

Krikri: Advancing Open Large Language Models for Greek

Anonymous EMNLP 2025 submission

Abstract

We introduce Llama-Krikri-8B, a cutting-edge Large Language Model tailored for the Greek language, built on Meta’s Llama 3.1-8B. Llama-Krikri-8B has been extensively trained on high-quality Greek data to ensure superior adaptation to linguistic nuances. With 8 billion parameters, it offers advanced capabilities while maintaining efficient computational performance. Llama-Krikri-8B supports both Modern Greek and English, and is also equipped to handle polytonic text and Ancient Greek. The chat version of Llama-Krikri-8B features a multi-stage post-training pipeline, utilizing both human and synthetic instruction and preference data, by applying techniques such as MAGPIE. In addition, for evaluation, we propose three novel public benchmarks for Greek. Our evaluation on existing as well as the proposed benchmarks shows notable improvements over comparable Greek and multilingual LLMs in both natural language understanding and generation as well as code generation.

1 Introduction

Recent advancements in AI have been largely driven by the development of large-scale foundation models. Meta’s Llama 3 (Grattafiori et al., 2024) fostered a new generation of open models, designed for strong multilingual capabilities, code generation, reasoning, and tool use. With extended context windows, and refined training strategies, models based on Llama 3 have achieved performance comparable to proprietary systems like GPT-4. A critical aspect in this evolution is the development of multilingual and language-specific models, democratizing access to AI technologies and preserving linguistic diversity.

While substantial progress has been made for widely spoken languages, low and medium resource languages remain underrepresented. Greek, in particular, has received limited attention despite its linguistic complexity, rich cultural her-

itage, and historical significance. Addressing this gap, we present Llama-Krikri-8B, a cutting-edge open Large Language Model tailored for the Greek language. Built on Meta’s Llama 3.1-8B architecture, Llama-Krikri has been continually pretrained on a diverse, high-quality Greek corpus. This allows the model to effectively capture the syntactic and semantic nuances of Greek, while retaining the multilingual strengths of the base model. Notably, Llama-Krikri also supports English and is capable of handling polytonic and Ancient Greek texts, addressing not only contemporary but also historical forms of the language.

Compared to Meltemi-7B (Voukoutis et al., 2024), the previous state-of-the-art open Greek LLM built on Mistral 7B (Jiang et al., 2023), Llama-Krikri-8B significantly increases the number of parameters, context length, and training data scale. Additionally, it features an enhanced post-training pipeline using both human and synthetic data. Following the MAGPIE methodology (Xu et al., 2024), we generate high-quality instruction-response pairs via prompting of aligned models, and apply multiple rounds of instruction tuning and alignment using Direct Preference Optimization (DPO) (Rafailov et al., 2024). This pipeline ensures that the model produces helpful, honest, and harmless outputs.

To evaluate Llama-Krikri-8B, we also introduce three novel public benchmarks specifically designed for Greek. These, alongside existing evaluation suites, show that Llama-Krikri outperforms comparable Greek and multilingual LLMs in both natural language understanding and generation, as well as code-related tasks. Moreover, it supports function calling and agentic behaviors, opening new application domains for Greek users. Llama-Krikri-8B is available under the Llama 3.1 Community License Agreement¹.

¹https://www.llama.com/llama3_1/license/

Our key contributions are:

- We present Llama-Krikri-8B, a state-of-the-art open Greek foundation model based on Llama 3.1, demonstrating strong capabilities in Modern and Ancient Greek, English, and code generation, along with support for function calling and agentic behavior.
- We implement a rigorous post-training pipeline incorporating synthetic instruction tuning via MAGPIE and alignment through DPO.
- We introduce three new benchmarks for evaluating Greek LLMs, covering language understanding, generation, and code tasks.
- We show that Llama-Krikri-8B significantly outperforms existing open Greek and multilingual models across several domains, while supporting advanced features such as function calling.

2 Background and Related Work

Large Language Models (LLMs) have achieved state-of-the-art performance across a wide variety of natural language processing (NLP) tasks. These models are typically trained on massive corpora dominated by English, leading to strong performance in English-language tasks but comparatively weaker capabilities in other languages (Devlin et al., 2019; Brown et al., 2020). As a result, the development of language-specific LLMs has become an active area of research, particularly for under-represented languages.

One prominent strategy for developing such models is continual pretraining, where a pretrained base model is further trained on data in the target language. This approach allows researchers to leverage the general capabilities of large base models while improving performance in specific linguistic domains, without the prohibitive cost of training from scratch (Gururangan et al., 2020).

Several recent projects have successfully applied continual pretraining to adapt existing models to new languages. BgGPT-GEMMA-2-27B-Instruct (Alexandrov et al., 2024) fine-tunes Google’s Gemma-2 model (Riviere et al., 2024) for Bulgarian, combining over 100B tokens of Bulgarian and English data and applying techniques such as Branch-and-Merge to mitigate catastrophic forgetting. Similarly, LeoLM (LAION, 2023) and the

Sabiá models (Pires et al., 2023) adapt LLaMA and Mistral-based architectures for German and Portuguese, respectively, through targeted continual pretraining.

For Greek, Meltemi-7B represents the first open generative LLM tailored to the language. It was developed by continually pretraining Mistral-7B on a substantial Greek corpus, followed by instruction fine-tuning. While effective, Meltemi’s performance is bounded by the size and capabilities of the base model, as well as the limited post-training alignment techniques employed at the time.

Beyond language adaptation, alignment of LLMs to generate helpful, harmless, and honest outputs has become increasingly central. Early approaches such as InstructGPT (Ouyang et al., 2022) and Constitutional AI (Bai et al., 2022) rely on multi-stage fine-tuning pipelines involving human feedback or rule-based constraints. More recently, DPO and data synthesis methods like MAGPIE have enabled scalable and effective instruction tuning. MAGPIE, in particular, leverages high-quality prompting of already-aligned models to generate large volumes of instruction-response pairs, demonstrating that synthetic data can rival or surpass human-curated datasets.

These advancements highlight a trend toward bootstrapping high-quality training data using strong base models, especially in low-resource languages. Our work builds on this foundation by employing Llama 3.1 as a base architecture, scaling up parameter count and context length, and applying a more rigorous post-training pipeline, including instruction tuning with MAGPIE-generated data and DPO for alignment.

3 Methodology

Llama-Krikri-8B is based on the Transformer architecture (Vaswani et al., 2023), which has become the de facto standard for large language models. The model inherits its architecture from Meta’s Llama 3.1-8B, leveraging the strong foundation in multilingual understanding, code generation, and reasoning provided by Llama 3.1.

Adapting an LLM for the Greek language requires addressing the lack of high-quality Greek data in the massive datasets typically used to train foundation models. Even though Llama 3.1’s pre-training corpus comprises trillions of tokens, it struggles to generate coherent Greek text, thus indicating that Greek data is only a tiny fraction of

Subcorpus	Original Tokens (B)	Percentage	Upsampled Tokens (B)	Percentage
Greek	56.7	62.3%	66.1	60.0%
English	21.0	23.1%	25.2	22.9%
Parallel	5.5	6.0%	8.8	8.0%
Math/Code	7.8	8.6%	10.1	9.1%
Total	91.0	100%	110.2	100%

Table 1: Composition of the pretraining corpus - original and upsampled

its training data; we should note however that there is limited information on the composition of its pretraining data. Our approach is to perform continual pretraining with Greek and parallel data to infuse the model with Greek knowledge. This training must be done carefully to avoid catastrophic forgetting (Luo et al., 2025) of the base model’s prior knowledge in other languages and domains. We tackle this by (a) constructing a high-quality, large-scale Greek corpus, (b) extending and tuning the tokenizer, (c) interleaving Greek training data with datasets containing English, mathematics, code, and parallel data in languages that the model has already been trained on, (d) employing a dataset sampling schedule during training that prefers data closer to the initial llama-3.1 distribution in the beginning, while shifting closer to our true dataset distribution as training continues and (d) re-warming and re-decaying the learning rate (Ibrahim et al., 2024).

Pretraining Data Collection & Cleaning: As a foundation for continual pretraining, we curated a large corpus of texts totalling approximately 91 billion tokens (after filtering and deduplication), which was upsampled to 110 billion tokens for the final pretraining mix. This corpus was constructed with a primary focus on Greek and aiming on retaining and enhancing the original model’s capabilities. The distribution included 56.7 billion monolingual Greek tokens (62.3%), 21 billion monolingual English tokens (23.1%), 5.5 billion parallel data tokens (6.0%), and 7.8 billion math and code tokens (8.6%). Table 1 presents the distribution of the pretraining data mix, with more details provided in Appendix A Pretrained data mix.

After corpus collection, we implemented a multi-stage preprocessing and filtering pipeline to ensure a high quality for the pretraining data. Various parts of our filtering methodology have been informed by approaches used in Voukoutis et al. (2024) and large-scale corpus creation efforts such as Zyda

(Tokpanov et al., 2024). However, we have adapted these approaches to cater for the peculiarities of the Greek language. We detail the preprocessing pipelines we used in Appendix B Pretraining data cleaning pipelines.

Tokenizer	Vocabulary Size	Fertility Greek	Fertility English
Mistral-7B	32,000	6.80	1.49
Meltemi-7B	61,362	1.52	1.44
Llama-3.1-8B	128,000	2.73	1.33
Llama-Krikri-8B	149,248	1.65	1.33

Table 2: Tokenizer statistics for Greek and English

Tokenizer and Embeddings Expansion: The original Llama 3 tokenizer comprises 128,000 tokens and is inefficient for Greek texts, as it generally performs character-level tokenization for Greek. This was determined, through the approximation of the llama-3.1 tokenizer’s fertility (Csaki et al., 2023), a metric of the average tokens per word produced. To determine the efficiency of the original Llama 3 tokenizer and compare with our approach, we conducted tests on diverse Greek and English corpora (each one containing 100,000 rows and totalling approximately 2M words) and calculated the difference in fertility, as can be seen in Table 2.

In order to develop an optimal tokenizer for Greek which is also efficient in historical dialects of the language, as well as in critical domains (such as legal and scientific texts), we extended the Llama 3 tokenizer with 20,992 new tokens through a multi-stage process which encompasses curating high-quality texts and allocating new tokens across five domains. This process is especially important during model inference, as it significantly reduces the input and output token cost during model use. Furthermore, more compact representations of input text help to improve model performance. We provide detail on the steps for the tokenizer and embeddings expansion in Appendix C, Tokenizer and embeddings expansion process.

Greek embeddings training: To effectively integrate newly introduced Greek tokens into the model, we implemented an initial, targeted training phase for their corresponding embeddings. By preparing the new token representations prior to full-scale pretraining, we mitigate potential disruptions to the existing model parameters.

We initialized the model with Llama 3.1-8B-Base weights, freezing all but the embeddings and output-projection weights for the 20,992 new tokens, allowing their initial training without large gradient updates to the rest of the model. The dataset for this step was comprised of 5B tokens sampled from the main pre-training dataset. This short, several-thousand-step training regimen ensured a smoother integration of the new vocabulary into the model’s existing knowledge representation.

Continual Pretraining Process: After embedding training, all parameters were unfrozen, and training continued on the 110B token corpus using a mixed-curriculum strategy. Training alternated between chunks of predominantly Greek text and supporting data (English, parallel, code) in a round-robin fashion. This interleaving improved Greek performance while maintaining/improving English validation perplexity, similar to findings in related work. The strategy involved curriculum learning and experience replay: starting with simpler Greek/more English, progressing to diverse Greek, and mixing in periodic English/code replays. Training was conducted on two machines, each equipped with 8 NVIDIA H200 GPUs, using DeepSpeed Zero 3 and bf16 mixed precision for ~ 50 days on the 110B token dataset at 128K context length. Training used packed sequences, cosine annealing LR, AdamW optimizer, gradient clipping, and weight decay.

Annealing Phase: Following pretraining, a short annealing pass used a curated 3.5B token dataset of very high-quality texts across all subcorpora. We used within-dataset normalized perplexity, calculated using KenLM (Heafield, 2011), to implement a dataset-aware fluency scoring method for document selection. To boost comprehension and reasoning, a synthetic question-answer dataset (189M tokens) was created by prompting Gemma-2-27B-IT to generate Q&A triplets with reasoning from curated documents. Annealing was tested with and without the synthetic QA data. Performance (Table 3) showed continual pretraining improved Greek (+8.7) but reduced English (-4) vs Llama-3.1.

Annealing with curated data gave modest gains. Most notably, adding synthetic QA significantly improved Greek (+2.1 vs continual pretraining) and enhanced English beyond original Llama-3.1 (+0.8). Liger kernels (Hsu et al., 2024) were used for efficiency.

Training Stage	Avg. Greek	Avg. English
Llama-3.1-8B	48.7	66.2
+ Continual Pretraining	57.4	62.2
+ Curated Corpora	58.0	63.4
+ Synthetic QA Dataset	59.5	67.0

Table 3: Average performance across training stages on Greek and English benchmarks

Instruction Tuning and Alignment: Llama-Krikri-8B-Instruct was created by fine-tuning the base model for instruction following and dialogue. Addressing Greek data scarcity and avoiding translation artifacts, the pipeline combined data synthesis, filtering, two-stage Supervised Fine-Tuning (SFT), and Direct Preference Optimization (DPO).

Data collection, synthesis, & curation involved collecting high-quality English instruction, e.g., Tulu 3 (Lambert et al., 2025) and SmolTalk (Allal et al., 2025), and preference data e.g., UltraFeedback (Cui et al., 2023). Additionally, Greek data was synthesized via translation (with post-editing), regenerating responses using LLMs (Gemma-2-27B-IT), and generating synthetic instructions directly in Greek using the MAGPIE technique. Curated corpora from annealing were reused for synthetic Q&A and dialogues. Data was scored and filtered using the Skyword-Reward-Gemma-2-27B-v0.2 (Liu et al., 2024) reward model, known for Greek accuracy. Rule-based filters ensured formatting and language verification.

SFT was done in two stages (~ 856 k pairs in Stage 1, ~ 638 k in Stage 2), with progressively higher data quality. Datasets included filtered original English, reward-model-filtered synthetic MAGPIE data (higher scores in Stage 2), translated/post-edited data (Stage 1), regenerated responses (including "thinking" in Stage 2), multi-language translation data, synthetic QA, synthetic multi-turn dialogues, and upsampled manual safety data. Training used linear decay LR, AdamW, and Liger kernels, masking prompt loss. SFT produced a model that followed instructions but needed alignment for helpfulness, precision, and safety.

DPO provided final alignment using ~ 92 k preference triplets. Data included selected/high-scored original preferences, preferences from MAGPIE-synthesized data (using reward model scores), pref-

ferences derived from translated data (contrasting regenerated responses or regenerated vs translated reference), and safety preferences. DPO maximized preferred response likelihood while minimizing dispreferred, using length normalization. Training used AdamW (Loshchilov and Hutter, 2017), linear decay LR, and Liger kernels. DPO significantly improved response quality, safety, and helpfulness compared to the SFT-only model. The DPO-tuned model is the final Llama-Krikri-8B-Instruct.

4 Evaluation

In this section, we present evaluation details for Llama-Krikri-8B-Base and Llama-Krikri-8B-Instruct, across six Greek and six English benchmarks. We compare our base model directly with the base model Llama-3.1-8B (Grattafiori et al., 2024) and the previous Greek state-of-the-art model Meltemi-7B-v1.5 (Voukoutis et al., 2024). Additionally, we evaluate our chat model, Llama-Krikri-8B-Instruct on three challenging English benchmarks, as well as three novel constructed Greek benchmarks which correspond to the English ones.

4.1 Base Model Evaluation: Krikri-8B-Base

We evaluated Llama-Krikri-8B-Base against Llama-3.1-8B and Meltemi-7B-v1.5 in a few-shot setting, consistent with the Open LLM Leaderboard².

Greek Benchmarks: The evaluation was carried out on a suite of six Greek-specific benchmarks³ used in Voukoutis et al. (2024), including machine-translated versions of established English datasets (ARC-Challenge Greek, Truthful QA Greek, Hel-laSwag Greek, MMLU Greek), the existing Belebele Greek benchmark (Bandarkar et al., 2024), and a novel medical QA benchmark (Medical MCQA).

Results in Table 4 demonstrate substantial improvements for Greek (+10.8%) compared to Llama-3.1-8B. Moreover, we observe that Llama-Krikri-8B-Base surpasses Meltemi-7B-v1.5 with a notable +11.6% average improvement across all benchmarks. On MMLU Greek, Llama-Krikri-8B-Base surpasses Llama-3.1-8B and Meltemi-7B-v1.5 by +9.4% and +10.8% respectively, while on

ARC-Challenge Greek, it achieves an accuracy of 49.4%, compared to Llama-3.1-8B’s and Meltemi-7B-v1.5’s 39.9% and 40.0%, respectively. Similar substantial gains are observed on the Belebele Greek dataset, where Llama-Krikri-8B-Base scores 82.7%, surpassing Meltemi-7B-v1.5 and Llama-3.1-8B by 21.7% and +9.9%, respectively. In the Greek Medical MCQA, Llama-Krikri-8B-Base reaches 53.8%, demonstrating clear advancements over Llama-3.1-8B (+20.4%) in a domain-specific Greek benchmark that was not translated from English.

English Benchmarks For the evaluation of base models on English, we utilized six benchmarks, with five of them being the original versions of those also used for Greek: ARC-Challenge (Clark et al., 2018), Truthful QA (Lin et al., 2022), Hel-laSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021), and Belebele (Bandarkar et al., 2024). Additionally, the Winogrande (Sakaguchi et al., 2021) test set was used as the sixth benchmark for English. In the results presented in Table 5 we see that our training methodology not only mitigates catastrophic forgetting effectively, but also improves average performance across all English test sets by +0.8%.

4.2 Chat Model Evaluation: Krikri-8B-Instruct

For evaluating the capabilities of Llama-Krikri-8B-Instruct as a conversational assistant, suitable for multi-turn dialogue, instruction-following and complex coding and math queries, we used a suite of benchmarks in both English and Greek. For English, we conducted evaluations across two paths:

- We submitted our model to the Open LLM Leaderboard (Fourrier et al., 2024) which automatically evaluates models on IFEval, BBH, MATH, GPQA, MUSR, and MMLU-Pro using the Eleuther AI Language Model Evaluation Harness (Gao et al., 2021), a unified framework to test generative language models on a large number of different evaluation tasks.
- We used the Arena Hard Auto v0.1 (Li et al., 2024; Chiang et al., 2024), IFEval (Zhou et al., 2023) (strict avg) and MT-Bench (Zheng et al., 2023) benchmarks. Although IFEval was already included in the Open LLM Leaderboard,

²https://huggingface.co/spaces/open-llm-leaderboard/open-llm_leaderboard

³<https://huggingface.co/collections/ilsp/ilsp-greek-evaluation-suite-6827304d5bf870d0346b02c>

Benchmark	Meltemi-7B-v1.5	Llama-3.1-8B	Krikri-8B-Base
Medical MCQA EL (15-shot)	42.2	33.4	53.8
Belebele EL (5-shot)	61.0	72.8	82.7
HellaSwag EL (10-shot)	53.8	52.1	64.6
ARC-Challenge EL (25-shot)	40.0	39.9	49.4
TruthfulQA MC2 EL (0-shot)	49.0	51.1	54.2
MMLU EL (5-shot)	41.2	42.6	52.0
Average	47.9	48.7	59.5

Table 4: Greek benchmark results (accuracy %) for base models.

Benchmark	Meltemi-7B-v1.5	Llama-3.1-8B	Krikri-8B-Base
Winogrande (5-shot)	73.4	74.6	72.6
Belebele EN (5-shot)	77.7	71.5	79.8
HellaSwag EN (10-shot)	79.6	82.0	80.7
ARC-Challenge EN (25-shot)	54.1	58.5	57.8
TruthfulQA MC2 EN (0-shot)	40.5	44.2	44.8
MMLU EN (5-shot)	56.9	66.2	65.1
Average	63.7	66.2	67.0

Table 5: English benchmark results (accuracy %) for base models.

we re-implemented it to enable accurate comparison with multiple models. In the evaluation of MT-Bench we used GPT-4o (2024-08-06) as the judge model, while in the evaluation of Arena Hard Auto v0.1 we used the standard approach with GPT-4-0314 as the baseline model (by default scoring 50%) and GPT-4-1106-Preview as the judge model, while also reusing the generations and judgments already computed by the authors.

For Greek, we created three novel evaluation benchmarks by translating three challenging, diverse, and widely used English benchmarks, ensuring high-quality translations through careful post-editing and validation:

- IFEval Greek (strict avg.): a manual translation of 541 prompts from the original Instruction-Following Evaluation benchmark (Zhou et al., 2023), featuring verifiable instructions such as "απάντησε με περισσότερες από 400 λέξεις" (answer with more than 400 words) and "ανάφερε τη λέξη ΤΝ τουλάχιστον 3 φορές" (mention the word AI at least 3 times), designed to assess the model's ability to follow specific instructions.
- MT-Bench Greek: a translated version of the Multi-turn Benchmark (Zheng et al., 2023) containing 80 high-quality, multi-turn conversations across eight diverse categories (e.g.,

STEM, humanities, roleplay, coding, etc.), carefully post-edited to ensure natural Greek phrasing and cultural appropriateness. MT-Bench is also used to evaluate the function-calling capabilities of LLMs (Chen et al., 2025). The performance of each model is calculated using LLM-as-Judge (Zheng et al., 2023) with GPT-4o (2024-08-06) serving as the scoring model.

- Arena-Hard-Auto Greek: a translated version of Arena-Hard-Auto v0.1, which originates from Chatbot Arena (Chiang et al., 2024) was included in m-ArenaHard (Dang et al., 2024) after translation with Google Translate API v3. We later post-edited using Claude Sonnet 3.5 (Anthropic, 2024) with 10-shot examples to address translation issues, particularly in coding-related prompts where some parts would best be left untranslated, as well as to retain the original style of the prompts, since some of them would be best left vaguely posed as in the original prompt. We used the version of the benchmark with style control methods for Markdown elements⁴. We set GPT-4o-Mini (2024-07-18) as the baseline model (by default 50% score) and GPT-4o (2024-08-06) as the judge model.

As shown in Table 6, Llama-Krikri-8B-Instruct

⁴<https://lmsys.org/blog/2024-08-28-style-control/>

Model	IFEval EL	IFEval EN	MT-Bench EL	MT-Bench EN
Qwen 2.5 7B	46.2	74.8	5.83	7.87
EuroLLM 9B	51.3	64.5	5.98	6.27
Aya Expanse 8B	50.4	62.2	7.68	6.92
Meltemi-7B-v1.5	32.7	41.2	6.25	5.46
Llama-3.1-8B	45.8	75.1	6.46	7.25
Llama-Krikri-8B	67.5	82.4	7.96	7.21
Gemma 2 27B IT	63.2	75.6	8.23	8.00
Aya Expanse 32B	60.3	70.2	8.27	7.40

Table 6: Greek and English evaluation results using IFEval and MT-Bench.

demonstrates exceptional performance across both Greek and English benchmarks, substantially outperforming not only its parent model Llama-3.1-8B-Instruct but also other competitive multilingual models in the 7-9B parameter range. It should be noted that the IFEval scores reported in this table reflect our own implementation of the benchmark, which may differ from the Open LLM Leaderboard implementation due to variations in prompt formatting and evaluation criteria. Despite these methodological differences, the relative performance comparisons remain valid within each implementation context.

On IFEval Greek, Llama-Krikri-8B-Instruct achieves a remarkable 67.5% accuracy, surpassing Llama-3.1-8B-Instruct by +21.7% and Meltemi-7B-v1.5 by +34.8%. Notably, our 8B model even outperforms much larger models like Gemma 2 27B IT (+4.3%) and Aya Expanse 32B (+7.2%) on this Greek instruction-following benchmark. As regards the original English IFEval, Llama-Krikri-8B-Instruct scores 82.4%, significantly higher than all other models, including those with 3-4 times more parameters. This dramatic improvement suggests that our data synthesis instruction tuning approach successfully addresses the unique challenges of following instructions in Greek, where naive translations of instruction data often fail to capture language-specific nuances.

For MT-Bench Greek, which evaluates multi-turn conversation quality, Llama-Krikri-8B-Instruct achieves a score of **7.96**, making it the top performer amongst other models in its size class. While larger models like Gemma 2 27B IT (8.23) and Aya Expanse 32B (8.27) achieve slightly higher scores on MT-Bench Greek, the margin is surprisingly small given the substantial difference in model size. On MT-Bench English, Llama-Krikri-8B maintains competitive

Task	Llama-3.1-Instr.	Krikri-8B-Instr.
IFEval	49.22	60.79
BBH	29.38	29.31
MATH	15.56	11.78
GPQA	8.72	7.05
MUSR	8.61	10.46
MMLU-PRO	31.09	25.70
Avg.	23.76	24.18

Table 7: Comparative evaluation on English benchmarks from the Open LLM Leaderboard.

performance at 7.21, essentially identical with Llama-3.1-8B-Instruct (-0.04), though understandably lower than the larger Gemma 2 27B IT (-0.79) and Aya Expanse (-0.19). This demonstrates that our instruction tuning approach can achieve a very high conversational performance on Greek conversational benchmarks, while also producing a competitive model for English benchmarks with a much more compact approach (Grattafiori et al., 2024).

As detailed in Table 7, Llama-Krikri-8B-Instruct’s official Open LLM Leaderboard submission shows an average score of 24.18% across all tests, slightly surpassing the 23.76% of Llama-3.1-8B-Instruct. The model shows particularly impressive gains on IFEval implementation (60.79% vs. 49.22%) and MUSR (10.46% vs. 8.61%), while closely matching performance on the Big Bench Hard (BBH) benchmark (29.31% vs. 29.38%). Although Llama-Krikri performs slightly below Meta-Llama-3.1-8B-Instruct in the MMLU-PRO category (25.70% vs. 31.09%), the overall performance indicates successful retention of English capabilities during the Greek-focused continual pre-training.

The results from our Arena Hard evaluations, presented in Table 8, reveal that, in the 8B parameter range, Llama-Krikri-8B-Instruct significantly outperforms its competitors, achieving a 31.8%

Model	ArenaHard EL	ArenaHard EN
Aya Expanse 8B	23.8	—
Llama 3.1 8B Instr.	4.0	19.7
Krikri 8B Instr.	31.8	35.1
Aya Expanse 32B	40.1	45.1
Gemma 2 27B IT	32.2	49.6
Llama 3.1 70B Instr.	27.4	53.9
GPT 4o Mini	50.0	65.0

Table 8: Arena Hard evaluation results (% win rate) for Greek and English.

win rate on Arena Hard Greek compared to Aya Expanse 8B’s 23.8% and Llama 3.1 8B Instruct’s 4.0% (+27.8% improvement). This demonstrates the effectiveness of our Greek-focused training approach. Even more impressively, Llama-Krikri-8B-Instruct achieves a 35.1% win rate on Arena Hard English, substantially outperforming the original Llama-3.1-8B-Instruct (19.7%) by +16.2%, despite our focus on Greek capabilities. While Aya Expanse 32B leads on Arena Hard Greek with 40.1%, our 8B model is on par with Gemma 2 27B IT (31.8% vs. 32.2%) and outperforms the 8.75 times larger Llama-3.1-70B-Instruct (27.4%) by +4.4% on the Greek evaluation data. On the original English Arena Hard, the larger models generally perform better, although it should be noted that Llama-Krikri-8B-Instruct outperforms Llama-3.1-8B-Instruct by +15.4% (35.1% vs. 19.7%).

Please note that while all models trail behind GPT-4o-Mini (used as baseline on the Greek Arena Hard), recent research (Li et al., 2025) has shown that judge models are biased towards student models, i.e., models finetuned on distilled data from the stronger & larger teacher model which also acts as a judge. While details on the post-training data of GPT-4o-Mini are undisclosed, it would be very reasonable to assume that it has been trained -at least partly- with GPT-4 and GPT-4o serving as teacher models and, therefore, that the judges that we utilized are biased towards it compared to all other evaluated models.

This performance comparison with much larger models highlights the efficiency of our approach since Llama-Krikri-8B-Instruct achieves comparable or even superior performance on Greek benchmarks compared to models with 3-4x more parameters, while maintaining strong English capabilities. This efficiency is particularly important for deployment scenarios where computational resources may be limited, demonstrating that a carefully trained smaller model can rival much larger ones for specific languages.

Apart from the comparative evaluations men-

tioned above, we have performed zeroshot MT experiments on an Ancient-Modern Greek (grc \leftrightarrow ell) translation dataset⁵ that includes 100 sentences of Ancient Greek texts manually translated into Modern Greek. Using Llama-Krikri-8B-Instruct we have observed a 54.66 BLEU score for the Ancient to Modern Greek (grc \rightarrow ell) translation direction, with the reverse direction (ell \rightarrow grc) being more challenging (20.41 BLEU).

5 Discussion and Conclusions

In this paper, we presented Llama-Krikri-8B, a new LLM that exhibits significant skills in understanding and generating Greek, while also showing highly accurate handling of text in English and historical Greek dialects. We achieved this by developing an efficient tokenizer that exhibits a low token/words fertility for Greek and by further training Llama 3.1-8B using a carefully constructed dataset that covered a wide variety of domains. In evaluation experiments on a benchmark suite comprising Greek and English datasets, we have observed that Llama-Krikri-8B performs significantly better in Greek (+10.8%) compared to its base model, while also showing gains in English (+0.8%). We then created Llama-Krikri-8B-Instruct, a version designed for following instructions and engaging in helpful conversations. This involved a multi-step process that comprised synthetic data generation in a multitude of domains, fine-tuning the model and then aligning it with human preferences. Evaluations revealed that Llama-Krikri-8B-Instruct significantly outperformed Llama-3.1-8B-Instruct in both Greek (+21.7%) and English (+7.3%) IFEval. Our model also demonstrated highly competitive chat abilities in both languages across several benchmarks.

6 Limitations

The quality and accessibility of Greek datasets are critical to the development of Krikri. Greek open-source corpora are becoming more numerous, but they might not be as large or varied as datasets for more extensively spoken languages, like English. This may result in biases in the model’s understanding of Greek, especially with regard to regional variances, dialects, and specialized fields like technical fields, law, or medicine.

⁵https://huggingface.co/datasets/ilsp/ancient-modern_greek_translations

As an 8B parameter model, our model shows a fairly high level of Greek fluency, but it is less effective than larger-class and commercial models at reasoning and instruction following, and is more likely to experience hallucinations.

In the future, our evaluation benchmarks should include more original Greek LLM datasets that are not the result of machine translation and post-editing. These datasets will help minimize the effect of machine translation on evaluation results and also better reflect the target language and culture.

7 Risks and ethical considerations

To mitigate potential risks, we took several steps to ensure the data used for training did not contain personally identifiable information, offensive, or otherwise inappropriate content. We sourced data from publicly available, licensed, or open-access datasets, ensuring compliance with their respective policies and any flagged data points were excluded. We did not collect data from private communications or data sources that could contain personally identifiable information. We have also given special care to align our model’s responses with safety guidelines followed by manual reviews.

However, we have not performed a systematic evaluation against LLM risks including risks related to discrimination, hate speech and exclusion; information hazards; and misinformation harms (Weidinger et al., 2022). As this was due to lack of relevant evaluation material in Greek, we want to contribute towards improving this situation in the future.

We recognize that these measures are not a substitute for more thorough evaluation protocols. Moving forward, we aim to contribute towards addressing these limitations by promoting the development of Greek-language evaluation resources for LLM risks. This will enable more robust and contextually appropriate assessments of ethical risks in future models.

References

Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang, Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, et al. 2024. A survey on data selection for language models. *arXiv preprint arXiv:2402.16827*.

Anton Alexandrov, Veselin Raychev, Dimitar I. Dimitrov, Ce Zhang, Martin Vechev, and Kristina

Toutanova. 2024. [Bggpt 1.0: Extending english-centric llms to other languages](#). *Preprint*, arXiv:2412.10893.

Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgren, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin Raffel, Leandro von Werra, and Thomas Wolf. 2025. [SmolLM2: When smol goes big – data-centric training of a small language model](#). *Preprint*, arXiv:2502.02737.

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

AI Anthropic. 2024. [The Claude 3 model family: Opus, sonnet, haiku](#). *Claude-3 Model Card*.

Mikel Artetxe and Holger Schwenk. 2018. Margin-based parallel corpus mining with multilingual sentence embeddings. *arXiv preprint arXiv:1811.01136*.

Mikel Artetxe and Holger Schwenk. 2019. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions of the association for computational linguistics*, 7:597–610.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022. [Constitutional ai: Harmlessness from ai feedback](#). *Preprint*, arXiv:2212.08073.

Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. 2024. The belebele benchmark: a parallel reading comprehension dataset in 122 language variants. In *Proceedings of the 62nd Annual*

765	<i>Meeting of the Association for Computational Lin-</i>	Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anasta-	823
766	<i>guistics (Volume 1: Long Papers)</i> , pages 749–775.	sios Nikolas Angelopoulos, Tianle Li, Dacheng Li,	824
767	Marta Bañón, Miquel Esplà-Gomis, Mikel L. For-	Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E	825
768	cada, Cristian García-Romero, Taja Kuzman, Nikola	Gonzalez, et al. 2024. Chatbot arena: An open plat-	826
769	Ljubešić, Rik van Noord, Leopoldo Pla Sempere,	form for evaluating llms by human preference. In <i>In-</i>	827
770	Gema Ramírez-Sánchez, Peter Rupnik, Vít Su-	<i>ternational Conference on Machine Learning</i> , pages	828
771	chomel, Antonio Toral, Tobias van der Werff, and	8359–8388. PMLR.	829
772	Jaume Zaragoza. 2022. MaCoCu: Massive collec-	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,	830
773	tion and curation of monolingual and bilingual data:	Ashish Sabharwal, Carissa Schoenick, and Oyvind	831
774	focus on under-resourced languages . In <i>Proceedings</i>	Tafjord. 2018. Think you have Solved Question An-	832
775	<i>of the 23rd Annual Conference of the European As-</i>	swering? Try ARC, the AI2 Reasoning Challenge .	833
776	<i>sociation for Machine Translation</i> , pages 303–304,	<i>arXiv:1803.05457v1</i> .	834
777	Ghent, Belgium. European Association for Machine	Zoltan Csaki, Pian Pawakapan, Urmish Thakker, and	835
778	Translation.	Qiantong Xu. 2023. Efficiently adapting pretrained	836
779	A Broder. 1997. On the resemblance and containment	language models to new languages. <i>arXiv preprint</i>	837
780	of documents. In <i>Proceedings of the Compression</i>	<i>arXiv:2311.05741</i> .	838
781	<i>and Complexity of Sequences 1997</i> , page 21.	Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao,	839
782	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and	840
783	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	Maosong Sun. 2023. Ultrafeedback: Boosting lan-	841
784	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	guage models with high-quality feedback . <i>Preprint</i> ,	842
785	Askeel, Sandhini Agarwal, Ariel Herbert-Voss,	<i>arXiv:2310.01377</i> .	843
786	Gretchen Krueger, Tom Henighan, Rewon Child,	John Dang, Shivalika Singh, Daniel D’souza, Arash	844
787	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	Ahmadian, Alejandro Salamanca, Madeline Smith,	845
788	Clemens Winter, Christopher Hesse, Mark Chen,	Aidan Peppin, Sungjin Hong, Manoj Govindassamy,	846
789	Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin	Terrence Zhao, Sandra Kublik, Meor Amer, Viraat	847
790	Chess, Jack Clark, Christopher Berner, Sam Mc-	Aryabumi, Jon Ander Campos, Yi-Chern Tan, Tom	848
791	Candlish, Alec Radford, Ilya Sutskever, and Dario	Kocmi, Florian Strub, Nathan Grinsztajn, Yannic	849
792	Amodei. 2020. Language models are few-shot learn-	Flet-Berliac, Acyr Locatelli, Hangyu Lin, Dwarak	850
793	ers . <i>Preprint</i> , arXiv:2005.14165.	Talupuru, Bharat Venkitesh, David Cairuz, Bowen	851
794	Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,	Yang, Tim Chung, Wei-Yin Ko, Sylvie Shang Shi,	852
795	Katherine Lee, Florian Tramèr, and Chiyuan Zhang.	Amir Shukayev, Sammie Bae, Aleksandra Piktus, Ro-	853
796	2023. Quantifying memorization across neural lan-	man Castagné, Felipe Cruz-Salinas, Eddie Kim, Lu-	854
797	guage models. In <i>The Eleventh International Confer-</i>	cas Crawhall-Stein, Adrien Morisot, Sudip Roy, Phil	855
798	<i>ence on Learning Representations</i> . OpenReview.	Blunsom, Ivan Zhang, Aidan Gomez, Nick Frosst,	856
799	Ilias Chalkidis, Manos Fergadiotis, and Ion Androut-	Marzieh Fadaee, Beyza Ermiş, Ahmet Üstün, and	857
800	sopoulos. 2021. Multieurlex – a multi-lingual and	Sara Hooker. 2024. Aya expanse: Combining re-	858
801	multi-label legal document classification dataset for	search breakthroughs for a new multilingual frontier .	859
802	zero-shot cross-lingual transfer . In <i>Proceedings of</i>	<i>Preprint</i> , arXiv:2412.04261.	860
803	<i>the 2021 Conference on Empirical Methods in Natu-</i>	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	861
804	<i>ral Language Processing</i> . Association for Computa-	Kristina Toutanova. 2019. Bert: Pre-training of deep	862
805	tional Linguistics.	bidirectional transformers for language understand-	863
806	Ilias Chalkidis, Manos Fergadiotis, Prodromos Malaka-	ing . <i>Preprint</i> , arXiv:1810.04805.	864
807	sotis, and Ion Androutsopoulos. 2019. Large-scale	Tomaž Erjavec, Maciej Ogródniczuk, Petya Osenova,	865
808	multi-label text classification on EU legislation . In	Nikola Ljubešić, Kiril Simov, Andrej Pančur, Michal	866
809	<i>Proceedings of the 57th Annual Meeting of the Asso-</i>	Rudolf, Matyáš Kopp, Starkađur Barkarson, Steinþór	867
810	<i>ciation for Computational Linguistics</i> , pages 6314–	Steingrímsson, Çağrı Çöltekin, Jesse de Does, Ka-	868
811	6322, Florence, Italy. Association for Computational	trien Depuydt, Tommaso Agnoloni, Giulia Ven-	869
812	Linguistics.	turi, María Calzada Pérez, Luciana D. de Macedo,	870
813	Yi-Chang Chen, Po-Chun Hsu, Chan-Jan Hsu, and Da-	Costanza Navarretta, Giancarlo Luxardo, Matthew	871
814	shan Shiu. 2025. Enhancing function-calling capa-	Coole, Paul Rayson, Vaidas Morkevičius, Tomas	872
815	bilities in LLMs: Strategies for prompt formats, data	Krilavičius, Roberts Darundėdis, Orsolya Ring,	873
816	integration, and multilingual translation . In <i>Proceed-</i>	Ruben van Heusden, Maarten Marx, and Darja Fišer.	874
817	<i>ings of the 2025 Conference of the Nations of the</i>	2022. The ParlaMint corpora of parliamentary pro-	875
818	<i>Americas Chapter of the Association for Computa-</i>	ceedings . <i>Lang. Resour. Eval.</i> , 57(1):415–448.	876
819	<i>tional Linguistics: Human Language Technologies</i>	Clémentine Fourrier, Nathan Habib, Alina Lozovskaya,	877
820	<i>(Volume 3: Industry Track)</i> , pages 99–111, Albu-	Konrad Szafer, and Thomas Wolf. 2024. Open	878
821	querque, New Mexico. Association for Computa-	llm leaderboard v2. https://huggingface.	879
822	tional Linguistics.	co/spaces/open-llm-leaderboard/open_llm_	880
		leaderboard .	881

882	Leo Gao, Jonathan Tow, Stella Biderman, Sid Black,	Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj	944
883	Anthony DiPofi, Charles Foster, Laurence Golding,	Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,	945
884	Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff,	Robert Stojnic, Roberta Raileanu, Rohan Maheswari,	946
885	Jason Phang, Laria Reynolds, Eric Tang, Anish Thite,	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-	947
886	Ben Wang, Kevin Wang, and Andy Zou. 2021. A	nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan	948
887	framework for few-shot language model evaluation.	Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-	949
		hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-	950
888	Maria Gavrilidou, Stelios Piperidis, Dimitrios Galanis,	hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-	951
889	Kanella Pouli, Penny Labropoulou, Juli Bakagianni,	ran Narang, Sharath Raparthy, Sheng Shen, Shengye	952
890	Iro Tsiouli, Miltos Deligiannis, Athanasia Kolovou,	Wan, Shruti Bhosale, Shun Zhang, Simon Van-	953
891	Dimitris Gkoumas, Leon Voukoutis, and Katerina	denhende, Soumya Batra, Spencer Whitman, Sten	954
892	Gkirtzou. 2023. The CLARIN:EL infrastructure:	Sootla, Stephane Collot, Suchin Gururangan, Syd-	955
893	Platform, Portal, K-Centre. In <i>Selected papers from</i>	ney Borodinsky, Tamar Herman, Tara Fowler, Tarek	956
894	<i>the CLARIN Annual Conference 2023.</i>	Sheasha, Thomas Georgiou, Thomas Scialom, Tobias	957
		Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal	958
895	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh	959
896	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-	960
897	Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-	ginie Do, Vish Vogeti, Vítor Albiero, Vladan Petro-	961
898	ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh	vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-	962
899	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-	ney Meers, Xavier Martinet, Xiaodong Wang, Xi-	963
900	tra, Archie Sravankumar, Artem Korenev, Arthur	aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-	964
901	Hinsvark, Arun Rao, Aston Zhang, Aurelien Ro-	feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-	965
902	driguez, Austen Gregerson, Ava Spataru, Baptiste	schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen,	966
903	Roziere, Bethany Biron, Binh Tang, Bobbie Chern,	Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,	967
904	Charlotte Caucheteux, Chaya Nayak, Chloe Bi,	Zacharie Delpierre Coudert, Zheng Yan, Zhengxing	968
905	Chris Marra, Chris McConnell, Christian Keller,	Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-	969
906	Christophe Touret, Chunyang Wu, Corinne Wong,	vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,	970
907	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,	971
908	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei	972
909	Danny Wyatt, David Esiobu, Dhruv Choudhary,	Baevski, Allie Feinstein, Amanda Kallet, Amit San-	973
910	Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,	gani, Amos Teo, Anam Yunus, Andrei Lupu, And-	974
911	Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,	res Alvarado, Andrew Caples, Andrew Gu, Andrew	975
912	Elina Lobanova, Emily Dinan, Eric Michael Smith,	Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-	976
913	Filip Radenovic, Francisco Guzmán, Frank Zhang,	dani, Annie Dong, Annie Franco, Anuj Goyal, Apar-	977
914	Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis An-	jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	978
915	derson, Govind Thattai, Graeme Nail, Gregoire Mi-	Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-	979
916	alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,	dan, Beau James, Ben Maurer, Benjamin Leonhardi,	980
917	Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	981
918	Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-	982
919	han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,	cock, Bram Wasti, Brandon Spence, Brani Stojkovic,	983
920	Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,	Brian Gamido, Britt Montalvo, Carl Parker, Carly	984
921	Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,	Burton, Catalina Mejia, Ce Liu, Changhan Wang,	985
922	Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,	Changkyu Kim, Chao Zhou, Chester Hu, Ching-	986
923	Jianfeng Chi, Jianyu Huang, Jiawen Li, Jie Wang,	Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-	987
924	Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,	ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,	988
925	Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-	Daniel Kreymer, Daniel Li, David Adkins, David	989
926	teng Jia, Kalyan Vasuden Alwala, Karthik Prasad,	Xu, Davide Testuggine, Delia David, Devi Parikh,	990
927	Kartikkeya Upasani, Kate Plawiak, Ke Li, Kenneth	Diana Liskovich, Didem Foss, Dingkan Wang, Duc	991
928	Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	992
929	Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal	Elaine Montgomery, Eleonora Presani, Emily Hahn,	993
930	Lakhotia, Lauren Rantala-Yearly, Laurens van der	Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-	994
931	Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,	ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	995
932	Louis Martin, Lovish Madaan, Lubo Malo, Lukas	Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat	996
933	Blecher, Lukas Landzaat, Luke de Oliveira, Madeline	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	997
934	Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar	Seide, Gabriela Medina Florez, Gabriella Schwarz,	998
935	Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew	Gada Badeer, Georgia Swee, Gil Halpern, Grant	999
936	Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-	Herman, Grigory Sizov, Guangyi, Zhang, Guna	1000
937	badur, Mike Lewis, Min Si, Mitesh Kumar Singh,	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	1001
938	Mona Hassan, Naman Goyal, Narjes Torabi, Niko-	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	1002
939	lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	1003
940	Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	1004
941	Alrassy, Pengchuan Zhang, Pengwei Li, Petar Va-	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,	1005
942	sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,	Irina-Elena Veliche, Itai Gat, Jake Weissman, James	1006
943	Praveen Krishnan, Punit Singh Koura, Puxin Xu,	Geboski, James Kohli, Janice Lam, Japhet Asher,	1007

1008	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783.	
1009	and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks . <i>Preprint</i> , arXiv:2004.10964.	1071
1010		1072
1011		1073
1012	Kenneth Heafield. 2011. KenLM: Faster and smaller language model queries . In <i>Proceedings of the Sixth Workshop on Statistical Machine Translation</i> , pages 187–197, Edinburgh, Scotland. Association for Computational Linguistics.	1074
1013		1075
1014		1076
1015		1077
1016		1078
1017		
1018	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding . In <i>International Conference on Learning Representations</i> .	1079
1019		1080
1020		1081
1021		1082
1022		1083
1023		
1024	Pin-Lun Hsu, Yun Dai, Vignesh Kothapalli, Qingquan Song, Shao Tang, Siyu Zhu, Steven Shimizu, Shivam Sahni, Haowen Ning, and Yanning Chen. 2024. Liger kernel: Efficient triton kernels for llm training . <i>arXiv preprint arXiv:2410.10989</i> .	1084
1025		1085
1026		1086
1027		1087
1028		1088
1029		
1030	Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L Richter, Quentin Anthony, Timothée Lesort, Eugene Belilovsky, and Irina Rish. 2024. Simple and scalable strategies to continually pre-train large language models. <i>arXiv preprint arXiv:2403.08763</i> .	1089
1031		1090
1032		1091
1033		1092
1034		1093
1035		
1036	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	1094
1037		1095
1038		1096
1039		1097
1040		1098
1041		
1042	Armand Joulin, Édouard Grave, Piotr Bojanowski, and Tomáš Mikolov. 2017. Bag of tricks for efficient text classification. In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers</i> , pages 427–431.	1099
1043		1100
1044		1101
1045		1102
1046		1103
1047		1104
1048	LAION. 2023. LeoLM: Igniting German-Language LLM Research. https://laion.ai/blog/leo-lm/ . Accessed: (12 July 2024).	1105
1049		1106
1050		1107
1051		
1052	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. 2025. Tulu 3: Pushing frontiers in open language model post-training . <i>Preprint</i> , arXiv:2411.15124.	1108
1053		1109
1054		1110
1055		1111
1056		1112
1057		1113
1058		1114
1059		1115
1060		1116
1061		1117
1062		
1063	Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. Deduplicating training data makes language models better. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8424–8445.	1118
1064		1119
1065		1120
1066		1121
1067		1122
1068		1123
		1124
1069	Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,	
1070		
		1125
		1126
		1127

1128	Dawei Li, Renliang Sun, Yue Huang, Ming Zhong,	cleaned, enormous, and multilingual dataset for large	1184
1129	Bohan Jiang, Jiawei Han, Xiangliang Zhang, Wei	language models in 167 languages. <i>arXiv preprint</i>	1185
1130	Wang, and Huan Liu. 2025. Preference leakage: A	<i>arXiv:2309.09400</i> .	1186
1131	contamination problem in llm-as-a-judge . <i>Preprint</i> ,		
1132	<i>arXiv:2502.01534</i> .		
1133	Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap,	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-	1187
1134	Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and	roll L. Wainwright, Pamela Mishkin, Chong Zhang,	1188
1135	Ion Stoica. 2024. From crowdsourced data to high-	Sandhini Agarwal, Katarina Slama, Alex Ray, John	1189
1136	quality benchmarks: Arena-hard and benchbuilder	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	1190
1137	pipeline. <i>arXiv preprint arXiv:2406.11939</i> .	Maddie Simens, Amanda Aspell, Peter Welinder,	1191
1138		Paul Christiano, Jan Leike, and Ryan Lowe. 2022.	1192
1139	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	Training language models to follow instructions with	1193
1140	TruthfulQA: Measuring how models mimic human	human feedback . <i>Preprint</i> , <i>arXiv:2203.02155</i> .	1194
1141	falsehoods . In <i>Proceedings of the 60th Annual Meet-</i>		
1142	<i>ing of the Association for Computational Linguistics</i>	Christos Papaloukas, Ilias Chalkidis, Konstantinos	1195
1143	<i>(Volume 1: Long Papers)</i> , pages 3214–3252, Dublin,	Athinaios, Despina-Athanasia Pantazi, and Manolis	1196
1144	Ireland. Association for Computational Linguistics.	Koubarakis. 2021. Multi-granular legal topic clas-	1197
1145	Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Ju-	sification on greek legislation . In <i>Proceedings of</i>	1198
1146	jie He, Chaojie Wang, Shuicheng Yan, Yang Liu,	<i>the Natural Legal Language Processing Workshop</i>	1199
1147	and Yahui Zhou. 2024. Skywork-reward: Bag of	2021, pages 63–75, Punta Cana, Dominican Republic.	1200
1148	tricks for reward modeling in llms. <i>arXiv preprint</i>	Association for Computational Linguistics.	1201
1149	<i>arXiv:2410.18451</i> .		
1150	Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kin-	Vassilis Papavassiliou, Sokratis Sofianopoulos, Prokopis	1202
1151	ney, and Daniel S Weld. 2020. S2orc: The semantic	Prokopidis, and Stelios Piperidis. 2018. The	1203
1152	scholar open research corpus. In <i>Proceedings of the</i>	ILSP/ARC submission to the WMT 2018 parallel	1204
1153	<i>58th Annual Meeting of the Association for Compu-</i>	corpus filtering shared task . In <i>Proceedings of the</i>	1205
1154	<i>tational Linguistics</i> , pages 4969–4983.	<i>Third Conference on Machine Translation: Shared</i>	1206
1155	Andrea Lösch, Valérie Mapelli, Khalid Choukri, Maria	<i>Task Papers</i> , pages 928–933, Belgium, Brussels. As-	1207
1156	Giagkou, Stelios Piperidis, Prokopis Prokopidis,	sociation for Computational Linguistics.	1208
1157	Vassilis Papavassiliou, Miltos Deligiannis, Aivars	Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov,	1209
1158	Berzins, Andrejs Vasiljevs, et al. 2021. Collection	Margaret Mitchell, Colin Raffel, Leandro Von Werra,	1210
1159	and Curation of Language Data within the European	Thomas Wolf, et al. 2024. The FineWeb Datasets:	1211
1160	Language Resource Coordination (ELRC) . In <i>Qura-</i>	Decanting the Web for the Finest Text Data at Scale.	1212
1161	<i>tor</i> .	<i>arXiv preprint arXiv:2406.17557</i> .	1213
1162	Ilya Loshchilov and Frank Hutter. 2017. Fixing	Ramon Pires, Hugo Abonizio, Thales Sales Almeida,	1214
1163	Weight Decay Regularization in Adam . <i>CoRR</i> ,	and Rodrigo Nogueira. 2023. Sabia: Portuguese	1215
1164	abs/1711.05101.	large language models . In <i>Intelligent Systems</i> , pages	1216
1165	Anton Lozhkov, Raymond Li, Loubna Ben Allal, Fed-	226–240, Cham. Springer Nature Switzerland.	1217
1166	erico Cassano, Joel Lamy-Poirier, Nouamane Tazi,	Ye Qi, Devendra Sachan, Matthieu Felix, Sarguna Pad-	1218
1167	Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei,	manabhan, and Graham Neubig. 2018. When and	1219
1168	et al. 2024. Starcoder 2 and the stack v2: The next	why are pre-trained word embeddings useful for neu-	1220
1169	generation. <i>arXiv preprint arXiv:2402.19173</i> .	ral machine translation? In <i>Proceedings of the 2018</i>	1221
1170	Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou,	<i>Conference of the North American Chapter of the</i>	1222
1171	and Yue Zhang. 2025. An empirical study of catas-	<i>Association for Computational Linguistics: Human</i>	1223
1172	trophic forgetting in large language models during	<i>Language Technologies, Volume 2 (Short Papers)</i> ,	1224
1173	continual fine-tuning . <i>Preprint</i> , <i>arXiv:2308.08747</i> .	pages 529–535, New Orleans, Louisiana. Associa-	1225
1174	Pedro Henrique Martins, Patrick Fernandes, João Alves,	tion for Computational Linguistics.	1226
1175	Nuno M Guerreiro, Ricardo Rei, Duarte M Alves,	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano	1227
1176	José Pombal, Amin Farajian, Manuel Faysse, Ma-	Ermon, Christopher D. Manning, and Chelsea Finn.	1228
1177	teusz Klimaszewski, et al. 2024. Eurollm: Multilin-	2024. Direct preference optimization: Your lan-	1229
1178	gual language models for europe. In <i>Proceedings of</i>	guage model is secretly a reward model . <i>Preprint</i> ,	1230
1179	<i>the Ninth Conference on Machine Translation</i> , pages	<i>arXiv:2305.18290</i> .	1231
1180	1393–1409. Association for Computational Linguis-		
1181	tics.	Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro,	1232
1182	Thuat Nguyen, Chien Van Nguyen, Viet Dac Lai, Hieu	Chrysoula Zerva, Ana C Farinha, Christine Maroti,	1233
1183	Man, Nghia Trung Ngo, Franck Dernoncourt, Ryan A	José G. C. de Souza, Taisiya Glushkova, Duarte	1234
	Rossi, and Thien Huu Nguyen. 2023. CulturaX: A	Alves, Luisa Coheur, Alon Lavie, and André F. T.	1235
		Martins. 2022. CometKiwi: IST-unbabel 2022 sub-	1236
		mission for the quality estimation shared task . In	1237
		<i>Proceedings of the Seventh Conference on Machine</i>	1238
		<i>Translation (WMT)</i> , pages 634–645, Abu Dhabi,	1239
		United Arab Emirates (Hybrid). Association for Com-	1240
		putational Linguistics.	1241

1242	Morgane Riviere, Shreya Pathak, Pier Giuseppe	arXiv:2408.00118.	1306
1243	Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard		
1244	Hussenot, Thomas Mesnard, Bobak Shahriari,	Richard J Roberts. 2001. Pubmed central: The genbank	1307
1245	Alexandre Ramé, Johan Ferret, Peter Liu, Pouya	of the published literature.	1308
1246	Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos,		
1247	Ravin Kumar, Charline Le Lan, Sammy Jerome, An-	Dimitrios Roussis and Vassilis Papavassiliou. 2022. The	1309
1248	ton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan	ARC-NKUA submission for the English-Ukrainian	1310
1249	Girgin, Nikola Momchev, Matt Hoffman, Shantanu	general machine translation shared task at WMT22.	1311
1250	Thakoor, Jean-Bastien Grill, Behnam Neyshabur,	In <i>Proceedings of the Seventh Conference on Ma-</i>	1312
1251	Olivier Bachem, Alanna Walton, Aliaksei Severyn,	<i>chine Translation (WMT)</i> , pages 358–365, Abu	1313
1252	Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin	Dhabi, United Arab Emirates (Hybrid). Association	1314
1253	Abdagic, Amanda Carl, Amy Shen, Andy Brock,	for Computational Linguistics.	1315
1254	Andy Coenen, Anthony Laforge, Antonia Pater-		
1255	son, Ben Bastian, Bilal Piot, Bo Wu, Brandon	Dimitrios Roussis, Vassilis Papavassiliou, Prokopis	1316
1256	Royal, Charlie Chen, Chintu Kumar, Chris Perry,	Prokopidis, Stelios Piperidis, and Vassilis Katsouris.	1317
1257	Chris Welly, Christopher A. Choquette-Choo, Danila	2022a. SciPar: A collection of parallel corpora	1318
1258	Sinopalnikov, David Weinberger, Dimple Vijayku-	from scientific abstracts. In <i>Proceedings of the Thir-</i>	1319
1259	mar, Dominika Rogozińska, Dustin Herbison, Elisa	<i>teenth Language Resources and Evaluation Confer-</i>	1320
1260	Bandy, Emma Wang, Eric Noland, Erica Moreira,	<i>ence</i> , pages 2652–2657, Marseille, France. European	1321
1261	Evan Senter, Evgenii Eltyshv, Francesco Visin,	Language Resources Association.	1322
1262	Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus		
1263	Martins, Hadi Hashemi, Hanna Klimczak-Plucińska,	Dimitrios Roussis, Vassilis Papavassiliou, Sokratis	1323
1264	Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda	Sofianopoulos, Prokopis Prokopidis, and Stelios	1324
1265	Mein, Jack Zhou, James Svensson, Jeff Stanway,	Piperidis. 2022b. Constructing parallel corpora from	1325
1266	Jetha Chan, Jin Peng Zhou, Joana Carrasqueira,	COVID-19 news using MediSys metadata. In <i>Pro-</i>	1326
1267	Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost	<i>ceedings of the Thirteenth Language Resources and</i>	1327
1268	van Amersfoort, Josh Gordon, Josh Lipschultz, Josh	<i>Evaluation Conference</i> , pages 1068–1072, Marseille,	1328
1269	Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya	France. European Language Resources Association.	1329
1270	Badola, Kat Black, Katie Millican, Keelin McDonell,		
1271	Kelvin Nguyen, Kiranbir Sodhia, Kish Greene,	Dimitris Roussis, Sokratis Sofianopoulos, and Stelios	1330
1272	Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre,	Piperidis. 2024. Enhancing scientific discourse: Ma-	1331
1273	Lena Heuermann, Leticia Lago, Lilly McNealus,	chine translation for the scientific domain. In <i>Pro-</i>	1332
1274	Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon,	<i>ceedings of the 25th Annual Conference of the Euro-</i>	1333
1275	Luciano Martins, Machel Reid, Manvinder Singh,	<i>pean Association for Machine Translation (Volume</i>	1334
1276	Mark Iverson, Martin Görner, Mat Velloso, Mateo	<i>1)</i> , pages 275–285.	1335
1277	Wirth, Matt Davidow, Matt Miller, Matthew Rahtz,		
1278	Matthew Watson, Meg Risdal, Mehran Kazemi,	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-	1336
1279	Michael Moynihan, Ming Zhang, Minsuk Kahng,	ula, and Yejin Choi. 2021. Winogrande: An adver-	1337
1280	Minwoo Park, Mofi Rahman, Mohit Khatwani, Na-	sarial winograd schema challenge at scale. <i>Commu-</i>	1338
1281	talie Dao, Nenshad Bardoliwalla, Nesh Devanathan,	<i>nications of the ACM</i> , 64(9):99–106.	1339
1282	Neta Dumai, Nilay Chauhan, Oscar Wahltinez,		
1283	Pankil Botarda, Parker Barnes, Paul Barham, Paul	Víctor M. Sánchez-Cartagena, Marta Bañón, Sergio	1340
1284	Michel, Pengchong Jin, Petko Georgiev, Phil Culli-	Ortiz-Rojas, and Gema Ramírez-Sánchez. 2018.	1341
1285	ton, Pradeep Kuppala, Ramona Comanescu, Ramona	Prompsit’s submission to WMT 2018 Parallel Cor-	1342
1286	Merhej, Reena Jana, Reza Ardeshtir Rokni, Rishabh	pus Filtering shared task. In <i>Proceedings of the</i>	1343
1287	Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc	<i>Third Conference on Machine Translation, Volume 2:</i>	1344
1288	Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R.	<i>Shared Task Papers</i> , Brussels, Belgium. Association	1345
1289	Arnold, Sebastian Krause, Shengyang Dai, Shruti	for Computational Linguistics.	1346
1290	Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Tim-		
1291	othy Jordan, Ting Yu, Tom Eccles, Tom Hennigan,	Liping Tang, Nikhil Ranjan, Omkar Pangarkar, Xuezhi	1347
1292	Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Ya-	Liang, Zhen Wang, Li An, Bhaskar Rao, Linghao Jin,	1348
1293	dav, Vilobh Meshram, Vishal Dharmadhikari, War-	Huijuan Wang, Zhoujun Cheng, Suqi Sun, Cun Mu,	1349
1294	ren Barkley, Wei Wei, Wenming Ye, Woohyun Han,	Victor Miller, Xuezhe Ma, Yue Peng, Zhengzhong	1350
1295	Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong,	Liu, and Eric P. Xing. 2024. Txt360: A top-quality	1351
1296	Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand	llm pre-training dataset requires the perfect blend.	1352
1297	Rao, Minh Giang, Ludovic Peran, Tris Warkentin,		
1298	Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia	Jörg Tiedemann. 2012. Parallel data, tools and inter-	1353
1299	Hadsell, D. Sculley, Jeanine Banks, Anca Dragan,	faces in OPUS. In <i>Proceedings of the Eighth In-</i>	1354
1300	Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hass-	<i>ternational Conference on Language Resources and</i>	1355
1301	abis, Koray Kavukcuoglu, Clement Farabet, Elena	<i>Evaluation (LREC’12)</i> , pages 2214–2218, Istanbul,	1356
1302	Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Ar-	Turkey. European Language Resources Association	1357
1303	mand Joulin, Kathleen Kenealy, Robert Dadashi,	(ELRA).	1358
1304	and Alek Andreev. 2024. Gemma 2: Improving		
1305	open language models at a practical size. <i>Preprint</i> ,	Yury Tokpanov, Beren Millidge, Paolo Glorioso,	1359
		Jonathan Pilault, Adam Ibrahim, James Whitting-	1360
		ton, and Quentin Anthony. 2024. ZydA: A 1.3 T	1361

1362	Dataset for Open Language Modeling. <i>arXiv preprint</i>	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	1418
1363	<i>arXiv:2406.01981</i> .	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	1419
1364	Santosh Tyss, Rashid Haddad, and Matthias Grabmair.	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023.	1420
1365	2024. Ecthr-pcr: A dataset for precedent understand-	Judging llm-as-a-judge with mt-bench and chatbot	1421
1366	ing and prior case retrieval in the european court	arena. <i>Advances in Neural Information Processing</i>	1422
1367	of human rights. In <i>Proceedings of the 2024 Joint</i>	<i>Systems</i> , 36:46595–46623.	1423
1368	<i>International Conference on Computational Linguis-</i>		
1369	<i>tics, Language Resources and Evaluation (LREC-</i>	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sid-	1424
1370	<i>COLING 2024</i>), pages 5473–5483.	dhārtha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,	1425
1371	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	and Le Hou. 2023. Instruction-following evalu-	1426
1372	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	ation for large language models. <i>arXiv preprint</i>	1427
1373	Kaiser, and Illia Polosukhin. 2023. Attention is all	<i>arXiv:2311.07911</i> .	1428
1374	you need . <i>Preprint</i> , arXiv:1706.03762.		
1375	Leon Voukoutis, Dimitris Roussis, Georgios	A Pretrained data mix	1429
1376	Paraskevopoulos, Sokratis Sofianopoulos, Prokopis	The pretraining data mix (Table 1) contains the	1430
1377	Prokopidis, Vassilis Papavasileiou, Athanasios	following four components:	1431
1378	Katsamanis, Stelios Piperidis, and Vassilis Katsouros.		
1379	2024. Meltemi: The first open large language model	1. Greek Texts (56.7B tokens): The Greek	1432
1380	for greek . <i>Preprint</i> , arXiv:2407.20743.	part of the dataset was sourced from publicly	1433
1381	Laura Weidinger, Jonathan Uesato, Maribeth Rauh,	available resources spanning a wide range	1434
1382	Conor Griffin, Po-Sen Huang, John Mellor, Amelia	of domains and sources: Wikipedia, ELRC-	1435
1383	Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh,	SHARE (Lösch et al., 2021), EUR-LEX &	1436
1384	Courtney Biles, Sasha Brown, Zac Kenton, Will	MultiEUR-LEX (Chalkidis et al., 2019, 2021),	1437
1385	Hawkins, Tom Stepleton, Abeba Birhane, Lisa Anne	MaCoCu (Bañón et al., 2022), CLARIN-EL	1438
1386	Hendricks, Laura Rimell, William Isaac, Julia Haas,	(Gavrilidou et al., 2023), EMEA ⁶ , parlia-	1439
1387	Sean Legassick, Geoffrey Irving, and Iason Gabriel.	mentary proceedings (Erjavec et al., 2022),	1440
1388	2022. Taxonomy of risks posed by language models .	⁷ , governmental and legal documents from	1441
1389	In <i>Proceedings of the 2022 ACM Conference on Fair-</i>	the Greek Government Gazette via the Na-	1442
1390	<i>ness, Accountability, and Transparency, FAccT ’22</i> ,	tional Printing House ⁸ , the Permanent Greek	1443
1391	page 214–229, New York, NY, USA. Association for	Legislation Code – Raptarchis dataset ⁹ (Pa-	1444
1392	Computing Machinery.	papaloukas et al., 2021), Greek School Books ¹⁰ ,	1445
1393	Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yun-	the Kallipos initiative of Greek open academic	1446
1394	tian Deng, Radha Poovendran, Yejin Choi, and	textbooks ¹¹ , full texts from publicly avail-	1447
1395	Bill Yuchen Lin. 2024. Magpie: Alignment data	able articles, theses, and dissertations from	1448
1396	synthesis from scratch by prompting aligned llms	academic repositories and the National Doc-	1449
1397	with nothing . <i>Preprint</i> , arXiv:2406.08464.	umentation Center ¹² , as well as pre-filtered	1450
1398	Jaume Zaragoza-Bernabeu, Gema Ramírez-Sánchez,	resources originally compiled from the web,	1451
1399	Marta Bañón, and Sergio Ortiz Rojas. 2022. Bi-	such as CulturaX (Nguyen et al., 2023) and	1452
1400	cleaner AI: Bicleaner goes neural . In <i>Proceedings</i>	CulturaY ¹³ . In addition to Modern Greek, we	1453
1401	<i>of the Thirteenth Language Resources and Evalua-</i>	incorporated a significant amount of Ancient	1454
1402	<i>tion Conference</i> , pages 824–831, Marseille, France.	Greek texts into our training corpus from Wik-	1455
1403	European Language Resources Association.	isource, school books, web pages, and Project	1456
1404	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali	Gutenberg ¹⁴ , which provides freely available	1457
1405	Farhadi, and Yejin Choi. 2019. HellaSwag: Can a ma-	Ancient Greek texts, including classical litera-	1458
1406	chine really finish your sentence? In <i>Proceedings of</i>	ture and historical documents. By including	1459
1407	<i>the 57th Annual Meeting of the Association for Com-</i>		
1408	<i>putational Linguistics</i> , pages 4791–4800, Florence,	⁶ https://www.ema.europa.eu/	
1409	Italy. Association for Computational Linguistics.	⁷ https://www.gutenberg.org/	
1410	Yifan Zhang, Yifan Luo, Yang Yuan, and Andrew Chi-	⁸ https://et.gr/	
1411	Chih Yao. 2024. Automathtext: Autonomous data se-	⁹ https://huggingface.co/datasets/AI-team-UoA/	
1412	lection with language models for mathematical texts.	greek_legal_code	
1413	<i>arXiv preprint arXiv:2402.07625</i> .	¹⁰ https://ebooks.edu.gr/ebooks/	
1414	Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji	¹¹ https://kallipos.gr/en/homepage/	
1415	Kawaguchi, and Lidong Bing. 2024. How do large	¹² https://www.ekt.gr/en	
1416	language models handle multilingualism? <i>arXiv</i>	¹³ https://huggingface.co/datasets/ontocord/	
1417	<i>preprint arXiv:2402.18815</i> .	CulturaY	
		¹⁴ https://www.gutenberg.org/	

Ancient Greek data, we ensured that Llama-Krikri-8B is able process polytonic Greek and engage with historical texts effectively. This enhances the model’s utility for classical studies, historical research, and philological applications.

2. **English Texts (21B tokens):** A subset of high-quality English data was mixed into the training corpus. This subset ensures that the model is continually trained on English data, and is drawn from sources that were also used for the Greek data, such as Wikipedia, Wikisource, Project Gutenberg (post-1900), EUR-LEX, EMEA, Greek academic repositories, etc. We also utilized additional English texts originating from abstracts and full texts of academic records found on multiple scientific repositories (Roussis et al., 2022a, 2024), ECtHR-PCR (Tyss et al., 2024), and pre-filtered datasets from TxT360 (Tang et al., 2024), like ArXiv, S2ORC (Lo et al., 2020), and PubMed Central (Roberts, 2001). By incorporating diverse and high-quality English texts, we mitigate the risk of catastrophic forgetting.

3. **Parallel Data (5.5B tokens):** We compiled a diverse parallel corpus with language pairs covering multiple languages: to Greek, English, French, Portuguese, German, Spanish, and Italian. The decision to add parallel data which covers other European languages (i.e., German, French, Italian, Portuguese, and Spanish) is informed from the languages that have been included in the multilingual instruction tuning of the original Llama-3.1. We utilized resources such as SciPar (Roussis et al., 2022a), MediSys (Roussis et al., 2022b), MultiEUR-LEX (Chalkidis et al., 2021), Europarl, TED Talk transcripts (Qi et al., 2018), and other sources with sentence pairs such as ELRC-SHARE (Lösch et al., 2021) & OPUS (Tiedemann, 2012). Our data include parallel documents and sentence pairs randomly sampled for each translation direction, e.g., EN-EL/EL-EN and EN-DE/DE-EN, as well as augmented training examples with concatenated parallel content across multiple languages (e.g., a Greek text followed by its English, German, and Spanish translations with appropriate prompt templates). The addition

of these documents has a twofold effect. It has been shown that parallel data boosts translation performance (Alves et al., 2024; Martins et al., 2024), while limited empirical evidence indicates that pretrained LLMs process multilingual queries by first translating the content into English, utilizing their English knowledge to answer the query and then translate the answer back to the original language (Zhao et al., 2024).

4. **Code and Math (7.8B tokens):** We also integrated datasets containing text with code and mathematics, leveraging Stack Overflow¹⁵, Python-Edu which is a subset of the SmolLM corpus (Allal et al., 2025) originating from The Stack V2 dataset (Lozhkov et al., 2024) and having been scored with an educational code classifier, and the AutoMathText dataset (Zhang et al., 2024), which is a collection of math-related documents originating from web data, papers on arXiv, and code/notebooks on GitHub. AutoMathText has undergone an automatic selection process using Qwen-72B (Bai et al., 2023) for relevancy to the mathematical domain and the educational value of each document. Code and Mathematics data, although not specific to Greek, were included to preserve and enhance the model’s ability to handle coding tasks, math problems and formal language. Maintaining these capabilities broadens the utility of Llama-Krikri beyond pure language tasks.

B Pretraining data cleaning pipelines

Our filtering processes began with format standardization in order to facilitate uniform processing across multiple heterogeneous datasets. We converted all textual content from various formats (e.g., PDF, HTML, plain text, etc.) into JSONL containing both the document text and relevant metadata such as identified language, word count, and source information (including source URLs).

For PDF documents such as academic records and laws, we implemented a specialized pipeline which integrated Marker¹⁶ for extraction and conversion into Markdown files, as it exhibits strong performance for Greek texts. Subsequently, the pipeline included language identification using

¹⁵<https://huggingface.co/datasets/code-rag-bench/stackoverflow-posts>

¹⁶<https://github.com/VikParuchuri/marker>

FastText (Joulin et al., 2017), removal of mark-down artifacts, and removal of lines with characters outside Unicode ranges for Greek, Latin, and other common and scientific symbols. Furthermore, we utilized document structure metrics (Marker also extracts various structural metadata) as quality indicators, such as the ratio of tables to pages and the fraction of removed lines in disallowed scripts.

Our main filtering pipeline used sequential rule-based and statistical filters to remove outlier documents across all data sources. First, we implemented URL-based filtering by removing content from several blacklisted domains known to contain low-quality or problematic content. This was particularly effective for web-crawled datasets like CulturaX (Nguyen et al., 2023) where source metadata was available. We then applied a set of minimal content-quality filters:

- Removal of documents containing multiple instances of profane or inappropriate terms from a curated list of Greek bad words
- Removal of short documents based on character and word counts
- Removal of documents containing multiple substrings like "lorem ipsum" which are indicative of content with low educational value
- Removal of documents containing extremely long words (>60 characters)
- Removal of documents with mean word length outside specified values.
- Removal of documents with a high fraction of non-alphanumeric characters.

Parallel datasets were filtered using a different pipeline featuring various steps from previous work (Papavassiliou et al., 2018; Roussis and Papavassiliou, 2022; Roussis et al., 2024) which include: (a) rule-based filters, such as length ratio, language identification verification, and (b) model-based alignment quality scores using tools like LASER (Artetxe and Schwenk, 2018, 2019), BiCleaner AI (Zaragoza-Bernabeu et al., 2022), and CometKiwi (Rei et al., 2022).

Additionally, in order to mitigate privacy concerns and protect sensitive information, we systematically identified and anonymized personally identifiable information (PII) with the use of regular expressions. In particular, we aimed to detect and

replace e-mail addresses with a generic placeholder ("email@example.gr") and mask IP addresses (replacing them with 0.0.0.0).

Finally, for Greek, English, and Mathematics/Code datasets we implemented intra-dataset deduplication, as well as cross-dataset deduplication. We utilized MinHashLSH near-deduplication (Broder, 1997; Leskovec et al., 2020) with 5-gram subsets, a MinHash signature of 128, and a Jaccard similarity threshold of 0.8, following parameter choices similar to those used in other works (Nguyen et al., 2023; Voukoutis et al., 2024). Regarding the deduplication of parallel datasets, we followed a different approach. All sentence pairs were normalized and cleaned, by converting them to lowercase and removing digits, punctuation. Pairs were then deduplicated based on the existence of either the source or target within the same dataset, thus ensuring that no sentence can be found multiple times in each parallel dataset (Roussis and Papavassiliou, 2022; Roussis et al., 2024). It should be noted that deduplication has consistently been shown to lead to higher performance, reduced training costs, as well as reduced model memorization; thus indirectly protecting sensitive information (Lee et al., 2022; Carlini et al., 2023; Grattafiori et al., 2024; Albalak et al., 2024). However, as we mentioned earlier, global deduplication may also remove documents of high quality and actually hurt performance (Tang et al., 2024; Penedo et al., 2024). For this reason, we decided to up-sample datasets of specific sources with important content. Table 1 summarizes the composition of the filtered and deduplicated pretraining corpus. In total, our collected dataset comprises roughly 91B tokens, of which 62.3% is Greek text. For the final training curriculum, we upsampled parts of the corpus to effectively train on an equivalent of 110B tokens. Upsampling was used to give higher relative importance to certain underrepresented but valuable segments and it also leads to higher memorization of important content (Carlini et al., 2023; Tang et al., 2024). For example, we assigned a slightly higher weight to datasets with long-context documents, Wikipedia-like sources, dialogue data, multi-parallel documents, and to certain important domains, such as legal, scientific, and medical. The decision to include a significant amount of English and parallel data (23.1% and 6% of tokens, respectively) was guided by prior work (Voukoutis et al., 2024) showing that mixed-language training can help retain the base model’s general knowledge and

prevent catastrophic forgetting.

C Tokenizer and embeddings expansion process

The tokenizer and embeddings expansion process involved the following steps:

- **Data Acquisition:** We acquired data by collecting sentences from high-quality sources of our pretraining mix in five domains:
 1. **General domain** which reuses a sampled portion of the data used to train the tokenizer of Meltemi (Voukoutis et al., 2024) and covers diverse domains,
 2. **Legal domain** which uses legal texts extracted from the Greek Government Gazette and are available via the National Printing House¹⁷, as well as the Permanent Greek Legislation Code – Raptarchis dataset¹⁸ (Papaloukas et al., 2021),
 3. **Scientific domain** which uses publicly available articles, theses, and dissertations found in the National Documentation Center¹⁹,
 4. **Literature domain** from public-domain books from Wikisource²⁰ and Project Gutenberg²¹ which contain literature, poetry, and other original writings across various variants of Greek (e.g., Koine Greek, Medieval Greek, Modern Greek, etc.),
 5. **Ancient Greek** which contains texts only in Ancient Greek sourced from Greek school books and various publicly available corpora.
- **Filtering and Preprocessing:** Each dataset underwent sequential processing and filtering including language identification verification with FastText (Joulin et al., 2017), application of regular expressions to remove URLs and other anomalies, symbol-to-word ratio filtering to remove outliers, and NFC normalization. We then performed sentence-level exact deduplication within each individual dataset.

¹⁷<https://et.gr/>

¹⁸https://huggingface.co/datasets/AI-team-UoA/greek_legal_code

¹⁹<https://www.ekt.gr/en>

²⁰<https://el.wikisource.org/>

²¹<https://gutenberg.org/>

To ensure text quality, we applied fluency scoring using Monocleaner (Sánchez-Cartagena et al., 2018) which leverages a 7-gram KenLM model for Greek, and setting a score threshold of 0.3 for non-polytonic text.

- **Creation of Training and Test Sets:** For the tokenizer training and test sets creation, we sampled 50% of the sentences from each source and divided it into training and testing splits (80%–20%).
- **Domain-specific Token Allocation:** New tokens were added sequentially for each domain until tokenizer fertility for this domain remained relatively stable, with most of the tokens being allocated to the General domain. This approach ensured that common Modern Greek patterns receive the largest coverage, while specialized terminology and older Greek variants are adequately represented.

Table 2 presents the fertility of several tokenizers on the original Greek and English corpora. Note that these test sets have not been used anywhere in the domain-specific multi-stage training process and could be considered as out-of-domain for the Llama-Krikri-8B tokenizer. We evaluated the Llama-3.1-8B tokenizer (128,000 vocabulary size) and our custom-trained Llama-Krikri-8B tokenizer (149,248 vocabulary size). We observe that the Llama-3.1-8B tokenizer exhibits a fertility of 2.73 for Greek and 1.33 for English. Our Llama-Krikri-8B tokenizer demonstrates a significantly lower fertility of 1.65 for Greek, while maintaining the same low fertility of 1.33 for English as the base Llama-3.1-8B tokenizer. This indicates that our Llama-Krikri-8B tokenizer is more efficient for Greek texts compared to the standard Llama-3.1-8B tokenizer. The table also includes the Mistral and the Meltemi-7B Greek tokenizer for comparison.

Furthermore, in Table 9, we list the token allocation per domain, as well as the tokenizer fertilities of Llama-3.1-8B and Llama-Krikri-8B for each of the five domains for which we created test sets during the tokenizer extension process. We can observe that tokenizer fertility has dropped significantly across all domains, with Ancient Greek benefiting from the highest reduction (-1.62) with just 492 added tokens, whereas fertility for Legal texts has almost halved (1.54 from 2.82) after the integration of a significant number of tokens (4,000).

Domain	Added Tokens	Llama-3.1-8B	Llama-Krikri-8B	Δ Fertility
General	15,000	2.65	1.59	-1.06
Legal	4,000	2.82	1.54	-1.28
Scientific	1,000	2.91	1.73	-1.18
Literature	500	2.90	1.89	-1.01
Ancient Greek	492	3.77	2.15	-1.62
Total	20,992	—	—	—

Table 9: Domain-specific token allocation and fertility comparison

Following the development of the Llama-Krikri-8B tokenizer with 149,248 tokens, containing the initial Llama 3 reserved tokens to reach a vocabulary size divisible by 256 for computational efficiency, we expanded the embeddings and lm head layer of our Llama 3 based model to match this vocabulary size.

The newly initialized embeddings and lm head weights for each new token were initialized by averaging the embeddings (or lm head weights respectively) associated with the tokens retrieved by the initial Llama-3.1 tokenizer when fed with that token. This provided us with a smoother initial Greek embeddings and lm head training phase of the frozen model.

D Examples of queries and responses

Ex. 1: Programming: Python Fibonacci function

User: Γράψε μια συνάρτηση Python που υπολογίζει την ακολουθία fibonacci.

Krikri: Ακολουθεί η υλοποίηση της αναδρομικής συνάρτησης Python για τον υπολογισμό των αριθμών Fibonacci:

```
def fibonacci(n):
    """
    Υπολογίζει τον n-οστό αριθμό Fibonacci.

    Args:
        n (int): Θετικός ακέραιος. n >= 0 (0-indexed προσέγγιση)

    Returns:
        int: ο n-οστός αριθμός Fibonacci

    Raises:
        ValueError: αν n είναι αρνητικός αριθμός
    """
    # Έλεγχος εγκυρότητας εισόδου
    if not isinstance(n, int) or n < 0:
        raise ValueError("Η είσοδος πρέπει να είναι μη αρνητικός ακέραιος")

    # Βάση αναδρομής:
    # F(0) = 0, F(1) = 1 (κλασικοί ορισμοί)
    if n in (0, 1):
        return n

    # Αναδρομικός τύπος: F(n) = F(n-1) + F(n-2)
    return fibonacci(n-1) + fibonacci(n-2)
```

Παράδειγμα χρήσης:

```
# Υπολογισμός μερικών όρων της
# ακολουθίας Fibonacci
for i in range(10):
    print(f"F({i}) = {fibonacci(i)}")

# Εκτύπωση αποτελεσμάτων
# F(0) = 0
# F(1) = 1
# F(2) = 1
# F(3) = 2
# F(4) = 3
# F(5) = 5
# F(6) = 8
# F(7) = 13
# F(8) = 21
# F(9) = 34
```

Ex. 2: Translation: Ancient to Modern Greek

User: Μετάφρασέ μου στα Νέα ελληνικά το κείμενο:

‘Οἱ μὲν δὴ Μήλιοι τοσαῦτα ἀπεκρίναντο· οἱ δὲ Ἀθηναῖοι διαλυόμενοι ἤδη ἐκ τῶν λόγων ἔφασαν «ἄλλ’ οὐν μόνον γε ἀπὸ τούτων τῶν βουλευμάτων, ὡς ἡμῖν δοκεῖτε, τὰ μὲν μέλλοντα τῶν ὀρωμένων σαφέστερα κρίνετε, τὰ δὲ ἄφανε τῷ βούλεσθαι ὡς γινόμενα ἤδη θεᾶσθε, καὶ Λακεδαιμονίους καὶ τύχῃ καὶ ἐλπίσι πλείστον δὴ παραβελημένοι καὶ πιστεύσαντες πλείστον καὶ σφαλῆσεσθε.»’

Krikri: Ακολουθεί η μετάφραση του κειμένου στα Νέα Ελληνικά:

Οι Μήλιοι απάντησαν με αυτά τα λόγια. Οι Αθηναίοι, καθώς αποχωρούσαν από τη συζήτηση, είπαν: ‘Όμως εσείς είστε οι μόνοι που από αυτές τις συζητήσεις συμπεραίνετε σωστά τόσο τα ορατά όσο και τα αόρατα, σαν να έχουν ήδη συμβεί. Είστε υπερβολικά αισιόδοξοι και έχετε εναποθέσει τις ελπίδες σας στους Λακεδαιμονίους, στην τύχη και στις προσδοκίες σας. Και όσο περισσότερο έχετε πιστέψει σε αυτά, τόσο περισσότερο θα απογοητευτείτε.’

Ex. 3: Translation to JSON Output: Greek to Multiple Languages

User: Μετάφρασε μου το παρακάτω κείμενο από τα νέα ελληνικά στα αγγλικά, στα ισπανικά και στα γαλλικά. Δώσε μου τις μεταφράσεις σε JSON.

Σύμφωνα με την οδηγία 2009/28/EK του Ευρωπαϊκού Κοινοβουλίου, ως ενέργεια από ανανεώσιμες μη ορυκτές πηγές θεωρείται η αιολική, ηλιακή, αεροθερμική, γεωθερμική, υδροθερμική και ενέργεια των ωκεανών, υδροηλεκτρική, από βιομάζα, από τα εκλυόμενα στους χώρους υγειονομικής ταφής αέρια, από αέρια μονάδων επεξεργασίας λυμάτων και από βιοαέρια.

Krikri:

```
{
  "en": "According to Directive 2009/28/EC of the European Parliament, energy from renewable non-fossil sources includes wind, solar, aerothermal, geothermal, hydrothermal and ocean energy, hydropower, biomass, landfill gas, sewage treatment plant gas and biogases.",
  "es": "De conformidad con la Directiva 2009/28/CE del Parlamento Europeo, la energía procedente de fuentes renovables no fósiles incluye la eólica, solar, aerotérmica, geotérmica, hidrotérmica y la energía del océano, hidroeléctrica, procedente de la biomasa, de los gases de vertedero, de los gases de estación depuradora de aguas residuales y de los biogás.",
  "fr": "Conformément à la directive 2009/28/CE du Parlement européen, les énergies renouvelables non fossiles comprennent l'éolien, le solaire, l'aérothermique, la géothermie, l'hydrothermie et l'énergie océanique, l'hydroélectricité, la biomasse, les gaz de décharge, les gaz de station d'épuration et les biogaz."
}
```

1767

Ex. 4: Summarization: Basketball game

User: Μπορείς να δημιουργήσεις μια περίληψη 5 προτάσεων για τον παρακάτω αγώνα μπάσκετ: (User provides the text of an article describing a 2025 basketball game)

Krikri: Ο Ολυμπιακός νίκησε την Μπαρτσελόνα με 90-88 σε ένα συναρπαστικό παιχνίδι, χάρη σε ένα νικητήριο καλάνι του Φουρνιέ στην εκπνοή. Ο Βεζένκοφ και ο Γκος σημείωσαν από 24 και 20 πόντους αντίστοιχα για τους νικητές, ενώ ο Μετού ήταν ο κορυφαίος των Καταλανών με 16 πόντους. Η ομάδα του Μπαρτζώκα παρέμεινε στην κορυφή της Ευρωλίγκας μετά από αυτή τη σημαντική νίκη.

1768

Ex. 5: History: Population exchange between Greece and Turkey

User: Συνόψισε σε μία παράγραφο την Ελληνοτουρκική ανταλλαγή πληθυσμών του 1923

Krikri: Η Ελληνοτουρκική ανταλλαγή πληθυσμών του 1923, που προβλέφθηκε από τη Συνθήκη της Λωζάνης, αποτέλεσε μια μαζική μετακίνηση πληθυσμών μεταξύ Ελλάδας και Τουρκίας με στόχο την ομογενοποίηση των εθνών-κρατών. Περίπου 1,5 εκατομμύριο Έλληνες χριστιανοί, κυρίως από την Ανατολική Θράκη, την Ιωνία, τον Πόντο και την Καππαδοκία, μετακινήθηκαν στην Ελλάδα, ενώ αντίστοιχα 500.000 Μουσουλμάνοι (κυρίως Τούρκοι αλλά και Πομάκοι και Αθίγγανοι) μετακινήθηκαν από την Ελλάδα στην Τουρκία. Αυτή η ανταλλαγή, παρά τις τραγικές της συνέπειες για τους ανθρώπους που εκτοπίστηκαν, αποτέλεσε ένα από τα σημαντικότερα γεγονότα που διαμόρφωσαν τη σύγχρονη δημογραφική και εθνοτική σύνθεση των δύο χωρών, οδηγώντας σε μια πιο ομοιογενή εθνική ταυτότητα για κάθε κράτος. Η διαδικασία αυτή ολοκληρώθηκε το 1924 και θεωρείται ορόσημο στις ελληνοτουρκικές σχέσεις.

1769

Ex. 6: High school biology: The process of photosynthesis

User: Δώσε μια σύνοψη της διαδικασίας της φωτοσύνθεσης σε μία παράγραφο.

Krikri: Η φωτοσύνθεση είναι η διαδικασία με την οποία τα φυτά, χρησιμοποιώντας την ηλιακή ενέργεια, μετατρέπουν το διοξείδιο του άνθρακα (CO₂) και το νερό (H₂O) σε γλυκόζη (C₆H₁₂O₆) και οξυγόνο (O₂). Συγκεκριμένα, η χλωροφύλλη στα φύλλα απορροφά φως, διασπά το νερό απελευθερώνοντας οξυγόνο, και χρησιμοποιώντας το υδρογόνο συνθέτει γλυκόζη από το διοξείδιο του άνθρακα. Η παραγόμενη γλυκόζη αποτελεί πηγή ενέργειας για το φυτό, ενώ το οξυγόνο απελευθερώνεται στην ατμόσφαιρα ως σημαντικό παραπροϊόν, υποστηρίζοντας την επιβίωση των αερόβιων οργανισμών.

1770

Ex. 7: Chain-of-Thought Reasoning: Math problem

System: (Appropriate system prompt used to enable thinking mode)

User: Η Αλίκη έχει ένα καλάθι με μήλα. Δίνει τα $\frac{2}{5}$ των μήλων στην αδερφή της και το $\frac{1}{5}$ στον αδερφό της. Τελικά της μένουν 6 μήλα. Πόσα μήλα είχε αρχικά στο καλάθι;

Krikri:

<thinking>

Ας λύσουμε αυτό το πρόβλημα βήμα προς βήμα:

1. Έστω x ο αρχικός αριθμός μήλων που είχε η Αλίκη
2. Γνωρίζουμε ότι:
 - * Δίνει $\frac{2}{5}$ του x στην αδερφή της
 - * Δίνει $\frac{1}{5}$ του x στον αδερφό της
 - * Μετά από όλα αυτά της μένουν 6 μήλα
3. Ας κάνουμε την εξίσωση:
 - * Αρχικά είχε: x μήλα
 - * Της μένουν: 6 μήλα
 - * Άρα: $x - (\frac{2}{5})x - (\frac{1}{5})x = 6$
 - * Απλοποιώντας: $x - (\frac{3}{5})x = 6$
 - * Συνδυάζοντας όρους: $(\frac{2}{5})x = 6$
 - * Πολλαπλασιάζοντας και τα δύο μέλη με $\frac{5}{2}$:
 - * $x = 15$

Ας επαληθεύσουμε:

- * Αρχικά είχε 15 μήλα
- * Δίνει: $\frac{2}{5} * 15 = 6$ μήλα στην αδερφή της
- * Δίνει: $\frac{1}{5} * 15 = 3$ μήλα στον αδερφό της
- * Μένουν: $15 - 6 - 3 = 6$ μήλα

Η λύση είναι σωστή!

</thinking>

<output>

Η Αλίκη είχε αρχικά 15 μήλα στο καλάθι της.

Μπορούμε να το επαληθεύσουμε:

- * Έδωσε 6 μήλα στην αδερφή της ($\frac{2}{5}$ του συνόλου)
- * Έδωσε 3 μήλα στον αδερφό της ($\frac{1}{5}$ του συνόλου)
- * Έτσι της έμειναν 6 μήλα ($15 - 6 - 3 = 6$)

</output>