). 2, 16, 19, 22, 28 [20] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al., , 35(2022) 27730 – 27744. 2, 7, 11, 16, 22 [21] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al., Llama 2: , arXiv arXiv:2307.09288 (2023). 2, 7, 10, 16, 25, 3 4 [22] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Y. O-gatama, M. Bosma, D. Zhou, D. Metzler, et al., , arXiv arXiv:2206.07682 (2022). 2 [23] T. Webb, K. J. Holyoak, H. Lu, , Nature Human Behaviour 7 (9) (2023) 1526 – 1541. 2 [24] D. A. Boiko, R. MacKnight, G. Gomes, , arXiv arXiv:2304.05332 (2023). 2 [25] G. Izacard, P. Lewis, M. Lomeli, L. Hosseini, F. Petroni, T. Schick, J. Dwivedi-Yu, A. Joulin, S. Riedel, E. Grave, Few-shot learning, arXiv arXiv:2208.03299 (2022). 2, 18, 19, 34 [26] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu , Palm-e: , arXiv arXiv:2303.03378 (2023). 2, 2 0, 22, 33 [27] A. Parisi, Y. Zhao, N. Fiedel, Talm: , arXiv arXiv:2205.12255 (2022). 2, 19, 20 [28] B. Zhang, H. Soh, , arXiv arXiv:2303.03548 (2023). 2, 33 [29] Q. Ye, H. Xu, G. Xu, J. Ye, M. Yan, Y. Zhou, J. Wang, A. Hu, P. Shi, Y. Shi, et al., mplug-owl: , arXiv arXiv:2304.14178 (2023). 2, 22 [30] W. Wang, Z. Chen, X. Chen, J. Wu, X. Zhu, G. Zeng, P. Luo, T. Lu, J. Zhou, Y. Qiao, et al., Visionllm: , arXiv arXiv:2305.11175 (2023). 2, 22 [31] R. Yang, L. Song, Y. Li, S. Zhao, Y. Ge, X. Li, Y. Shan, Gpt4tools: , arXiv arXiv:2305.18752(2023). 2, 19, 22, 23 [32] E. Saravia, Prompt Engineering Guide, <https://github.com/dair-ai/Prompt-Engineering-Guide>(12 2022). 2, 7, 18, 34 [33] A. Zeng, X. Liu, Z. Du, Z. Wang, H. Lai, M. Ding, Z. Yang, Y. Xu, W. Zheng, X. Xia, et al., Glm-130b: , arXiv arXiv:2210.02414 (2022). 2, 10, 23, 24, 25 [34] Y. Wang, H. Le, A. D. Gotmare, N. D. Bui, J. Li, S. C. Hoi, Codet5+: , arXiv arXiv:2305.07922 (2023). 2, 11, 24, 25 [35] S. Wang, Y. Sun, Y. Xiang, Z. Wu, S. Ding, W. Gong, S. Feng, J. Shang, Y. Zhao, C. Pang , Ernie 3.0 titan: , arXiv arXiv:2112.12731(2021). 2, 8, 24, 2 5 [36] J. Rasley, S. Rajbhandari, O. Ruwase, Y. He, Deepspeed: 1,000 가 , 2020 26 ACM SIGKDD , 3505-3506 . 2, 5 [37] S. Rajbhandari, J. Rasley, O. Ruwase, Y. He, Zero: 1 , SC20: , IEEE, 2020, 1-16 . 2, 4, 24 [38] J. He, C. Zhou, X. Ma, T. Berg-Kirkpatrick, G. Neubig, , arXiv arXiv:2110.04366 (2021). 2, 20, 21 [39] Z. Hu, Y. Lan, L. Wang, W. Xu, E.-P. Lim, R. K.-W. Lee, L. Bing, S. Poria, Llm-adapters: , arXiv arXiv:2304.01933 (2023). 2, 20 [40] B. Lester, R. Al-Rfou, N. Constant, , arXiv arXiv:2104.08691(2021). 2, 8, 20, 21 [41] X. L. Li, P. Liang, : , arXiv arXiv:2101.00190(2021). 2, 20, 21 [42] X. Ma, G. Fang, X. Wang, Llm-pruner: , arXiv arXiv:2305.11627(2023). 2, 22 [43] R. Xu, F. Luo, C. Wang, B. Chang, J. Huang, S. Huang, F. Huang, : , AAAI , 36 , 2022 , 11547~11555 . 2, 22 [44] G. Xiao, J. Lin, M. Seznec, H. Wu, J. Demouth, S. Han, Smoothquant: , ICML, 202 , PMLR, 2023 , 38087~38099 . 2, 21 [45] C. Tao, L. Hou, W. Zhang, L. Shang, X. Jiang, Q. Liu, P. Luo, N. Wong, , arXiv arXiv:2203.10705 (2022). 2, 21 [46] A. Pal, D. Karkhanis, M. Roberts, S. Dooley, A. Sundararajan, S. Naidu, Giraffe: llms , arXiv arXiv:2308.10882(2023). 2, 17 [47] B. Peng, J. Quesnelle, H. Fan, E. Shippole, Yarn: , arXiv arXiv:2309.00071(2023). 2, 17 [48] M. Guo, J. Ainslie, D. Uthus, S. Ontanon, J. Ni, Y.-H. Sung, Y. Yang,

- LongT5: , arXiv  
arXiv:2112.07916 (2021). 2, 18 [49] S. Chen, S. Wong, L. Chen, Y. Tian, , arXiv  
arXiv:2306.15595 (2023). 2, 17 [50] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al., , arXiv  
arXiv:2303.18223 (2023). 2, 3, 7 [51] U. Naseem, I. Razzak, S. K. Khan, M. Prasad, :  
20(5) (2021) 1  
– 35. 2, 3 [52] B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heinz, D. Roth, : , arXiv  
arXiv:2111.01243 (2021). 2, 3 [53] C. Zhou, Q. Li, C. Li, J. Yu, Y. Liu, G. Wang, K. Zhang, C. Ji, Q. Yan, L. He , : bert chat  
gpt , arXiv  
arXiv:2302.09419(2023). 2, 3 [54] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, Z. Sui, , arXiv  
arXiv:2301.00234(2022). 2, 7, 18 [55] J. Huang, K. C.-C. Chang, :  
 , arXiv  
arXiv:2212.10403 (2022). 2, 7, 18 [56] Y. Wang, W. Zhong, L. Li, F. Mi, X. Zeng, W. Huang, L. Shang, X. Jiang, Q. Liu, : , arXiv  
arXiv:2307.12966 (2023). 2 [57] X. Zhu, J. Li, Y. Liu, C. Ma, W. Wang, , arXiv  
arXiv:2308.07633 (2023). 2 [58] S. Yin, C. Fu, S. Zhao, K. Li, X. Sun, T. Xu, E. Chen, , arXiv  
arXiv:2306.13549(2023). 2, 22, 23 [59] J. J. Webster, C. Kit, NLP  
COL-ING 1992 4 : 14 , 1992. 4 [60] T. Kudo, :  
56 ( 1 :  
) , 2018, 66-75 . 4 [61] R. Sennrich, B. Haddow, A. Birch, 7  
 , 2016 54  
( 1 : ) , pp. 1715-1725. 4 [62] M. Schuster, K. Nakajima, , 2012 IEEE  
(ICASSP), IEEE, 2012, pp. 5149-5152. 4 [63] S. J. Mielke, Z. Alyafeai, E. Salessky, C. Raffel, M. Dey, M. Gallé, A. Raja, C. Si, W. Y. Lee, B. Sagot, et al., : nlp  
 , arXiv  
arXiv:2112.10508 (2021). 4 [64] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, .  
30(2017). 4, 7 [65] O. Press, N. Smith, M. Lewis, , :  
7 , International Conference on Learning Representations, 2022 . URL <https://openreview.net/forum?id=R8sQPpGCv0> 4, 17 [66] J. Su, Y. Lu, S. Pan, A. Murtadha, B. Wen, Y. Liu, Roformer: , arXiv  
arXiv:2104.09864(2021). 4, 9, 17 [67] R. Child, S. Gray, A. Radford, I. Sutskever, , arXiv  
arXiv:1904.10509(2019). 4, 7, 23 [68] T. Dao, D. Fu, S. Ermon, A. Rudra, C. Ré, Flashattention: io-awareness  
35(2022) 16344 – 16359. 4 [69] K. H. Hornik, M. Stinchcombe, H. White, ,  
2(5)(1989) 359 – 366. 4 [70] V. Nair, G. E. Hinton, , 27  
(ICML-10), 2010, 807-814 . 4 [71] D. Hendrycks, K. Gimpel, 7  
(gelus), arXiv  
arXiv:1606.08415(2016). 4 [72] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, :  
15(1)(2014) 1929-1958. 4 [73] D. Krueger, T. Maharaj, J. Krámár, M. Pezeshki, N. Ballas, N. R. Ke, A. Goyal, Y. Bengio, A. Courville, C. Pal, Zoneout: rms , arXiv  
arXiv:1606.01305(2016). 4 [74] N. Shazeer, Glu  
 , arXiv  
arXiv:2002.05202(2020). 4 [75] Y. N. Dauphin, A. Fan, M. Auli, D. Grangier, , PMLR, 2017, 933-941 . 4 [76] J. L. Ba, J. R. Kiros, G. E. Hinton, , arXiv  
arXiv:1607.06450(2016). 4 [77] B. Zhang, R. Sennrich, ,  
32(2019). 4 [78] A. Baevski, M. Auli, , arXiv  
arXiv:1809.10853(2018). 4 [79] H. Wang, S. Ma, L. Dong, S. Huang, D. Zhang, F. Wei, Deepnet: 1,000 , arXiv  
arXiv:2203.00555(2022). 4 [80] M. Shoyebi, M. Patwary, R. Puri, P. LeGresley, J. Casper, B. Catanzaro, Megatron-lm: , arXiv  
arXiv:1909.08053(2019). 4, 5 [81] "bmtrain: ". URL <https://github.com/OpenBMB/BMTrain> 4, 5 [82] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. M. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al., Transformers: State-of-the-art natural language processing, in: Proceedings of the 2020 conference on empirical methods in natural language processing: system demons, 2020 , pp. 38 – 45. 5 [83] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. L. Eary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, et al., Jax: python+ numpy 7 (2018). 5 [84] S. Li, J. Fang, Z. Bian, H. Liu, Y. Liu, H. Huang, B. Wang, Y. You, Colossal -ai: , arXiv  
arXiv:2110.14883(2021). 5 [85] J. He, J. Qiu, A. Zeng, Z. Yang, J. Zhai, J. Tang, Fastmoe: 7 , arXiv  
arXiv:2103.13262(2021). 5 [86] L. Huawei Technologies Co., Huawei mindspore ai , Springer, 2022, pp. 137 – 162. 5 [87] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: 32(2019). 5 [88] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al., Tensorflow: , Osd, Vol. 16, Savannah, GA, USA, 2016, pp. 265 – 283. 5 [89] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, Z. Zhang, Mxnet: , a  
arXiv  
arXiv:1512.01274(2015). 5 [90] W. Fedus, B. Zoph, N. S. Shazeer, :  
 , The Journal of Machine Learning Research 23(1)(2022) 5232 – 5270. 5, 9 [91] N. Du, Y. Huang, A. M. Dai, S. Tong, D. Lepikhin, Y. Xu, M. Krikun, Y. Zhou, A. W. Yu, O. Firat , Glam: 7  
 , PMLR, 2022, 5547-5569 . 5, 9, 23, 24, 25 [92] X. Ren, P. Zhou, X. Meng, X. Huang, Y. Wang, W. Wang, P. Li, X. Zhang, A. Podolskiy, G. Arshinov , Pangu-Σ: , arXiv  
arXiv:2303.10845(2023). 5, 10, 16, 23, 24, 25 [93] T. Wang, A. Roberts, D. Hesslow, T. Le Scao, H. W. Chung, I. Beltagy, J. Launay, C. Raffel,

: ing objective works best for zero-shot generalization?, in: International Conference on Machine Learning, PMLR, 2022, pp. 22964 – 22984. 5 [94] L. Dong, N. Yang, W. Wang, F. Wei, X. Liu, Y. Wang, J. Gao, M. Zhou, H.-W. Hon, 32(2019). 6 [95] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, D. Amodei, , arXiv arXiv:2001.08361(2020). 6 [96] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark , arXiv arXiv:2203.15556(2022). 6, 9, 25, 29 [97] S. Iyer, X. V. Lin, R. Pasunuru, T. Mihaylov, D. Simig, P. Yu, K. Shuster, T. Wang, Q. Liu, P. S. Koura, et al., Opt-impl: , arXiv arXiv:2212.12017(2022). 7, 11, 16, 17, 22, 25, 28 [98] Z. Sun, Y. Shen, Q. Zhou, H. Zhang, Z. Chen, D. Cox, Y. Yang, C. Gan, , arXiv arXiv:2305.03047(2023). 7, 17 [99] A. Askeel, Y. Bai, A. Chen, D. Drain, D. Ganguli, T. Henighan, A. Jones, N. Joseph, B. Mann, N. DasSarma, et al., , arXiv arXiv:2112.00861(2021). 7 [100] D. M. Ziegler, N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. Christiano, G. Irving, , arXiv arXiv:1909.08593(2019). 7 [101] S. Kim, S. J. Joo, D. Kim, J. Jang, S. Ye, J. Shin, M. Seo, Cot : , arXiv arXiv:2305.14045(2023). 7, 16 [102] Q. Liu, F. Zhou, Z. Jiang, L. Dou, M. Lin, 0 : , arXiv arXiv:2304.07995(2023). 7, 16 [103] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al., 35(2022) 24824 – 24837. 7, 20, 23 [104] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowhery, D. Zhou, , arXiv arXiv:2203.11171(2022). 7, 20 [105] S. Yao, D. Yu, J. Zhao, I. Shafraan, T. L. Griffiths, Y. Cao, K. Narasimhan, : , arXiv arXiv:2305.10601(2023). 7, 20 [106] N. Housley, A. Giurghi, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, S. Gelly, NLP : , PMLR, 2019, pp. 2790 – 2799. 7, 20 [107] S. McCandlish, J. Kaplan, D. Amodei, O. D. Team, , arXiv arXiv:1812.06162(2018). 7 [108] W. Zeng, X. Ren, T. Su, H. Wang, Y. Liao, Z. Wang, X. Jiang, Z. Yang, K. Wang, X. Zhang , Pangu- $\alpha$ : , arXiv arXiv:2104.12369(2021). 8, 23, 24, 25 [109] S. Yuan, H. Zhao, Z. Du, M. Ding, X. Liu, Y. Cen, X. Zou, Z. Yang, J. Tang, Wudaocorpora: , AI Open 2(2021) 65 – 68. 8, 30 [110] Y. Sun, S. Wang, S. Feng, S. Ding, C. Pang, J. Shang, J. Liu, X. Chen, Y. Zhao, Y. Lu, et al., Ernie 3.0: , arXiv arXiv:2107.02137(2021). 8, 25 [111] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, R. Salakhutdinov, Transformer-xl: , arXiv arXiv:1901.02860(2019). 8 [112] O. Lieber, O. Sharir, B. Lenz, Y. Shoham, Jurassic-1: , AI21 Labs 1(2021). 8, 24, 25 [113] Y. Levine, N. Wies, O. Sharir, H. Bata, A. Shashua, , 33(2020) 22640 – 22651. 8, 11 [114] B. Kim, H. Kim, S.-W. Lee, G. Lee, D. Kwak, D. H. Jeon, S. Park,

S. Kim, S. Kim, D. Seo, et al., 가 ? Hyperclova : , arXiv arXiv:2109.04650(2021). 8, 25 [115] S. Wu, X. Zhao, T. Yu, R. Zhang, C. Shen, H. Liu, F. Li, H. Zhu, J. Luo, L. Xu, et al., Yuan 1.0: , arXiv arXiv:2110.04725(2021). 8, 24, 25 [116] J. W. Rae, S. Borgeaud, T. Cai, K. Millican, J. Hoffmann, F. Song, J. Aslanides, S. Henderson, R. Ring, S. Young, et al., : Gopher , arXiv arXiv:2112.11446(2021). 8, 9, 25, 28 [117] S. Smith, M. Patwary, B. Norick, P. LeGresley, S. Rajbhandari, J. Casper, Z. Liu, S. Prabhumoye, G. Zerveas, V. Korthikanti , deepspeed megatron megatron-turing nlG 530b , arXiv arXiv:2201.11990(2022). 8, 9, 24, 25 [118] S. Black, S. Biderman, E. Hallahan, Q. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonnell, J. Phang, et al., Gpt-neox-20b: , arXiv arXiv:2204.06745(2022). 9, 23, 24, 25 [119] W. Ben, K. Aran, Gpt-j-6b: 60 (2021). 9 [120] P. Micikevicius, S. Narang, J. Alben, G. Diamos, E. Elsen, D. Garcia, B. Ginsburg, M. Houston, O. Kuchaiev, G. Venkatesh, et al., , arXiv arXiv:1710.03740(2017). 9, 23 [121] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hin-ton, J. Dean, : , arXiv arXiv:1701.06538(2017). 9, 23 [122] S. Soltan, S. Ananthakrishnan, J. FitzGerald, R. Gupta, W. Hamza, H. Khan, C. Peris, S. Rawls, A. Rosenbaum, A. Rumshisky, et al., Alexatm 20b: seq2seq Few-shot , arXiv preprint arXiv:2208.01448(2022). 9, 23, 24, 25 [123] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen, et al., Palm 2 , arXiv arXiv:2305.10403(2023). 9, 25 [124] Y. Tay, J. Wei, H. W. Chung, V. Q. Tran, D. R. So, S. Shakeri, X. Garcia, H. S. Zheng, J. Rao, A. Chowdhery, et al., 0.1% 가 , arXiv arXiv:2210.11399(2022). 9, 24, 25 [125] Y. Tay, M. Dehghani, V. Q. Tran, X. Garcia, J. Wei, X. Wang, H. W. Chung, D. Bahri, T. Schuster, S. Zheng, et al., UL2: , 11 , 2022. 9, 10, 24, 25 [126] Z. Du, Y. Qian, X. Liu, M. Ding, J. Qiu, Z. Yang, J. Tang, Gln: , 60 ( 1 : ), 2022, 3 20-335 . 10 [127] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al., Llama: , arXiv arXiv:2302.13971(2023). 10, 23, 25 [128] M. N. Rabe, C. Staats, o(n<sup>2</sup>) 가 , arXiv arXiv:2112.05682(2021). 10 [129] V. A. Korthikanti, J. Casper, S. Lym, L. McAfee, M. Andersch, M. Shoenybi, B. C. atanzaro, , Machine Learning and Systems 5(2023). 10 [130] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al., 3 , arXiv arXiv:2407.21783(2024). 10, 25 [131] https://mistral.ai/news/mixtral-8x22b/. 10, 25 [132] https://github.com/Snowflake-Labs/snowflake-arctic. 10, 25 [133] https://github.com/xai-org/grok-1. 10 [134] https://x.ai/blog/grok-1.5. 10 [135] G. Team, R. Anil, S. Borgeaud, Y. Wu, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, et al., Gemini: , arXiv arXiv:2312.11805(2023). 10 [136] M. Reid, N. Savinov, D. Teplyashin, D. Lepikhin, T. Lillicrap, J.-b.

Alayrac, R. Soricut, A. Lazaridou, O. Firat, J. Schrittwieser, et al., Gemini 1.5: , arXiv:2403.05530(2024). 10 [137] B. Adler, N. Agarwal, A. Ait hal, D. H. Anh, P. Bhattacharya, A. Brundyn, J. Casper, B. Catanzaro, S. Clay, J. Cohen , Nemotron-4 340b , arXiv:2406.11704(2024). 10, 25 [138] X. Bi, D. Chen, G. Chen, S. Chen, D. Dai, C. Deng, H. Ding, K. Dong, Q. Du, Z. Fu, et al., Deepseek llm: , arXiv:2401.02954(2024). 10, 25 [139] DeepSeek-AI, A. Liu, B. Feng, B. Wang, B. Wang, B. Liu, C. Zhao, C. Deng, C. Ruan, D. Dai, D. Guo, D. Yang, D. Chen, D. Ji, E. Li, F. Lin, F. Luo, G. Hao, G. Chen, G. Li, H. Zhang, H. Xu, H. Yang, H. Zhang, H. Ding, H. Xin, H. Gao, H. Li, H. Qu, J. L. Cai, J. Liang, J. Guo, J. Ni, J. Li, J. Chen, J. Yuan , J. Qiu, J. Song, K. Dong, K. Gao, K. Guan, L. Wang, L. Zhang, L. Xu, L. Xia, L. Zhao, L. Zhang, M. Li, M. , M. Zhang, M. Zhang, M. Tang, M. Li, N. Tian, P. Huang, P. Wang, P. Zhang, Q. Zhu, Q. Chen, Q. Du, R. J. Chen, R. L. Jin, R. Ge, R. Pan, R. Xu, R. Chen, S. S. Li, S. Lu, S. Zhou, S. Chen, S. Wu, S. Ye, S. Ma, S. Wang, S. Zhou , S. Yu, S. Zhou, S. Zheng, T. Wang, T. Pei, T. Yuan, T. Sun, W. L. Xiao, W. Zeng, W. An, W. Liu, W. Liang, W. Gao , W. Zhang, X. Q. Li, X. Jin, X. Wang, X. Bi, X. Liu, X. Wang, X. Shen, X. Chen, X. Chen, X. Nie, X. Sun, Deepseek-v2: , CoRR abs/2405.04434(2024). 10, 25 [140] E. Nijkamp, B. Pang, H. Hayashi, L. Tu, H. Wang, Y. Zhou, S. Savarese, C. Xiong, Codegen : , arXiv:2203.13474(2022). 11, 23, 25, 28 [141] M. Chen, J. Tworek , H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman , , arXiv:2107.03374(2021). 11 , 25, 29, 31 [142] Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. Dal Lago, et al., alphacode , Science 378(6624)(2022) 1092 – 1097. 11, 23, 25, 29 [143] N. Shazeer, : , arXiv:1911.02150(2019). 11 [144] R. Y. Pang, H. He, , arXiv:2009.07839(2020). 11 [145] R. Dabre, A. Fujita, Softmax , arXiv:2009.09372(2020). 11 [146] Y. Wang, W. Wang, S. Joty, S. C. Hoi, Codet5: , arXiv:2109.00859(2021). 11 [147] R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chim, et al., Starcoder: !, arXiv:2305.06161(2023). 11, 25 [148] R. Taylor, M. Kardas, G. Cucurull, T. Scialom, A. Hartshorn, E. Saravia, A. Poulton, V. Kerkez, R. Stojnic, Galactica: , arXiv:2211.09085(2022). 11, 24, 25, 29 [149] FairScale , FairScale: , https://github.com/facebookresearch/fairscale(2021). 11 [150] R. Thoppilan, D. De Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H.-T. Cheng, A. Jin, T. Bos, L. Baker, Y. Du, et al., Lamda: , arXiv:2201.08239(2022). 11, 25 [151] S. Wu, O. Irsoy, S. Lu, V. Dabrovolski, M. Dredze, S. Gehrmann, P. Kambadur, D. Rosenberg, G. Mann, Bloomberggpt: , arXiv:2303.17564(2023). 11, 25, 33 [152] X. Zhang, Q. Yang, D. Xu, Xuanyuan 2.0: 7, arXiv:2305.12002(2023). 11, 17, 25 [153] W. Ben, Mesh-transformer-jax (2021). 12, 24 [154] N. Muennighoff, T. Wang, L. Sutawika, A. Roberts, S. Biderman, T. L. Scao, M. S. Bari, S. Shen, Z.-X. Yong, H. Schoelkopf, et al., , arXiv:2211.01786(2022). 16, 25, 28, 31 [155] D. Yin, X. Liu, F. Yin, M. Zhong, H. Bansal, J. Han, K.-W. Chang, Dynosaur: ration, arXiv:2305.14327(2023). 16 [156] P. Gao, J. Han, R. Zhang, Z. Lin, S. Geng, A. Zhou, W. Zhang, P. Lu, C. He, X. Yue, et al., Llama-adapter v2: , arXiv:2304.15010(2023). 16, 24 [157] Openai. gpt-4 (2023). 16, 35 [158] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, T. B. Hashimoto, Stanford : , https://github.com/tatsu-lab/stanford\_alpaca(2023). 16, 25, 28 [159] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez, I. Stoica, E. P. Xing, Vicuna: 90%\* chatgpt gpt-4 (2023 3 ). URL https://lmsys.org/blog/2023-03-30-vicuna/ 16, 22, 25, 28 [160] B. Peng, C. Li, P. He, M. Galley, J. Gao, gpt-4 , arXiv:2304.03277(2023). 16, 28 [161] T. Liu, B. K. H. Low, : , arXiv:2305.14201(2023). 16 [162] H. Wang, C. Liu, N. Xi, Z. Qiang, S. Zhao, B. Qin, T. Liu, Huatuo: , arXiv:2304.06975(2023). 16 [163] C. Xu, Q. Sun, K. Zheng, X. Geng, P. Zhao, J. Feng, C. Tao, D. Jiang, Wizardlm: , arXiv:2304.12244(2023). 16 [164] Z. Luo, C. Xu, P. Zhao, Q. Sun, X. Geng, W. Hu, C. Tao, J. Ma, Q. Lin, D. Jiang, Wizardcoder: evol-instruct , arXiv:2306.08568(2023). 16, 25 [165] J. Menick, M. Trebacz, V. Mikulik, J. Aslanides, F. Song, M. Chadwick, M. Glaese, S. Young, L. Campbell-Gillingham, G. Irving , , arXiv:2203.11147(2022). 17 [166] R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang , C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders , Webgpt: , arXiv:2112.09332(2021). 17, 19, 20, 25, 31 [167] A. Glaese, N. McAleese, M. Třebacz, J. Aslanides, V. Firoiu, T. Ewalds, M. Rauh, L. Weidinger, M. Chadwick, P. Thacker , , arXiv:2209.14375(2022). 17, 20, 25 [168] R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, C. Finn, : , arXiv:2305.18290(2023). 17 [169] H. Dong, W. Xiong, D. Goyal, R. Pan, S. Diao, J. Zhang, K. Shum, T. Zhang, Raft: , arXiv:2304.06767(2023). 17 [170] Z. Yuan, H. Yuan, C. Tan, W. Wang , S. Huang, F. Huang, Rhf: , arXiv:2304.05302(2023). 17 [171] F. Song, B. Yu, M. Li, H. Yu, F. Huang, Y. Li, H. Wang, , arXiv:2306.17492(2023). 17 [172] H. Liu, C. Sferrazza, P. Abbeel, : , arXiv:2302.02676(2023). 17 [173] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon , ai: ai , arXiv:2212.08073(2022). 17 [174] Y. Dubois, X. Li, R. Taori, T. Zhang, I. Gulrajani, J. Ba, C. Guestrin, P. Liang, T. B. Hashimoto, AlpacaFarm: , arXiv:2305.14387(2023). 17 [175] C. Si, Z. Gan, Z. Yang, S. Wang, J. Wang, J. Boyd-Graber, L. Wang, gpt-3 , arXiv:2210.09150(2022). 17 [176] D. Ganguli, A. Askell, N. Schiefer, T. Liao, K. Lukošiuaitė, A. Chen, A. Goldie, A. Mirhoseini, C. Olsson, D. Hernandez, et al., , arXiv:2302.07459(2023). 17 [177] A. Wei, N. Haghtalab, J. Steinhardt, : LLM M ?, arXiv:2307.02483(2023). 17

- [178] D. Ganguli, L. Lovitt, J. Kernion, A. Askell, Y. Bai, S. Kadavath, B. Mann, E. Perez, N. Schiefer, K. Ndousse, [arXiv:2209.00000](#), 2022. 17 [179] S. Casper, J. Lin, J. Kwon, G. Culp, D. Hadfield-McEll, [arXiv:2306.09442\(2023\)](#). 17 [180] E. Perez, S. Huang, F. Song, T. Cai, R. Ring, J. Aslanides, A. Glaese, N. McAleese, G. Irving, [arXiv:2202.03286\(2022\)](#). 17 [181] T. Scialom, T. Chakrabarty, S. Muresan, [arXiv:2202.03286\(2022\)](#). 2022, 6107-6122. 17 [182] Z. Shi, A. Lipani, [arXiv:2305.01711\(2023\)](#). 17 [183] H. Gupta, S. A. Sawant, S. Mishra, M. Nakamura, A. Mitra, S. Mashtty, C. Baral, [arXiv:2306.05539\(2023\)](#). 17 [184] H. Chen, Y. Zhang, Q. Zhang, H. Yang, X. Hu, X. Ma, Y. Yanggong, J. Zhao, [arXiv:2305.09246\(2023\)](#). 17 [185] C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al., [arXiv:2305.11206\(2023\)](#). 17, 25, 28 [186] C. Han, Q. Wang, W. Xiong, Y. Chen, H. Ji, S. Wang, Lm-infinite: [arXiv:2308.16137\(2023\)](#). 17, 18 [187] J. Ainslie, T. Lei, M. de Jong, S. Ontañón, S. Brahma, Y. Zemlyan-skiy, D. Uthus, M. Guo, J. Lee-Thorp, Y. Tay, et al., [arXiv:2303.09752\(2023\)](#). 18 [188] J. Ding, S. Ma, L. Dong, X. Zhang, S. Huang, W. Wang, F. Wei, Longnet: [arXiv:2307.02486\(2023\)](#). 18 [189] Y. Chen, S. Qian, H. Tang, X. Lai, Z. Liu, S. Han, J. Jia, Longlora: [arXiv:2309.12307\(2023\)](#). 18 [190] N. Ratner, Y. Levine, Y. Belinkov, O. Ram, I. Magar, O. Abend, E. Karpas, A. Shashua, K. Leyton-Brown, Y. Shoham, [arXiv:2308.16137\(2023\)](#). 18 [191] W. Wang, L. Dong, H. Cheng, X. Liu, X. Yan, J. Gao, F. Wei, [arXiv:2306.07174\(2023\)](#). 18 [192] X. Xu, Z. Gou, W. Wu, Z.-Y. Niu, H. Wu, H. Wang, S. Wang, [arXiv:2203.05797\(2022\)](#). 18 [193] S. Borgeaud, A. Mensch, J. Hoffmann, T. Cai, E. Rutherford, K. Milli-cian, G. B. Van Den Driessche, J.-B. Lespiau, B. Damoc, A. Clark, [arXiv:2305.10250\(2023\)](#). 18, 19, 34 [194] W. Zhong, L. Guo, Q. Gao, Y. Wang, Memorybank: [arXiv:2303.11366\(2023\)](#). 18, 20 [196] C. Hu, J. Fu, C. Du, S. Luo, J. Zhao, H. Zhao, Chatdb: [arXiv:2306.03901\(2023\)](#). 18 [197] Z. Jiang, F. F. Xu, L. Gao, Z. Sun, Q. Liu, J. Dwivedi-Yu, Y. Yang, J. Callan, G. Neubig, [arXiv:2305.06983\(2023\)](#). 18 [198] O. Ram, Y. Levine, I. Dalmedigos, D. Muhlgay, A. Shashua, K. Leyton-Brown, Y. Shoham, [arXiv:2302.00083\(2023\)](#). 18, 34 [199] X. Li, X. Qiu, Mot: [arXiv:2305.05181\(2023\)](#). 18 [200] D. Schuurmans, [arXiv:2301.04589\(2023\)](#). 18 [201] A. Modarressi, A. Imani, M. Fayyaz, H. Schütze, Ret-llm: [arXiv:2305.14322\(2023\)](#). 18 [202] S. Robertson, H. Zaragoza, [arXiv:2305.14322\(2023\)](#). 18 [203] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, D. Zhou, [arXiv:2305.14322\(2023\)](#). 18 [204] F. Zhang, B. Chen, Y. Zhang, J. Liu, D. Zan, Y. Mao, J.-G. Lou, W. Chen, Repocoder: [arXiv:2303.12570\(2023\)](#). 18 [205] B. Wang, W. Ping, P. Xu, L. McAfee, Z. Liu, M. Shoenybi, Y. Dong, O. Kuchaiev, B. Li, C. Xiao, et al., [arXiv:2304.06762\(2023\)](#). 19 [206] L. Wang, N. Yang, F. Wei, [arXiv:2307.07164\(2023\)](#). 19 [207] J. Liu, D. Shen, Y. Zhang, B. Dolan, L. Carin, W. Chen, gpt-3: [arXiv:2101.06804\(2021\)](#). 19 [208] O. Rubin, J. Herzig, J. Berant, [arXiv:2112.08633\(2021\)](#). 19 [209] W. Shi, S. Min, M. Yasunaga, M. Seo, R. James, M. Lewis, L. Zettlemoyer, W.-t. Yih, Replug: [arXiv:2301.12652\(2023\)](#). 19 [210] O. Rubin, J. Berant, [arXiv:2306.13421\(2023\)](#). 19 [211] K. Guu, K. Lee, Z. Tung, P. Pasupat, M. Chang, [arXiv:2306.13421\(2023\)](#). 19 [212] S. Hofstätter, J. Chen, K. Raman, H. Zamani, Fidelity: [arXiv:2306.13421\(2023\)](#). 19 [213] M. Komeili, K. Shuster, J. Weston, [arXiv:2107.07566\(2021\)](#). 19 [214] A. Lazaridou, E. Gribovskaya, W. Stokowiec, N. Grigorev, [arXiv:2203.05115\(2022\)](#). 19 [215] D. Gao, L. Ji, L. Zhou, K. Q. Lin, J. Chen, Z. Fan, M. Z. Shou, Assist-gpt: [arXiv:2306.08640\(2023\)](#). 19 [216] P. Lu, B. Peng, H. Cheng, M. Galley, K.-W. Chang, Y. N. Wu, S.-C. Zhu, J. Gao, [arXiv:2304.09842\(2023\)](#). 19, 20, 23 [217] B. Paranjape, S. Lundberg, S. Singh, H. Hajishirzi, L. Zettlemoyer, M. T. Ribeiro, [arXiv:2303.09014\(2023\)](#). 19 [218] C.-Y. Hsieh, S.-A. Chen, C.-L. Li, Y. Fujii, A. Ratner, C.-Y. Lee, R. Krishna, T. Pfister, [arXiv:2308.00675\(2023\)](#). 19 [219] Y. Song, W. Xiong, D. Zhu, C. Li, K. Wang, Y. Tian, S. Li, Restgpt: Restful API [arXiv:2306.06624\(2023\)](#). 19 [220] S. Hao, T. Liu, Z. Wang, Z. Hu, Toolkengpt: [arXiv:2305.11554\(2023\)](#). 19 [221] S. G. Patil, T. Zhang, X. Wang, J. E. Gonzalez, Gorilla: [arXiv:2305.15334\(2023\)](#). 19 [222] Q. Xu, F. Hong, B. Li, C. Hu, Z. Chen, J. Zhang, [arXiv:2305.16504\(2023\)](#). 19 [223] Y. Qin, S. Liang, Y. Ye, K. Zhu, L. Yan, Y. Lu, Y. Lin, X. Cong, X. Tang, B. Qian, et al., Toolllm: 16000+ [arXiv:2307.16789\(2023\)](#). 19,



- 20 [224] Y. Shen, K. Song, X. Tan, D. Li, W. Lu, Y. Zhuang, Hugging gpt: huggingface chatgpt ai , arXiv arXiv:2303.17580 (2023). 19, 20, 33 [225] Y. Liang, C. Wu, T. Song, W. Wu, Y. Xia, Y. Liu, Y. Ou, S. Lu, L. Ji, S. Mao, et al., Taskmatrix.ai: api , arXiv arXiv:2303.16434 (2023). 19 [226] D. Surís, S. Menon, C. Vondrick, Vipergpt: Python , arXiv arXiv:2303.08128 (2023). 20 [227] A. Maedche, S. Morana, S. Schacht, D. Werth, J. Krumeich, , 58 (2016) 367 – 370. 20 [228] M. Campbell, A. J. Hoane Jr, F.-h. Hsu, , 134 (1-2) (2002) 57 – 83. 20 [229] S. Hong, X. Zheng, J. Chen, Y. Cheng, J. Wang, C. Zhang, Z. Wang, S. K. S. Yau, Z. Lin, L. Zhou , Metagpt: , arXiv arXiv:2308.00352(2023). 20 [230] Z. Xi, W. Chen, X. Guo, W. He, Y. Ding, B. Hong, M. Zhang, J. Wang, S. Jin, E. Zhou , , arXiv arXiv:2309.07864(2023). 20 [231] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin, et al., , arXiv arXiv:2308.11432(2023). 20 [232] W. Huang, P. Abbeel, D. Pathak, I. Mordatch, , PMLR, 2022, pp. 9118 – 9147. 20 [233] S. Hao, Y. Gu, H. Ma, J. J. Hong, Z. Wang, D. Z. Wang, Z. Hu, , arXiv arXiv:2305.14992 (2023). 20, 33 [234] W. Yao, S. Heinecke, J. C. Nieves, Z. Liu, Y. Feng, L. Xue, R. Murthy, Z. Chen, J. Zhang, D. Arpit, et al., Retroformer: , arXiv arXiv:2308.02151 (2023). 20, 33 [235] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch, Y. Chebotar, P. Sermanet, T. Jackson, N. Brown, L. Luu, S. Levine, K. Hausman, Brian Ichter, , in: 6 , 2022. URL https://openreview.net/forum?id=3R3Pz5i0tye 20 [236] C. Jin, W. Tan, J. Yang, B. Liu, R. Song, L. Wang, J. Fu, Alphablock: , arXiv arXiv:2305.18898 (2023). 20, 33 [237] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomas on, A. Garg, Progprompt: , 2023 IEEE (ICRA), IEEE, 2023, pp. 11523 – 11530. 20, 33 [238] W. Yu, N. Gileadi, C. Fu, S. Kirmani, K.-H. Lee, M. G. Arenas, H.-T. L. Chiang, T. Erez, L. Hasenclever, J. Humpalik , arXiv arXiv:2306.08647(2023). 20 [239] X. Tang, A. Zou, Z. Zhang, Y. Zhao, X. Zhang, A. Cohen, M. Gerstein, Medagents: , arXiv arXiv:2311.10537(2023). 20 [240] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian , 가 , PMLR, 2023, 287-318 . 20, 33 [241] H. Ha, P. Florence, S. Song, , : , arXiv arXiv:2307.14535(2023). 20 [242] A. Rajvanshi, K. Sikka, X. Lin, B. Lee, H.-P. Chiu, A. Velasquez, Say-nav: , arXiv arXiv:2309.04077(2023). 20 [243] C. H. Song, J. Wu, C. Washington, B. M. Sadler, W.-L. Chao, Y. Su, Llm-planner: Few-shot , arXiv arXiv:2212.04088(2022). 20 [244] V. S. Dorbala, J. F. Mullen Jr, D. Manocha, 가 , arXiv arXiv:2303.03480(2023). 20 [245] C. Huang, O. Mees, A. Zeng, W. Burgard, , 2023 IEEE (ICRA), IEEE, 2023, pp. 10608 – 10615. 20 [246] Y. Ding, X. Zhang, C. Paxton, S. Zhang, , arXiv arXiv:2303.06247(2023). 20, 33 [247] X. Liu, Y. Zheng, Z. Du, M. Ding, Y. Qian, Z. Yang, J. Tang, Gpt arXiv arXiv:2103.10385(2021) . 20, 21 [248] G. Chen, F. Liu, Z. Meng, S. Liang, , : , arXiv arXiv:2202.07962(2022). 20 [249] Y. Wang, S. Mukherjee, X. Liu, J. Gao, A. H. Awadallah, J. Gao, Adamix: , arXiv arXiv:2205.12410 1 (2) (2022) 4. 20 [250] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, Lora: , arXiv arXiv:2106.09685 (2021). 21, 22, 23 [251] X. Liu, K. Ji, Y. Fu, W. Tam, Z. Du, Z. Yang, J. Tang, P-tuning: Prompt tuning can be comparable to fine-tuning across scales and task, in: Pro-ceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), 2022, pp. 61 – 68. 21 [252] A. Razdaibiedina, Y. Mao, R. Hou, M. Khabza, M. Lewis, A. Almahairi, Progressive prompts: Continual learning for language models, arXiv preprint arXiv:2301.12314 (2023). 21 [253] Z.-R. Zhang, C. Tan, H. Xu, C. Wang, J. Huang, S. Huang, Towards adaptive prefix tuning for parameters-efficient language model fine-tuning, arXiv preprint arXiv:2305.15212 (2023). 21 [254] E. B. Zaken, S. Ravfogel, Y. Goldberg, Bitfit: , arXiv preprint arXiv:2106.10199 (2021). 21 [255] T. Dettmers, M. Lewis, Y. Belkada, L. Zettlemoyer, Llm.int8(): 8 , arXiv preprint arXiv:2208.07339 (2022). 21, 22 [256] E. Frantar, S. Ashkboos, T. Hoefler, D. Alistarh, Gptq: , arXiv arXiv:2210.17323 (2022). 21 [257] X. Wei, Y. Zhang, Y. Li, X. Zhang, R. Gong, J. Guo, X. Liu, +: , arXiv arXiv:2304.09145 (2023). 21 [258] E. Frantar, D. Alistarh, : 가 , 35(2022) 4475 – 4488. 21 [259] C. Lee, J. Jin, T. Kim, H. Kim, E. Park, Owq: 가 , arXiv arXiv:2306.02272(2023). 21 [260] S. J. Kwon, J. Kim, J. Bae, K. M. Yoo, J.-H. Kim, B. Park, B. Kim, J.-W. Hwang, N. Sung, D. Lee, Alpatuning: Quantization-aware parameters-efficient adaptation of large-scale pre-trained language models, arXiv preprint arXiv:2210.03858 (2022). 21 [261] T. Dettmers, A. Pagnoni, A. Holtzman, L. Zettlemoyer, Qlora: Efficient finetuning of quantized llms, arXiv preprint arXiv:2305.14314 (2023). 21, 22 [262] Z. Liu, B. Oguz, C. Zhao, E. Chang, P. Stock, Y. Mehdad, Y. Shi, R. Krishnamoorthi, V. Chandra, Llm-qat: , arXiv arXiv:2305.17888(2023). 21, 22 [263] Y. Guo, A. Yao, H. Zhao, Y. Chen, : CNN , IEEE , 2017, 5955-5963 . 21 [264] J. Kim, J. H. Lee, S. Kim, J. Park, K. M. Yo, S. J. Kwon, D. Lee, 4 , arXiv arXiv:2305.14152(2023). 22 [265] M. Sun, Z. Liu, A. Bair, J. Z. Kolter, 가 , arXiv arXiv:2306.11695(2023). 22 [266] Z. Wang, J. Wohlwend, T. Lei, 가 , arXiv arXiv:1910.04732(2019). 22

- [267] L. Yin, Y. Wu, Z. Zhang, C.-Y. Hsieh, Y. Wang, Y. Jia, M. Pechenizkiy, Y. Liang, Z. Wang, S. Liu, 가 ( ): llms 가 , arXiv arXiv:2310.05175(2023). 22 [268] C. Tao, L. Hou, H. Bai, J. Wei, X. Jiang, Q. Liu, P. Luo, N. Wong, 가 , in: Findings of the Association for Computational Linguistics: ACL 2023, 2023, pp. 10880 – 10895. 22 [269] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, et al., Flamingo: a visual language model for few-shot learning, 35(2022) 23716 – 23736. 22 [270] J. Li, D. Li, S. Savarese, S. Hoi, Blip-2: , arXiv arXiv:2301.12597(2023). 22 [271] H. Liu, C. Li, Q. Wu, Y. J. Lee, Visual instruction tuning, arXiv arXiv:2304.08485(2023). 22 [272] K. Li, Y. He, Y. Wang, Y. Li, W. Wang, P. Luo, Y. Wang, L. Wang, Y. Qiao, Videochat: , arXiv arXiv:2305.06355(2023). 22 [273] M. Maaz, H. Rasheed, S. Khan, F. S. Khan, Video-chatgpt: , arXiv arXiv:2306.05424(2023). 2 [274] H. Zhang, X. Li, L. Bing, Video-llama: , arXiv arXiv:2306.02858(2023). 22 [275] X. Mei, C. Meng, H. Liu, Q. Kong, T. Ko, C. Zhao, M. D. Plumbley, Y. Zou, W. Wang, Wavcaps: - chatgpt , arXiv arXiv:2303.17395(2023). 22 [276] C. Lyu, M. Wu, L. Wang, X. Huang, B. Liu, Z. Du, S. Shi, Z. Tu, Macaw-llm: , arXiv arXiv:2306.09093(2023). 22 [277] D. Zhu, J. Chen, X. Shen, X. Li, M. Elhoseiny, Minigpt-4: , arXiv arXiv:2304.10592(2023). 22 [278] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., 16x16 가 가 : , arXiv arXiv:2010.11929(2020). 22 [279] W. Dai, J. Li, D. Li, A. M. H. Tiong, J. Zhao, W. Wang, B. Li, P. Fung, S. Hoi, Instructblip: , arXiv arXiv:2305.06500(2023). 22 [280] Z. Xu, Y. Shen, L. Huang, Multiinstruct: , arXiv arXiv:2212.10773(2022). 22 [281] Z. Zhao, L. Guo, T. Yue, S. Chen, S. Shao, X. Zhu, Z. Yuan, J. Liu, Chatbridge: , arXiv arXiv:2305.16103 (2023). 22 [282] L. Li, Y. Yin, S. Li, L. Chen, P. Wang, S. Ren, M. Li, Y. Yang, J. Xu, X. Sun, et al., M3 it: , arXiv arXiv:2306.04387 (2023). 22 [283] R. Pi, J. Gao, S. Diao, R. Pan, H. Dong, J. Zhang, L. Yao, J. Han, H. Xu, L. K. T. Zhang, Detgpt: , arXiv arXiv:2305.14167(2023). 22 [284] G. Luo, Y. Zhou, T. Ren, S. Chen, X. Sun, R. Ji, : , arXiv arXiv:2305.15023(2023). 22 [285] R. Zhang, J. Han, A. Zhou, X. Hu, S. Yan, P. Lu, H. Li, P. Gao, Y. Qiao, Llama-adapter: zero-init attention , arXiv arXiv:2303.16199(2023). 22 [286] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, I. Sutskever, , PMLR, 2023, pp. 28492 – 28518. 22 [287] Z. Zhang, A. Zhang, M. Li, H. Zhao, G. Karypis, A. Smola, , arXiv arXiv:2302.00923(2023). 23 [288] J. Ge, H. Luo, S. Qian, Y. Gan, J. Fu, S. Zhan, , arXiv arXiv:2304.07919(2023). 23 [289] C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, N. Duan, chatgpt: , arXiv arXiv:2303.04671 (2023). 23 [290] Z. Yang, L. Li, J. Wang, K. Lin, E. Azarnasab, F. Ahmed, Z. Liu, C. Liu, M. Zeng, L. Wang, Mm-react: chatgpt , arXiv arXiv:2303.11381 (2023). 2 [291] T. Wang, J. Zhang, J. Fei, Y. Ge, H. Zheng, Y. Tang, Z. Li, M. Gao, S. Zhao, Y. Shan , : , arXiv arXiv:2305.02677(2023). 23 [292] X. Zhu, R. Zhang, B. He, Z. Zeng, S. Zhang, P. Gao, Pointclip v2 : 3D , arXiv arXiv:2211.11682(2022). 2 [293] T. Gupta, A. Kembhavi, : , IEEE/CVF , 2023, p. 14953 – 14962. 23 [294] P. Gao, Z. Jiang, H. You, P. Lu, S. C. Hoi, X. Wang, H. Li, , IEEE /CVF , 2019, 6639-6648 . 23 [295] Z. Yu, J. Yu, Y. Cui, D. Tao, Q. Tian, , IEEE/CVF , 2019, 6281-6290 . 23 [296] H. You, R. Sun, Z. Wang, L. Chen, G. Wang, H. A. Ayyubi, K.-W. Chang, S.-F. Chang, Idealgpt: , arXiv arXiv:2305.14985(2023). 23 [297] R. Zhang, X. Hu, B. Li, S. Huang, H. Deng, Y. Qiao, P. Gao, H. Li, , : -shot , IEEE/CVF , 2023, pp. 15211 – 15222. 23 [298] T. Q. Nguyen, J. Salazar, : , CoRR abs/1910.05895(2019). 24 [299] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, V. Stoyanov, Roberta: bert , arXiv arXiv:1907.11692(2019). 24, 30 [300] X. Geng, A. Gudibande, H. Liu, E. Wallace, P. Abbeel, S. Levine, D. Song, Koala: , (2023 4 ). URL <https://bair.berkeley.edu/blog/2023/04/03/koala/> 25 [301] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite, N. Nabeshima, et al., The pile: 800gb , arXiv arXiv:2101.00027(2020). 28, 30 [302] H. Laurençon, L. Saulnier, T. Wang, C. Akiki, A. Villanova del Moral, T. Le Scao, L. Von Werra, C. Mou, E. González Ponferrada, H. Nguyen , et al., : 1.6TB 35 (2022) 31809-31826. 28 [303] . URL [https://en.wikipedia.org/wiki/Main\\_Page](https://en.wikipedia.org/wiki/Main_Page) 28 [304] Together Computer, Redpajama: (2023 4 ). URL <https://github.com/togethercomputer/RedPajama-Data> 28 [305] O. Honovich, T. Scialom, O. Levy, T. Schick, : , arXiv arXiv:2212.09689 (2022). 28 [306] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan , et al., , arXiv arXiv:2204.05862(2022). 2 [307] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, J. Steinhardt, , arXiv arXiv:2009.03300(2020). 26, 29 [308] A. Srivastava, A. Rastogi, A. Rao, A. A. M. Shole, A. Abid, A. Fisch, A. R. Brown, A. Santoro, A. Gupta, A. Garriga-Alonso , Beyond

- , arXiv arXiv:2206.04615 (2022). 26, 29 [309] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, S. R. Bowman, Glue: , arXiv arXiv:1804.07461 (2018). 26, 29 [310] Y. Yao, Q. Dong, J. Guan, B. Cao, Z. Zhang, C. Xiao, X. Wang, F. Qi, J. Bao, J. Nie, et al., Cuge: , arXiv arXiv:2112.13610 (2021). 29 [311] L. Xu, H. Hu, X. Zhang, L. Li, C. Cao, Y. Li, Y. Xu, K. Sun, D. Yu, C. Yu, et al., Clue: , arXiv arXiv:2004.05986 (2020). 29 [312] L. Xu, X. Lu, C. Yuan, X. Zhang, H. Xu, H. Yuan, G. Wei, X. Pan, X. Tian, L. Qin, et al., Fewclue: few-shot , arXiv arXiv:2107.07498 (2021). 29 [313] E. M. Smith, M. Williamson, K. Shuster, J. Weston, Y.-L. Boureau, Can you put it all together: Evaluating conversational agents' ability to blend skills, arXiv preprint arXiv:2004.08449 (2020). 29 [314] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, et al., Holistic evaluation of language models, arXiv preprint arXiv:2211.09110 (2022). 29 [315] S. Park, J. Moon, S. Kim, W. I. Cho, J. Han, J. Park, C. Song, J. Kim, Y. Song, T. Oh, et al., Klue: , arXiv preprint arXiv:2105.09680 (2021). 29 [316] S. Reddy, D. Chen, C. D. Manning, Coqa: , 7(2019) 249–266. 27, 29 [317] M. T. Pilehvar, J. Camacho-Collados, Wic: , 10,000 , arXiv arXiv:1808.09121 (2018). 27, 29 [318] S. Merity, C. Xiong, J. Bradbury, R. Socher, , arXiv arXiv:1609.07843(2016). 28, 29 [319] J. W. Rae, A. Potapenko, S. M. Jayakumar, T. P. Lillicrap, , arXiv arXiv:1911.05507(2019). 28, 29 [320] X. Liu, Q. Chen, C. Deng, H. Zeng, J. Chen, D. Li, B. Tang, Lcqm: , 27 , 2018, pp. 1952–1962. 28, 29 [321] S. Iyer, N. Dandekar, K. Csernai, quora : , https://quoradata.quora.com/ First-Quora-Dataset-Release-Question-Pairs. 29 [322] R. Rudinger, J. Naradowsky, B. Leonard, B. Van Durme, , arXiv arXiv:1804.09301(2018). 29 [323] M.-C. De Marneffe, M. Simons, J. Tonhause r, : , 23 , 2019 , 107-124 . 29 [324] Z. Li, N. Ding, Z. Li u, H. Zheng, Y. Shen, , 2019 57 , 4377-4386 . 29 [325] J. Xu, J. Wen, X. Sun, Q. Su, , arXiv arXiv:1711.07010(2017). 29 [326] J. Chen, Q. Chen, X. Liu, H. Yang, D. Lu, B. Tang, bq : , 2018 , 4946-4951 . 29 [327] B. Liu, D. Niu, H. Wei, J. Lin, Y. He, K. Lai, Y. Xu, , arXiv arXiv:1802.07459(2018) . 29 [328] P. Li, W. Li, Z. He, X. Wang, Y. Cao, J. Zhou, W. Xu, , arXiv arXiv:1607.06275(2016). 29 [329] N. Peng, M. Dredze, , 2015 , 2015, 548-554 . 29 [330] W. Ling, D. Yogatama, C. Dyer, P. Blunsom, : , arXiv arXiv:1705.04146(2017). 29 [331] R. Weischedel, S. Pradhan, L. Ramshaw, M. Palmer, N. Xue, M. Marcus, A. Taylor, C. Greenberg, E. Hovy, R. Belvin , Ontonotes 4.0, LDC2011T03, , : (2011). 29 [332] D. Vilarés, C. Gómez-Rodríguez, Head-qa: , arXiv arXiv:1906.04701(2019). 29 [333] S. L. Blodgett, L. Green, B. O'Connor, , arXiv arXiv:1608.08868(2016). 29 [334] N. Mostafazadeh, N. Chambers, X. He, D. Parikh, D. Batra, L. Van-derwende, P. Kohli, J. Allen, , arXiv arXiv:1604.01696(2016). 28, 29 [335] D. Paperno, G. Kruszewski, A. Lazaridou, Q. N. Pham, R. Bernardi, S. Pezzelle, M. Baroni, G. Boleda, R. Fernández, : , arXiv arXiv:1606.06031(2016). 28, 29 [336] B. Hu, Q. Chen, F. Zhu, Lcsts: , arXiv arXiv:1506.05865(2015). 29 [337] Z. Shao, M. Huang, J. Wen, W. Xu, X. Zhu, , arXiv arXiv:1908.06605(2019). 29 [338] J. Novikova, O. Dušek, V. Rieser, e2e : , arXiv arXiv:1706.09254(2017). 29 [339] C. Zheng, M. Huang, A. Sun, Chid: , arXiv arXiv:1906.01265(2019). 29 [340] Y. Bisk, R. Zellers, J. Gao, Y. Choi , Piqa: , AAAI , 34 , 2020 , 7432~7439 . 28, 29 [341] M. Joshi, E. Choi, D. S. Weld, L. Zettlemoyer, Triviaqa: , arXiv arXiv:1705.03551(2017). 28, 29, 31 [342] P. Clark, I. Cowhey, O. Etzioni, T. Khot, A. Sabharwal, C. Schoenick, O. Tafjord, ? try arc, ai2 , arXiv arXiv:1803.05457(2018). 28, 29, 31 [343] S. Aroca-Ouellette, C. Paik, A. Roncone, K. Kann, Prost: , arXiv arXiv:2106.03634(2021). 29 [344] T. Mihaylov, P. Clark, T. Khot, A. Sabharwal, ? , arXiv arXiv:1809.02789(2018). 29 [345] T. C. Ferreira, C. Gardent, N. Ilinykh, C. Van Der Lee, S. Mille, D. Moussallem, A. Shimorina, 2020 , webnlg+ , 3 (webnlg+ WebN LG+), 2020 . 29 [346] C. Xu, W. Zhou, T. Ge, K. Xu, J. McAuley, F. Wei, : , arXiv arXiv:2104.02704(2021). 29 [347] G. Lai, Q. Xie, H. Liu, Y. Yang, E. Hovy, : , arXiv arXiv:1704.04683(2017). 29 [348] E. Choi, H. He, M. Iyyer, M. Yatskar, W.-t. Yih, Y. Choi, P. Liang, L. Zettlemoyer, Quac: , arXiv arXiv:1808.07036(2018). 29 [349] M. Geva, D. Khashabi, E. Segal, T. Khot, D. Roth, J. Berant, ? , 9(2021) 346–361. 29, 31 [350] J. Boyd-Graber, B. Sattinoff, H. He, H. Daumé III, : , 2012 , 2012, pp. 1290–1301. 29 [351] S. Zhang, X. Zhang, H. Wang, J. Cheng, P. Li, Z. Ding, CNN 7(8)(2017) 767. 29 [352] S. Zhang, X. Zhang, H. Wang, L. Guo, S. Liu, , IEEE

- Access 6(2018) 74061 – 74071. 29 [353] C. Xu, J. Pei, H. Wu, Y. Liu, C. Li, Matinf: , arXiv arXiv:2004.12302(2020). 29 [354] K. Sa kaguchi, R. L. Bras, C. Bhagavatula, Y. Choi, Winogrande: Winograd , ACM 64(9)(2021) 99 – 106. 27, 29 [355] R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, Y. Choi, Hellaswag: 가 , arXiv arXiv:1905.07830(2019). 29 [356] M. Roemmele, C. A. Bejan, A. S. Gordon, 가 : 가, AAAI : , 2011, 90-95 . 29 [357] H. Levesque, E. Davis, L. Morge nstern, Winograd , 13 , 2012. 27, 29 [358] A. Talmor, J. Herzig, N. Lourie, J. Berant, Commonsenseqa: , arXiv arXiv:1811.00937(2018). 29, 31 [359] M. Sap, H. Rashkin, D. Che n, R. LeBras, Y. Choi, Socialqa: , ar Xiv arXiv:1904.09728(2019). 29 [360] K. Sun, D. Yu, D. Yu, C. Cardie, , 8(2020) 141 – 155. 29 [361] S. Zhang, X. Liu, J. Liu, J. Ga o, K. Duh, B. Van Durme, : , arXiv arXiv:1810.12885(2018). 29 [362] P. Rajpurkar, J. Zhang, K. Lopyrev, P. Liang, Squad: 100,0 00+ , arXiv arXiv:1606.05250(2016). 29, 31 [363] C. Clark, K. Lee, M.-W. Chang, T. Kwiatkowski, M. Collins, K. Toutanova, Boolq: / , arXiv arXiv:1905.10044(2019). 29, 31 [364] P. Rajpurkar, R. Jia, P. Liang, : Squad , arXiv arXiv:1806.03822(2018). 29, 31 [365] D. Dua, Y. Wang, P. Dasigi, G. Stanovsky, S. Singh, M. Gardner, Drop: , arXiv arXiv:1903.00161(2019). 29, 31 [366] I. Dagan, O. Glickman, B. Magnini, , Springer, 2005, pp. 177 – 190. 29, 31 [367] Y. Chang, M. Narang, H. Suzuki, G. Cao, J. Gao, Y. Bisk, Webqa: Multi-tihop multi modal qa, IEEE/CVF , 2022, pp. 16495 – 16504. 29, 31 [368] Y. Cui, T. Liu, Z. Chen, W. Ma, S. Wang, G. Hu, 가 , arXiv arXiv:1709.08299(2017). 29 [369] Y. Cui, T. Liu, W. Che, L. Xiao, Z . Chen, W. Ma, S. Wang, G. Hu, , arXiv arXiv:1810.07366(2018). 29, 31 [370] Y. Cui, T. Liu, Z. Yang, Z. Chen, W. Ma, W. Che, S. Wang, G. Hu, , arXiv arXiv:2004 .03116(2020). 29 [371] Y. Li, T. Liu, D. Li, Q. Li, J. Shi, Y. Wang, pos bilstm-crf, Asian Confer ence on Machine Learning, PMLR, 2018, 518-533 . 29 [372] D. Khashabi, S . Chaturvedi, M. Roth, S. Upadhyay, D. Roth, : , 2018 : , 1 ( ), 2018, 252-262 . 29 [373] T. Kwi atkowski, J. Palomaki, O. Redfield, M. Collins, A. Parikh, C. Alberti, D. Epste in, I. Polosukhin, J. Devlin, K. Lee, et al., : 7(2019) 453 – 466. 29 [374] C. C. Shao, T. Liu, Y. Lai, Y. Tseng, S. Tsai, Dred: , arXiv arXiv:1806.00920(2018). 29 [375] W. He , K. Liu, J. Liu, Y. Lyu, S. Zhao, X. Xiao, Y. Liu, Y. Wang, H. Wu, Q. She, et al., Dureader: , arXiv arXiv:1711.05073(2017). 29 [376] H. Tang, J. Liu, H. Li, Y. Hong, H. Wu, H. Wang, Dureaderrobust: 가 , arXiv arXiv:2004.11142(20 20). 29 [377] J. Welbl, N. F. Liu, M. Gardner, , arXiv arXiv:1707.06209(2017). 29 [378] C. Xiong, Z. Dai, J. Callan, Z. Liu, R. Power, , 40 ACM SIGIR , 201 7, 55-64 . 29 [379] A. Peñas, E. Hovy, P. Forner, Á. Rodrigo, R. Sutcliffe, R . Morante, Qa4mre 2011-2013: 가 , Info rmation Access Evaluation , CLEF 2013, , 2013 9 23-26 . Pro-ceedings 4, Springer, 2013, pp. 303-320. 29 [380] S. Lim, M. Ki m, J. Lee, Korquad1.0: qa , arXi v arXiv:1909.07005(2019). 29 [381] C. Xiao, H. Zhong, Z. Guo, C. Tu, Z. Liu, M. Sun, Y. Feng, X. Han, Z. Hu, H. Wang, et al., Cail2018: , arXiv arXiv:1807.02 478(2018). 29 [382] D. Hendrycks, S. Basart, S. Kadavath, M. Mazeika, A. Ar ora, E. Guo, C. Burns, S. Puranik, H. He, D. Song, et al., , arXiv arXiv:2105.09938(2021). 29, 31 [383] Y . Wang, X. Liu, S. Shi, 2017 , 2017, 845-854 . 29, 31 [38 4] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plapp ert, J. Tworek, J. Hilton, R. Nakano, et al., , arXiv arXiv:2110.14168(2021). 29, 31 [385] J. Austin, A. Odena, M. I. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. J. Ca i, M. Terry, Q. V. Le, C. Sutton, , CoRR abs/2108.07732(2021). 29 [386] F. Shi, M. Suzgun, M. Freitag, X. Wa ng, S. Srivats, S. Vosoughi, H. W. Chung, Y. Tay, S. Ruder, D. Zhou, et al., , arXiv arXiv:2210.03057( 2022). 29 [387] S. Roy, D. Roth, , arXiv arXiv:1608.01413(2016). 29 [388] S.-Y. Miao, C.-C. Liang, K.-Y. Su, 가 , arXi v arXiv:2106.15772(2021). 29 [389] R. Koncel-Kedziorski, S. R oy, A. Amini, N. Kushman, H. Hajishirzi, Mawps: , 2 016 : , 2016, 1152-1157 . 29 [390] A. Patel, S. Bhattamishra, N. Goyal, NLP ? , arXiv arXiv:2103.07191(2021). 29 [391] Y. Lai, C. Li, Y. Wang, T. Zhang, R. Zhong, L. Zettlemoyer, W.-t. Y ih, D. Fried, S. Wang, T. Yu, Ds-1000: , International Conference on Machine Lear ning, PMLR, 2023, pp. 18319 – 18345. 29 [392] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le, et al. , , arXiv arXiv:210 8.07732(2021). 29 [393] Y. Nie, A. Williams, E. Dinan, M. Bansal, J. Weston, D. Kiela, nli: , arXiv arXiv:1910.14599(2019). 29, 31 [394] A. Williams, N. Nangia, S. R. Bo wman, , arXiv arXiv:1704.05426(2017). 29 [395] R. T. McCoy, E. Pavlick, T. Lin zen, :

- , arXiv  
arXiv:1902.01007 (2019). 29 [396] J. Liu, L. Cui, H. Liu, D. Huang, Y. Wang, Y. Zhang, Logiqa:  
, arXiv arXiv:2007.08124 (2020). 29 [397] P. Lewis, B. Oğuz, R. Rinott, S. Riedel, H. Schwenk, Mlqa:  
7}, arXiv arXiv:1910.07475 (2019). 29 [398] A. Conneau, G. Lample, R. Rinott, A. Williams, S. R. Bowman, H. Schwenk, V. Stoyanov, Xnli:  
7}, arXiv arXiv:1809.05053(2018). 29, 31  
[399] Y. Yang, Y. Zhang, C. Tar, J. Baldridge, Paws-x:  
, arXiv arXiv:1908.11828(2019). 29, 31  
[400] S. Narayan, S. B. Cohen, M. Lapata,  
!, ArXiv, abs (1808). 29 [401] E. M. Ponti, G. Glavaš, O. Majewska, Q. Liu, I. Vulić, A. Korhonen, Xcopa:  
arXiv:2005.00333 (2020). 29 [402] A. Tikhonov, M. Ryabinin,  
:  
, arXiv arXiv:2106.12066 (2021). 29 [403] J. H. Clark, E. Choi, M. Collins, D. Garrette, T. Kwiatkowski, V. Nikolaev, J. Palomaki, Tydi qa:  
, 8(2020) 454 – 470. 29 [404] T. Scialom, P.-A. Dray, S. Lamprier, B. Piwowarski, J. Staiano, Msum:  
, arXiv arXiv:2004.14900(2020). 29 [405] S. Lin, J. Hilton, O. Evans, Truthfulqa:  
, arXiv arXiv:2109.07958(2021). 29, 32 [406] I. Augenstein, C. Lioma, D. Wang, L. C. Lima, C. Hansen, C. Hansen, J. G. Simonsen, Multific:  
, arXiv arXiv:1909.03242(2019). 29 [407] J. Thorne, A. Vlachos, C. Christodoulopoulos, A. Mittal, Fever:  
, arXiv arXiv:1803.05355(2018). 29 [408] I. Mollas, Z. Chrysopoulou, S. Karlos, G. Tsoumakas, Ethos:  
, arXiv arXiv:2006.08328(2020). 29, 32 [409] M. Nadeem, A. Bethke, S. Reddy, StereoSet:  
, arXiv arXiv:2004.09456(2020). 29, 32 [410] A. Parrish, A. Chen, N. Nangia, V. Padmakumar, J. Phang, J. Thompson, P. M. Htut, S. R. Bowman, Bbq:  
arXiv arXiv:2110.08193(2021). 29 [411] J. Zhao, T. Wang, M. Yatskar, V. Ordonez, K.-W. Chang,  
, arXiv arXiv:1804.06876(2018). 29 [412] N. Nangia, C. Vania, R. Bhalerao, S. R. Bowman, Crows-pairs: 7}, arXiv  
arXiv:2010.00133(2020). 29 [413] S. Gehman, S. Gururangan, M. Sap, Y. Choi, N. A. Smith, Realtocixityprompts:  
7}, arXiv arXiv:2009.11462(2020). 29 [414] D. Borkan, L. Dixon, J. Sorensen, N. Thain, L. Vasserman,  
, 2019  
, 2019, 491-500. 29 [415] O. Bojar, R. Chatterjee, C. Federmann, Y. Graham, B. Haddow, M. Huck, A. J. Yepes, P. Koehn, V. Logacheva, C. Monz, et al., 2016  
, 2016 1  
: 2, , 131-198. 29 [416] B. Loïc, B. Magdalena, B. Ondřej, F. Christian, G. Yvette, G. Roman, H. Barry, H. Matthias, J. Eric, K. Tom, et al., 2020 (wmt20), in: 5, Association for Computational Linguistics „ 2020, pp. 1 – 55. 29
- [417] W. Li, F. Qi, M. Sun, X. Yi, J. Zhang, Ccpm:  
, arXiv arXiv:2106.01979(2021). 29 [418] E. Dinan, S. Roller, K. Shuster, A. Fan, M. Auli, J. Weston,  
, arXiv arXiv:1811.01241(2018). 29 [419] H. Rashkin, E. M. Smith, M. Li, Y.-L. Boureau,  
:  
, arXiv arXiv:1811.00207(2018). 29 [420] E. Dinan, V. Logacheva, V. Malykh, A. Miller, K. Shuster, J. Urbanek, D. Kiela, A. Szlam, I. Serban, R. Lowe, et al., (convai2), NeurIPS'18  
:  
, Springer, 2020, pp. 187 – 208. 29 [421] H. Zhou, C. Zheng, K. Huang, M. Huang, X. Zhu, Kdconv:  
, arXiv arXiv:2004.04100(2020). 29 [422] L. CO, Iflytek:  
, (2019). 29 [423] J. Baumgartner, S. Zannettou, B. Keegan, M. Squire, J. Blackburn, pushshift reddit  
, 14, 2020, 830-839. 30 [424] A. Fan, Y. Jernite, E. Perez, D. Grangier, J. Weston, M. Auli, Eli5:  
, arXiv arXiv:1907.09190(2019). 31 [425] Y. Wang, S. Mishra, P. Alipoormolabashi, Y. Kordi, A. Mirzaei, A. Arunkumar, A. Ashok, A. S. Dhanasekaran, A. N. Aik, D. Stap, et al., 1,600+  
, arXiv arXiv:2204.07705 (2022). 31 [426] T. Xie, C. H. Wu, P. Shi, R. Zhong, T. Scholak, M. Yasunaga, C.-S. Wu, M. Zhong, P. Yin, S. I. Wang, Unifiedskg:  
, arXiv arXiv:2201.05966(2022). 31 [427] Q. Ye, B. Y. Lin, X. Ren, Crossfit: NLP  
7}, arXiv arXiv:2104.08835(2021). 31 [428] V. Aribandi, Y. Tay, T. Schuster, J. Rao, H. S. Zheng, S. V. Mehta, H. Zhuang, V. Q. Tran, D. Bahri, J. Ni, Ext5:  
, arXiv arXiv:2111.10952(2021). 31 [429] A. Williams, N. Nangia, S. Bowman,  
, 2018  
:  
, 1 ( ),  
2018, 1112-1122. doi:10.18653/v1/N18-1101. URL https://aclanthology.org/N18-1101 31 [430] Y. Zhang, J. Baldridge, L. He, PAWS:  
, 2019  
:  
, 1 ( ),  
2019, 1298-1308. doi:10.18653/v1/N19-1131. URL https://aclanthology.org/N19-1131 32 [431] C. Qin, A. Zhang, Z. Zhang, J. Chen, M. Yasunaga, D. Yang, GPT7  
7}, 2023  
2023, 2023. URL https://openreview.net/forum?id=u03xn1COsO 32 [432] M. U. Hadi, R. Qureshi, A. Shah, M. Irfan, A. Zafar, M. B. Shaikh, N. Akhtar, J. Wu, S. Mirjalili, et al.,  
:  
, TechRxiv(2023). 32 [433] X. L. Dong, S. Moon, Y. E. Xu, K. Malik, Z. Yu, LLM  
, 2023 29 ACM SIGKDD  
, 5792-5793. 32 [434] K. Pandya, M. Holia, Langchain  
:  
GPT  
, arXiv arXiv:2310.05421(2023). 32 [435] J. Li, B. Hui, G. Qu, B. Li, J. Yang, B. Li, B. Wang, B. Qin, R. Cao, R. Geng,  
LLM  
?

- Xiv : arXiv:2305.03111 (2023). 32 [436] A. Rao, J. Kim, M. Kamineni, M. Pang, W. Lie, M. D. Succi, chatgpt 7}, medRxiv (2023) 2023 – 02. 32 [437] M. Benary, X. D. Wang, M. Schmidt, D. Soll, G. Hilfenhaus, M. Nas-sir, C. Sigler, M. Knödler, U. Keller, D. Beule, et al., JAMA Network Open 6 (11) (2023) e2343689 – e2343689. 32 [438] C. M. Chiesa-Estomba, J. R. Lechien, L. A. Vaira, A. Brunet, G. C. amaroto, M. Mayo-Yanez, A. Sanchez-Barrueco, C. Saga-Gutierrez, chat-gpt (2023) 1 – 6. 32 [439] S. Montagna, S. Ferretti, L. C. Klopfenstein, A. Florio, M. F. Pengo, 7} llm , 2023 ACM , 2023, 205 – 212 . 32 [440] D. Bill, T. Eriksson, LLM (2023). 32 [441] M. Abbasian, I. Azimi, A. M. Rahmani, R. Jain, : LLM , arXiv arXiv:2310.02374(2023). 32 [442] K. V. Lemley, chatgpt7} ?, Journal of the American Society of Nephrology(2023) 10 – 1681. 32 [443] S. Pal, M. Bhattacharya, S.-S. Lee, C. Chakraborty, (llm) chatgpt7} . Annals of Biomedical Engineering(2023) 1 – 4. 32 [444] Y. Du, S. Zhao, Y. Chen, R. Bai, J. Liu, H. Wu, H. Wang, B. Qin, calla : llm , arXiv arXiv:2309.04198(2023). 32 [445] A. Abd-Alrazaq, R. AlSaad, D. Alhuwail, A. Ahmed, P. M. Healy, S. Latifi, S. Aziz, R. Damseh, S. A. Alrazak, J. Sheikh , JMIR 9(1)(2023) e48291. 32 [446] A. B. Mbakwe, I. Lourentzou, L. A. Celi, O. J. Mechanic, A. Dagan, Chatgpt (2023). 32 [447] S. Ahn, 35(1) (2023) 103. 32 [448] E. Waisberg, J. Ong, M. Masalkhi, A. G. Lee, (llm) , Eye (2023) 1 – 3. 32 [449] G. Deiana, M. Dettori, A. Arghittu, A. Azara, G. Gabutti, P. Castiglia, : chatgpt 7}, Vaccines 11(7)(2023) 1217. 32 [450] L. De Angelis, F. Baglivo, G. Arzilli, G. P. Privitera, P. Ferragina, A. E. Tozzi, C. Rizzo, Chatgpt : AI , Frontiers in Public Health 11(2023) 1166120. 32 [451] N. L. Rane, A. Tawde, S. P. Choudhary, J. Rane, chatgpt (llm) 5(10) (2023) 875 – 899. 32 [452] W. Dai, J. Lin, H. Jin, T. Li, Y.-S. Tsai, D. Gašević, G. Chen, ? chatgpt , 2023 IEEE (ICALT), IEEE E, 2023, pp. 323 – 325. 32 [453] E. Kasneci, K. Seßler, S. Küchemann, M. Bannert, D. Dementieva, F. Fischer, U. Gasser, G. Groh, S. Günemann, E. Hüllermeier, et al., Chatgpt 7}? , 103(2023) 102274. 32 [454] N. Rane, chatgpt : (2023 9 15 )(2023). 32 [455] J. C. Young, M. Shishido, openai chatgpt , 14(6)(2023). 32 [456] J. Irons, C. Mason, P. Cooper, S. Sidra, A. Reeson, C. Paris, c hatgpt , SocArXiv(2023). 32 [457] P. G. Schmidt, A. J. Meir, AI : ? , arXiv arXiv:2311.06981(2023). 32 [458] Y. Zheng, H. Y. Koh, J. Ju, A. T. Nguyen, L. T. May, G. I. Webb, S. Pan, , arXiv arXiv:2310.07984(2023). 33 [459] B. Aczel, E.-J. Wagenmakers, chatgpt , PsyArXiv(2023). 33 [460] S. Altmäe, A. Sola-Leyva, A. Salumets, : 7} 7}?, Reproductive BioMedicine Online(2023). 33 [461] S. Imani, L. Du, H. Shrivastava, Mat hprompter: , arXiv arXiv:2303.05398(2023). 33 [462] Z. Yuan, H. Yuan, C. Li, G. Dong, C. Tan, C. Zhou, , arXiv arXiv:2308.01825(2023). 33 [463] K. Yang, A. M. Swope, A. Gu, R. Chalamala, P. Song, S. Yu, S. Godil, R. Prenger, A. Anandkumar, L. eandojo: , arXiv arXiv:2306.15626(2023). 33 [464] K. M. Collins, A. Q. Jiang, S. Frieder, L. Wong, M. Zilka, U. Bhatt, T. Lukasiewicz, Y. Wu, J. B. Tenenbaum, W. Hart, et al., 7}, arXiv arXiv:2306.01694(2023). 33 [465] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, H. He, A. Li, M. He, Z. Liu , chatgpt , Meta-Radiology (2023) 100017. 33 [466] J. Drápal, H. Westermann, J. Savelka, , arXiv arXiv:2310.18729 (2023). 33 [467] J. Savelka, K. D. Ashley, M. A. Gray, H. Westermann, H. Xu, (gpt-4) , arXiv arXiv:2306.09525 (2023). 33 [468] N. Guha, J. Nyarko, D. E. Ho, C. Ré, A. Chilton, A. Narayana, A. Chohlas-Wood, A. Peters, B. Waldon, D. N. Rockmore , Legal-bench: , arXiv arXiv:2308.11462(2023). 33 [469] J. Cui, Z. Li , Y. Yan, B. Chen, L. Yuan, Chatlaw: , arXiv arXiv:2306.16092(2023). 33 [470] H. Yang, X.-Y. Liu, C. D. Wang, Fingpt: , arXiv arXiv:2306.06031(2023). 33 [471] Y. Li, S. Wang, H. Ding, H. Chen, : , 4 ACM AI , 2023, 374-382 . 33 [472] A. Lykov, D. Tsetserukou, Llm-b rain: Ai , arXiv arXiv:2305.19352(2023). 33 [473] E. Billing, J. Rosén, M. Lamb, - , ACM/IEEE - , 2023 3 13-16 , , ACM , 2023, 905-906 . 33 [474] Y. Ye, H. You, J. Du, chatgpt , IEEE Access(2023). 33 [475] Y. Ding, X. Zhang, C. Paxton, S. Zhang, , RSS 2023 , 2023. 33 [476] J. Wu, R. Antonova, A. Kan, M. Lepert, A. Zeng, S. Song, J. Bohg, S. Rusinkiewicz, T. Funkhouser, Tidybot: , arXiv arXiv:2305.05658(2023). 33 [477] E. Strubell, A. Ganesh, A. McCallum, NLP , arXiv arXiv:1906.02243(2019). 34 [478] E. M. Bender, T. Gebru, A. McMillan-Major, S. Shmitchell,

, 2021, 610-623 . 34 [479] C. Zhang, S. Bengio, M. Har  
 dt, B. Recht, O. Vinyals, ( ), 가 , AC  
 M 64(3)(2021) 107-115. 34 [480] M. Tănzer, S. Ruder, M. Rei,  
 , arXiv arXiv:2105.00828(2021). 34 [481] S.  
 M. West, M. Whittaker, K. Crawford, , AI Now(2019) 1 – 33. 34 [482] K. Valmee  
 kam, A. Olmo, S. Sreedharan, S. Kambhampati,  
 ( llms ), arXiv arXiv:2206.10498(2022). 34  
 [483] Y. Zhang, Y. Li, L. Cui, D. Cai, L. Liu, T. Fu, X. Huang, E. Zhao, Y. Zhang, Y. Chen, et  
 al., AI : , arXiv arXiv:  
 2309.01219(2023). 34 [484] A. Webson, E. Pavlick,  
 ? , arXiv arXiv:2109.01247(2021). 34 [485] O. Shaikh, H.  
 Zhang, W. Held, M. Bernstein, D. Yang, !  
 , arXiv arXiv:2212.08061(2022). 34 [486] B. C. Das, M.  
 H. Amini, Y. Wu, : , arXiv  
 arXiv:2402.00888(2024). 34 [487] X. Liu, H. Cheng, P. He, W. Chen, Y. Wang, H. Po  
 on, J. Gao, , ArXiv(2020 4 ). URL https://w  
 ww.microsoft.com/en-us/research/ /adversarial-training-for-large-neural- -models/ 34  
 [488] E. Shayegani, M. A. A. Mamun, Y. Fu, P. Zaree, Y. Dong, N. Abu-Ghazaleh,  
 (2023). arXiv:2310.10844. 34 [489] X. Xu,  
 K. Kong, N. Liu, L. Cui, D. Wang, J. Zhang, M. Kankanhalli, llm :  
 (2023). arXiv: 2310.13345. 34 [490] H. Zhao, H. Chen, F. Yan  
 g, N. Liu, H. Deng, H. Cai, S. Wang, D. Yin, M. Du, 가 :  
 (2023). arXiv:2309.01029. 35 [491] S. Huang, S. Mamidanna, S. Jangam, Y. Zhou, L. H.  
 Gilpin, ? llm  
 (2023). arXiv:2310.11207. 35 [492] H. Brown, K. Lee, F. Mireshghallah, R. Shokri, F.  
 Tramèr, ?2022 ACM  
 , 2022, 2280-2292 .35 [493] R. Plant, V. Giuffrida,  
 D. Gkatzia, : , ar  
 Xiv arXiv:2204.09391(2022).35 [494] W. Niu, Z. Kong, G. Yuan, W. Jiang, J. G  
 uan, C. Ding, P. Zhao, S. Liu, B. Ren, Y. Wang,  
 (2020). arXiv:2009.06823. 35 [495] C. Guo, J. Tang, W. Hu, J. Leng, C. Zhang, F. Yang, Y.  
 Liu, M. Guo, Y. Zhu, Olive: -  
 가 , 2023 50 , 1-15 . 35 [496]  
 B. Meskó, E. J. Topol, ( AI)  
 , npj Digital Medicine 6(1)(2023) 120 . 35 [497] J. Zhang, X. Ji, Z. Zhao, X. Hei, K.-  
 K. R. Choo, :  
 , arXiv arXiv:2308.02678(2023). 35 [498] J. Mökander, J. Schuett, H.  
 R. Kirk, L. Floridi, : 3 , AI (2023) 1 – 31. 35