

HaloRAG: Towards Mitigating LLM Hallucinations with Low-Cost Real-Time Retrieval

Anonymous ACL submission

Abstract

Large Language Models (LLMs) often struggle to stay up-to-date due to their reliance on static datasets, leading to outdated responses and hallucinations. We introduce HaloRAG, a cost-efficient agentic wrapper that enhances LLMs with real-time information retrieval using advanced web scraping technologies. Leveraging semantic searches and Retrieval-Augmented Generation (RAG), this wrapper fetches, validates, and summarizes up-to-date web data, extending the LLM’s knowledge base without retraining. This method significantly enhances the accuracy and relevance of LLM responses, particularly for queries requiring the latest information. Comparative analysis indicates that the wrapper-enhanced LLM outperforms models like GPT-3.5 and Claude on queries involving recent events and emerging technologies. This work advocates for integrating real-time data retrieval techniques to significantly reduce hallucinations and extend the practical applicability of LLMs across various domains.

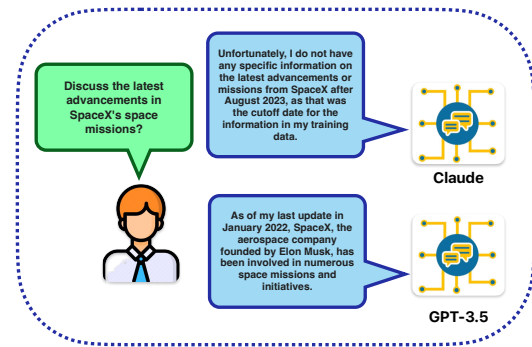


Figure 1: LLMs trained on older datasets are not capable of generating responses to queries involving recent recent developments.

1 Introduction

Large Language Models (LLMs) have revolutionized natural language processing (NLP) with their ability to perform tasks such as text generation, translation, summarization, and question-answering. Models like GPT-3.5 by OpenAI (Brown et al., 2020) and Claude by Anthropic (Bai et al., 2022) demonstrate impressive capabilities. However, LLMs still face significant challenges, particularly the issue of hallucinations—generating information that appears plausible but is factually incorrect or outdated (Bender et al., 2021; Liu et al., 2021; Maynez et al., 2020; Ji et al., 2023). This issue often arises from limitations in their training datasets, which may not cover the latest information or be sufficiently comprehensive.

The lack of efficient mechanisms to integrate external data exacerbates these issues, as current

methods frequently rely on expensive APIs (Thop-pilan et al., 2022). Moreover, running LLMs locally can lead to high RAM consumption, making it impractical for users with limited computational resources. To address these limitations, we propose HaloRAG, a web-based Retrieval-Augmented Generation (RAG) (Lewis et al., 2021a) approach, which enhances the interaction between user queries and LLMs by integrating real-time information from the web. Our method significantly reduces hallucinations by grounding responses to updated information, thus enhancing accuracy and relevance by retrieving the most pertinent documents, and offers a cost-effective, scalable solution through automated web scraping techniques. By continually incorporating fresh data, HaloRAG ensures that the underlying LLMs remain current and capable of addressing queries related to recent developments and emerging topics, extending their utility across various dynamic informational contexts. This approach represents a low-cost step towards making LLMs more reliable, accurate, and applicable in real-time scenarios. To mitigate high resource consumption, the approach is designed to be efficient, allowing it to run on local computers with minimal resources and enabling users to

067	benefit from enhanced LLM capabilities without	is critical for global applicability (Lample and Con-	114
068	requiring server infrastructure.	neau, 2019).	115
069	The key contribution of our work is a wrapper	Finally, the high computational costs associated	116
070	for LLMs with the following key features:	with large models highlight the need for efficient	117
071		scaling strategies. These strategies should balance	118
072	• Low computational overhead	model size, computational efficiency, and perfor-	119
073	• Ability to handle multilingual prompts	mance to allow for scalable and cost-effective AI	120
074	• Real-time data retrieval and integration	systems (Kaplan et al., 2020). The approach of	121
075	• Compatibility with various LLM architectures	refreshing LLMs with search engine augmentation	122
076	• Reduced hallucinations in generated re-	presents a promising direction, although it currently	123
077	sponses without the need for re-training with	faces challenges related to cost and language limi-	124
078	new data	tations (Vu et al., 2023).	125
079	2 Related Work	3 Methodology	126
080	Hallucination in natural language generation has	HaloRAG consists of five different modules and we	127
081	been identified as a major challenge, attributed to	describe each of it in this section. Figure 2 shows	128
082	static training datasets that fail to reflect the evol-	the architecture of our approach.	129
083	ving real-world knowledge (Ji et al., 2023). Dy-	3.1 Multilingual Support	130
084	namic updates to models are essential to mitigate	The inclusion of non-English speaking users is	131
085	this problem, with strategies that enhance model	achieved through the integration of a multilingual	132
086	architectures and improve training data quality	support system. To this end, the user queries are	133
087	being crucial (Liu et al., 2023). Real-time data	first translated into the primary operating language	134
088	retrieval and the use of adversarial training are	of the LLM and subsequently the LLM’s responses	135
089	proposed to enhance factual accuracy in dialogue	back into the user’s language. Bidirectional trans-	136
090	systems. Similar to our work, FreshLLMs (Vu et al., 2023)	lation ensures effective processing and response to	137
091	explores augmenting LLMs with web search	queries from a diverse global audience. The com-	138
092	capabilities. However, the approach diverges in	ponent is designed to create a seamless interaction	139
093	critical ways. The FreshLLMs framework relies on	loop, thereby minimizing language barriers. This	140
094	the Google search API, which incurs costs and	design enhances the usability and accessibility of	141
095	may not be accessible to all users. Moreover,	LLM technologies across various linguistic demo-	142
096	FreshLLM is limited to English language prompts,	graphics.	143
097	whereas HaloRAG supports multilingual queries,	3.2 Integration with LLMs	144
098	offering a more versatile solution.	The integration of large language models (LLMs)	145
099	The limitations of static datasets in LLMs, which	is designed to be straightforward and minimally	146
100	prevent access to up-to-date information and	invasive, allowing for the incorporation of various	147
101	result in outdated responses, have been address-	models without significant modifications to their	148
102	ed by integrating real-time data sources to	architectures. This approach ensures broad appli-	149
103	maintain model currency (Komeili et al., 2021). The	cability across different models and configura-	150
104	Retrieval-Augmented Generation (RAG) approach	tions, facilitating widespread adoption. The in-	151
105	has shown promise in improving LLM responses	tegration process supports any LLM, including	152
106	by using relevant documents to inform the	those hosted on platforms such as HuggingFace ,	153
107	generation process, which enhances both	and is compatible with both general-purpose	154
108	accuracy and contextual appropriateness	and specialized models.	155
109	(Lewis et al., 2021b).	For LLMs fine-tuned for specific applications,	156
110	Concerns about the dominance of the English	such as financial forecasting or medical re-	157
111	language in LLMs have led to calls for the	search, this integration provides access to	158
112	development of multilingual models to	current external information, enhancing per-	159
113	ensure inclusivity and accessibility (Bender et al., 2021). Cross-	formance and accuracy. This allows LLMs	160
	lingual model training has shown that LLMs	to respond to queries based not only on	161
	can effectively operate across different lan-	their pre-trained knowledge but also on the	

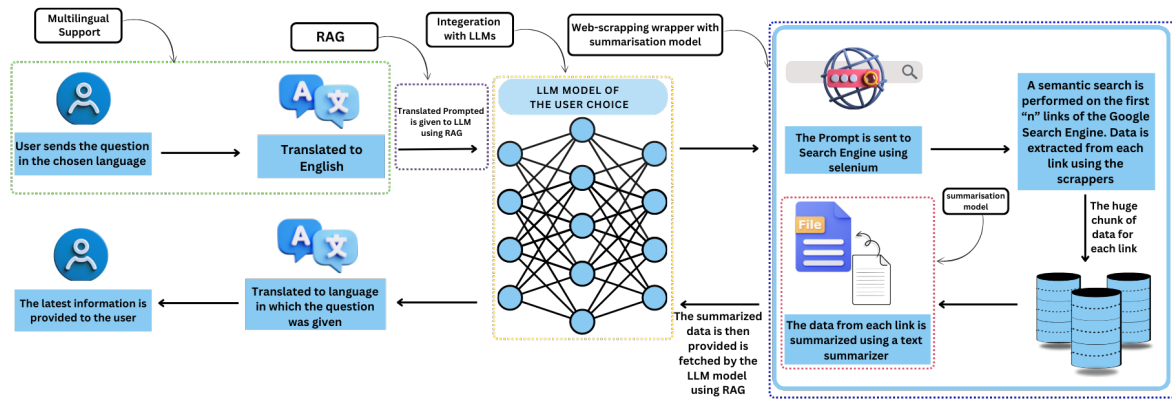


Figure 2: Architecture Diagram. Method for enhancing language model accuracy using web scraping and RAG. Starts with user queries in native languages, translation, web scraping, data summarisation, and re-integration into the LLM for refined responses.

latest available data. This capability is particularly valuable in dynamic and rapidly evolving information landscapes. The integration strategy interfaces with the model’s input and output processes, enabling real-time data enhancement before response generation. This ensures that LLMs remain current and reliable, addressing the limitations of static training datasets and expanding the utility of LLMs across various applications and industries.

3.3 Retrieval-Augmented Generation (RAG)

The Retrieval-Augmented Generation (RAG) component is designed to enhance the interaction between a user’s query and the LLM by integrating external data. This approach combines retrieval-based and generation-based methods to improve the accuracy and relevance of the LLM’s responses (Lewis et al., 2021b). Upon receiving a user query, the LLM processes the input to understand its context and intent. The RAG component then formulates a search query to retrieve relevant information from various online sources. This retrieved data is evaluated for relevance and credibility, ensuring that only authoritative and reliable information is used. The integration of retrieval and generation in RAG allows the LLM to provide responses that are not only based on static knowledge but are also enriched with real-time, contextually appropriate information. This dynamic interaction significantly reduces the incidence of generating hallucinations—responses that are plausible but factually incorrect—by continually updating the LLM with the latest information (Izcard and Grave, 2021). By continually integrating fresh data from the web, the RAG component ensures that the LLM remains current and capable of addressing queries about

recent developments and emerging topics.

3.4 Web Scraping Module

The web scraping module is designed to access recent and relevant information from the internet. Upon receiving a query, the module structures the input for search engines. The number of links processed is adjustable based on research depth and computational resource availability. The semantic analysis assesses the relevance of each link’s metadata to ensure the accuracy and substance of the information. Each identified link is open in a controlled browser session, extracting textual information. Semantic processing is applied to ensure the extracted text is contextually relevant. The scraped data from each link is aggregated into unified data, undergoing preliminary processing to eliminate duplicates, correct formatting issues, and prepare the data. This method ensures the retrieval of the most current and relevant information, enhancing the system’s overall effectiveness.

3.5 Summarization Models

The summarization models in the system are crucial for distilling extensive information retrieved from the web into concise, essential content that can be effectively utilized by large language models (LLMs). HaloRAG employs the FalconAI summarization model (Almazrouei et al., 2023) which is known for its efficiency and accuracy in processing large texts. FalconAI summarization model utilizes a trained network to extract key facts and themes, ensuring the summaries are both coherent and focused on the most relevant details. This process significantly reduces the informational load on LLMs, enabling quicker and more accurate re-

sponses. The system is designed with flexibility, allowing for the integration of other summarization models if required. This adaptability ensures continuous updating the approach to summarisation based on evolving technological advancements or specific application needs.

4 Experiments

To assess the performance of HaloRAG, a sample dataset consisting of factual prompts was curated. These prompts were specifically selected to cover recent developments (specifically from 2023 and 2024) in technology, politics, and related fields, hence allowing us to focus on the model’s ability to handle up-to-date content. More details can be found in Appendix A. We compare the responses from several models, including GPT-3.5 (Brown et al., 2020), GPT-4o (OpenAI et al., 2024), Claude (Bai et al., 2022), Phind (Rozière et al., 2024), and HaloRAG with Qwen (Bai et al., 2023) as the underlying LLM. For our experiments we used Google Colab with a RAM of 12.5 GB and GPU of 15 GB.

5 Results

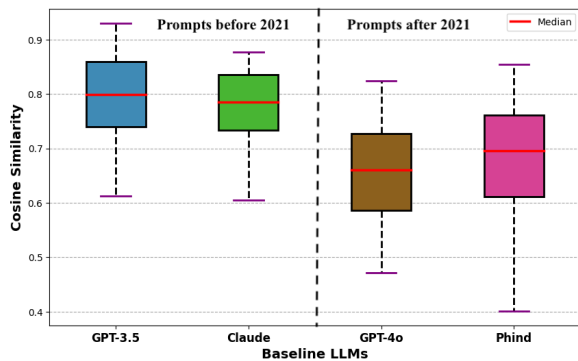


Figure 3: Cosine similarity between the responses generated by GPT-3.5, Claude, GPT-4o and Phind

Figure 3 shows the cosine similarity between the responses generated by our method (HaloRAG). The results show a median cosine similarity of 0.6599 with GPT-4o and 0.6956 with Phind, demonstrating textual and conceptual alignment between the responses generated. The highest cosine similarity recorded is 0.8240 with GPT-4o and 0.8545 with Phind whereas models like GPT-3.5 and Claude were unable to generate the responses for these same prompts (Fig 1). This shows that our model is able to generate responses to queries similar to those generated by state-of-the-art LLMs

without the requirement of re-training with new data. In order to measure the effectiveness of HaloRAG with pre-2021 queries, we evaluate the cosine similarity between our approach and GPT-3.5 and Claude on another dataset consisting of 30 prompts centered around information before 2021 (more details can be found in Appendix B). The results indicate a median cosine similarity of 0.7997 with GPT-3.5 and 0.7848 with Claude, demonstrating textual and conceptual alignment between the responses generated.

6 Conclusion and Future Work

We presented HaloRAG, an approach to enhancing Large Language Models (LLMs) by developing a cost-efficient, no-cost agentic wrapper that leverages web scraping technologies to perform semantic searches and retrieve real-time information from the web. This wrapper utilizes RAG to extend the knowledge base of any LLM it is paired with, enabling it to provide accurate and current answers without needing to be retrained on new datasets. The approach demonstrates significant improvements over existing models like GPT-3.5 and Claude, particularly in handling prompts about recent events or emerging technologies. Future work includes improving the reliability and computational speed of the wrapper and performing a comprehensive comparison with other baseline LLMs on standard datasets.

7 Limitations

While the wrapper-enhanced LLM offers significant advancements, it is essential to acknowledge certain limitations. One notable limitation is the dependency on web scraping technologies, specifically Selenium and BeautifulSoup (bs4). If the structure of the Google search results page changes, the web scraping components may fail to function correctly, leading to incomplete or inaccurate data extraction. This reliance necessitates frequent updates to the scraping logic to adapt to changes in the search engine’s HTML structure. The wrapper ensures that information is retrieved from verified sources, there remains a risk of encountering biased or misleading information. Users must critically evaluate the responses and cross-reference with other sources when necessary. In conclusion, while the wrapper-enhanced LLM reduces hallucinations and provides up-to-date information, it has limitations such as dependency on web scraping.

References

- 316 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Al-
317 shamsi, Alessandro Cappelli, Ruxandra Cojocaru,
318 M rouane Debbah,  tienne Goffinet, Daniel Hesslow,
319 Julien Launay, Quentin Malartic, Daniele Mazzotta,
320 Badreddine Noune, Baptiste Pannier, and Guilherme
321 Penedo. 2023. [The falcon series of open language
322 models](#). *Preprint*, arXiv:2311.16867.
- 323 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
324 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
325 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,
326 Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,
327 Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,
328 Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong
329 Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang
330 Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian
331 Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi
332 Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang,
333 Yichang Zhang, Zhenru Zhang, Chang Zhou, Jin-
334 gren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023.
335 [Qwen technical report](#). *Preprint*, arXiv:2309.16609.
- 336 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda
337 Askell, Anna Chen, Nova DasSarma, Dawn Drain,
338 Stanislav Fort, Deep Ganguli, Tom Henighan,
339 Nicholas Joseph, Saurav Kadavath, Jackson Kernion,
340 Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac
341 Hatfield-Dodds, Danny Hernandez, Tristan Hume,
342 Scott Johnston, Shauna Kravec, Liane Lovitt, Neel
343 Nanda, Catherine Olsson, Dario Amodei, Tom
344 Brown, Jack Clark, Sam McCandlish, Chris Olah,
345 Ben Mann, and Jared Kaplan. 2022. [Training
346 a helpful and harmless assistant with reinforce-
347 ment learning from human feedback](#). *Preprint*,
348 arXiv:2204.05862.
- 349 Emily Bender, Timnit Gebru, Angelina McMillan-
350 Major, and Shmargaret Shmitchell. 2021. [On the
351 dangers of stochastic parrots: Can language models
352 be too big?](#) pages 610–623.
- 353 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
354 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
355 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
356 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
357 Gretchen Krueger, Tom Henighan, Rewon Child,
358 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
359 Clemens Winter, Christopher Hesse, Mark Chen,
360 Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin
361 Chess, Jack Clark, Christopher Berner, Sam Mc-
362 Candlish, Alec Radford, Ilya Sutskever, and Dario
363 Amodei. 2020. [Language models are few-shot learn-
364 ers](#). *Preprint*, arXiv:2005.14165.
- 365 Gautier Izacard and Edouard Grave. 2021. [Lever-
366 aging passage retrieval with generative models
367 for open domain question answering](#). *Preprint*,
368 arXiv:2007.01282.
- 369 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan
370 Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea
371 Madotto, and Pascale Fung. 2023. [Survey of halluci-
372 nation in natural language generation](#). *ACM Comput-
373 ing Surveys*, 55(12):1–38.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B.
Brown, Benjamin Chess, Rewon Child, Scott Gray,
Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.
[Scaling laws for neural language models](#). *Preprint*,
arXiv:2001.08361.
- Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2021.
[Internet-augmented dialogue generation](#). *Preprint*,
arXiv:2107.07566.
- Guillaume Lample and Alexis Conneau. 2019. [Cross-
lingual language model pretraining](#). *Preprint*,
arXiv:1901.07291.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
rich K ttler, Mike Lewis, Wen tau Yih, Tim Rock-
t schel, Sebastian Riedel, and Douwe Kiela. 2021a.
[Retrieval-augmented generation for knowledge-
intensive nlp tasks](#). *Preprint*, arXiv:2005.11401.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
rich K ttler, Mike Lewis, Wen tau Yih, Tim Rock-
t schel, Sebastian Riedel, and Douwe Kiela. 2021b.
[Retrieval-augmented generation for knowledge-
intensive nlp tasks](#). *Preprint*, arXiv:2005.11401.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan,
Lawrence Carin, and Weizhu Chen. 2021. [What
makes good in-context examples for gpt-3?](#) *Preprint*,
arXiv:2101.06804.
- Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu
Chen, and Jiawei Han. 2023. [Understanding
the difficulty of training transformers](#). *Preprint*,
arXiv:2004.08249.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and
Ryan McDonald. 2020. [On faithfulness and fac-
tuality in abstractive summarization](#). *Preprint*,
arXiv:2005.00661.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,
Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
man, Diogo Almeida, Janko Altschmidt, Sam Alt-
man, Shyamal Anadkat, Red Avila, Igor Babuschkin,
Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-
ing Bao, Mohammad Bavarian, Jeff Belgum, Ir-
wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,
Christopher Berner, Lenny Bogdonoff, Oleg Boiko,
Madeline Boyd, Anna-Luisa Brakman, Greg Brock-
man, Tim Brooks, Miles Brundage, Kevin Button,
Trevor Cai, Rosie Campbell, Andrew Cann, Brittany
Carey, Chelsea Carlson, Rory Carmichael, Brooke
Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully
Chen, Ruby Chen, Jason Chen, Mark Chen, Ben
Chess, Chester Cho, Casey Chu, Hyung Won Chung,
Dave Cummings, Jeremiah Currier, Yunxing Dai,
Cory Decareaux, Thomas Degry, Noah Deutsch,
Damien Deville, Arka Dhar, David Dohan, Steve
Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,
Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,
Sim n Posada Fishman, Juston Forte, Isabella Ful-
ford, Leo Gao, Elie Georges, Christian Gibson, Vik

431	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeef Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Bar-		
	ret Zoph. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.		495 496
	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code llama: Open foundation models for code . <i>Preprint</i> , arXiv:2308.12950.		497 498 499 500 501 502 503 504 505 506
	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Agueras-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. 2022. Lamda: Language models for dialog applications . <i>Preprint</i> , arXiv:2201.08239.		507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527
	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. 2023. Freshllms: Refreshing large language models with search engine augmentation . <i>Preprint</i> , arXiv:2310.03214.		528 529 530 531 532

A Prompts based on information in 2023-2024	533
• Tell me about the current status of Palestine war?	534
• Discuss the latest developments in the Ukraine conflict?	535
• What are the current issues surrounding the U.S. immigration policy?	536
• Explain the controversy over the 2024 U.S. Presidential election.	537
• What are the recent advancements in renewable energy technology?	538
• What is the current status of the COVID-19 pandemic worldwide?	539
• Explain the controversy surrounding the recent tech layoffs?	540
• Discuss the impact of the recent stock market fluctuations due to recession?	541
• What are the latest trends in cryptocurrency regulations in 2024?	542
• Explain the current issues in global supply chain disruptions?	543
• Discuss the latest findings in climate change research from 2023 onwards?	544
• Explain the controversy over the new AI regulations in the EU since 2024?	545
• Discuss the recent developments in space exploration by India since 2023?	546
• Explain the impact of the recent data breaches occurred in AIMS India?	547
• Discuss the current debates on healthcare reform in the U.S. from January 2024?	548
• What is the current status of the trade war between the U.S. and China?	549
• Explain the controversy surrounding the recent Supreme Court decisions about the farmers protest?	550
• Explain the controversy revolving around NVIDIA and Jensen Huang in 2024?	551
• What are the latest advancements in cancer research?	552
• Explain the current issues in the global refugee crisis?	553
• Discuss the implications of the latest tech company mergers?	554
• Discuss the controversy of Google using Reddit for search AI?	555
• Tell me about the fight going between Elon and Mark?	556
• Tell me about the fight going between Elon and Yann LeCun?	557
• Explain the impact recent cyclone that hit the east coast of India?	558
• Discuss the current debates on WhatsApp shutting down its service in India?	559
• What is the current status of the opioid crisis in the U.S.?	560
• Discuss the latest advancements in SpaceX's space missions?	561
• Discuss the current debates on gun control laws in the U.S.?	562
• Tell me about the controversy of Devin-AI?	563

B Prompts based on information before 2021

- What were the key elements of the Paris Climate Agreement?
- Describe the impact of Brexit on the European Union's trade policies.
- How did the US-China trade war begin, and what were its major impacts?
- What are the principles of the Green New Deal proposed in the United States?
- How has artificial intelligence impacted healthcare in the last decade?
- Discuss the role of NATO in the 21st century.
- What were the main issues during the verbal confrontations between the US and North Korea under the Trump administration?
- How did the Fukushima nuclear disaster affect energy policies worldwide?
- What are the ongoing challenges in managing global plastic pollution?
- Describe the rise of electric vehicles and their impact on the global oil industry.
- What strategies have been effective in combating deforestation in the Amazon rainforest?
- How has the gig economy transformed traditional employment models?
- What were the significant outcomes of the COP26 summit?
- How do cryptocurrency regulations differ around the world?
- Describe the impact of remote work on urban and suburban development.
- What are the challenges and successes of the Mars Rover missions?
- How did the Arab Spring reshape politics in the Middle East?
- What are the main goals of the United Nations Sustainable Development Goals (SDGs)?
- How has online education evolved apart from the COVID-19 pandemic?
- What are the implications of quantum computing for data security?
- How have drones been integrated into commercial and military operations?
- What are the ethical considerations of gene editing technologies?
- How has social media influenced political campaigns in the 21st century?
- What were the causes and consequences of the global financial crisis of 2008?
- Discuss the impact of the MeToo movement on global workplace policies.
- What are the technological advancements in renewable energy in the last five years?
- How did the Venice Flood Barriers help combat rising sea levels?
- What are the major factors driving urban sprawl in major cities worldwide?
- How have international space laws evolved with the advent of private space travel?
- What are the diplomatic challenges faced by countries in the Arctic as ice melts?

C Disclaimer

We used ChatGPT to rephrase some of the sentences in this paper. But we ensured (to the best of our extent) that all the content in the paper was original.