

ACER: AUTOMATIC LANGUAGE MODEL CONTEXT EXTENSION VIA RETRIEVAL

Anonymous authors

Paper under double-blind review

ABSTRACT

Long-context modeling is one of the critical capabilities of language AI for digesting and reasoning over complex information pieces. In practice, long-context capabilities are typically built into a pre-trained language model (LM) through a carefully designed context extension stage, with the goal of producing generalist long-context capabilities. In our preliminary experiments, however, we discovered that the current open-weight generalist long-context models are still lacking in practical long-context processing tasks. While this means perfectly effective long-context modeling demands task-specific data, the cost can be prohibitive. In this paper, we draw inspiration from how humans process a large body of information: a lossy **retrieval** stage ranks a large set of documents while the reader ends up reading deeply only the top candidates. We build an **automatic** data synthesis pipeline that mimics this process using short-context LMs. The short-context LMs are further tuned using these self-generated data to obtain task-specific long-context capabilities. Similar to how pre-training learns from imperfect data, we hypothesize and further demonstrate that the short-context model can bootstrap over the synthetic data, outperforming not only long-context generalist models but also the retrieval and read pipeline used to synthesize the training data.

1 INTRODUCTION

The field of Artificial Intelligence (AI) and Natural Language Processing (NLP) have made substantial progress in building and teaching neural language models (LMs) to understand and generate language (Radford et al., 2019; Brown et al., 2020; OpenAI, 2023; Anthropic, 2023; 2024; Touvron et al., 2023a;b; MetaAI et al., 2024). Large-scale deep learning has enabled large LMs to learn from massive amounts of human-generated text (Radford et al., 2019; Brown et al., 2020). However, new challenges emerge as researchers consider building capabilities beyond those of humans. One popular example, also the focus of this paper, is understanding long-contexts of text (Dai et al., 2019; Su et al., 2024; Xiong et al., 2023). Despite the advancements in modeling (Dao et al., 2022; Su et al., 2024; Xiong et al., 2023), data keeps being lacking simply because humans no longer produce them naturally. Specifically, while longer contexts give more degrees of freedom in forming possible language sequences and long-interaction/complex-information tasks, the existing long texts are limited to organized and compact ones like novels and code (Fu et al., 2024).

A common methodology in building long-context LM is to incorporate a context extension phase into the model building cycle after the general pre-training phase and before the task-aware post-training phase (Xiong et al., 2023; Rozière et al., 2024; MetaAI et al., 2024). As the standard pre-training stage builds into the model the general capabilities over diverse text distributions, the long-context extension phase is designed with the hope to extend generalist capabilities further to long-context patterns. The subsequent post-training is supposed to activate and align these extra long-context capabilities to task instructions.

Despite the exciting promise of this context extension scheme, the limited portion of model training dedicated to long-context contradicts the aforementioned increasingly more complex space language and task patterns when the text gets longer (Xiong et al., 2023; MetaAI et al., 2024). In other words, using a context extension phase to cover all long-context understanding tasks can be mathematically intractable, and consequently, building a long-context generalist may not succeed. To investigate this, in this paper, we conduct experiments and demonstrate empirically that practical tasks like

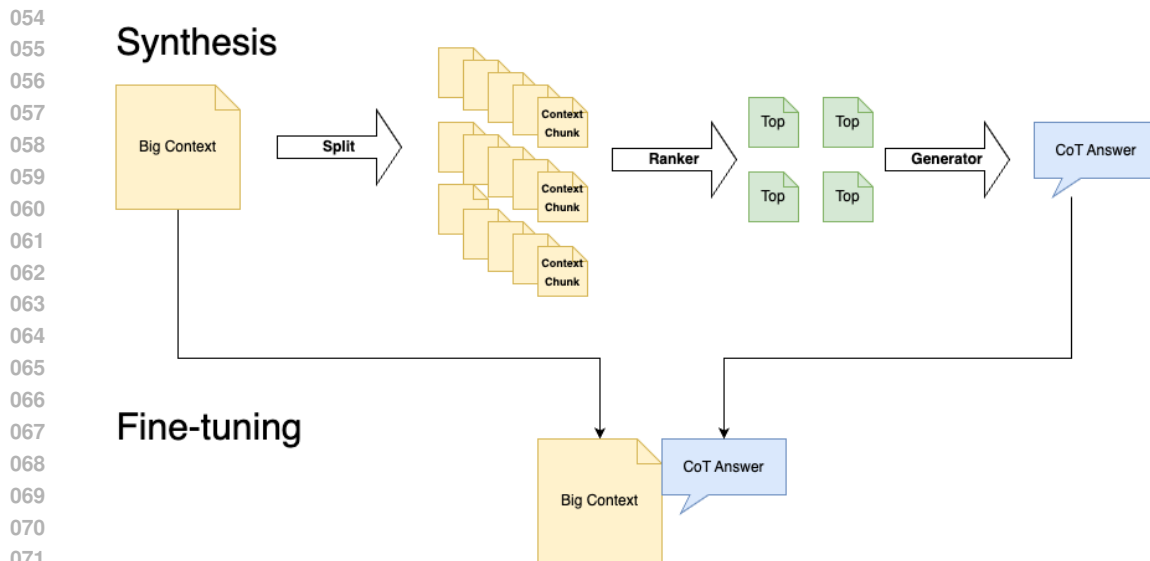


Figure 1: The full process of ACER involves a data synthesis stage and a fine-tuning stage. (top) The data synthesis stage splits and retrieves a set of relevant text chunks for a problem and use a short-context model to generate an answer with CoT (Wei et al., 2022). (bottom) The fine-tuning stage use the original long-context data and the synthetic CoT answer to fine-tune a long-context model.

long-context retrieval augmented generation can easily break existing long-context models. While creating domain-specific supervised training may help fix this. Unlike short-context training data that can be collected from everyday people (Kopf et al., 2023), long-context training can be much harder to curate. Reading long pieces of text is inherently hard for humans, which could make the annotation process not only costly but also very demanding and, therefore, less reliable.

In this paper, to alleviate this problem and provide an intermediate solution, we propose a new approach which Automates Context Extension via Retrieval (ACER; Figure 1). Overall, ACER is a two-stage method. We start with synthesizing *imperfect* data by combining retrieval with an LM *excels in short context*. In the subsequent stage, we will fine-tune a large LM to bootstrap over this data. ACER will start with some pairs of question and its long context, while no labeling is required. In the synthesis pipeline, the long context is broken into chunks and a retrieval model will score and rank the chunks. A small set of top-ranked chunks will be fed into the short-context LM to produce answer *with chain-of-thought* (CoT; Wei et al. (2022)) reasoning. Then, in the fine-tuning stage, we train the model using the context, question, and the CoT. On the other hand, the retrieval-based data synthesis process will be hidden from the model. We desire that the deep model will learn a generic long-context understanding function by fitting over the full context and the CoT reasoning. We hypothesize that the model may discover a better latent function that transcends the original ranking mechanism used in the first stage. This also shares some spirit with LM pre-training. Whereas typically pre-training bootstrap from existing human data, due to the apparent data scarcity, ACER will bootstrap from synthetic long-context data generated using retrieval.

Our experiments demonstrated the effectiveness of ACER. We found model trained with ACER without supervision can outperform contemporary generalist long-context models. It also outperforms its own retrieval-based answering pipeline if applied to the test sets.

2 ACER

In this section, we will give an overview of ACER. The ACER method consists of two major stages: 1) automatic data synthesis, and 2) self training, as illustrated in Figure 1. In this chapter, we will describe how each of these stages work.

2.1 AUTOMATIC DATA SYNTHESIS

Our data synthesis process combines a heuristic-based retrieval pipeline and a short-context generator. We use the following ingredients,

- **Prompts:** a set of prompts/problems, $\{p_1, p_2, \dots, p_N\}$, each of which consists of a pair of context and question $p_i = (c_i, q_i)$
- **Short-Context Ranker:** A ranker model $r(q, t) \rightarrow \mathbb{R}$ which takes a question q and a piece of *short* text t . This ranker determines a relevance score corresponding to how helpful the text t in answering the question q .
- **Short-Context Generator:** A generator model $g(x) \rightarrow a$ that takes some short prompt x and returns an answer a . This will be an aligned instruction-tuned model like many existing contemporary LM.

We start the data synthesis process from the set of prompts. For one prompt $p_i = (c_i, q_i)$, we break the context into a set of text chunks $\{t_{i1}, t_{i2}, t_{i3}, \dots\}$. The ranking model will assign chunk t_{ij} a relevance score $s_{ij} = r(q_i, t_{ij})$. Based on these scores, we can produce a ranking of the text chunk indices using the estimated helpfulness, $[r_{i1}, r_{i2}, r_{i3}, \dots]$. We collect the top M ranked text chunks while making sure that their concatenation can still fit in the short generator. These “most helpful” chunks will together be fed into a generator to produce an answer,

$$\hat{a}_i = g(\text{INST} \circ q_i \circ t_{r_{i1}} \circ t_{r_{i2}} \circ \dots \circ t_{r_{iM}}) \quad (1)$$

Here `INST` denotes an instruction to the model to elicit an explicit reasoning in a chain-of-thought like form. This makes sure that the model can describe how it compares the information in the provided text pieces to identify the most useful ones, as well as how it combines them to arrive at the final answer. This reasoning process will be used in the next stage to help train the model in the fine-tuning stage.

This data synthesis process leverages two critical capabilities in the original LM, relevance/usefulness analysis and understanding of short pieces of text. We combine them heuristically to build a surrogate long-context pipeline. Readers familiar with the concept of map-reduce may recognize our data synthesis as such a process.

2.2 FINE-TUNING

During data synthesis, we produce for each prompt, a full document ranking, a small set of helpful documents and an answer. For fine-tuning, we discard the ranking and the document set and use only the generated answer, because we do not want our model to be exposed to and learn from the *lossy* retrieval pipeline. We fine-tune an LM f_θ to produce the CoT answer \hat{a}_i , i.e.,

$$\hat{\theta} = \operatorname{argmin}_\theta \sum_i f(\hat{a}_i | (c_i, q_i)) \quad (2)$$

In this setup, we provide the model a lossless, unfiltered access to the full context. In addition, we pair it with an answer with extra training signals, a CoT describing an information extraction process of picking up useful pieces of information and sorting through. As with any other deep learning application, we desire that the large over-parametrized LM can learn during the optimization process to fit a long-context understanding function that will lead to similar reasoning and answer. To keep it simple, we fine-tune the model with teacher forcing using a log likelihood loss. Only the answer tokens participate in loss computation with the context token loss masked out during training.

3 EXPERIMENTAL SETUP

3.1 TASKS

We consider two realistic tasks, long-context retrieval augmented generation (Jiang et al. (2024)) and long-context reading comprehension. We pick long-context RAG as our major evaluation since the task, while being very useful, has only very recently be carefully studied. This is more

aligned with our desired setup, distinct yet useful long context tasks with little supervision. We consider the following two datasets: **Natural Question (NQ)** (Kwiatkowski et al., 2019) and **TriviaQA (TQA)** (Joshi et al., 2017) We use Wikipedia as the knowledge source, and BGE (Xiao et al., 2024b) which is a dense retriever Karpukhin et al. (2020). We report the exact match (EM) metric.

For long-context reading comprehension, we used the **NarrativeQA** dataset Kočiský et al. (2018). The dataset has been out for several years at the time of writing and has since been a standard evaluation for long context. We expect many long-context models being trained on some form of its train set. Nevertheless, we still use it as reference to understand how our self-supervised approach compares to the supervised ones. We used the curated test set from LongBench Bai et al. (2023). We report the token-level F1 metric.

While these datasets have short gold answer, modern LMs tend to produce long answers. In order to evaluate EM or F1 scores, we employ an additional short answer extraction step by few-shot prompting a Llama-3-8B-Instruct model using the prompt introduced by Jiang et al. (2024).

3.2 DATA SYNTHESIS

For data synthesis, in order to keep it simple, we implement the ranker and the short context generator by prompting the same LM, Llama-3-8B-Instruct. This is the short-context version of the latest generation of Llama models (MetaAI et al., 2024). We use the 8B variant instead of the 70B or 405B to demonstrate the applicability of ACER in a cost efficient setup.

Ranker Model We show the prompt of the ranker model in Figure 2. We instruct the model to read the question and context chunk (referred to as passage in the prompt) and *think step-by-step* to decide the helpfulness of the passage for answering the question. We formulate it as a multiple choice question where the model is instructed to choose from 5 options from the most a) “providing exact answer” to e) “not related”.

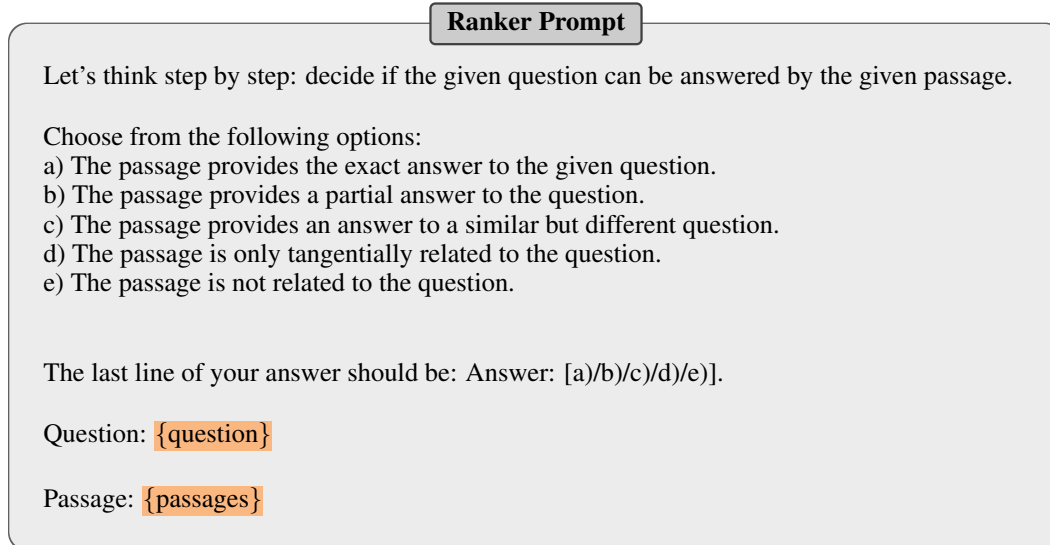


Figure 2: Prompts given to the LM to produce relevance judgement.

Generator Model We show the prompt for the generator model in Figure 3. The model is instructed to read the given passages and extract *several* pieces of potentially helpful information before it makes decision on what evidence to use and then use it to answer the question.

We use prompts in the corresponding dataset’s training set. For RAG tasks, we obtain the context by using a dense retriever to retrieve 100 passages from Wikipedia. We use the DPR version of Wikipedia dump, where the documents are splitted into chunks of 100 words. For NarrativeQA, the context is simply the original full book. We do not make any additional edit over it. We apply

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

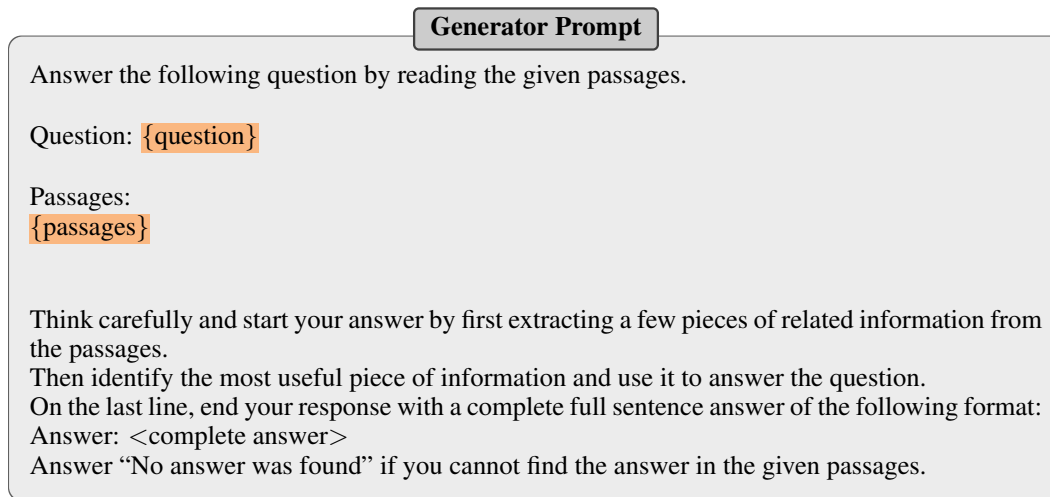


Figure 3: Prompts given to the LM to produce the final CoT answer.

an addition data augmentation technique to the RAG data: we perform a sliding window shuffle of candidates passages. We use window size of 30 and a stride of 20. This is used to efficiently emulate a ranking noise where the absolute change in position decays with the change size. One copy of the noised data is used as well as one with original ranking.

For the creation of the final data, we reuse the generator prompt template (Figure 3) with the full context as question. We combine it with the generated answer using a chat template.

3.3 FINE-TUNING

We directly fine-tune the ACER models over long context without approximation. This section describes the training details as well as contemporary training technologies/techniques we adopted.

3.3.1 IMPLEMENTATION

We train over the synthesized data for 1 epoch with a batch size of 128 examples. We use a max sequence length of 32k for RAG and 64k for reading comprehension, and therefore this amounts up to 4M tokens and 8M tokens batch, respectively. These lengths are picked to cover most of the test lengths while at the same time keeping training on our accelerators viable and efficient. At the test time, RoPE extrapolation (Su et al., 2024) is utilized for sequences of longer length. We use a $3e - 6$ learning rate with Adam optimizer (Kingma & Ba, 2015). For model initialization, we still require a model that has a sufficiently long context length. This can be achieved unsupervisedly by training on long-context corpora. To save cost, we borrow Llama-3-8B-ProLong-Base model (Gao et al., 2024)¹, an open model unsupervisedly context-extended from our data synthesizer model Llama-3-8B-Instruct (The authors made decision to build a base model from an instruction tuned version.)

3.3.2 INFRASTRUCTURE

We train our models using JAX (Bradbury et al., 2018) on cloud TPUv4s (Jouppi et al., 2023) provided by TPU Research Cloud. These 4th generation TPUs are built with relatively small-size 32GB HBM per chip while interconnected with high bandwidth fibers. We therefore opt to partition computation and perform our training with 3D parallelism of data parallelism, sequence parallelism, and tensor parallelism. We train on 32 TPUv4 chips in a (2, 4, 4) configuration mesh and map these axes to data, sequence, and tensor parallelism axes respectively. Optimizer states as well as the full-precision copy of model weights are kept fully sharded over the mesh, over the entire course of training. We manually shard model, optimizer, and critical activation in attention and feed-forward

¹princeton-nlp/Llama-3-8B-ProLong-512k-Base

network. We use JAX’s `shard_map` to shard flash attention (Dao et al., 2022; Dao, 2024) TPU kernels² across model axis by heads. The rest of the shardings are decided by the XLA SPMD compiler (Xu et al., 2021).

We compute loss only over the answer, with the long context receiving only implicit gradients. Taking advantage of this to save memory and flops, on top of using a loss mask, we (dynamically) slice the last layer of hidden states at each batch instance into a smaller answer tensor. Only this smaller tensor is passed to the LM head for out-projection and loss computation. In the backward pass, this will translate into a (dynamic) update of the gradient tensor.

3.4 COMPARED MODELS

For comparison, we consider models that are relatively open and closely related to the base model we use for ACER, Llama-3.1-Instruct. We acknowledge the effectiveness of proprietary models such as GPT4 OpenAI (2023). We do not include them here, because our focus is to evaluate a method of data synthesis and the models and data they use are not directly comparable.

Specifically, we compare ACER with long-context LMs fine-tuned on open data as well as those fine-tuned on closed data. Specifically, we consider the following,

- **Together-Llama-2-7B-32K-Instruct (Together, 2023):** A model by TogetherAI by extending Llama2 to a 32k context size. It is fine-tuned on an open mixture of long-context question answering and summarization data.
- **Llama-3-8B-ProLong-512k-Instruct (Gao et al., 2024):** This is a Llama3-8B-Instruct context extended on open long context dataset. It is then fine-tuned on UltraChat (Ding et al., 2023), a large fine-tuning dataset generated by GPT-4. The creator found it most helpful amongst open instruction-tuning datasets. This model closely relates to ours and share the same base model, offering a straightforward comparison.
- **Llama3-8B-Instruct (Truncation) (MetaAI et al., 2024):** Llama3-8B-Instruct is a 8k context model pre-trained and instruction-tuned by MetaAI with closed data. In this setup, we truncate the input to fit it into the context.
- **Llama3-8B-Instruct (RAG) (MetaAI et al., 2024):** This is Llama3-8B-Instruct reading a small set of input chunks generated with the same retrieval pipeline used in the ACER data synthesis.
- **Llama3.1-8B-Instruct (MetaAI et al., 2024):** This is the second iteration of the Llama3 model by MetaAI. Its context is systematically extended by Meta GenAI team.
- **Mistral-Nemo-Instruct-2407 (MistralAI, 2024):** A larger 12B model trained by MistralAI and Nvidia using closed data. It has a native 128k context length.

We perform inference with vLLM (Kwon et al., 2023) with greedy decoding. Models always read the full context and extrapolate when reading longer context than the original training-time length except in Llama3-8B-Instruct (Truncation).

4 EXPERIMENTAL RESULTS

In Table 1, we show the performance of the compared systems as well as ACER on the evaluation datasets. We see a general trend that ACER outperforms the compared systems with decent margins especially on the novel long-context RAG tasks. We observe a general trend that the closed data model performing better than models trained on open data. Specifically, the Together-Llama-2 model, which is based on the previous generation Llama model, significantly underperforms all other model. This is likely due to two facts that the model essentially stems from an earlier generation, and its shorter context requires it to do more extrapolation. Now, to remind our reader, the rest of the 8B models, all came from the same base model Llama-3-8B-Base. Despite this fact, we see vastly different performance. The ProLong model performs decently on NarrativeQA but falls behind on the RAG tasks. This is not surprising as the UltraChat supervision it uses relates more

²https://github.com/jax-ml/jax/blob/main/jax/experimental/pallas/ops/tpu/flash_attention.py

Model	NQ (EM)	TQA (EM)	NarrativeQA (F1)
<i>Supervised Open Data</i>			
Together-Llama-2-7B-32K-Instruct	0.172	0.299	0.013
Llama-3-8B-ProLong-512k-Instruct	0.260	0.457	0.150
<i>Supervised Closed Data</i>			
Llama3-8B-Instruct (Truncation)	0.388	0.567	0.076
Llama3-8B-Instruct (RAG)	0.408	0.611	0.161
Llama3.1-8B-Instruct	0.312	0.518	0.217
Mistral-Nemo-Instruct-2407 ^{12B}	0.365	0.534	0.072
<i>Self-Supervised</i>			
ACER (ours)	0.446	0.648	0.189

Table 1: Performance comparison of ACER and baselines on Natural Question (NQ), TriviaQA (TQA), and NarrativeQA. The best performance on each dataset is boldfaced.

closely to standard-form question answering but can be very different from tasks in the wild like long-context RAG here. On the other hand, we see with the more involved long-context extension done by MetaAI, the Llama-3.1 model significantly outperforms Prolong. It actually also attains the best performance on NarrativeQA, likely because the data it uses may contain long-context QA data. Our ACER model is the second best on NarrativeQA, still out-performing all other systems, suggesting that our self-supervised method is still competitive when no specific supervised data is present. Nevertheless, ACER is the best-performing on the RAG datasets, again confirming the usefulness of our self-supervised method.

We do also note that the 12B Mistral-Nemo does not always show decisive advantage over the 8B models. This demonstrates further that the model size and capability advantage do not always translate into long-context performance. As we discussed before, long context means more task diversity and a model’s task fitness becomes as critical; here we see the Mistral model performs decently on RAG but is lacking on Narrative QA. This again shows the usefulness of ACER which is able to self-supervisedly teach a model a new task.

When compared with Llama3-8B-Instruct (RAG), which uses Llama-3 with the ACER retrieval pipeline in a test-time ad-hoc manner, we still see a decent performance boost in ACER. This suggests that learning and bootstrapping from ACER generated data can transcend the original data synthesis pipeline, as we desired.

5 ANALYSIS

5.1 COMPARING LLAMA-3.1 AND ACER AT DIFFERENT CONTEXT SIZES

To get a more fine-grained understanding of how context extension affects each model, in this section, we compare models at a variety of context lengths. Specifically, we varies the number of retrieved passages fed into the model for generating the final answer. Recall that Llama-3.1 and ACER both derive from the same pre-trained Llama-3 base model using different context extension processes. In Figure 4, we plot the two models’ performance against number of passages on both Natural Question and TriviaQA (1k subset to reduce inference cost.)

We observed a very interesting yet intuitive phenomenon. Here the two models perform very similarly when reading a small context. This means they possesses similar capability and alignment behavior at short context, as *sibling models*. However, as context increases, the performance keeps diverge. While ACER can keep digesting more useful information from the context, Llama-3.1 seemingly suffers from the longer contexts with performance decaying. This shows one other example that ACER achieved the goal we set for it and successfully **extended** the model’s capability presented in a shorter context onto a much longer context.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

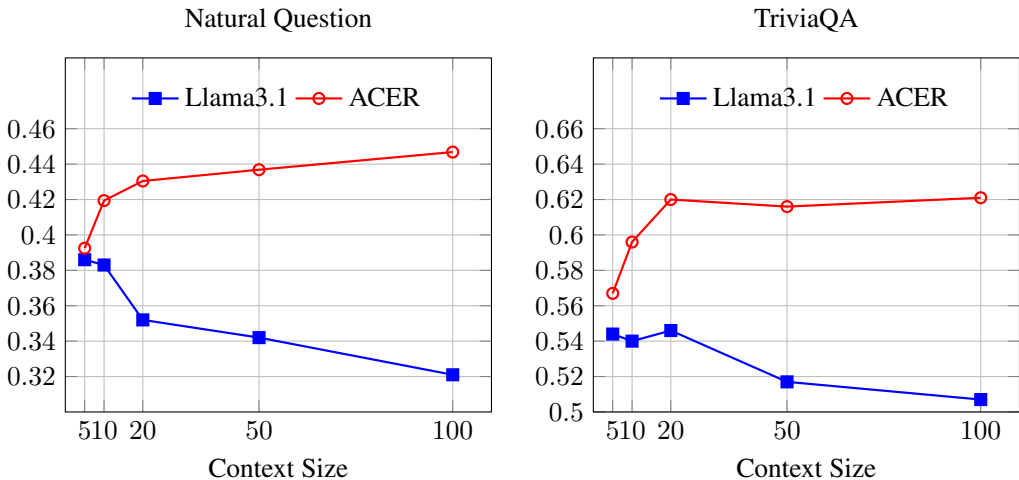


Figure 4: Performance comparison between Llama3.1 and ACER when reading different context sizes on the Natural Question and Trivia QA datasets.

5.2 USING AN UNSUPERVISED RETRIEVER

Retriever	Model	NQ (EM)	TQA (EM)
BM25	Llama3 (Truncation)	0.260	0.542
	Llama3 (RAG)	0.354	0.576
	Llama3.1	0.305	0.507
	ACER	0.383	0.632
Dense	Llama3 (Truncation)	0.388	0.567
	Llama3 (RAG)	0.408	0.611
	Llama3.1	0.312	0.518
	ACER	0.446	0.648

Table 2: Performance of ACER and baselines with different base retrievers.

The previously discussed ACER pipeline for RAG tasks adopts a supervisedly trained dense retriever. In certain situations where even retrieval data is scarce, obtaining such a retriever may not even be possible. People instead need to fall back to the classical BM25 retrievers. In this section, we consider such a situation by running ACER and some of the baseline systems with BM25. This means the top-100 candidates will be different and of lower quality. In Table 2, we show the results of this BM25 setup compared with the original dense setup. While using BM25, all systems take a hit in performance because of lower-quality candidates, the general performance order of the systems remain the same. Our ACER method still outperforms the other systems by decent margins. In fact, when looking at the absolute numbers, we can observe that ACER with BM25 attains close or even better performance than best-performing baseline systems. The results suggest that ACER is relatively agnostic to the retriever used and attains good performance to warm-start a system without requiring extra supervision.

6 RELATED WORKS

Language Modeling Recent advancements of large language models have shown strong language understanding ability and ace a wide range of natural language processing tasks. Based on the Transformer structure (Vaswani et al., 2017), LMs can be broadly categorized as decoder-only models (e.g., GPT (Radford et al., 2018; 2019)), encoder-only models (e.g., BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019)), and encoder-decoder models (e.g., BART Lewis et al. (2020a)

and T5 Raffel et al. (2020)). Most recently, the large language models such as the GPT family (Brown et al., 2020; OpenAI, 2023), Gemini (GeminiTeam et al., 2024a;b), Llama Touvron et al. (2023a;b); MetaAI et al. (2024), and Claude Anthropic (2023; 2024) have attract great attention by achieving human-level performance on various tasks and showing structure-following ability. These models are typically trained in several stages, including pre-training on web-scale text data with auto-regressive language modeling, supervised fine-tuning on specific applications, and human preference alignment, where the last step plays the pivotal role to steer the LMs as a dialogue system to answer human’s instruction and generate responses that are in desired quality and style, safe, and ethical. Some recent LM alignment methods include reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022), proximal policy optimization (PPO) (Schulman et al., 2017), direct policy optimization (DPO) (Rafailov et al., 2023), and Kahneman-Tversky optimization (KTO) (Ethayarajh et al., 2024).

Long Context Modeling Transformer (Vaswani et al., 2017) is the foundation architecture for recent advancements in language models. However, its quadratic temporal and computation complexity to the sequence length poses great challenge to scale to long input sequence, and the lack of robust position embeddings also degrades the performance of long context understanding. Therefore, tremendous efforts are made to improve the long-context modeling of language models from different aspects, including more efficient attention mechanism for long context Beltagy et al. (2020); Kwon et al. (2023); Dao et al. (2022); Dao (2024); Xiao et al. (2024a), a (recurrent) internal or external memory bank (Dai et al., 2019; Wu et al., 2022a;b), and length extrapolation through positional encoding (Press et al., 2022; Su et al., 2024; Peng et al., 2024; Zhu et al., 2024). To evaluate the long-context modeling ability of LLMs, several synthetic benchmarks are proposed, such as Lost in the Middle (Liu et al., 2024), Needle in a Haystack (Kamradt, 2023), and Ruler (Hsieh et al., 2024). Because we are proposing new data synthesis method for long-context LM training, we do not include the synthetic benchmarks.

Retrieval Augmented Generation Retrieval-augmented generation, or RAG, enhances LLMs’ ability in knowledge intensive tasks. Originated from open-domain question answering (Chen et al., 2017; Xiong et al., 2021b), RAG will enrich the prompt (i.e., a question) with additional relevant context retrieved from an external corpus to help a reader better answer the question (Lewis et al., 2020b; Guu et al., 2020; Izacard et al., 2024). The retrieval methods range from sparse features (Robertson & Zaragoza, 2009; Roberts et al., 2020) and deep-learning based dense representations (Karpukhin et al., 2020; Lewis et al., 2020b) to direct generation by LLMs (Sun et al., 2023; Yu et al., 2023). Recent studies have explored various ways to enhance RAG, such as better query understanding (Kim et al., 2023; Chan et al., 2024), a better retriever (Karpukhin et al., 2020; Xiong et al., 2021a; Yao et al., 2023), and a better reading model (Izacard & Grave, 2021; Cheng et al., 2021; Borgeaud et al., 2021). Specifically, Self-RAG (Asai et al., 2024) trains a critique model that reflects the retrieved passages and the generation to adaptively decide whether the retrieval is necessary. SuRe (Kim et al., 2024) proposes summarized retrieval which generate multiple summaries of retrieved context based on answer candidates. RAPTOR Sarthi et al. (2024) builds a hierarchical structure by iteratively cluster text chunks and generate summaries for the clusters, which can aggregate text chunks to better answer summarizing questions.

7 CONCLUSION

In this paper, we propose ACER, a method for automatically extending language model capabilities to longer contexts without using human generated supervision. ACER pivots through a retrieval pipeline to help a short-context model to heuristically process long pieces of text and produce a synthetic *imperfect* answer. We demonstrate that this model can then bootstrap over this synthetic answer to gain even stronger long context capabilities, often outperforms carefully built long-context generalist models. We believe ACER demonstrate an effective unsupervised approach to extend short context capabilities into longer contexts. It can be a useful tool to improve specific long-context task performance where little training data exists. Users only need to write a few prompts to run the pipeline. Broadly, we introduce automating capabilities extension onto longer context as a new research problem; future works may consider better data synthesis processes and paths to produce better extension results.

REFERENCES

- Anthropic. Introducing 100K Context Windows, 2023. URL <https://www.anthropic.com/index/100k-context-windows>.
- Anthropic. Introducing the next generation of claude, 2024. URL <https://www.anthropic.com/news/claude-3-family>.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-RAG: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=hSyW5go0v8>.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. Longbench: A bilingual, multitask benchmark for long context understanding, 2023.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv:2004.05150*, 2020.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, T. W. Hennigan, Saffron Huang, Lorenzo Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and L. Sifre. Improving language models by retrieving from trillions of tokens. In *International Conference on Machine Learning*, 2021. URL <https://api.semanticscholar.org/CorpusID:244954723>.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/jax-ml/jax>.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020.
- Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. RQ-RAG: Learning to refine queries for retrieval augmented generation. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=tzE7VqsaJ4>.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer open-domain questions. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1870–1879, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1171. URL <https://aclanthology.org/P17-1171>.
- Hao Cheng, Yelong Shen, Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Unit-QA: A hybrid approach for open domain question answering. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 3080–3090, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.240. URL <https://aclanthology.org/2021.acl-long.240>.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2978–2988, Florence, Italy, July 2019. Association for

- 540 Computational Linguistics. doi: 10.18653/v1/P19-1285. URL <https://aclanthology.org/P19-1285>.
- 541
- 542
- 543 Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *International Conference on Learning Representations (ICLR)*, 2024.
- 544
- 545 Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 546
- 547
- 548 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, pp. 4171–4186, 2019.
- 549
- 550
- 551 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*, 2023.
- 552
- 553
- 554 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Model alignment as prospect theoretic optimization. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 12634–12651. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/ethayarajh24a.html>.
- 555
- 556
- 557
- 558
- 559
- 560 Yao Fu, Rameswar Panda, Xinyao Niu, Xiang Yue, Hanna Hajishirzi, Yoon Kim, and Hao Peng. Data engineering for scaling language models to 128k context. *ArXiv*, abs/2402.10171, 2024.
- 561
- 562
- 563
- 564
- 565
- 566
- 567 Tianyu Gao, Alexander Wettig, Howard Yen, and Danqi Chen. Prolong long-context language model series, 2024. URL <https://huggingface.co/princeton-nlp/Llama-3-8B-ProLong-64k-Instruct>.
- 568
- 569
- 570
- 571
- 572
- 573
- 574
- 575
- 576
- 577
- 578
- 579
- 580
- 581
- 582
- 583
- 584
- 585
- 586
- 587
- 588
- 589
- 590
- 591
- 592
- 593
- GeminiTeam, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillcrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan, Manaal Faruqui, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah, Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan, Jeremiah Liu, Andras Orban, Fabian Gura, Hao Zhou, Xinying Song, Aurelien Boffy, Harish Ganapathy, Steven Zheng, HyunJeong Choe, goston Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Merey, Martin Baeuml, Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anas White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rustemi, Natalie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette, Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh Tomar, Evan Senter, Martin Chadwick, Ilya Kornakov, Nithya Attaluri, Iñaki Iturrate, Ruibo Liu, Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Xavier Garcia, Thanumalayan Sankaranarayanan Pillai, Jacob Devlin, Michael Laskin, Diego de Las Casas, Dasha Valter, Connie Tao, Lorenzo Blanco, Adrià Puigdomènech Badia, David Reitter, Mianna Chen, Jenny Brennan, Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski, Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yiming Gu, Kate Olszewska, Ravi Addanki, Antoine Miech, Annie Louis, Denis Teplyashin, Geoff Brown, Elliot Catt, Jan Balaguer, Jackie Xiang, Pidong Wang, Zoe Ashwood, Anton Briukhov, Albert Webson, Sanjay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-Wei Chang, Axel Stjerngren, Josip Djolonga, Yuting

594 Sun, Ankur Bapna, Matthew Aitchison, Pedram Pejman, Henryk Michalewski, Tianhe Yu, Cindy
595 Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter Humphreys, Thibault
596 Sellam, James Bradbury, Varun Godbole, Sina Samangooei, Bogdan Damoc, Alex Kaskasoli,
597 Sébastien M. R. Arnold, Vijay Vasudevan, Shubham Agrawal, Jason Riesa, Dmitry Lepikhin,
598 Richard Tanburn, Srivatsan Srinivasan, Hyeontaek Lim, Sarah Hodgkinson, Pranav Shyam, Johan
599 Ferret, Steven Hand, Ankush Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Giang, Alexander
600 Neitz, Zaheer Abbas, Sarah York, Machel Reid, Elizabeth Cole, Aakanksha Chowdhery, Dipan-
601 jan Das, Dominika Rogozińska, Vitaliy Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas Zilka,
602 Flavien Prost, Luheng He, Marianne Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan, Dawei
603 Jia, Miltiadis Allamanis, Clara Huiyi Hu, Raoul de Liedekerke, Justin Gilmer, Carl Saroufim,
604 Shruti Rijhwani, Shaobo Hou, Disha Shrivastava, Anirudh Baddepudi, Alex Goldin, Adnan
605 Ozturel, Albin Cassirer, Yunhan Xu, Daniel Sohn, Devendra Sachan, Reinald Kim Amplayo,
606 Craig Swanson, Dessie Petrova, Shashi Narayan, Arthur Guez, Siddhartha Brahma, Jessica Lan-
607 don, Miteyan Patel, Ruizhe Zhao, Kevin Vilella, Luyu Wang, Wenhao Jia, Matthew Rahtz, Mai
608 Giménez, Legg Yeung, James Keeling, Petko Georgiev, Diana Mincu, Boxi Wu, Salem Haykal,
609 Rachel Saputro, Kiran Vodrahalli, James Qin, Zeynep Cankara, Abhanshu Sharma, Nick Fer-
610 nando, Will Hawkins, Behnam Neyshabur, Solomon Kim, Adrian Hutter, Priyanka Agrawal, Alex
611 Castro-Ros, George van den Driessche, Tao Wang, Fan Yang, Shuo yiin Chang, Paul Komarek,
612 Ross McIlroy, Mario Lučić, Guodong Zhang, Wael Farhan, Michael Sharman, Paul Natsev, Paul
613 Michel, Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shakeri, Christina Butterfield, Justin
614 Chung, Paul Kishan Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar Soparkar, Karel Lenc,
615 Timothy Chung, Aedan Pope, Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo Wang, Joshua
616 Maynez, Mary Phuong, Taylor Tobin, Andrea Tacchetti, Maja Trebacz, Kevin Robinson, Yash
617 Katariya, Sebastian Riedel, Paige Bailey, Kefan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose
618 Slone, Neil Houlsby, Xuehan Xiong, Zhen Yang, Elena Gribovskaia, Jonas Adler, Mateo Wirth,
619 Lisa Lee, Music Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay
620 Ghemawat, Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay Bolina, Mariko Iinuma, Polina
621 Zablotskaia, James Besley, Da-Woon Chung, Timothy Dozat, Ramona Comanescu, Xiance Si,
622 Jeremy Greer, Guolong Su, Martin Polacek, Raphaël Lopez Kaufman, Simon Tokumine, Hexi-
623 ang Hu, Elena Buchatskaya, Yingjie Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad Toma-
624 sevic, Jinwei Xing, Christina Greer, Helen Miller, Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada
625 Ma, Angelos Filos, Milos Besta, Rory Blevins, Ted Klimenko, Chih-Kuan Yeh, Soravit Chang-
626 pinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir, Vered Cohen, Charline Le Lan,
627 Krishna Haridasan, Amit Marathe, Steven Hansen, Sholto Douglas, Rajkumar Samuel, Mingqiu
628 Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo, Lars Lowe
629 Sjösund, Sébastien Cevey, Zach Gleicher, Thi Avrahami, Anudhyan Boral, Hansa Srinivasan,
630 Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard Hussenot, Livio Baldini Soares, Kate
631 Baumli, Michael B. Chang, Adrià Recasens, Ben Caine, Alexander Pritzel, Filip Pavetic, Fabio
632 Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan Horgan, Kartikeya Badola, Nora Kass-
633 ner, Subhrajit Roy, Ethan Dyer, Víctor Campos Campos, Alex Tomala, Yunhao Tang, Dalia El
634 Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal, Sharad Vikram, Zhitao
635 Gong, Sergi Caelles, Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng, Wojciech Stokowiec,
636 Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh, James Svensson,
637 Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsihlias, Arpi Vezer, Marco
638 Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran Rong, Allan
639 Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks, Marie Pel-
640 lat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi Hashemi,
641 Richard Ives, Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze Wang,
642 Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Moufarek, Samer Has-
643 san, Kaushik Shivakumar, Joost van Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal,
644 Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević,
645 Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot,
646 Matthew Lamm, Nicola De Cao, Charlie Chen, Sidharth Mudgal, Romina Stella, Kevin Brooks,
647 Gautam Vasudevan, Chenxi Liu, Mainak Chain, Nivedita Melinker, Aaron Cohen, Venus Wang,
648 Kristie Seymore, Sergey Zubkov, Rahul Goel, Summer Yue, Sai Krishnakumaran, Brian Albert,
649 Nate Hurley, Motoki Sano, Anhad Mohanane, Jonah Joughin, Egor Filonov, Tomasz Kepa,
650 Yomna Eldawy, Jiawern Lim, Rahul Rishi, Shirin Badiezadegan, Taylor Bos, Jerry Chang, Sanil
651 Jain, Sri Gayatri Sundara Padmanabhan, Subha Puttagunta, Kalpesh Krishna, Leslie Baker, Nor-
652 bert Kalb, Vamsi Bedapudi, Adam Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin, Xiang Zhou,

648 Zhichun Wu, Sam Sobell, Andrea Siciliano, Alan Papir, Robby Neale, Jonas Bragagnolo, Tej
649 Toor, Tina Chen, Valentin Anklin, Feiran Wang, Richie Feng, Milad Gholami, Kevin Ling, Lijuan
650 Liu, Jules Walter, Hamid Moghaddam, Arun Kishore, Jakub Adamek, Tyler Mercado, Jonathan
651 Mallinson, Siddhinita Wandekar, Stephen Cagle, Eran Ofek, Guillermo Garrido, Clemens Lom-
652 briser, Maksim Mukha, Botu Sun, Hafeezul Rahman Mohammad, Josip Matak, Yadi Qian, Vikas
653 Peswani, Pawel Janus, Quan Yuan, Leif Schelin, Oana David, Ankur Garg, Yifan He, Oleksii
654 Duzhyi, Anton Algmyr, Timothée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex Chinien, Rakesh
655 Shivanna, Aleksandr Chuklin, Josie Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed, Sub-
656 habrata Das, Zihang Dai, Kyle He, Daniel von Dincklage, Shyam Upadhyay, Akanksha Mau-
657 rya, Luyan Chi, Sebastian Krause, Khalid Salama, Pam G Rabinovitch, Pavan Kumar Reddy M,
658 Aarush Selvan, Mikhail Dektiarev, Golnaz Ghiasi, Erdem Guven, Himanshu Gupta, Boyi Liu,
659 Deepak Sharma, Idan Heimlich Shtacher, Shachi Paul, Oscar Akerlund, François-Xavier Aubet,
660 Terry Huang, Chen Zhu, Eric Zhu, Elico Teixeira, Matthew Fritze, Francesco Bertolini, Liana-
661 Eleonora Marinescu, Martin Bölle, Dominik Paulus, Khyatti Gupta, Tejasi Latkar, Max Chang,
662 Jason Sanders, Roopa Wilson, Xuewei Wu, Yi-Xuan Tan, Lam Nguyen Thiet, Tulsee Doshi,
663 Sid Lall, Swaroop Mishra, Wanming Chen, Thang Luong, Seth Benjamin, Jasmine Lee, Ewa
664 Andrejczuk, Dominik Rabiej, Vipul Ranjan, Krzysztof Styrz, Pengcheng Yin, Jon Simon, Mal-
665 colm Rose Harriott, Mudit Bansal, Alexei Robsky, Geoff Bacon, David Greene, Daniil Mirylenka,
666 Chen Zhou, Obaid Sarvana, Abhimanyu Goyal, Samuel Andermatt, Patrick Siegler, Ben Horn,
667 Assaf Israel, Francesco Pongetti, Chih-Wei "Louis" Chen, Marco Selvatici, Pedro Silva, Kathie
668 Wang, Jackson Tolins, Kelvin Guu, Roey Yogev, Xiaochen Cai, Alessandro Agostini, Maulik
669 Shah, Hung Nguyen, Noah Ó Donnaile, Sébastien Pereira, Linda Friso, Adam Stambler, Adam
670 Kurzrok, Chenkai Kuang, Yan Romanikhin, Mark Geller, ZJ Yan, Kane Jang, Cheng-Chun Lee,
671 Wojciech Fica, Eric Malmi, Qijun Tan, Dan Banica, Daniel Balle, Ryan Pham, Yanping Huang,
672 Diana Avram, Hongzhi Shi, Jasjot Singh, Chris Hidey, Niharika Ahuja, Pranab Saxena, Dan Do-
673 oley, Srividya Pranavi Potharaju, Eileen O'Neill, Anand Gokulchandran, Ryan Foley, Kai Zhao,
674 Mike Dusenberry, Yuan Liu, Pulkit Mehta, Ragha Kotikalapudi, Chalence Safranek-Shrader, An-
675 drew Goodman, Joshua Kessinger, Eran Globen, Prateek Kolhar, Chris Gorgolewski, Ali Ibrahim,
676 Yang Song, Ali Eichenbaum, Thomas Brovelli, Sahitya Potluri, Preethi Lahoti, Cip Baetu, Ali
677 Ghorbani, Charles Chen, Andy Crawford, Shalini Pal, Mukund Sridhar, Petru Gurita, Asier Mu-
678 jika, Igor Petrovski, Pierre-Louis Cedo, Chenmei Li, Shiyuan Chen, Niccolò Dal Santo, Sid-
679 dharth Goyal, Jitesh Punjabi, Karthik Kappaganthu, Chester Kwak, Pallavi LV, Sarmishta Velury,
680 Himadri Choudhury, Jamie Hall, Premal Shah, Ricardo Figueira, Matt Thomas, Minjie Lu, Ting
681 Zhou, Chintu Kumar, Thomas Jurdi, Sharat Chikkerur, Yenai Ma, Adams Yu, Soo Kwak, Victor
682 Åhdel, Sujevan Rajayogam, Travis Choma, Fei Liu, Aditya Barua, Colin Ji, Ji Ho Park, Vincent
683 Hellendoorn, Alex Bailey, Taylan Bilal, Huanjie Zhou, Mehrdad Khatir, Charles Sutton, Woj-
684 ciech Rzadkowski, Fiona Macintosh, Konstantin Shagin, Paul Medina, Chen Liang, Jinjing Zhou,
685 Pararth Shah, Yingying Bi, Attila Dankovics, Shipra Banga, Sabine Lehmann, Marissa Bredesen,
686 Zifan Lin, John Eric Hoffmann, Jonathan Lai, Raynald Chung, Kai Yang, Nihal Balani, Arthur
687 Bražinskis, Andrei Sozanschi, Matthew Hayes, Héctor Fernández Alcalde, Peter Makarov, Will
688 Chen, Antonio Stella, Liselotte Snijders, Michael Mandl, Ante Kärrman, Paweł Nowak, Xinyi
689 Wu, Alex Dyck, Krishnan Vaidyanathan, Raghavender R, Jessica Mallet, Mitch Rudominer, Eric
690 Johnston, Sushil Mittal, Akhil Udathu, Janara Christensen, Vishal Verma, Zach Irving, Andreas
691 Santucci, Gamaleldin Elsayed, Elnaz Davoodi, Marin Georgiev, Ian Tenney, Nan Hua, Geoffrey
692 Cideron, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy Zheng, Dylan
693 Scandinaro, Heinrich Jiang, Jasper Snoek, Mukund Sundararajan, Xuezhi Wang, Zack Ontiveros,
694 Itay Karo, Jeremy Cole, Vinu Rajashekhar, Lara Tumeh, Eyal Ben-David, Rishub Jain, Jonathan
695 Uesato, Romina Datta, Oskar Bunyan, Shimu Wu, John Zhang, Piotr Stanczyk, Ye Zhang, David
696 Steiner, Subhajit Naskar, Michael Azzam, Matthew Johnson, Adam Paszke, Chung-Cheng Chiu,
697 Jaime Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin Miao, Andrew Lee, Nino Vieil-
698 lard, Jane Park, Jiageng Zhang, Jeff Stanway, Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong
699 Lee, Aviral Kumar, Luowei Zhou, Jonathan Evens, William Isaac, Geoffrey Irving, Edward Loper,
700 Michael Fink, Isha Arkatkar, Nanxin Chen, Izhak Shafraan, Ivan Petrychenko, Zhe Chen, John-
701 son Jia, Anselm Levskaya, Zhenkai Zhu, Peter Grabowski, Yu Mao, Alberto Magni, Kaisheng
Yao, Javier Snaider, Norman Casagrande, Evan Palmer, Paul Suganthan, Alfonso Castaño, Irene
Giannoumis, Wooyeol Kim, Mikofaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David So-
ergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Di-
ana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li,
Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Mar-

702 cus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Ginger Perng, Elena Allica Abellan,
703 Mingyang Zhang, Ishita Dasgupta, Nate Kushman, Ivo Penchev, Alena Repina, Xihui Wu, Tom
704 van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse,
705 Fan Yang, Jeff Piper, Nathan Ie, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Daniel
706 Andor, Pedro Valenzuela, Minnie Lui, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan
707 Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Cassidy Hardin, Lucas Dixon, Lili
708 Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Dayou Du, Dan McKinnon,
709 Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi
710 Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Ken Franko, Anna Bu-
711 lanova, Rémi Leblond, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer,
712 Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, Han-
713 nah Forbes, Dylan Banarse, Zora Tung, Mark Omernick, Colton Bishop, Rachel Sterneck, Rohan
714 Jain, Jiawei Xia, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Daniel J.
715 Mankowitz, Alex Polozov, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Den-
716 nis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Matthieu Geist, Ser tan Girgin, Hui Li,
717 Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Danila
718 Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Kathy Wu, David Miller, Nico-
719 las Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian
720 Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu
721 Zheng, Frederick Liu, Fan Yang, Rui Zhu, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko,
722 Nejc Trdin, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver
723 Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina,
724 Rishika Sinha, Alice Talbert, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton,
725 Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Saaber Fatehi, John Wieting, Omar Ajmeri, Be-
726 nigno Urias, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Yeqing
727 Li, Nir Levine, Ariel Stolovich, Rebeca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim,
728 Robin Strudel, Ali Elqursh, Charlie Deck, Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoffmann,
729 Dan Holtmann-Rice, Olivier Bachem, Sho Arora, Christy Koh, Soheil Hassas Yeganeh, Siim
730 Pöder, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu
731 Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun
732 Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Vinod Koverkathu, Christo-
733 pher A. Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash Shroff, Mani Varadarajan,
734 Sanaz Bahargam, Rob Willoughby, David Gaddy, Guillaume Desjardins, Marco Cornero, Brona
735 Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson,
736 Alireza Ghaffarkhah, Morgane Rivièrre, Alanna Walton, Clément Crepy, Alicia Parrish, Zongwei
737 Zhou, Clement Farabet, Carey Radebaugh, Praveen Srinivasan, Claudia van der Salm, Andreas
738 Fidljeland, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson,
739 Dario de Cesare, Tom Hudson, Piermaria Mendolicchio, Lexi Walker, Alex Morris, Matthew
740 Mauger, Alexey Guseynov, Alison Reid, Seth Odoom, Lucia Loher, Victor Cotruta, Madhavi
741 Yenugula, Dominik Grewe, Anastasia Petrushkina, Tom Duerig, Antonio Sanchez, Steve Yad-
742 lowsky, Amy Shen, Amir Globerson, Lynette Webb, Sahil Dua, Dong Li, Surya Bhupatiraju, Dan
743 Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj Khare, Shreyas Ram-
744 mohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin Sang, Brennan
745 Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Manish Reddy
746 Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke, Xiao Ma, Evgenii Eltyshev, Nina Martin,
747 Hardie Cate, James Manyika, Keyvan Amiri, Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier,
748 Nilesh Tripuraneni, David Madras, Mandy Guo, Austin Waters, Oliver Wang, Joshua Ainslie,
749 Jason Baldridge, Han Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Riham Mansour, Jason
750 Gelman, Yang Xu, George Polovets, Ji Liu, Honglong Cai, Warren Chen, XiangHai Sheng,
751 Emily Xue, Sherjil Ozair, Christof Angermueller, Xiaowei Li, Anoop Sinha, Weiren Wang, Julia
752 Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand Iyer, Madhu Gurumurthy, Mark Gold-
753 enson, Parashar Shah, MK Blake, Hongkun Yu, Anthony Urbanowicz, Jennimaria Palomaki,
754 Chrisantha Fernando, Ken Durden, Harsh Mehta, Nikola Momchev, Elahe Rahimtoroghi, Maria
755 Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw, Jinhyuk Lee, Denny Zhou, Komal
Jalan, Dinghua Li, Blake Hechtman, Parker Schuh, Milad Nasr, Kieran Milan, Vladimir Mikul-
lik, Juliana Franco, Tim Green, Nam Nguyen, Joe Kelley, Aroma Mahendru, Andrea Hu, Joshua
Howland, Ben Vargas, Jeffrey Hui, Kshitij Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang,
Ke Ye, Jean Michel Sarr, Melanie Moranski Preston, Madeleine Elish, Steve Li, Aakash Kaku,
Jigar Gupta, Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi M., Xinyun Chen, Aida Amini,

- 756 Alex Fabrikant, Eric Chu, Xuanyi Dong, Amruta Muthal, Senaka Buthpitiya, Sarthak Jauhari, Nan
757 Hua, Urvashi Khandelwal, Ayal Hitron, Jie Ren, Larissa Rinaldi, Shahar Drath, Avigail Dabush,
758 Nan-Jiang Jiang, Harshal Godhia, Uli Sachs, Anthony Chen, Yicheng Fan, Hagai Taitelbaum,
759 Hila Noga, Zhuyun Dai, James Wang, Chen Liang, Jenny Hamer, Chun-Sung Ferng, Chenel
760 Elkind, Aviël Atias, Paulina Lee, Vít Listík, Mathias Carlen, Jan van de Kerkhof, Marcin Pikuś,
761 Krunoslav Zaher, Paul Müller, Sasha Zykova, Richard Stefanec, Vitaly Gatsko, Christoph Hirn-
762 schall, Ashwin Sethi, Xingyu Federico Xu, Chetan Ahuja, Beth Tsai, Anca Stefanoiu, Bo Feng,
763 Keshav Dhandhanian, Manish Katyal, Akshay Gupta, Atharva Parulekar, Divya Pitta, Jing Zhao,
764 Vivaan Bhatia, Yashodha Bhavnani, Omar Alhadlaq, Xiaolin Li, Peter Danenberg, Dennis Tu,
765 Alex Pine, Vera Filippova, Abhipso Ghosh, Ben Limonchik, Bhargava Urala, Chaitanya Krishna
766 Lanka, Derik Clive, Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak, Ianna Li, Kalind Thakkar,
767 Kuanysh Omarov, Kushal Majmundar, Michael Alverson, Michael Kucharski, Mohak Patel, Mu-
768 dit Jain, Maksim Zabelin, Paolo Pelagatti, Rohan Kohli, Saurabh Kumar, Joseph Kim, Swetha
769 Sankar, Vineet Shah, Lakshmi Ramachandruni, Xiangkai Zeng, Ben Bariach, Laura Weidinger,
770 Tu Vu, Alek Andreev, Antoine He, Kevin Hui, Sheleem Kashem, Amar Subramanya, Sissie Hsiao,
771 Demis Hassabis, Koray Kavukcuoglu, Adam Sadovsky, Quoc Le, Trevor Strohman, Yonghui Wu,
772 Slav Petrov, Jeffrey Dean, and Oriol Vinyals. Gemini: A family of highly capable multimodal
773 models, 2024a. URL <https://arxiv.org/abs/2312.11805>.
- 774 GeminiTeam, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,
775 Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng,
776 Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin,
777 Andrea Tacchetti, Colin Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love,
778 Paul Voigtlaender, Rohan Jain, Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn,
779 Tao Zhu, Kornraphop Kawintiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz,
780 Manaal Faruqui, Natalie Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki
781 Morioka, Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Samer
782 Hassan, Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth Goyal,
783 Paul Barham, DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad Vikram, Dmitry
784 Lepikhin, Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, Kiran Vo-
785 drahalli, Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Sid-
786 dhartha Brahma, David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo
787 Xu, Raphael Lopez Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den
788 Driessche, Seth Odoom, Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, San-
789 tiago Ontanon, Luheng He, Denis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis
790 Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran
791 Kazemi, Ehsan Amid, Anastasia Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris
792 Larkin, Mario Pinto, Geng Yan, Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave
793 Orr, Sebastien M. R. Arnold, Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas
794 Yadav, Xi Chen, Elena Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel, Paul Komarek,
795 Sophia Austin, Sebastian Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, Kris Cao, Da-
796 Woon Chung, Matthew Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia Chen,
797 Kaushik Shivakumar, Rishabh Agarwal, Harshal Godhia, Ravi Rajwar, Javier Snaider, Xerxes
798 Dotiwalla, Yuan Liu, Aditya Barua, Victor Ungureanu, Yuan Zhang, Bat-Orgil Batsaikhan, Ma-
799 teo Wirth, James Qin, Ivo Danihelka, Tulsee Doshi, Martin Chadwick, Jilin Chen, Sanil Jain,
800 Quoc Le, Arjun Kar, Madhu Gurusurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lam-
801 prou, Rich Munoz, Nathan Lintz, Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo,
802 Quan Wang, Lorenzo Blanco, Albin Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee, Jakub
803 Sygnowski, Zach Fisher, James Besley, Richard Powell, Zafarali Ahmed, Dominik Paulus, David
804 Reitter, Zalan Borsos, Rishabh Joshi, Aedan Pope, Steven Hand, Vittorio Selo, Vihan Jain, Nikhil
805 Sethi, Megha Goel, Takaki Makino, Rhys May, Zhen Yang, Johan Schalkwyk, Christina Butter-
806 field, Anja Hauth, Alex Goldin, Will Hawkins, Evan Senter, Sergey Brin, Oliver Woodman, Mar-
807 vin Ritter, Eric Noland, Minh Giang, Vijay Bolina, Lisa Lee, Tim Blyth, Ian Mackinnon, Machel
808 Reid, Obaid Sarvana, David Silver, Alexander Chen, Lily Wang, Loren Maggiore, Oscar Chang,
809 Nithya Attaluri, Gregory Thornton, Chung-Cheng Chiu, Oskar Bunyan, Nir Levine, Timothy
810 Chung, Evgenii Eltyshev, Xiance Si, Timothy Lillicrap, Demetra Brady, Vaibhav Aggarwal, Boxi
811 Wu, Yuanzhong Xu, Ross McIlroy, Kartikeya Badola, Paramjit Sandhu, Erica Moreira, Wojciech
812 Stokowiec, Ross Hemsley, Dong Li, Alex Tudor, Pranav Shyam, Elahe Rahimtoroghi, Salem
813 Haykal, Pablo Sprechmann, Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki, Kalpesh Krishna,

810 Xiao Wu, Alexandre Frechette, Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang, Thi Avrahami,
811 Ye Zhang, Emanuel Taropa, Hanzhao Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano, Hyun-
812 Jeong Choe, Alex Tomala, Chalence Safranek-Shrader, Nora Kassner, Mantas Pajarskas, Matt
813 Harvey, Sean Sechrist, Meire Fortunato, Christina Lyu, Gamaleldin Elsayed, Chenkai Kuang,
814 James Lottes, Eric Chu, Chao Jia, Chih-Wei Chen, Peter Humphreys, Kate Baumli, Connie Tao,
815 Rajkumar Samuel, Cicero Nogueira dos Santos, Anders Andreassen, Nemanja Rakićević, Do-
816 minik Grewe, Aviral Kumar, Stephanie Winkler, Jonathan Caton, Andrew Brock, Sid Dalmia,
817 Hannah Sheahan, Iain Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Feryal Behbahani, Flavien
818 Prost, Yanhua Sun, Artiom Myaskovsky, Thanumalayan Sankaranarayana Pillai, Dan Hurt, An-
819 geliki Lazaridou, Xi Xiong, Ce Zheng, Fabio Pardo, Xiaowei Li, Dan Horgan, Joe Stanton,
820 Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu Wang, Basil Mustafa, Albert Webson, Hyo
821 Lee, Rohan Anil, Martin Wicke, Timothy Dozat, Abhishek Sinha, Enrique Piqueras, Elahe Dabir,
822 Shyam Upadhyay, Anudhyan Boral, Lisa Anne Hendricks, Corey Fry, Josip Djolonga, Yi Su,
823 Jake Walker, Jane Labanowski, Ronny Huang, Vedant Misra, Jeremy Chen, RJ Skerry-Ryan,
824 Avi Singh, Shruti Rijhwani, Dian Yu, Alex Castro-Ros, Beer Changpinyo, Romina Datta, Sumit
825 Bagri, Arnar Mar Hrafnkelsson, Marcello Maggioni, Daniel Zheng, Yury Sulsky, Shaobo Hou,
826 Tom Le Paine, Antoine Yang, Jason Riesa, Dominika Rogozinska, Dror Marcus, Dalia El Badawy,
827 Qiao Zhang, Luyu Wang, Helen Miller, Jeremy Greer, Lars Lowe Sjos, Azade Nova, Heiga Zen,
828 Rahma Chaabouni, Mihaela Rosca, Jiepu Jiang, Charlie Chen, Ruiho Liu, Tara Sainath, Maxim
829 Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh Newlan, Zeynep Cankara, Soo Kwak, Yun-
830 han Xu, Phil Chen, Andy Coenen, Clemens Meyer, Katerina Tsihlias, Ada Ma, Juraj Gottweis,
831 Jinwei Xing, Chenjie Gu, Jin Miao, Christian Frank, Zeynep Cankara, Sanjay Ganapathy, Ishita
832 Dasgupta, Steph Hughes-Fitt, Heng Chen, David Reid, Keran Rong, Hongmin Fan, Joost van
833 Amersfoort, Vincent Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohanane, Anastasija
834 Ilic, Taylor Tobin, John Wieting, Anna Bortsova, Phoebe Thacker, Emma Wang, Emily Caveness,
835 Justin Chiu, Eren Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed Elhawaty,
836 Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun Dai, Wenhao Jia, Matthew Wiethoff, El-
837 naz Davoodi, Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel Gao, Golan Pundak, Su-
838 san Zhang, Michael Azzam, Khe Chai Sim, Sergi Caelles, James Keeling, Abhanshu Sharma,
839 Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes Bostock, Yamini Bansal, Zachary Nado,
840 Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar, Lev Proleev, Abe Ittycheriah, Soheil Has-
841 sas Yeganeh, George Polovets, Aleksandra Faust, Jiao Sun, Alban Rrustemi, Pen Li, Rakesh
842 Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh Baddepudi, Sebastian Krause,
843 Emilio Parisotto, Radu Soriccut, Zheng Xu, Dawn Bloxwich, Melvin Johnson, Behnam Neyshabur,
844 Justin Mao-Jones, Renshen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal,
845 Duc Dung Nguyen, James Svensson, Le Hou, Sarah York, Kieran Milan, Sophie Bridgers, Wiktor
846 Gworek, Marco Tagliasacchi, James Lee-Thorp, Michael Chang, Alexey Guseynov, Ale Jakse
847 Hartman, Michael Kwong, Ruizhe Zhao, Sheleem Kashem, Elizabeth Cole, Antoine Miech,
848 Richard Tanburn, Mary Phuong, Filip Pavetic, Sebastien Cevey, Ramona Comanescu, Richard
849 Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, Mariko Iinuma, Clara Huiyi Hu, Aurko Roy,
850 Shaan Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel Saputro, Anita Gergely, Steven Zheng,
851 Dawei Jia, Ioannis Antonoglou, Adam Sadovsky, Shane Gu, Yingying Bi, Alek Andreev, Sina
852 Samangooei, Mina Khan, Tomas Kocisky, Angelos Filos, Chintu Kumar, Colton Bishop, Adams
853 Yu, Sarah Hodgkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, Yong Cheng, Adam Blo-
854 niarz, Jaehoon Lee, Pedram Pejman, Paul Michel, Stephen Spencer, Vladimir Feinberg, Xuehan
855 Xiong, Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary Chesus, Bernd
856 Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa, Ambrose
857 Slone, Kedar Soparkar, Disha Shrivastava, James Cobon-Kerr, Michael Sharman, Jay Pavagadhi,
858 Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin, David Barker, Qijia Li, An-
859 ton Briukhov, Neil Houlsby, Mia Glaese, Balaji Lakshminarayanan, Nathan Schucher, Yunhao
860 Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria Recasens, Guangda Lai, Alberto
861 Magni, Nicola De Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay, Mostafa Dehghani, Jenny
862 Brennan, Yifan He, Kelvin Xu, Yang Gao, Carl Saurouf, James Molloy, Xinyi Wu, Seb Arnold,
863 Solomon Chang, Julian Schrittwieser, Elena Buchatskaya, Soroush Radpour, Martin Polacek,
Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, Sarah
Cogan, Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan Qiao,
Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kishan Ruben-
stein, Gan Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao, Parker Schuh, Axel
Stjerngren, Golnaz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Felipe Tiengo Ferreira, Aish-

864 warya Kamath, Ted Klimenko, Ken Franko, Kefan Xiao, Indro Bhattacharya, Miteyan Patel, Rui
865 Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica
866 Landon, Mara Finkelstein, Priya Jhakra, Justin Frye, Megan Barnes, Matthew Mauger, Dennis
867 Daun, Khuslen Baatarsukh, Matthew Tung, Wael Farhan, Henryk Michalewski, Fabio Viola, Fe-
868 lix de Chaumont Quitry, Charline Le Lan, Tom Hudson, Qingze Wang, Felix Fischer, Ivy Zheng,
869 Elspeth White, Anca Dragan, Jean baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam Iwan-
870 icki, Michael Isard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, Srivatsan Srini-
871 vasan, Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae, Gary
872 Wang, Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Alberti, Dan Gar-
873 rette, Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel Sohn, Josip Matak, Inaki
874 Iturrate, Michael B. Chang, Jackie Xiang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian Hut-
875 ter, Vahab Mirrokni, Nanxin Chen, Kaisheng Yao, Zoltan Egyed, Francois Galilee, Tyler Liechty,
876 Praveen Kallakuri, Evan Palmer, Sanjay Ghemawat, Jasmine Liu, David Tao, Chloe Thornton,
877 Tim Green, Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel, Hongkun
878 Yu, Ed Chi, Alexander Neitz, Jens Heitkaemper, Anu Sinha, Denny Zhou, Yi Sun, Charbel
879 Kaed, Brice Hulse, Swaroop Mishra, Maria Georgaki, Sneha Kudugunta, Clement Farabet, Izhak
880 Shafran, Daniel Vlasic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin, Guolong Su,
881 Pei Sun, Shashank V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina, William Wong,
882 Warren Weilun Chen, Peter Hawkins, Egor Filonov, Lucia Loher, Christoph Hirsenschall, Weiye
883 Wang, Jingchen Ye, Andrea Burns, Hardie Cate, Diana Gage Wright, Federico Piccinini, Lei
884 Zhang, Chu-Cheng Lin, Ionel Gog, Yana Kulizhskaya, Ashwin Sreevatsa, Shuang Song, Luis C.
885 Cobo, Anand Iyer, Chetan Tekur, Guillermo Garrido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven
886 Zheng, Hui Li, Ananth Agarwal, Christel Ngani, Kati Goshvadi, Rebeca Santamaria-Fernandez,
887 Wojciech Fica, Xinyun Chen, Chris Gorgolewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali
888 Eslami, Nan Hua, Jon Simon, Pratik Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen
889 Thiet, Quan Yuan, Florian Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srini-
890 vasan, Minmin Chen, Vinod Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith
891 Anderson, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan,
892 Michael Quinn, Amit Raul, Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard,
893 Achintya Singhal, Thang Luong, Boyu Wang, Sujeewan Rajayogam, Julian Eisenschlos, Johnson
894 Jia, Daniel Finchelstein, Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing Li,
895 Dj Dvijotham, Shalini Pal, Kai Kang, Jaelyn Konzelmann, Jennifer Beattie, Olivier Dousse, Di-
896 ane Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-
897 Rice, Krystal Kallarackal, Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen
898 Jafari, Huanjie Zhou, Lilly Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya
899 Kopparapu, Francoise Beaufays, Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hi-
900 lal Dib, Jeff Stanway, Frank Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li,
901 Xihui Wu, Sonam Goenka, David Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou, Yaxin
902 Liu, Yannie Liang, Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy
903 Basu, Li Lao, Adnan Ozturel, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna
904 Walton, Lucas Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew Bolt, Olivia Wiles, Mi-
905 lad Nasr, Iliia Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara McCarthy,
906 Misha Khalman, Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Villela, Haibin Zhang,
907 Harry Richardson, James Martens, Matko Bosnjak, Shreyas Rammohan Belle, Jeff Seibert, Mah-
908 moud Alnahlawi, Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici,
909 Lenin Simicich, Laura Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi,
910 Jovana Mitrovic, Alex Grills, Joseph Pagadora, Dessie Petrova, Danielle Eisenbud, Zhishuai
911 Zhang, Damion Yates, Bhavishya Mittal, Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli,
912 Prateek Jain, Mihajlo Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar,
913 Jun Xu, Haroon Qureshi, Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock,
914 Matthias Bauer, Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Ros-
915 gen, Ashish Shenoy, Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei
916 Xia, Jiri Simsa, Andrea Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey
917 Radebaugh, Andre Elisseeff, Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert Cui, Eri
Latorre-Chimoto, Solomon Kim, William Zeng, Ken Durden, Priya Ponnappalli, Tiberiu Sosea,
Christopher A. Choquette-Choo, James Manyika, Brona Robenek, Harsha Vashisht, Sebastien
Pereira, Hoi Lam, Marko Velic, Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia
Parrish, Shawn Lu, Jane Park, Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung
Cheng, Andrei Sozanschi, Adam Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama,

- 918 Wooyeol Kim, Nandita Dukkupati, Anthony Baryshnikov, Christos Kaplanis, XiangHai Sheng,
919 Yuri Chervonyi, Caglar Unlu, Diego de Las Casas, Harry Askham, Kathryn Tunyasuvunakool,
920 Felix Gimeno, Siim Poder, Chester Kwak, Matt Miecnikowski, Vahab Mirrokni, Alek Dimitriev,
921 Aaron Parisi, Dangyi Liu, Tomy Tsai, Toby Shevlane, Christina Kouridi, Drew Garmon, Adrian
922 Goedeckemeyer, Adam R. Brown, Anitha Vijayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang,
923 Sara Mc Carthy, Jay Hoover, Lucy Kim, Sandeep Kumar, Wei Chen, Courtney Biles, Garrett
924 Bingham, Evan Rosen, Lisa Wang, Qijun Tan, David Engel, Francesco Pongetti, Dario de Ce-
925 sare, Dongseong Hwang, Lily Yu, Jennifer Pullman, Srini Narayanan, Kyle Levin, Siddharth
926 Gopal, Megan Li, Asaf Aharoni, Trieu Trinh, Jessica Lo, Norman Casagrande, Roopali Vij, Loic
927 Matthey, Bramandia Ramadhana, Austin Matthews, CJ Carey, Matthew Johnson, Kremena Gora-
928 nova, Rohin Shah, Shereen Ashraf, Kingshuk Dasgupta, Rasmus Larsen, Yicheng Wang, Man-
929 ish Reddy Vuyyuru, Chong Jiang, Joana Ijazi, Kazuki Osawa, Celine Smith, Ramya Sree Bop-
930 pana, Taylan Bilal, Yuma Koizumi, Ying Xu, Yasemin Altun, Nir Shabat, Ben Bariach, Alex Ko-
931 rchemniy, Kiam Choo, Olaf Ronneberger, Chimezie Iwuanyanwu, Shubin Zhao, David Soergel,
932 Cho-Jui Hsieh, Irene Cai, Shariq Iqbal, Martin Sundermeyer, Zhe Chen, Elie Bursztein, Chai-
933 tanya Malaviya, Fadi Biadsy, Prakash Shroff, Inderjit Dhillon, Tejasi Latkar, Chris Dyer, Hannah
934 Forbes, Massimo Nicosia, Vitaly Nikolaev, Somer Greene, Marin Georgiev, Pidong Wang, Nina
935 Martin, Hanie Sedghi, John Zhang, Praseem Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang,
936 Jiageng Zhang, Viorica Patraucean, Dayou Du, Igor Mordatch, Ivan Jurin, Lewis Liu, Ayush
937 Dubey, Abhi Mohan, Janek Nowakowski, Vlad-Doru Ion, Nan Wei, Reiko Tojo, Maria Abi Raad,
938 Drew A. Hudson, Vaishakh Keshava, Shubham Agrawal, Kevin Ramirez, Zhichun Wu, Hoang
939 Nguyen, Ji Liu, Madhavi Sewak, Bryce Petrini, DongHyun Choi, Ivan Philips, Ziyue Wang,
940 Ioana Bica, Ankush Garg, Jarek Wilkiewicz, Priyanka Agrawal, Xiaowei Li, Danhao Guo, Emily
941 Xue, Naseer Shaik, Andrew Leach, Sadh MNM Khan, Julia Wiesinger, Sammy Jerome, Ab-
942 hishek Chakladar, Alek Wenjiao Wang, Tina Ornduff, Folake Abu, Alireza Ghaffarkhah, Marcus
943 Wainwright, Mario Cortes, Frederick Liu, Joshua Maynez, Andreas Terzis, Pouya Samangouei,
944 Riham Mansour, Tomasz Kepa, François-Xavier Aubet, Anton Algyrn, Dan Banica, Agoston
945 Weisz, Andras Orban, Alexandre Senegés, Ewa Andrejczuk, Mark Geller, Niccolo Dal Santo,
946 Valentin Anklin, Majd Al Merey, Martin Baeuml, Trevor Strohman, Junwen Bai, Slav Petrov,
947 Yonghui Wu, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. Gemini
1.5: Unlocking multimodal understanding across millions of tokens of context, 2024b. URL
<https://arxiv.org/abs/2403.05530>.
- 948 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Realm: retrieval-
949 augmented language model pre-training. In *Proceedings of the 37th International Conference on*
950 *Machine Learning*, ICML'20. JMLR.org, 2020.
- 951 Cheng-Ping Hsieh, Simeng Sun, Samuel Krizan, Shantanu Acharya, Dima Rekeshe, Fei Jia, Yang
952 Zhang, and Boris Ginsburg. Ruler: What's the real context size of your long-context language
953 models? *arXiv preprint arXiv:2404.06654*, 2024.
- 954
955 Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open
956 domain question answering. In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty (eds.), *Proceed-*
957 *ings of the 16th Conference of the European Chapter of the Association for Computational Lin-*
958 *guistics: Main Volume*, pp. 874–880, Online, April 2021. Association for Computational Linguis-
959 tics. doi: 10.18653/v1/2021.eacl-main.74. URL [https://aclanthology.org/2021.](https://aclanthology.org/2021.eacl-main.74)
960 [eacl-main.74](https://aclanthology.org/2021.eacl-main.74).
- 961
962 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane
963 Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: few-shot learning
964 with retrieval augmented language models. *J. Mach. Learn. Res.*, 24(1), March 2024. ISSN
965 1532-4435.
- 966
967 Ziyang Jiang, Xueguang Ma, and Wenhui Chen. Longrag: Enhancing retrieval-augmented gen-
968 eration with long-context llms. *ArXiv*, abs/2406.15319, 2024. URL [https://api.](https://api.semanticscholar.org/CorpusID:270688725)
969 [semanticscholar.org/CorpusID:270688725](https://api.semanticscholar.org/CorpusID:270688725).
- 970
971 Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly
supervised challenge dataset for reading comprehension. In Regina Barzilay and Min-Yen Kan
(eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*

- 972 (Volume 1: Long Papers), pp. 1601–1611, Vancouver, Canada, July 2017. Association for Com-
 973 putational Linguistics. doi: 10.18653/v1/P17-1147. URL [https://aclanthology.org/
 974 P17-1147](https://aclanthology.org/P17-1147).
 975
- 976 Norman P. Jouppi, George Kurian, Sheng Li, Peter C. Ma, Rahul Nagarajan, Lifeng Nai, Nishant
 977 Patil, Suvinay Subramanian, Andy Swing, Brian Towles, Cliff Young, Xiaoping Zhou, Zongwei
 978 Zhou, and David A. Patterson. Tpu v4: An optically reconfigurable supercomputer for machine
 979 learning with hardware support for embeddings. *Proceedings of the 50th Annual International
 980 Symposium on Computer Architecture*, 2023. URL [https://api.semanticscholar.
 981 org/CorpusID:257921908](https://api.semanticscholar.org/CorpusID:257921908).
- 982 Gregory Kamradt. Needle in a haystack - pressure testing llms, 2023. URL [https://github.
 983 com/gkamradt/LLMTest_NeedleInAHaystack/tree/main](https://github.com/gkamradt/LLMTest_NeedleInAHaystack/tree/main).
 984
- 985 Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi
 986 Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In Bonnie
 987 Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on
 988 Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6769–6781, Online, Novem-
 989 ber 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.550.
 990 URL <https://aclanthology.org/2020.emnlp-main.550>.
 991
- 992 Gangwoo Kim, Sungdong Kim, Byeongguk Jeon, Joonsuk Park, and Jaewoo Kang. Tree of clar-
 993 ifications: Answering ambiguous questions with retrieval-augmented large language models.
 994 In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference
 995 on Empirical Methods in Natural Language Processing*, pp. 996–1009, Singapore, December
 996 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.63. URL
 997 <https://aclanthology.org/2023.emnlp-main.63>.
- 998 Jaehyung Kim, Jaehyun Nam, Sangwoo Mo, Jongjin Park, Sang-Woo Lee, Minjoon Seo, Jung-Woo
 999 Ha, and Jinwoo Shin. Sure: Summarizing retrievals using answer candidates for open-domain
 1000 QA of LLMs. In *The Twelfth International Conference on Learning Representations*, 2024. URL
 1001 <https://openreview.net/forum?id=w4DW6qkRmt>.
 1002
- 1003 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR (Poster)*,
 1004 2015. URL <http://arxiv.org/abs/1412.6980>.
- 1005 Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis,
 1006 and Edward Grefenstette. The NarrativeQA reading comprehension challenge. *Transactions of
 1007 the Association for Computational Linguistics*, 6:317–328, 2018. doi: 10.1162/tacl.a.00023.
 1008 URL <https://aclanthology.org/Q18-1023>.
 1009
- 1010 Andreas Kopf, Yannic Kilcher, Dimitri von Rutte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens,
 1011 Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Rich’ard Nagyfi, ES Shahul, Sameer
 1012 Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen,
 1013 and Alexander Mattick. Openassistant conversations - democratizing large language model
 1014 alignment. *ArXiv*, abs/2304.07327, 2023. URL [https://api.semanticscholar.org/
 1015 CorpusID:258179434](https://api.semanticscholar.org/CorpusID:258179434).
- 1016 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
 1017 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion
 1018 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav
 1019 Petrov. Natural questions: A benchmark for question answering research. *Transactions of the
 1020 Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl.a.00276. URL
 1021 <https://aclanthology.org/Q19-1026>.
 1022
- 1023 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
 1024 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
 1025 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating
 Systems Principles*, 2023.

- 1026 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer
1027 Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: denoising sequence-to-sequence pre-
1028 training for natural language generation, translation, and comprehension. In *ACL*, pp. 7871–7880,
1029 2020a.
- 1030 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
1031 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe
1032 Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the*
1033 *34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook,
1034 NY, USA, 2020b. Curran Associates Inc. ISBN 9781713829546.
- 1035 Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and
1036 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the*
1037 *Association for Computational Linguistics*, 12:157–173, 2024. doi: 10.1162/tacl_a.00638. URL
1038 <https://aclanthology.org/2024.tacl-1.9>.
- 1039 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
1040 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining
1041 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 1042 MetaAI, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,
1043 Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
1044 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
1045 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
1046 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
1047 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
1048 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
1049 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
1050 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
1051 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
1052 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
1053 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
1054 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
1055 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
1056 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
1057 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
1058 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
1059 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der
1060 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
1061 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
1062 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
1063 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
1064 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
1065 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
1066 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
1067 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
1068 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
1069 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
1070 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
1071 Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
1072 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
1073 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
1074 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
1075 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
1076 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
1077 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Ji, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,
1078 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
1079 Couderc, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay
Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda

- 1080 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew
 1081 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita
 1082 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh
 1083 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De
 1084 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
 1085 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
 1086 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
 1087 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
 1088 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
 1089 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
 1090 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
 1091 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
 1092 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
 1093 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
 1094 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
 1095 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
 1096 man, Ibrahim Damraj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
 1097 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
 1098 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
 1099 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
 1100 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
 1101 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
 1102 Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
 1103 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
 1104 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
 1105 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
 1106 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
 1107 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
 1108 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
 1109 ata Bawa, Nayan Singhal, Nick Gebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
 1110 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
 1111 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
 1112 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
 1113 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
 1114 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
 1115 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
 1116 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
 1117 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
 1118 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
 1119 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
 1120 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
 1121 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
 1122 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
 1123 Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
 1124 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
 1125 stable, Xiaocheng Tang, Xiaofang Wang, Xiaojuan Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
 1126 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
 1127 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef
 1128 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024.
- 1126 MistralAI. Mistral nemo, 2024. URL <https://mistral.ai/news/mistral-nemo/>.
- 1127
- 1128 OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 1129
- 1130 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,
 1131 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser
 1132 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan
 1133 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback.
 In *NeurIPS*, 2022.

- 1134 Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. YaRN: Efficient context win-
1135 dows extension of large language models. In *The Twelfth International Conference on Learning*
1136 *Representations*, 2024. URL <https://openreview.net/forum?id=wHBfxhZulu>.
- 1137 Ofir Press, Noah Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables
1138 input length extrapolation. In *International Conference on Learning Representations*, 2022. URL
1139 <https://openreview.net/forum?id=R8sQPpGCv0>.
- 1140 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language under-
1141 standing by generative pre-training. 2018.
- 1142 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
1143 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 1144 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and
1145 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
1146 In *NeurIPS*, 2023.
- 1147 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
1148 Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text
1149 transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020.
- 1150 Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the
1151 parameters of a language model? In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu
1152 (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Process-*
1153 *ing (EMNLP)*, pp. 5418–5426, Online, November 2020. Association for Computational Linguistics.
1154 doi: 10.18653/v1/2020.emnlp-main.437. URL [https://aclanthology.org/2020.](https://aclanthology.org/2020.emnlp-main.437)
1155 [emnlp-main.437](https://aclanthology.org/2020.emnlp-main.437).
- 1156 Stephen Robertson and Hugo Zaragoza. The probabilistic relevance framework: Bm25 and beyond.
1157 *Found. Trends Inf. Retr.*, 3(4):333–389, apr 2009.
- 1158 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
1159 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Ev-
1160 timov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafori, Wenhan Xiong,
1161 Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier,
1162 Thomas Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2024.
1163 URL <https://arxiv.org/abs/2308.12950>.
- 1164 Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D Man-
1165 ning. RAPTOR: Recursive abstractive processing for tree-organized retrieval. In *The Twelfth*
1166 *International Conference on Learning Representations*, 2024. URL [https://openreview.](https://openreview.net/forum?id=GN921JHCRw)
1167 [net/forum?id=GN921JHCRw](https://openreview.net/forum?id=GN921JHCRw).
- 1168 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
1169 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 1170 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-
1171 hanced transformer with rotary position embedding. *Neurocomput.*, 568(C), March 2024. ISSN
1172 0925-2312. doi: 10.1016/j.neucom.2023.127063. URL [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.neucom.2023.127063)
1173 [neucom.2023.127063](https://doi.org/10.1016/j.neucom.2023.127063).
- 1174 Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. Recitation-augmented lan-
1175 guage models. In *International Conference on Learning Representations*, 2023. URL [https://openreview.](https://openreview.net/forum?id=-cqvvvb-NkI)
1176 [net/forum?id=-cqvvvb-NkI](https://openreview.net/forum?id=-cqvvvb-NkI).
- 1177 Together. OpenChatKit: An Open Toolkit and Base Model for Dialogue-style Applications, 3 2023.
1178 URL <https://github.com/togethercomputer/OpenChatKit>.
- 1179 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
1180 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez,
1181 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation lan-
1182 guage models. *ArXiv*, abs/2302.13971, 2023a. URL [https://api.semanticscholar.](https://api.semanticscholar.org/CorpusID:257219404)
1183 [org/CorpusID:257219404](https://api.semanticscholar.org/CorpusID:257219404).

- 1188 Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
1189 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas
1190 Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes,
1191 Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S.
1192 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian
1193 Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut
1194 Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,
1195 Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta,
1196 Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh
1197 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov,
1198 Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert
1199 Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat
1200 models. *ArXiv*, abs/2307.09288, 2023b. URL [https://api.semanticscholar.org/
1201 CorpusID:259950998](https://api.semanticscholar.org/CorpusID:259950998).
- 1202 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
1203 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st Inter-
1204 national Conference on Neural Information Processing Systems, NIPS'17*, pp. 6000–6010, Red
1205 Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- 1206 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,
1207 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language
1208 models. In *NeurIPS*, 2022.
- 1209 Qingyang Wu, Zhenzhong Lan, Kun Qian, Jing Gu, Alborz Geramifard, and Zhou Yu. Memformer:
1210 A memory-augmented transformer for sequence modeling. In Yulan He, Heng Ji, Sujian Li,
1211 Yang Liu, and Chua-Hui Chang (eds.), *Findings of the Association for Computational Linguistics:
1212 ACL-IJCNLP 2022*, pp. 308–318, Online only, November 2022a. Association for Computational
1213 Linguistics. URL <https://aclanthology.org/2022.findings-acl.29>.
- 1214 Yuhuai Wu, Markus Norman Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing trans-
1215 formers. In *International Conference on Learning Representations*, 2022b. URL [https:
1216 //openreview.net/forum?id=TrjbxzRcnf-](https://openreview.net/forum?id=TrjbxzRcnf-).
- 1217 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
1218 language models with attention sinks. In *The Twelfth International Conference on Learning Rep-
1219 resentations*, 2024a. URL <https://openreview.net/forum?id=NG7sS51zVF>.
- 1220 Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. C-
1221 pack: Packed resources for general chinese embeddings. In *Proceedings of the 47th Interna-
1222 tional ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR
1223 '24*, pp. 641–649, New York, NY, USA, 2024b. Association for Computing Machinery. ISBN
1224 9798400704314. doi: 10.1145/3626772.3657878. URL [https://doi.org/10.1145/
1226 3626772.3657878](https://doi.org/10.1145/
1225 3626772.3657878).
- 1227 Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed,
1228 and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense
1229 text retrieval. In *International Conference on Learning Representations*, 2021a. URL [https:
1230 //openreview.net/forum?id=zeFrfgYzln](https://openreview.net/forum?id=zeFrfgYzln).
- 1231 Wenhan Xiong, Xiang Lorraine Li, Srinivasan Iyer, Jingfei Du, Patrick Lewis, William Yang Wang,
1232 Yashar Mehdad, Wen-tau Yih, Sebastian Riedel, Douwe Kiela, and Barlas Oğuz. Answering
1233 complex open-domain questions with multi-hop dense retrieval. In *International Conference on
1234 Learning Representations (ICLR)*, 2021b.
- 1235 Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin,
1236 Rashi Rungta, Karthik Abinav Sankararaman, Barlas Oğuz, Madian Khabsa, Han Fang, Yashar
1237 Mehdad, Sharan Narang, Kshitiz Malik, Angela Fan, Shruti Bhosale, Sergey Edunov, Mike Lewis,
1238 Sinong Wang, and Hao Ma. Effective long-context scaling of foundation models. In *North
1239 American Chapter of the Association for Computational Linguistics*, 2023. URL [https://
1241 api.semanticscholar.org/CorpusID:263134982](https://
1240 api.semanticscholar.org/CorpusID:263134982).

1242 Yuanzhong Xu, HyoukJoong Lee, Dehao Chen, Blake A. Hechtman, Yanping Huang, Rahul Joshi,
1243 Maxim Krikun, Dmitry Lepikhin, Andy Ly, Marcello Maggioni, Ruoming Pang, Noam M.
1244 Shazeer, Shibo Wang, Tao Wang, Yonghui Wu, and Zhifeng Chen. Gspmd: General and scal-
1245 able parallelization for ml computation graphs. *ArXiv*, abs/2105.04663, 2021. URL <https://api.semanticscholar.org/CorpusID:234357958>.
1246
1247 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
1248 ReAct: Synergizing reasoning and acting in language models. In *International Conference on*
1249 *Learning Representations (ICLR)*, 2023.
1250
1251 Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu,
1252 Michael Zeng, and Meng Jiang. Generate rather than retrieve: Large language models are strong
1253 context generators. In *The Eleventh International Conference on Learning Representations, 2023*.
1254 URL <https://openreview.net/forum?id=fB0hRu9GZUS>.
1255
1256 Dawei Zhu, Nan Yang, Liang Wang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. PoSE:
1257 Efficient context window extension of LLMs via positional skip-wise training. In *The Twelfth*
1258 *International Conference on Learning Representations, 2024*. URL <https://openreview.net/forum?id=3Z1gXuAQRa>.
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295