

Large Language Models for Mental Health: A Multilingual Evaluation

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have demonstrated impressive capabilities across various NLP tasks, but their performance in multilingual settings is often underexplored. This study evaluates proprietary and open-source LLMs on eight mental health datasets of various languages. We compare their performance in zero-shot, few-shot, and fine-tuned settings against traditional methods. Results show that LLMs achieve competitive or superior F1 scores across several datasets, with fine-tuned models often surpassing state-of-the-art results. However, performance varies across languages, highlighting both the strengths and limitations of LLMs in this critical application. These findings provide actionable insights into the application of LLMs for multilingual mental health text classification.

1 Introduction

While LLMs have been transforming research in NLP, caution must be exercised when adopting these models in sensitive domains such as mental health (Hua et al., 2024). Due to the potential risks and ethical considerations, experts are cautious about the integration of LLMs in applications in sensitive domains. These concerns are further amplified in multilingual settings where studies have demonstrated that LLMs tend to perform worse when prompted in languages other than English (Jin et al., 2024; Raihan et al., 2024a).

Most mental health datasets are curated from social media mining platforms such as Reddit and X. A recent survey by Kumar et al. (2024) shows that the clear majority of such datasets (Mariappan et al., 2024; Turcan and Mckeown, 2019; Raihan et al., 2024b) are in English. Recent efforts are being made to curate similar

resources in other languages such as Arabic (Baghdadi et al., 2022; Helmy et al., 2024), Russian (Narynov et al., 2020), and Thai (Hämäläinen et al., 2021). However, none of the aforementioned studies evaluate the performance of LLMs on non-English datasets, leaving an important gap in understanding their performance in multilingual mental health settings.

A few recent studies explore the performance of LLMs in English mental health datasets. Xu et al. (2024) examines LLMs performance on multiple datasets compared to statistical and traditional encoder-only models (Devlin et al., 2019; Alsentzer et al., 2019). Similarly, Kuzmin et al. (2024), Yang et al. (2023), and Wei et al. (2022) investigate LLM performance exploring various prompting strategies. Finally, Yang et al. (2024) presents a fine-tuning approach with the release of MentaLLaMA, a task-specific model for the domain. While these approaches achieve competitive results, they are limited to English, leaving a significant gap in research in other languages.

To address this gap, we present the first multilingual evaluation of state-of-the-art LLMs on mental health datasets. We consider mental health datasets in six languages - Arabic, Bengali, Spanish, Portuguese, Russian, and Thai - and two tasks, namely depression and suicide ideation detection. We address the following Research Questions (RQs):

- **RQ1:** How does the performance of LLMs compare to previously proposed models (e.g., statistical, neural, BERT-based)?
- **RQ2:** What are the best prompting strategies for LLMs on mental health?
- **RQ3:** What is the impact of instruction fine-tuning on the performance of open-source LLMs?

Dataset	Language (ISO code)	Mental Disorder	Platform	Expert Labeling	Size
Narynov et al. (2020)	Russian (ru)	Depression	VKontakte	Yes	32,018
Hämäläinen et al. (2021)	Thai (tha)	Depression	Blogs	Yes	33,436
Boonyarat et al. (2024)	Thai (tha)	Suicidal Ideation	X	No	2,400
Uddin et al. (2019)	Bengali (ben)	Depression	X	Yes	3,914
de Oliveira et al. (2022)	Portuguese (por)	Suicidal Ideation	X	Yes	3,788
Baghdadi et al. (2022)	Arabic (ar)	Suicidal Ideation	X	N/A	14,576
Helmy et al. (2024)	Arabic (ar)	Depression	X	No	10,000
Valeriano et al. (2020)	Spanish (es)	Suicidal ideation	X	N/A	1,068

Table 1: Overview of the eight mental disorder datasets across different languages.

2 Datasets

For our study, we acquire eight publicly available datasets presented in Table 1. The datasets include posts in Russian from VKontakte and posts in Thai from online blogs that have been annotated for depression (Narynov et al., 2020; Hämäläinen et al., 2021). Additionally, we acquire multiple datasets containing data sourced from X, annotated for depression in Bengali (Uddin et al., 2019) and Arabic (Helmy et al., 2024). Furthermore, there are posts from X annotated for suicidal ideation in Thai (Boonyarat et al., 2024), Portuguese (de Oliveira et al., 2022), Arabic (Baghdadi et al., 2022) and Spanish (Valeriano et al., 2020).

3 Experiments and Results

We evaluate a diverse set of LLMs, encompassing both proprietary and open-source architectures. Our evaluation includes multiple prompting strategies and also fine-tuning open-source models to gather better insights.

3.1 LLMs

We evaluate seven state-of-the-art LLMs spanning both proprietary and open-source architectures, as listed in Table 2.

LLMs	OS?	Size	Ref.
GPT4-omni	✗	–	OpenAI
Claude3.5-Sonnet	✗	–	Anthropic
Gemini2-Flash	✗	–	Team et al.
LLaMA3.2	✓	11B	Dubey et al.
Gemma2	✓	27B	Team et al.
Mistral	✓	8B	MistralAI
R1	✓	14B	Guo et al.

Table 2: List of LLMs used for the experiments. (OS - Open-Source).

Our selection includes three proprietary models—GPT-4 Omni, Claude 3.5 Sonnet, and Gemini 2 Flash—as well as four open-source models: LLaMA 3.2, Gemma 2, Mistral, and R1. The proprietary models remain closed-source, with limited architectural details, while the open-source models offer greater transparency and adaptability for fine-tuning. These models have demonstrated strong performance across multiple tasks and domains, making them well-suited for our multilingual evaluation. We analyze their capabilities in both zero-shot and few-shot settings, leveraging their diverse architectures and parameter sizes to assess their effectiveness in multilingual tasks.

3.2 Prompting

We evaluate three prompting methods: zero-shot, few-shot (5 examples), and Chain-of-Thought (CoT) prompting (Wei et al., 2022). For the 5-shot setting, we randomly select five examples from the respective datasets. For CoT prompting, we adopt the SOTA prompting method for the mental health tasks, CoT_{Emo_FS} , introduced by Yang et al. (2023).

Table 3 presents a comprehensive comparison of F1 scores obtained via different prompting methods across eight multilingual depression and suicide ideation datasets. Our analysis reveals that CoT prompting generally improves performance, with models such as GPT-4 and Claude3.5 often achieving the highest scores—for example, GPT-4 increases its F1 from 0.76 to 0.87 on the Russian dataset and from 0.75 to 0.84 on the Spanish dataset. However, the gains are not uniform across all settings, as seen with the Bengali dataset where few-shot and CoT strategies yield comparable results. Moreover, while some baseline methods

Baseline - Reported results				Our Results (LLMs)							
Dataset	lang	Models	F1	Prompting	GPT4	Claude3.5	Gemini2	LLaMA 3.2	Gemma2	Ministral	R1
	ISO		reported	method	omni	Sonnet	Flash	11B	27B	8B	14B
Narynov et al.	ru	-	-	zero	0.76	0.74	0.68	0.56	0.69	0.41	0.71
				few	0.79	0.83	0.73	0.62	0.71	0.53	0.73
				CoT	0.87	0.85	0.80	0.59	0.73	0.44	0.79
Hämäläinen et al.	tha	Thai-BERT	0.78	zero	0.77	0.77	0.66	0.45	0.68	0.20	0.76
				few	0.84	0.81	0.69	0.58	0.66	0.31	0.75
				CoT	0.85	0.80	0.70	0.40	0.69	0.40	0.81
Boonyarat et al.	tha	LFBERT	0.93	zero	0.83	0.81	0.83	0.63	0.76	0.26	0.69
				few	0.87	0.85	0.86	0.71	0.72	0.39	0.71
				CoT	0.91	0.95	0.87	0.77	0.84	0.47	0.84
Uddin et al.	ben	GRU	0.76	zero	0.78	0.85	0.79	0.73	0.73	0.36	0.66
				few	0.86	0.91	0.88	0.59	0.71	0.43	0.64
				CoT	0.86	0.91	0.88	0.59	0.71	0.43	0.64
Oliveira et al.	por	Random Forest	0.94	zero	0.86	0.86	0.81	0.71	0.80	0.56	0.61
				few	0.89	0.93	0.85	0.73	0.63	0.67	0.69
				CoT	0.94	0.95	0.89	0.71	0.80	0.51	0.82
Baghdadi et al.	ar	AraElectra	0.96	zero	0.80	0.85	0.81	0.58	0.73	0.34	0.77
				few	0.87	0.92	0.89	0.67	0.82	0.47	0.79
				CoT	0.89	0.91	0.87	0.61	0.81	0.47	0.83
Helmy et al.	ar	LR (TF-IDF)	0.95	zero	0.87	0.91	0.79	0.56	0.62	0.50	0.84
				few	0.93	0.95	0.86	0.73	0.79	0.61	0.86
				CoT	0.95	0.95	0.82	0.83	0.67	0.50	0.87
Valeriano et al.	es	LR (W2V)	0.79	zero	0.75	0.69	0.62	0.37	0.41	0.23	0.67
				few	0.81	0.76	0.69	0.46	0.51	0.31	0.67
				CoT	0.84	0.79	0.70	0.43	0.60	0.21	0.76

Table 3: F1 score comparison for **Zero-Shot**, **Few-Shot**, and **Chain-of-Thought** prompting across the eight (8) multilingual depression and suicide ideation datasets. We compare the reported best methods and results in the original papers with the proprietary and open-source LLMs with different prompting strategies. The highest F1 score for each dataset is shown in orange. For all other F1 scores (in blue) - the darker the shade, the higher the score. For the language names, ISO-639 codes are used. ('LR' - Logistic Regression, 'W2V' - Word2Vec, 'CoT' - Chain-of-Thought).

(e.g., Random Forest and AraElectra) achieve competitive performance in certain languages, the results underscore the potential of advanced prompting techniques to narrow the gap with or even surpass traditional approaches. These observations motivate further investigation into model- and language-specific factors that influence the efficacy of prompt engineering.

3.3 Fine-tuning

Due to the intrinsic black box nature of proprietary models and their high costs, we wanted to explore models that could be fully-customization to this task. Therefore, we experiment with fine-tuning the open-source models. The hyperparameters are chosen empirically as we run a set of experiments with different combinations of parameters and report the best results. The final selection of hyper-parameters is presented in Appendix A.

Table 4 presents a comparative analysis of F1 scores before and after fine-tuning on eight multilingual depression and suicide ideation

datasets. The results indicate that fine-tuning generally enhances model performance, with Gemma2 and R1 often reaching the highest scores. While LLaMA 3.2 and Ministral show notable improvements in several datasets, their performance gains are not uniform—for instance, LLaMA 3.2 exhibits a decrease in the Bengali dataset. These findings underscore the potential of fine-tuning to optimize multilingual performance while also revealing the need for further investigation into model- and dataset-specific factors that modulate the benefits of fine-tuning.

4 Observation and Analysis

We now revisit the 3 RQs posed in the introduction (see Section 1):

RQ₁ How does the performance of LLMs compare to previously proposed models (e.g., statistical, neural, BERT-based)?

Dataset Info		Before Fine-Tuning (Zero-Shot)				After Fine-Tuning			
Dataset	lang	LLaMA 3.2	Gemma2	Ministral	R1	LLaMA 3.2	Gemma2	Ministral	R1
Narynov et al.	ru	0.56	0.69	0.41	0.71	0.79	0.83	0.62	0.79
Hämäläinen et al.	tha	0.45	0.68	0.20	0.76	0.62	0.73	0.43	0.82
Boonyarat et al.	tha	0.63	0.76	0.26	0.69	0.70	0.75	0.51	0.74
Uddin et al.	ben	0.73	0.73	0.36	0.66	0.65	0.77	0.63	0.64
Oliveira et al.	por	0.71	0.80	0.56	0.61	0.72	0.86	0.64	0.70
Baghdadi et al.	ar	0.58	0.73	0.34	0.77	0.80	0.88	0.58	0.81
Helmy et al.	ar	0.56	0.62	0.50	0.84	0.70	0.81	0.71	0.93
Valeriano et al.	es	0.37	0.41	0.23	0.67	0.55	0.62	0.48	0.76

Table 4: F1 score comparison before and after fine-tuning across eight multilingual depression and suicide ideation datasets. The columns under **Before Fine-Tuning (Zero-Shot)** report the initial prompting results, while those under **After Fine-Tuning** display the fine-tuned performance. The highest F1 score in the fine-tuned setting is highlighted with an orange cell. For the language names, ISO-639 codes are used.

Our analysis indicates that LLMs, when equipped with effective prompting strategies, achieve performance that is competitive with or superior to traditional approaches. While statistical, neural, and BERT-based models demonstrate strong performance in certain linguistic scenarios, LLMs exhibit robust and consistent F1 scores across diverse multilingual mental health datasets, highlighting their capacity for broad generalization and adaptability.

ment underscores the value of targeted fine-tuning in adapting LLMs to domain-specific tasks, thereby enhancing their overall effectiveness in mental health applications while also mitigating performance variability observed in zero-shot configurations.

RQ₂ What are the best prompting strategies for LLMs on mental health?

The results reveal that Chain-of-Thought (CoT) prompting is the most effective strategy for mental health applications, consistently yielding higher F1 scores compared to zero-shot and few-shot methods. This structured approach to prompting enhances reasoning capabilities, enabling LLMs to better extract nuanced signals from text data, which is critical in sensitive domains such as mental health.

RQ₃ What is the impact of instruction fine-tuning on the performance of open-source LLMs?

Instruction fine-tuning markedly improves the performance of open-source LLMs, as evidenced by substantial increases in F1 scores across all evaluated datasets. This improve-

5 Conclusion and Future Work

This work represents the first investigation of LLMs in the multilingual mental health domain. Our findings show that advanced prompting strategies—particularly chain-of-thought prompting—and targeted instruction fine-tuning substantially enhance model performance, often surpassing traditional statistical, neural, and BERT-based approaches. While our results demonstrate considerable promise, the variability in performance across languages and models indicates that further research is required to optimize these techniques for sensitive mental health applications.

Overall, this study lays an important foundation for future efforts aimed at refining LLM-based methodologies in complex, multilingual settings. In future work, we would like to include more tasks and languages to broaden our understanding and gain more insights. Additionally, we plan to adapt open-source models to the domain with methods like Continual Pretraining and Synthetic Fine-tuning to potentially increase their performance.

216	Limitations		
217	While our approach is limited by the inherent		
218	variability in data sources, evaluation protocols,		
219	and reporting standards across the literature, it		
220	also represents a significant strength: we are the		
221	first to systematically synthesize and critically		
222	evaluate LLM performance in this sensitive and		
223	underexplored area. The exclusive reliance on		
224	publicly available data restricts the diversity		
225	and depth of our analysis, and the absence of		
226	direct model development or human subject		
227	involvement means that practical deployment		
228	challenges remain unaddressed. These limita-		
229	tions notwithstanding, our work lays a founda-		
230	tional framework for future research that can		
231	leverage standardized benchmarks and broader		
232	datasets to further validate and enhance the util-		
233	ity of LLMs in mental health applications.		
234	Ethical Considerations		
235	This work is entirely analytical and does		
236	not involve the collection of new data, the		
237	development of new models, or engage-		
238	ment with human subjects. All analyses		
239	are based solely on previously published		
240	and publicly available data. We adhere to		
241	the ethical guidelines outlined in the ACL		
242	Code of Ethics (https://www.aclweb.org/portal/content/acl-code-ethics), and we		
243	emphasize that any research in the mental health		
244	domain must be conducted with utmost sensi-		
245	tivity to privacy and ethical considerations. Al-		
246	though our study is retrospective in nature, we		
247	recognize the critical importance of safeguard-		
248	ing vulnerable populations, and we advocate		
249	for strict adherence to ethical standards in any		
250	practical applications derived from our findings.		
251			
252	References		
253	Emily Alsentzer, John Murphy, William Boag, Wei-		
254	Hung Weng, Di Jindi, Tristan Naumann, and		
255	Matthew McDermott. 2019. Publicly available		
256	clinical bert embeddings. In <i>Proceedings of the</i>		
257	<i>2nd Clinical Natural Language Processing Work-</i>		
258	<i>shop</i> .		
259	Anthropic. 2023. Claude: The anthropic ai lan-		
260	guage model. <i>Online documentation</i> . Available		
261	at: https://www.anthropic.com .		
	Nadiah A Baghdadi, Amer Malki, Hossam Magdy	262	
	Balaha, Yousry AbdulAzeem, Mahmoud Badawy,	263	
	and Mostafa Elhosseini. 2022. An optimized	264	
	deep learning approach for suicide detection	265	
	through arabic tweets. <i>PeerJ Computer Science</i> .	266	
	Panchanit Boonyarat, Di Jie Liew, and Yung-Chun	267	
	Chang. 2024. Leveraging enhanced bert models	268	
	for detecting suicidal ideation in thai social media	269	
	content amidst covid-19. <i>Information Processing</i>	270	
	<i>& Management</i> .	271	
	Adonias C de Oliveira, Evandro JS Diniz, Silmar	272	
	Teixeira, and Ariel S Teles. 2022. How can	273	
	machine learning identify suicidal ideation from	274	
	user’s texts? towards the explanation of the bo-	275	
	amente system. <i>Procedia Computer Science</i> .	276	
	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	277	
	Kristina Toutanova. 2019. Bert: Pre-training of	278	
	deep bidirectional transformers for language un-	279	
	derstanding. In <i>Proceedings of the 2019 Con-</i>	280	
	<i>ference of the North American Chapter of the</i>	281	
	<i>Association for Computational Linguistics</i> .	282	
	Abhimanyu Dubey, Abhinav Jauhri, Abhinav	283	
	Pandey, Abhishek Kadian, et al. 2024. The	284	
	llama 3 herd of models. <i>arXiv preprint</i>	285	
	<i>arXiv:2407.21783</i> .	286	
	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao	287	
	Song, et al. 2025. Deepseek-r1: Incentivizing	288	
	reasoning capability in llms via reinforcement	289	
	learning. <i>arXiv preprint arXiv:2501.12948</i> .	290	
	Mika Härmäläinen, Pattama Patpong, Khalid Alna-	291	
	jjar, Niko Partanen, and Jack Rueter. 2021. De-	292	
	tecting depression in thai blog posts: a dataset	293	
	and a baseline. In <i>Proceedings of the Seventh</i>	294	
	<i>Workshop on Noisy User-generated Text (W-NUT</i>	295	
	<i>2021)</i> .	296	
	AbdelMoniem Helmy, Radwa Nassar, and Nagy	297	
	Ramdan. 2024. Depression detection for twit-	298	
	ter users using sentiment analysis in english and	299	
	arabic tweets. <i>Artificial intelligence in medicine</i> .	300	
	Yining Hua, Fenglin Liu, Kailai Yang, Zehan Li,	301	
	et al. 2024. Large language models in mental	302	
	health care: a scoping review. <i>arXiv preprint</i>	303	
	<i>arXiv:2401.02984</i> .	304	
	Yiqiao Jin, Mohit Chandra, Gaurav Verma, Yibo Hu,	305	
	Munmun De Choudhury, and Srijan Kumar. 2024.	306	
	Better to ask in english: Cross-lingual evaluation	307	
	of large language models for healthcare queries.	308	
	In <i>Proceedings of the ACM on Web Conference</i>	309	
	<i>2024</i> .	310	
	Puneet Kumar, Alexander Vedernikov, and Xiaobai	311	
	Li. 2024. Measuring non-typical emotions for	312	
	mental health: A survey of computational ap-	313	
	proaches. <i>arXiv preprint arXiv:2403.08824</i> .	314	

315	Gleb Kuzmin, Petr Strepetov, Maksim Stankevich,	Jason Wei, Xuezhi Wang, Dale Schuurmans,	369
316	Artem Shelmanov, and Ivan Smirnov. 2024. Men-	Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,	370
317	tal disorders detection in the era of large language	Denny Zhou, et al. 2022. Chain-of-thought	371
318	models. <i>arXiv preprint arXiv:2410.07129</i> .	prompting elicits reasoning in large language	372
		models. <i>Advances in neural information process-</i>	373
319	Umasree Mariappan, D Balakrishnan, G Merline,	<i>ing systems</i> .	374
320	M Sandhia, Dubba Saitej Reddy, and Satti-		
321	neni Gagan Teja. 2024. Mental health disorder	Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saa-	375
322	prediction using recurrent neural network algo-	dia Gabriel, Hong Yu, James Hendler, Marzyeh	376
323	rithm. In <i>2024 Asia Pacific Conference on Inno-</i>	Ghassemi, Anind K Dey, and Dakuo Wang. 2024.	377
324	<i>vation in Technology (APCIT)</i> .	Mental-llm: Leveraging large language models	378
		for mental health prediction via online text data.	379
325	Sergazy Narynov, Daniyar Mukhtarkhanuly, and	<i>Proceedings of the ACM on Interactive, Mobile,</i>	380
326	Batyrkhan Omarov. 2020. Dataset of depressive	<i>Wearable and Ubiquitous Technologies</i> .	381
327	posts in russian language collected from social		
328	media. <i>Data in brief</i> .	Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian	382
		Xie, Ziyang Kuang, and Sophia Ananiadou. 2023.	383
329	OpenAI. 2023. Gpt-4 technical report. <i>arXiv</i>	Towards interpretable mental health analysis with	384
330	<i>preprint arXiv:2303.08774</i> .	large language models. In <i>Proceedings of the</i>	385
		<i>2023 Conference on Empirical Methods in Natu-</i>	386
331	Nishat Raihan, Antonios Anastasopoulos, and Mar-	<i>ral Language Processing</i> .	387
332	cos Zampieri. 2024a. mhumaneval—a multi-		
333	lingual benchmark to evaluate large language	Kailai Yang, Tianlin Zhang, Ziyang Kuang, Qianqian	388
334	models for code generation. <i>arXiv preprint</i>	Xie, Jimin Huang, and Sophia Ananiadou. 2024.	389
335	<i>arXiv:2410.15037</i> .	Mentallama: interpretable mental health analy-	390
		sis on social media with large language models.	391
336	Nishat Raihan, Sadiya Sayara Chowdhury Puspo,	In <i>Proceedings of the ACM on Web Conference</i>	392
337	Shafkat Farabi, Ana-Maria Bucur, Tharindu	<i>2024</i> .	393
338	Ranasinghe, and Marcos Zampieri. 2024b. Men-		
339	talhelp: A multi-task dataset for mental health		
340	in social media. In <i>Proceedings of the 2024</i>		
341	<i>Joint International Conference on Computational</i>		
342	<i>Linguistics, Language Resources and Evaluation</i>		
343	<i>(LREC-COLING 2024)</i> .		
344	Gemini Team, Rohan Anil, Sebastian Borgeaud,		
345	Jean-Baptiste Alayrac, et al. 2023. Gemini: a		
346	family of highly capable multimodal models.		
347	<i>arXiv preprint arXiv:2312.11805</i> .		
348	Gemma Team, Morgane Riviere, Shreya Pathak,		
349	Pier Giuseppe Sessa, Cassidy Hardin, et al. 2024.		
350	Gemma 2: Improving open language models at a		
351	practical size. <i>arXiv preprint arXiv:2408.00118</i> .		
352	Elsbeth Turcan and Kathleen Mckeown. 2019.		
353	Dreaddit: A reddit dataset for stress analysis in		
354	social media. In <i>Proceedings of the Tenth Inter-</i>		
355	<i>national Workshop on Health Text Mining and</i>		
356	<i>Information Analysis (LOUHI 2019)</i> .		
357	Abdul Hasib Uddin, Durjoy Bapery, and Abu		
358	Shamim Mohammad Arif. 2019. Depression		
359	analysis of bangla social media data using gated		
360	recurrent neural network. In <i>2019 1st Interna-</i>		
361	<i>tional conference on advances in science, en-</i>		
362	<i>gineering and robotics technology (ICASERT)</i> ,		
363	pages 1–6. IEEE.		
364	Kid Valeriano, Alexia Condori-Larico, and José		
365	Sulla-Torres. 2020. Detection of suicidal in-		
366	intent in spanish language social networks using		
367	machine learning. <i>International Journal of Ad-</i>		
368	<i>vanced Computer Science and Applications</i> .		

394

A Experimental Details

395

The fine-tuning stage is performed on a single
 396 NVIDIA A100 GPU with 40 GB of memory,
 397 accessed via Google Colab¹. The system is
 398 further equipped with 80 GB of RAM and 256
 399 GB of disk storage to support computational
 400 efficiency.

Parameter	Value
Max Sequence Length	2048
Batch Size (Train/Eval)	8
Gradient Accumulation Steps	4
Number of Epochs	3
Learning Rate	5e-5
Weight Decay	0.02
Warmup Steps	10%
Optimizer	AdamW (8-bit)
LR Scheduler	Cosine
Precision	BF16
Evaluation Strategy	Steps
Evaluation Steps	50
Save Strategy	Steps
Save Steps	Varies
Seed	42

Table 5: Final set of hyperparameters, chosen empirically after several iterations of trial and error, for fine-tuning.

¹<https://colab.research.google.com/>