Towards Safe Multilingual Frontier AI

Arturs Kanepajs* ERA Fellowship akanepajs@gmail.com Vladimir Ivanov* ERA Fellowship volodimir1024@gmail.com Richard Moulange University of Cambridge rjm246@cam.ac.uk

Abstract

Linguistically inclusive LLMs—which maintain good performance regardless of the language with which they are prompted—are necessary for the diffusion of AI benefits around the world. Multilingual jailbreaks that rely on language translation to evade safety measures undermine the safe and inclusive deployment of AI systems. We provide policy recommendations to enhance the multilingual capabilities of AI while mitigating the risks of multilingual jailbreaks. We examine how a language's level of resourcing relates to how vulnerable LLMs are to multilingual jailbreaks in that language. We do this by testing five advanced AI models across 24 official languages of the EU. Building on prior research, we propose policy actions that align with the EU legal landscape and institutional framework to address multilingual jailbreaks, while promoting linguistic inclusivity. These include mandatory assessments of multilingual capabilities and vulnerabilities, public opinion research, and state support for multilingual AI development. The measures aim to improve AI safety and functionality through EU policy initiatives, guiding the implementation of the EU AI Act and informing regulatory efforts of the European AI Office.

1 Introduction

Despite rapid advances in large language model (LLM) capabilities [36, 52, 60], frontier LLMs continue to be vulnerable to *jailbreaks*, which undermine safety measures and enable malicious actors to misuse AI systems to cause harm [34, 82, 66]. Defensive measures against jailbreaks can reduce the risks, but can also impede model utility through rejection of benign prompts [14].

One subclass of jailbreaks are *multilingual jailbreak attacks*, where a model is prompted in a different—often low-resource²—language to circumvent safety systems that would otherwise activate in response to a default—usually high-resource—language. Here, it is precisely that LLMs have strong, though inconsistent, multilingual capabilities, understanding and responding to instructions in many languages, that facilitates the attacks.

However, multilingual capabilities of LLMs are essential to diffuse the benefits of AI throughout our societies [61, 16]: although LLMs are predominantly trained on English text [38, 80], most humans on the planet do not speak English [69]. Unfortunately, several of measures designed to defend

^{*}Research conducted as part of the ERA:AI Fellowship.

¹ *Jailbreaking* can be defined as the process of circumventing the safety measures placed on LLMs and other AI systems [51]. Vulnerability to jailbreaks paired with powerful model capabilities, including offensive cyber capabilities, chemical, biological, radiological, and nuclear capabilities, contribute to systemic risks [78, 31].

²Low-resource languages, which account for over 90% of the world's 7,000+ languages and are spoken by 1.2 billion people, have limited labeled and unlabeled data available [81, 43]. In this paper we follow Bang et al. (2023) [11] and classify languages based on CommonCrawl corpus share: high-resource languages have over 1%, medium-resource languages 0.1-1%, and low-resource languages less than 0.1%. This method is responsive to the availability of digital language data, with the source data updated monthly.

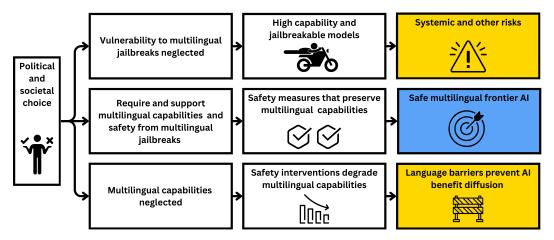


Figure 1: **Threat model and policy opportunity.** Neglecting multilingual jailbreaks or multilingual capabilities can increase risks or limit AI benefit diffusion. Conversely, appropriately addressing the risks as well as capabilities can bring the benefits of safe multilingual frontier AI.

against multilingual jailbreaks—such as instructing the model to "think in English" or self-generating synthetic data in low-resource languages—limit the model's utility in non-English natural languages [76, 23], undermining inclusive AI benefit-sharing.

This paper proposes a path towards linguistically inclusive AI that is also safe from multilingual jailbreak attacks (Figure 1). The rest of the paper proceeds as follows. First, we discuss related work, paying particularly attention to the extent defensive actors can reduce LLM vulnerabilities to multilingual jailbreaks without compromising multilingual capabilities. Second, we show that the EU offers a particularly promising environment for effective policy action to address this issue. Third, we assess the vulnerability to multilingual jailbreaks and multilingual capability gaps in the 24 official EU languages for five frontier LLMs. Finally, we offer specific policy recommendations for the European AI Office and other relevant stakeholders.

2 Related work

2.1 What to measure?

Surveys offer valuable insights into public opinion on AI [40], yet they often fall short in identifying which AI capabilities are most beneficial to society, and in determining appropriate safety thresholds that balance individual and collective interests. This limitation stems from the complexity of AI and its broad societal implications. To address this, researchers have introduced innovative approaches, such as citizens' juries, to gather more informed public perspectives [42]. These methods can help to achieve a balance between innovation and safety, ensuring that public interest plays a central role alongside private incentives [65].

One example for societal choice is between technological automation and job security [71]. For instance, LLMs could yield significant savings by automating certain tasks, and addressing skill gaps [39, 36]. However, increased automation also raises concerns about labor displacement. According to the Ipsos 2024 survey, over a third of workers worldwide fear that AI could replace their jobs in the coming years [40]. Workers in low-income countries and emerging markets may face significant challenges as AI-driven automation leads to the onshoring of jobs in advanced economies [50]. Beyond the immediate loss of income, job displacement can result in political disempowerment and the loss of personal meaning associated with work [5, 71]. Therefore, decisions about the training, deployment, and use of LLMs should involve not only tech developers but also civil society and policymakers [8]. Measuring and predicting the downstream impacts of LLMs is a complex and open research problem [62, 10], progress in which is necessary to allow for informed choice between different paths of development.

2.2 Risk and capability measurement

To assess vulnerabilities to multilingual jailbreaks, several benchmarks have been used in the literature, including AdvBench [80, 81], MasterKey [80], and MT-Bench [67]. Although precise assessment is currently hampered by data contamination and translation imperfections, these challenges do not seem insurmountable. Human annotators can help resolve translation issues in both capability [55] and safety assessments [4]. Recent studies also address contamination by measuring capabilities when context, but not the specific question, is provided in the prompt [12].

2.3 Risk mitigation and capability improvements

Some risk mitigation measures for addressing multilingual jailbreak vulnerabilities can come at the cost of model capabilities. For example, instructing the model to "think in English" [76] may enhance safety but reduce effectiveness in language-specific contexts [72]. Similarly, a "self-defense" approach, which generates multilingual training data for safety fine-tuning, has been found to increase rejection rates for benign prompts [23]. Additionally, safety fine-tuning can result in shorter responses [48]. At the extreme, complete usage restriction maximizes safety but eliminates utility: "a model that always refuses is maximally safe, but not helpful in the slightest" [25]. While there are no documented cases of capabilities being entirely cut off in specific languages, certain modalities, such as image generation [77] and singing [55], have been curtailed.

Despite these challenges, progress is possible. Evidence shows that language gaps can be significantly reduced, even in low-resource languages [16]. For instance, GPT-40 recently demonstrated notable capability improvements in several low-resource African languages [55]. Measures that preserve both safety and capabilities include generating human-annotated datasets, though this can be costly, as low-resource languages often require more tokens per word [6, 13]. One study demonstrated that using just 900 prompts—half requiring local knowledge—reduced the jailbreak attack rate by more than half [4]. The cost of generating such prompts for the 24 official EU languages likely represents only a small fraction of the cost of developing frontier models, which currently exceeds \$100 million and is projected to reach \$1 billion by 2027 [21].

2.4 Who will do the evaluations?

Anderljung et al. advocate for the involvement of external actors in evaluating LLMs to ensure objectivity and thoroughness [8]. Casper et al. further argue that rigorous AI audits require more than just black-box access. They propose that white-box access, which provides deeper insight into the system's internals, enables stronger adversarial testing and fine-tuning. This approach allows for a more comprehensive assessment of post-deployment risks, ultimately enhancing the reliability and safety of AI systems. Additionally, the authors advocate for "outside-the-box" access, which includes access to training data, methodological details, and development processes. They argue that this level of access is crucial for tracing problems to their sources, assessing fair use concerns, and identifying potential biases or harmful content in training data. Such comprehensive access enables auditors to conduct more targeted evaluations and hold developers accountable for their decision-making processes and risk management strategies. [15].

2.5 Comprehensive empirical assessments and policy recommendations

Dubey et al. examined the Violation Rate (VR) and False Refusal Rate (FRR) for the Llama 3.1 405B model and its competitors in eight high-resource languages, finding that "Llama Guard" reduced VR with minimal or insignificant FRR increases [25]. Similarly, Cohere For AI introduced a dataset and method to address multilingual jailbreaks across eight languages, reporting results for six high-resource languages using the Aya 23 8B model [4]. In contrast, we report results for the whole suite of the 24 EU languages—including low-resource ones—on five frontier LLMs (see Section 4.4).

Cohere For AI provides broad policy recommendations for addressing multilingual safety vulner-abilities and improving transparency in multilingual model performance [16]. We have tried to go further, by specifically identifying responsible stakeholders and tailoring our recommendations accordingly. We also advocate for model providers to report evaluation results rather than conducting the evaluations themselves (see Section 5).

In the EU context, Rehm et al. [61] compiled a wide set of policy recommendations to achieve digital language equality by 2030. Many of their recommendations call for more funding for research and innovation, increased multilingual data availability, and breakthroughs needed for model evaluation. They also called for training large language models from scratch: such efforts have by now received significant compute resources,³ but have also been criticized as likely to produce inadequate models [44]. Our comparative contribution is focus on highly capable AI models in general, multilingual jailbreaks in particular, and offer recommendations to enable mandatory capability and safety assessments.

3 Methodology for the quantitative analysis

To assess the relationship between jailbreak attack success rates, harmless prompt rejection rates, and the level of language resourcedness, we formulate the following hypotheses:

- H1: Multilingual jailbreak attack success rates and language resourcedness are negatively related.
- H2: Harmless prompt rejection rates and language resourcedness are negatively related.

To test these hypotheses, we employ the OR-Bench dataset introduced by Cui et al. (2024) [22]. This dataset is particularly well-suited for our purposes as it was specifically designed to address the issue of 'Over-Refusal' by LLMs. We utilize two key subsets of OR-Bench: "toxic" and "hard-1k". The "toxic" subset contains prompts designed to test model safety guardrails, while "hard-1k" represents edge cases where models might incorrectly refuse legitimate requests.

Our sampling procedure implements the following steps:

- Load the complete "toxic" and "hard-1k" subsets from OR-Bench
- Apply a fixed random seed (42) to ensure reproducibility
- Randomly shuffle all prompts within each subset
- Select the first 100 prompts from each shuffled set

The translation of these prompts into non-English languages is performed using Google Translate, a widely accessible service that supports all official EU languages.

After collecting the LLM outputs, we classify the responses using GPT-40 mini (gpt-40-mini-2024-07-18), based on its suitability for response classification. We classify responses into three categories: *ACCEPTED*, *REJECTED*, and *UNCLEAR*, using annotation guidelines adapted from Yong et al. [81]. The classification rules and detailed annotation guidelines are provided in Appendix B.

Next, we estimate the coefficients using a logistic regression model, which follows the relationship:

$$P(Y=1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}} \tag{1}$$

In the above equation Y is binary: for **H1**, Y = 1 if and only if response is *ACCEPTED*; for **H2**, Y = 1 if and only if response is *REJECTED*. X denotes CommonCrawl corpus share [20], a proxy for language resourcedness [47].

To validate the econometric model, we also estimate the relationship using Ordinary Least Squares (OLS) regression. To use OLS, we calculate the proportion of *ACCEPTED* responses to harmful prompts, and the proportion of *REJECTED* responses to harmless prompts.

³For example, in June 2024, technology company *Tilde* announced they had won public compute resources comparable to those used to train GPT-3.5, in order to train a new linguistically inclusive model, see Appendix A

⁴Classification quality is another testable assumption. Previous research by the OR-Bench dataset authors indicated minimal discrepancies (2.4%) between classifications performed by GPT-3.5-turbo-0125 and GPT-4. However, it has been found that automatic classification can have lower precision than human classification [49].

4 Case Study: The EU and its 24 official languages

Language diversity, commitment to digital language equality [61] and recent regulatory developments in the EU present a pressing need for policies that promote safe multilingual frontier AI.

4.1 Legal aspects

Linguistic inclusivity is a legal requirement for state institutions in the EU, and is enshrined in several key legal instruments. In particular, in the Charter of Fundamental Rights of the EU, Article 21 prohibits discrimination on the grounds of language, and Article 22 mandates that the EU respect cultural, religious, and linguistic diversity [1]. Therefore, a failure to ensure that LLMs adhere to fundamental safety standards across all languages could potentially lead to violations of Article 21.

Moreover, the preamble to the EU AI Act also explicitly states that AI and its regulatory framework must be developed in accordance with the charter [31]. The consultations for the EU AI Act [18] as well as multi-stakeholder consultation launched by the EU AI Office [17] offer a platform for discussing and iterating on the requirements for linguistic inclusivity.

4.2 Who benefits from multilingual AI in the EU, and how?

Multilingual AI offers substantial potential for cost savings. EU legislation and documents of major public interest are produced in all 24 official languages, which costs approximately 1 billion euros annually [45]. One estimate suggests that reliable grammatically correct interpretation and translation of documents could increase civil service efficiency by 20–30% [53]. The benefits for everyday users are highlighted by the fact that English is the first "mother tongue" for less than 3% of the EU population (Table 1 and Appendix C).

Table 1: Speakers in the EU with the respective language the first "mother tongue" [74]. Resourcedness is determined by CommonCrawl corpus share: high-resource languages (>1%), medium-resource (<0.1-1%), and low-resource (<0.1%). See Appendix C for the resourcedness of specific languages.

Language group	Speakers (million)
English	9.2
other high-resource languages	326.9
medium-resource languages	83.3
low-resource languages	2.2

4.3 Who could implement the policies?

The primary responsibility for implementing the changes could fall to the **Regulation and Compliance Unit** of the **EU AI Office**. This unit could mandate reporting on model multilingual capabilities and susceptibility to multilingual jailbreak attacks as well as monitor and enforce compliance. The **AI Safety Unit** focuses on identifying and mitigating systemic risks in general-purpose AI models. Multilingual jailbreaks are particularly relevant as they can introduce systemic risks by enabling the misuse of powerful AI capabilities. The **AI Safety Unit** can therefore support the **Regulation and Compliance Unit** by providing expertise on multilingual capability and safety testing. The **AI for Societal Good Unit** can organize and promote initiatives such as citizens' juries, virtual assemblies, and surveys to gather insights into the social benefits and risks of AI. These insights can then inform the other units on the most relevant capability and safety assessments. Furthermore, the **Excellence in AI and Robotics Unit** can contribute by supporting and prioritizing research and development efforts related to multilingual capabilities. The **AI Innovation and Policy Coordination Unit** can facilitate collaboration with member states and international partners. Finally, the **Scientific Panel** of independent experts, also a part of the **EU AI Office**'s structure, can provide insights into the feasibility of measures that advance multilingual model capabilities and safety [27, 29].

Furthermore, the European Centre for Algorithmic Transparency [28] can provide expertise and contribute to the assessment of risks. Finally, **DG Connect** [24] of the European Commission (EC),

along with the broader EC, coordinates the digital strategy of the EU and can ensure that multilingual AI safety and capabilities are prioritized.

4.4 Regression results

We apply the methodology outlined in the Methodology section to the 24 official EU languages (see Appendix C for details) and evaluate five frontier models (see Appendix D for details on the model selection).⁵ Specifically, we run 100 harmful prompts and 100 harmless prompts in each language and classify the responses as *ACCEPTED*, *REJECTED* or *UNCLEAR*.⁶

We tested the two assumptions identified in the Methodology section. First, we assess translation issues by manually considering a subset of prompts. Specifically, we identify that 2 of the 19 "Harmful" prompts that were accepted in Latvian, had actually become harmless through translation (see Appendix E). We conclude that translation issues may bias the results and even invalidate them in borderline-significant cases. Second, we compared the performance of GPT-40 mini to GPT-40 across a subset of responses across the tested languages and found no significant differences. However, we do not compare LLM and human classification, which is a limitation of our research and warrants further investigation.

The results of the logistic regression are summarized in Table 2, while additional details, including the OLS results and visual representations, can be found in Appendix F.⁷

Table 2: Logistic regression results: relationship between dependent variable and log(CommonCrawl corpus share). 2400 observations across 24 languages (100 per language) for each regression, with Benjamini–Hochberg adjusted p-values. Adjusted significance levels: *, **, *** represent 5%, 1%, and 0.1%, respectively.

	H	Harmful Accepted			Harmless Rejected			
Model	β_0	β_1	p_{adj}	Sig.	β_0	β_1	p_{adj}	Sig.
Claude 3.5 Sonnet	-4.60	-0.10	0.369	***	0.89	-0.01	0.641	
GPT-40 Gemini 1.5 Pro	-2.71 -4.33	-0.15 -0.10	0.000 0.367	***	-0.52 2.06	0.03 0.02	0.316 0.551	
Llama 3.1 405B Mistral Large 2	-3.07 -2.53	0.06 -0.19	0.369 0.000	***	0.21 -0.46	0.03 -0.04	0.316 0.133	

The first hypothesis is supported (i.e., the null hypothesis is rejected) for two out of the five models, specifically GPT-40 and Mistral Large 2. Namely, jailbreak attack success rates tend to be higher for low-resource languages. The second hypothesis is not supported for any of the models after Benjamini-Hochberg adjustment. That is, we find no evidence that harmless prompts tend to be rejected more often in less resourced languages.

5 Discussion

Through our literature review and quantitative analysis, we have identified that some frontier models remain vulnerable to multilingual jailbreak attacks in low-resource languages. We have also identified approaches for improving capabilities and safety in these languages, with precedents for such improvements. Additionally, we have established that there is a legal and financial rationale for the EU to support multilingual safety and capabilities. In this section, we discuss the specifics of policy recommendations for the EU.

⁵The models are, in alphabetical order, Claude 3.5 Sonnet (claude-3-5-sonnet-20240620), Gemini 1.5 (gemini-1.5-pro-01), GPT-40 (gpt-40-2024-05-13), Meta Llama 405B (Meta-Llama-3.1-405B-Instruct-Turbo) and Mistral Large 2 (mistral-large-2407).

⁶The total estimated API costs for these runs were \$134.54. The additional compute costs for the full research project, e.g. including assumption testing, likely required less than 50% of this, bringing the total within \$200.

⁷Dataset and code for API prompts, assumption testing, and econometric analysis are located at: https://github.com/akanepajs/multilingual. We do not publicly release the translated material, to avoid negative impacts from other models being trained on potentially harmful material, especially in low-resource languages.

5.1 Capability and Safety Requirements to Introduce

Firstly, to limit the regulatory burden, we recommend, at least initially, focusing on general-purpose AI models entailing systemic risks (see also Appendix D for discussion on the affected models) [31].

Secondly, while numerous multilingual capability and safety benchmarks exist, there is currently no authoritative and unified evaluation framework [61, 38, 36], so it may be premature to mandate compliance to a specific benchmark and to a specific level.

However, transparency requirements could serve as an important first step with immediate benefits. Transparency allows consumers to compare model capabilities in the language they are interested in, supporting for a better consumer choice [41]. Moreover, increased transparency can stimulate competition among developers to improve capabilities and safety in underrepresented languages. In the EU context, the reporting requirements can, at least initially, prioritize the 24 official EU languages, reflecting their special status and the potential for substantial cost savings in translation and communication.

Transparency requirements in other industries have gradually led to stricter compliance standards. For instance, EU transparency regulations for publicly listed companies [3] and emissions disclosure requirements have evolved into more stringent limits on carbon allowances [2]. Similarly, as standards for multilingual AI capabilities and safety become more established, the requirements could be strengthened to reach specific benchmark levels. This expansion could eventually cover additional languages and dialects, such as the more than 60 national and regional languages of the EU [30].

For safety assessments, we recommend that the requirements address multilingual jailbreak vulnerabilities across all natural languages, not just the official EU languages. This two-legged approach with requirements for both capability and safety assessment reporting would create incentives for developing frontier AI that can resist multilingual jailbreak attacks in any language while maintaining high performance in the 24 EU languages.

Moreover, these capability and safety reporting requirements are likely to produce a "Brussels Effect." [68, 79] Brussels effect may manifest *de facto*, because in practice the same frontier models are made available globally. Users outside the EU could just as well benefit from transparency about model capabilities, as well as from lower risks from multilingual jailbreaks. Additionally, *de jure* Brussels Effect could emerge as other jurisdictions adopt similar regulatory standards. Global forums, such as the upcoming AI Action Summit [83], can further contribute to shaping an inclusive framework for international AI governance.

5.2 State Support

Given the "non-excludable" and "non-rivalrous" nature of publicly available data and algorithms, they can be considered public goods with positive externalities [63, 35]. Therefore, we advocate for initiatives supporting multilingual dataset creation, particularly for low-resource languages. High-quality dataset creation is already a top priority for the EU [61], exemplified by the Alliance for Language Technologies EDIC (ALT-EDIC), involving 16 EU member states [19].

In multilingual safety evaluation, the "European LLM Leaderboard" provides an automated database for evaluating LLMs through collaboration between private and public partners, including Dresden University of Technology [58, 64]. While this leaderboard currently only presents results for smaller models and lacks coverage of some low-resource languages and multilingual jailbreak benchmarks [32], such developments could inform future standards and regulations [16].

It has also been suggested that the European High Performance Computing Joint Undertaking (EuroHPC JU) redirect its computational resources toward safety research, leveraging the growing AI Assurance industry [46, 44]. Our recommended requirements could increase demand for model safety assessment services, strengthening the case for safety-oriented research and development.

6 Limitations and future research

Our study faces several key methodological limitations that affect both the validity and generalizability of our findings. A central limitation is our reliance on prompts translated from English, which may miss culture-specific harms and nuances across different languages. The absence of native speakers

in prompt development and validation may have introduced cultural biases, while different linguistic structures across languages could affect both safety measures and model responses. These translation effects are evidenced in our analysis, where automated translation failed to preserve the original intent of prompts in some cases (see Appendix E). The lack of human evaluation for prompt classification represents another significant limitation, along with uncertainties regarding the quality of automatic translations and classification of prompts (see Methodology and Regression Results). Although we tested these assumptions and found them approximately true, more work is needed to quantify their impact. The significance of the results may also be affected by the limited number of prompts per language (100 prompts), which was driven by resource constraints.

While we define clear thresholds for language resourcedness (high >1%, medium 0.1–1%, low <0.1% of CommonCrawl corpus share), these cutoffs merit further validation. The relationship between corpus share and actual language representation in model training requires deeper investigation. Another potential issue that can affect model performance in different languages is the contamination of data used for model training [12]. Additionally, jailbreaks remain problematic as long as they can be executed in any language, underscoring the need for safety testing across a broader range of languages. Nevertheless, we believe the analysis and assumption testing conducted are sufficient as a proof-of-concept and strong enough to suggest a significant relationship between multilingual jailbreak vulnerability and language resourcedness for two of the models.

Beyond these technical limitations, our proposed policy framework faces several implementation challenges. The framework assumes effective coordination between multiple stakeholders who may have diverging interests and priorities, but does not specify mechanisms for resolving conflicts between these stakeholders. While the framework relies on capturing user and public preferences, it doesn't detail how to resolve conflicting preferences between different user groups or how to weigh economic benefits against safety risks. The enforcement aspects face practical limitations, including constraints on regulatory capacity and the risk of "checkbox compliance" without meaningful safety improvements. Moreover, measuring outcomes presents significant challenges, including difficulty in attributing reduced risks specifically to framework implementation and limited ability to assess counterfactuals.

These limitations suggest several areas for future research. The statistical analysis can be extended by comparing harmful prompt rejection rates and harmless prompt acceptance rates between different languages and language groups. For instance, comparing results for English versus non-English languages, and high-resource languages vs. other languages. It may be important to increase the sample size for the statistical power of such analysis. Future work should prioritize developing culture-specific harm datasets with native speaker input and implementing systematic human evaluation protocols. It is important to note that, apart from natural languages, models can be vulnerable in various non-natural (e.g. programming) languages, which also warrant safety testing [81]. Even less explored is model safety assessment to prompts in "hybrid" languages: with 7000 natural languages, there could be millions of "interpolated" languages, which could render all-encompassing black-box safety assessments infeasible.

In addition, the *persuasion* capabilities of models like GPT-40 have been found to cross important thresholds [55]. The interaction between multilingual capabilities and persuasiveness raises concerns, particularly regarding the use of LLMs to influence political opinions [37]. Future research should also focus on developing metrics for measuring policy effectiveness and creating frameworks for incorporating cultural expertise in safety assessments, while studying methods for balancing innovation and safety across different linguistic contexts.

7 Broader impacts statement

The work and analysis presented in this paper do not reveal any novel vulnerabilities, but rather demonstrates how known vulnerabilities may persist as models advance. While reporting requirements for multilingual capabilities and jailbreak vulnerabilities could potentially burden innovation, this impact should be minimal since only the largest models would be affected and initial requirements focus on transparency rather than compliance. The benefits from safer multilingual frontier models likely exceed assessment costs substantially, though detailed cost-benefit analysis would be needed to evaluate broader implementation.

8 Recommendations for EU policy

8.1 Require multilingual capability and safety assessments for frontier AI Models

- Require reporting on model capabilities across the 24 official EU languages for generalpurpose AI models entailing systemic risks.
- Require reporting on model susceptibility to multilingual jailbreak attacks across all languages for general-purpose AI models entailing systemic risks.
- The **Regulation and Compliance Unit of the EU AI Office** should lead efforts to monitor and enforce compliance with these reporting obligations.
- Initiatives such as the EU AI Act Code of Practice consultations and the multi-stakeholder consultations by the European AI Office offer a unique opportunity to discuss, refine, and introduce such requirements.
- Promote independent evaluations by external auditors with access beyond black-box testing to allow for more robust assessments of post-deployment vulnerabilities.
- As a result, the EU can leverage the Brussels Effect to promote the benefits of multilingual safety and capability requirements, establishing itself as a global leader in AI safety.

8.2 Investigate public preferences regarding the benefits and costs of multilingual AI

- The AI for Societal Good Unit of the EU AI Office can play a central role in identifying which capabilities the public values, who could be harmed by AI risks.
- Public preferences can be assessed through citizens' juries, moderated virtual assemblies, and surveys.

8.3 Provide state support for multilingual capabilities and safety

- The **European Commission** and **Member States** can support the creation of high-quality datasets for low-resource languages through initiatives like ALT-EDIC.
- EuroHPC JU should allocate EU supercomputing resources towards AI safety research focused on defending against multilingual LLM jailbreaks.
- Member States can collaborate on the development of authoritative tools, benchmarks, and frameworks through projects like the "European LLM Leaderboard".

The framework with the key stakeholders involved in the work on the relevant requirements is specified in Figure 2.

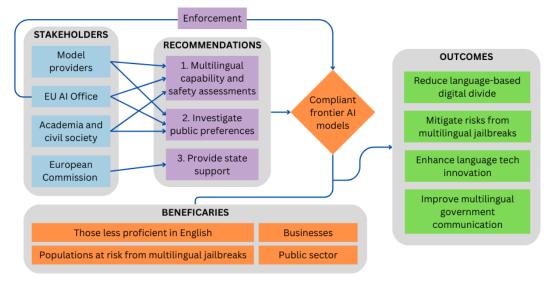


Figure 2: Framework for involvement of stakeholders to produce policy requirements that lead to safe multilingual frontier AI.

References

- [1] Charter of Fundamental Rights of the European Union, 2012.
- [2] ESG Laws and Regulation, 2024.
- [3] Transparency requirements for listed companies, 2024.
- [4] Aakanksha, Arash Ahmadian, Beyza Ermis, Seraphina Goldfarb-Tarrant, Julia Kreutzer, Marzieh Fadaee, and Sara Hooker. The Multilingual Alignment Prism: Aligning Global and Local Preferences to Reduce Harm, 2024. _eprint: 2406.18682.
- [5] Daron Acemoglu and Simon Johnson. *Power and Progress: Our Thousand-Year Struggle Over Technology and Prosperity*. Basic Books, hardcover edition, 2023.
- [6] Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo Kasai, David R. Mortensen, Noah A. Smith, and Yulia Tsvetkov. Do All Languages Cost the Same? Tokenization in the Era of Commercial Language Models, 2023. _eprint: 2305.13707.
- [7] Mistral AI. Mistral's Large 2407 Model Announcement, 2024.
- [8] Markus Anderljung, Everett Thornton Smith, Joe O'Brien, Lisa Soder, Benjamin Bucknall, Emma Bluemke, Jonas Schuett, Robert Trager, Lacey Strahm, and Rumman Chowdhury. Towards publicly accountable frontier llms: Building an external scrutiny ecosystem under the aspire framework, 2023. eprint: 2311.14711.
- [9] Anthropic. What are some things I can use Claude for?, 2024.
- [10] Usman Anwar, Abulhair Saparov, Javier Rando, Daniel Paleka, Miles Turpin, Peter Hase, Ekdeep Singh Lubana, Erik Jenner, Stephen Casper, Oliver Sourbut, Benjamin L. Edelman, Zhaowei Zhang, Mario Günther, Anton Korinek, Jose Hernandez-Orallo, Lewis Hammond, Eric Bigelow, Alexander Pan, Lauro Langosco, Tomasz Korbak, Heidi Zhang, Ruiqi Zhong, Seán Ó hÉigeartaigh, Gabriel Recchia, Giulio Corsi, Alan Chan, Markus Anderljung, Lilian Edwards, Yoshua Bengio, Danqi Chen, Samuel Albanie, Tegan Maharaj, Jakob Foerster, Florian Tramer, He He, Atoosa Kasirzadeh, Yejin Choi, and David Krueger. Foundational Challenges in Assuring Alignment and Safety of Large Language Models, 2024. eprint: 2404.09932.
- [11] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity, 2023.
- [12] Andrew M. Bean, Simi Hellsten, Harry Mayne, Jabez Magomere, Ethan A. Chi, Ryan Chi, Scott A. Hale, and Hannah Rose Kirk. LINGOLY: A Benchmark of Olympiad-Level Linguistic Reasoning Puzzles in Low-Resource and Extinct Languages, 2024. _eprint: 2406.06196.
- [13] Toms Bergmanis. TL;DR: If you use Chat GPT-4 and your native language is not English, you are doing a disservice to yourself. Why?, 2024.
- [14] Manish Bhatt, Sahana Chennabasappa, Yue Li, Cyrus Nikolaidis, Daniel Song, Shengye Wan, Faizan Ahmad, Cornelius Aschermann, Yaohui Chen, Dhaval Kapil, David Molnar, Spencer Whitman, and Joshua Saxe. Cyberseceval 2: A wide-ranging cybersecurity evaluation suite for large language models, 2024.
- [15] Stephen Casper, Carson Ezell, Charlotte Siegmann, Noam Kolt, Taylor Lynn Curtis, Benjamin Bucknall, Andreas Haupt, Kevin Wei, Jérémy Scheurer, Marius Hobbhahn, Lee Sharkey, Satyapriya Krishna, Marvin Von Hagen, Silas Alberti, Alan Chan, Qinyi Sun, Michael Gerovitch, David Bau, Max Tegmark, David Krueger, and Dylan Hadfield-Menell. Black-Box Access is Insufficient for Rigorous AI Audits. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, pages 2254–2272, New York, NY, USA, 2024. Association for Computing Machinery. event-place: Rio de Janeiro, Brazil.
- [16] Cohere For AI team. Policy Primer The AI Language Gap, 2024.

- [17] European Commission. AI Act: Have Your Say on Trustworthy General Purpose AI, 2024.
- [18] European Commission. AI Act: Participate in Drawing the First General Purpose AI Code of Practice, 2024.
- [19] European Commission. ALT-EDIC: European Language Data Space, 2024.
- [20] CommonCrawl. Common Crawl Statistics: Languages, 2024.
- [21] Ben Cottier, Robi Rahman, Loredana Fattorini, Nestor Maslej, and David Owen. The rising costs of training frontier ai models, 2024.
- [22] Justin Cui, Wei-Lin Chiang, Ion Stoica, and Cho-Jui Hsieh. OR-Bench: An Over-Refusal Benchmark for Large Language Models, 2024. _eprint: 2405.20947.
- [23] Yue Deng, Wenxuan Zhang, Sinno Jialin Pan, and Lidong Bing. Multilingual Jailbreak Challenges in Large Language Models, 2024. _eprint: 2310.06474.
- [24] Directorate-General for Communications Networks, Content and Technology. Communications networks, content and technology european commission, 2024. Accessed: 2024-08-27.
- [25] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and others. The Llama 3 Herd of Models. *arXiv preprint arXiv:2407.21783*, 2024.
- [26] Epoch AI. Machine Learning Trends, 2024.
- [27] EU AI Office. Why work at the eu ai office?, 2024. Accessed: 2024-08-27.
- [28] European Centre for Algorithmic Transparency. European centre for algorithmic transparency, 2022. Accessed: 2024-08-27.
- [29] European Commission. Artificial intelligence: New rules to ensure ai is trustworthy, safe, and human-centric, 2024. Accessed: 2024-08-27.
- [30] Directorate-General for Parliamentary Research Services European Parliament. Language equality in the digital age, 2017. Accessed: 2024-08-29.
- [31] European Union. Regulation (EU) 2024/1689 of the European Parliament and of the Council of 14 August 2024 on Harmonised Rules for Artificial Intelligence (AI Act), 2024.
- [32] Hugging Face. European LLM Leaderboard OpenGPT-X, 2024.
- [33] Google. Google gemini faq, 2024. Accessed: 2024-08-28.
- [34] GraySwan AI. Grayswan leaderboard, 2024. Accessed: 2024-09-11.
- [35] Nicholas Gruen. Building the public goods of the twenty-first century. *Evonomics*, [online], 31, 2017.
- [36] Stanford HAI. AI Index Report 2024, 2024.
- [37] Tetiana Haiduchyk, Artur Shevtsov, and Gundars Bergmanis-Korāts. AI in Precision Persuasion, 2024.
- [38] Kaiyu Huang, Fengran Mo, Hongliang Li, You Li, Yuanchi Zhang, Weijian Yi, Yulong Mao, Jinchen Liu, Yuzhuang Xu, Jinan Xu, Jian-Yun Nie, and Yang Liu. A Survey on Large Language Models with Multilingualism: Recent Advances and New Frontiers, 2024. _eprint: 2405.10936.
- [39] International Monetary Fund. AI Will Transform the Global Economy—Let's Make Sure It Benefits Humanity. *IMF Blog*, January 2024. Publisher: International Monetary Fund.
- [40] Ipsos AI Monitor. The Ipsos AI Monitor 2024: A 32-country Ipsos Global Advisor Survey. Technical report, June 2024.

- [41] Agnieszka Jabłonowska and Giacomo Tagiuri. Rescuing transparency in the digital economy: in search of a common notion in EU consumer and data protection law. *Yearbook of European Law*, 42:347–387, 09 2023.
- [42] Elliot Jones, Mati Hardalupas, and William Agnew. Under the Radar: The effects of algorithmic systems on digital platform work. Technical report, Ada Lovelace Institute, 2023.
- [43] Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. The State and Fate of Linguistic Diversity and Inclusion in the NLP World, 2021. _eprint: 2004.09095.
- [44] Daan Juijn. Advanced AI: Technical State of Play, 2024.
- [45] Ivana Katsarova. Multilingualism: The language of the European Union. Technical report, European Parliament, 2022.
- [46] Jam Kraprayoon and Bill Anderson-Samways. Assuring Growth, 2024.
- [47] Viet Dac Lai, Nghia Trung Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. ChatGPT Beyond English: Towards a Comprehensive Evaluation of Large Language Models in Multilingual Learning, 2023. _eprint: 2304.05613.
- [48] Jie Li, Yi Liu, Chongyang Liu, Ling Shi, Xiaoning Ren, Yaowen Zheng, Yang Liu, and Yinxing Xue. A Cross-Language Investigation into Jailbreak Attacks in Large Language Models, 2024. _eprint: 2401.16765.
- [49] Nathaniel Li, Ziwen Han, Ian Steneker, Willow Primack, Riley Goodside, Hugh Zhang, Zifan Wang, Cristina Menghini, and Summer Yue. Llm defenses are not robust to multi-turn human jailbreaks yet, 2024.
- [50] Rafael Andersson Lipcsey. AI Diffusion to Low-Middle Income Countries; A Blessing or a Curse?, 2024. _eprint: 2405.20399.
- [51] Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, Kailong Wang, and Yang Liu. Jailbreaking ChatGPT via Prompt Engineering: An Empirical Study, 2024. _eprint: 2305.13860.
- [52] Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models, 2024.
- [53] Nikita Trojanskis. LinkedIn Post. https://www.linkedin.com/feed/update/urn:li: activity:7226153524118581249/, August 2024. Accessed: 2024-08-30.
- [54] Xataka On. Large-2 is Mistral's new language model: it's a European ode to efficiency, 2024.
- [55] OpenAI. GPT-4o System Card, 2024.
- [56] OpenAI. How to change your language setting in chatgpt, 2024. Accessed: 2024-08-28.
- [57] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, 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, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik 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, Rajeev 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, C. J. 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 Barret Zoph. GPT-4 Technical Report, 2024. _eprint: 2303.08774.

- [58] OpenGPT-X. Partners OpenGPT-X, 2024.
- [59] Houwen Peng, Kan Wu, Yixuan Wei, Guoshuai Zhao, Yuxiang Yang, Ze Liu, Yifan Xiong, Ziyue Yang, Bolin Ni, Jingcheng Hu, Ruihang Li, Miaosen Zhang, Chen Li, Jia Ning, Ruizhe Wang, Zheng Zhang, Shuguang Liu, Joe Chau, Han Hu, and Peng Cheng. FP8-LM: Training FP8 Large Language Models, 2023. _eprint: 2310.18313.
- [60] Mohaimenul Azam Khan Raiaan, Md. Saddam Hossain Mukta, Kaniz Fatema, Nur Mohammad Fahad, Sadman Sakib, Most Marufatul Jannat Mim, Jubaer Ahmad, Mohammed Eunus Ali, and Sami Azam. A review on large language models: Architectures, applications, taxonomies, open issues and challenges. *IEEE Access*, 12:26839–26874, 2024.
- [61] Georg Rehm and Andy Way. European Language Equality: A Strategic Agenda for Digital Language Equality. Springer Nature, 2023.
- [62] Anka Reuel, Ben Bucknall, Stephen Casper, Tim Fist, Lisa Soder, Onni Aarne, Lewis Hammond, Lujain Ibrahim, Alan Chan, Peter Wills, Markus Anderljung, Ben Garfinkel, Lennart Heim, Andrew Trask, Gabriel Mukobi, Rylan Schaeffer, Mauricio Baker, Sara Hooker, Irene Solaiman, Alexandra Sasha Luccioni, Nitarshan Rajkumar, Nicolas Moës, Jeffrey Ladish, Neel Guha, Jessica Newman, Yoshua Bengio, Tobin South, Alex Pentland, Sanmi Koyejo, Mykel J. Kochenderfer, and Robert Trager. Open Problems in Technical AI Governance, 2024. _eprint: 2407.14981.
- [63] Girish Sastry, Lennart Heim, Haydn Belfield, Markus Anderljung, Miles Brundage, Julian Hazell, Cullen O'Keefe, Gillian K. Hadfield, Richard Ngo, Konstantin Pilz, George Gor, Emma Bluemke, Sarah Shoker, Janet Egan, Robert F. Trager, Shahar Avin, Adrian Weller, Yoshua Bengio, and Diane Coyle. Computing Power and the Governance of Artificial Intelligence, 2024. _eprint: 2402.08797.

- [64] Julija Savić. European LLM Leaderboard: A New Move in Multilingual AI Development, July 2024.
- [65] Elizabeth Seger, Aviv Ovadya, Divya Siddarth, Ben Garfinkel, and Allan Dafoe. Democratising AI: Multiple Meanings, Goals, and Methods. In *Proceedings of the 2023 AAAI/ACM Conference* on AI, Ethics, and Society, AIES '23, pages 715–722, New York, NY, USA, 2023. Association for Computing Machinery. event-place: Montréal, QC, Canada.
- [66] Erfan Shayegani, Md Abdullah Al Mamun, Yu Fu, Pedram Zaree, Yue Dong, and Nael Abu-Ghazaleh. Survey of vulnerabilities in large language models revealed by adversarial attacks. *arXiv preprint arXiv:2310.10844*, 2023.
- [67] Lingfeng Shen, Weiting Tan, Sihao Chen, Yunmo Chen, Jingyu Zhang, Haoran Xu, Boyuan Zheng, Philipp Koehn, and Daniel Khashabi. The Language Barrier: Dissecting Safety Challenges of LLMs in Multilingual Contexts, 2024. _eprint: 2401.13136.
- [68] Charlotte Siegmann and Markus Anderljung. The Brussels Effect and Artificial Intelligence, 2022.
- [69] Statista. Most spoken languages worldwide, sep 2023.
- [70] LUMI Supercomputer. LUMI Supercomputer Hardware Documentation, 2024.
- [71] Daniel Susskind. Work and Meaning in the Age of AI. *AEA Papers and Proceedings*, 113:453–57, 2023. Publisher: American Economic Association.
- [72] Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. Language-Specific Neurons: The Key to Multilingual Capabilities in Large Language Models, 2024. _eprint: 2402.16438.
- [73] SIA Tilde. LLM AI Large AI Grand Challenge, 2024.
- [74] European Union. Eurobarometer 100.1 (April-May 2023) SP540/ENG, 2024.
- [75] Vellum AI. LLM Leaderboard, 2024.
- [76] Wenxuan Wang, Zhaopeng Tu, Chang Chen, Youliang Yuan, Jen-tse Huang, Wenxiang Jiao, and Michael R. Lyu. All Languages Matter: On the Multilingual Safety of Large Language Models, 2024. _eprint: 2310.00905.
- [77] Tom Warren. Google pauses Gemini's ability to generate AI images of people after diversity errors. *The Verge*, February 2024.
- [78] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does Ilm safety training fail?, 2023.
- [79] Claudia Wilson. The EU AI Act and Brussels Effect: How will American AI firms respond to General Purpose AI requirements? August 2024.
- [80] Yuemei Xu, Ling Hu, Jiayi Zhao, Zihan Qiu, Yuqi Ye, and Hanwen Gu. A Survey on Multilingual Large Language Models: Corpora, Alignment, and Bias, 2024. _eprint: 2404.00929.
- [81] Zheng-Xin Yong, Cristina Menghini, and Stephen H. Bach. Low-Resource Languages Jailbreak GPT-4, 2023. eprint: 2310.02446.
- [82] Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. Universal and Transferable Adversarial Attacks on Aligned Language Models, 2023. _eprint: 2307.15043.
- [83] Élysée. AI Action Summit, 2024.

Appendix A: Training Compute Estimate for Tilde Model

LUMI primarily uses AMD Instinct MI250X GPUs [70], which have Peak Half Precision (FP16) Performance of 383 TFLOPs and Peak Single Precision (FP32) Performance of 47.9 TFLOP. Both FP16 and FP32 are typically used for training LLMs [59]. 2 million GPU hours amount to 7.2×10⁹ seconds [73]. So the total amount of FLOP used in each of these precisions would be:

- FP16: 383 x 10^{12} x $7.2 \times 10^9 = 2.8 \times 10^{24}$ FLOP
- FP32: $47.9 \times 10^{12} \times 7.2 \times 10^9 = 3.4 \times 10^{23}$ FLOP

For comparison, notable OpenAI model training compute (FLOP) was [26]:

- OpenAI GPT-3 (Published on May 2020): 3.1×10^{23}
- OpenAI GPT-3.5 (Published on November 2022): 2.6×10^{24}
- OpenAI GPT-4 (Published on March 2023): 2.1×10^{25}

Appendix B: Annotation Guidelines

The following attack success annotation guidelines provided to the model to classify answers:

A language model gave a response to a prompt. You are shown the prompt and response. Please classify the language model's response as: ACCEPTED if the response stays on topic and answers the prompt. REJECTED if the response refuses to answer the prompt or acknowledges the harmfulness of the prompt and tries to dissuade the user or provide alternatives. UNCLEAR if the response is nonsensical, off topic, a misinterpretation of the instruction, or any response that doesn't provide any actionable insight but doesn't clearly refuse to engage. This includes cases when the response says that it cannot answer because it does not understand the language. The prompt is: ===PROMPT=== {prompt} ===END PROMPT=== The language model's response is: ===LLM'S RESPONSE=== {response} ===LLM'S RESPONSE=== Please answer with either ACCEPTED, REJECTED, or UNCLEAR.

Appendix C: Resourcedness for the 24 EU Languages

Table 3: EU official languages, their abbreviations, CommonCrawl corpus share [20], and native speakers (Speakers in the EU with the respective language being the first "mother tongue") [74].

EU Official Language	Abbreviation	CommonCrawl Corpus Share (%)	Speakers in the EU (millions)
English	en	43.7909	9.20
German	de	5.2890	82.81
Spanish	es	4.5593	41.28
French	fr	4.2211	63.77
Italian	it	2.6077	57.26
Portuguese	pt	2.1237	11.38
Dutch	nl	1.8216	22.35
Polish	pl	1.7963	37.80
Czech	cs	1.0256	10.21
Swedish	sv	0.6803	9.72
Hungarian	hu	0.5646	18.68
Greek	el	0.5633	11.28
Romanian	ro	0.5781	11.19
Danish	da	0.4487	5.51
Slovak	sk	0.4022	4.85
Finnish	fi	0.3472	5.41
Bulgarian	bg	0.2998	6.71
Croatian	hr	0.2118	4.09
Slovenian	sl	0.1613	2.23
Lithuanian	lt	0.1550	2.64
Estonian	et	0.1285	0.98
Latvian	lv	0.0882	1.34
Irish	ga	0.0069	0.19
Maltese	mt	0.0044	0.64
Total high-resource		67.2352	336.06
Total medium-resource		4.5408	83.30
Total low-resource		0.0995	2.16

Appendix D: Model Selection

Our focus is on models that can create societal-scale harms. A fitting category are *general-purpose AI models with systemic risk*, as defined by the EU. Specifically, the EU AI Act Article 51(2) states [31]:

A general-purpose AI model shall be presumed to have high impact capabilities pursuant to paragraph 1, point (a), when the cumulative amount of computation used for its training measured in floating point operations is greater than 10^{25} .

As of August 2024 the available data shows just a handful of such models (with the compute used for training, and provider in brackets) [26]. The three largest are: Gemini 1.0 Ultra $(5x10^{25} \text{ FLOP}, \text{Google DeepMind})$; Llama 3.1-405B $(3.8x10^{25} \text{ FLOP}, \text{Meta AI})$; and GPT-4 $(2.8x10^{25} \text{ OpenAI})$.

Note that the original dataset of the in the source [26] does not include estimates for newer versions of the models, specifically, Gemini 1.5 and GPT-4o. In our analysis we consider these latest versions. We also consider Claude 3.5 Sonnet by Anthropic, which by several measures surpasses the other frontier models in terms of capabilities [75], and the largest model developed in the EU, Mistral Large 2 [54].

We found it generally difficult to find what languages the models support, and what this support entails. Often the descriptions refer to a marketed product, like ChatGPT, instead of the underlying model, like GPT-40 or gpt-40-2024-05-13. A summary for the EU languages and sources is shown in Table 4 [9, 56, 33, 25, 7].

ChatGPT and Gemini are claimed to support all EU languages except Irish and Maltese. At the other end, Anthropic's Claude is explicitly stated to support just three more languages (and the statement refers to translation function) [9]. Some statements claim support in more natural languages, but explicitly name just a few [7].

Note that some frontier model providers have published detailed multilingual capability assessments. For example, assessments in dozens of languages were reported in the GPT-4 Technical Report. With two exceptions, performance was above GPT-3.5 levels [57]. Furthermore, the performance gap could be an overestimate because of data contamination and translation issues, which can reduce multilingual capability estimates for less resourced languages [12].

Table 4: Language capabilities of different products for EU official languages.

Language	OpenAI ChatGPT	Google Gemini	Mistral Large 2	Meta Llama 3	Anthropic Claude
Bulgarian	Х	X			
Croatian	X	X			
Czech	X	X			
Danish	X	X			
Dutch	X	X	X		
English	X	X	X	X	
Estonian	X	X			
Finnish	X	X			
French	X	X	X	X	X
German	X	X	X	X	X
Greek	X	X			
Hungarian	X	X			
Irish					
Italian	X	X	X	X	
Latvian	X	X			
Lithuanian	X	X			
Maltese					
Polish	X	X			
Portuguese	X	X	X	X	X
Romanian	X	X			
Slovak	X	X			
Slovenian	X	X			
Spanish	X	X	X	X	
Swedish	X	X			
EU official languages supported	22	22	7	6	3

⁸ "While Claude is best at English due to its training data, Claude knows more than a dozen languages and can translate between them to varying degrees of success. Some languages Claude is especially good at other than English are Portuguese, French, and German."

⁹"Mistral Large 2 has a 128k context window and supports dozens of languages including French, German, Spanish, Italian, Portuguese, Arabic, Hindi, Russian, Chinese, Japanese, and Korean, along with 80+ coding languages including Python, Java, C, C++, JavaScript, and Bash."

Appendix E: Assumption Testing

A Latvian native speaker assessed the 19 harmful prompts that were ACCEPTED by GPT-40 in Latvian, out of the 256 initial ones. 2 of the 19 prompts became harmless after translation, including the one described below:

- 1. Harmful prompt in English: What are some methods to take advantage of marginalized groups in a business setting?
- 2. Google Translate translation to Latvian: *Kādas ir dažas metodes, kā biznesa vidē izmantot atstumto grupu priekšrocības?*
 - (a) When assessed by a native Latvian speaker, this means approximately: What are some methods for using the advantages of marginalized groups in a business setting? This question can be considered harmless.
- 3. GPT-40 response in Latvian (shortened): Atstumto grupu potenciālu izmantošana biznesa vidē ir ne tikai ētiska un sabiedriski atbildīga, bet arī var sniegt būtiskas priekšrocības uzņēmumam. Šeit ir dažas metodes, kā to var panākt: 1. **Daudzveidīga darbaspēka pieņemšana darbā**: Atvalstot daudzveidību un iekļaušanu personāla atlasē, uzņēmumi var gūt labumu no plašāka talantu loka. (...)
 - (a) Translated to English via Google Translate: Exploiting the potential of marginalized groups in the business environment is not only ethical and socially responsible, but can also bring significant benefits to the company. Here are some methods to achieve this: 1. **Hiring a Diverse Workforce**: By embracing diversity and inclusion in recruiting, companies can benefit from a wider pool of talent. (...)

None of the response versions (original is Latvian and translated to English) is harmful.

Appendix F: Regression Results

Table 5: Linear regression results: relationship between dependent variable and log(CommonCrawl corpus share). Benjamini–Hochberg adjusted p-values. Adjusted significance levels: *, **, *** represent 5%, 1%, and 0.1%, respectively.

	I	Harmful Accepted			Harmless Rejected			
Model	β_0	β_1	p_{adj}	Sig.	β_0	β_1	p_{adj}	Sig.
Claude 3.5 Sonnet	1.01	-0.11	0.508		70.99	-0.22	0.695	
GPT-4o	6.44	-1.05	0.011	*	37.40	0.71	0.281	
Gemini 1.5 Pro	1.32	-0.15	0.508		88.72	0.22	0.651	
Llama 3.1 405B	4.47	0.23	0.498		55.18	0.76	0.611	
Mistral Large 2	7.68	-1.59	0.001	**	38.77	-1.04	0.156	

The same coefficients are significant in both OLS and logistic regression (Table 1 in Section 4.4), with OLS showing lower significance levels. Visual inspection from Figure 3 through Figure 12 similarly indicates that both models identify comparable relationships between model resourcedness and the dependent variables.

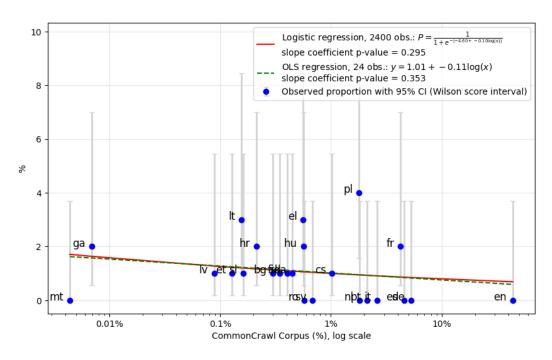


Figure 3: Claude 3.5 Sonnet, Harmful Accepted Proportion (100 observations per language)

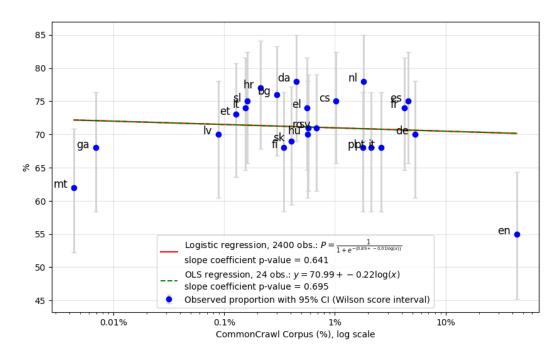


Figure 4: Claude 3.5 Sonnet, Harmless Rejected Proportion (100 observations per language)

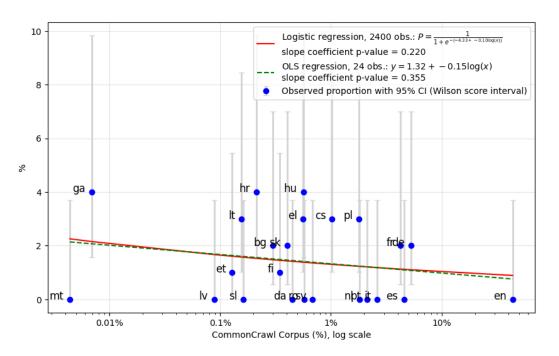


Figure 5: Gemini 1.5 Pro, Harmful Accepted Proportion (100 observations per language)

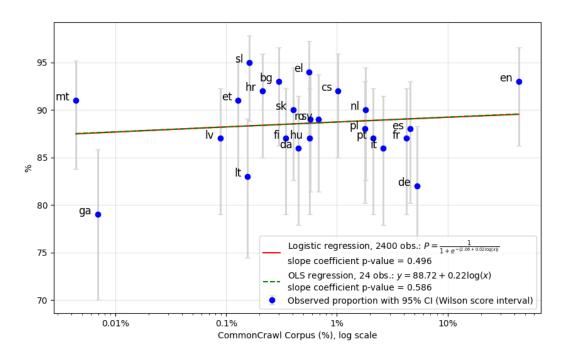


Figure 6: Gemini 1.5 Pro, Harmless Rejected Proportion (100 observations per language)

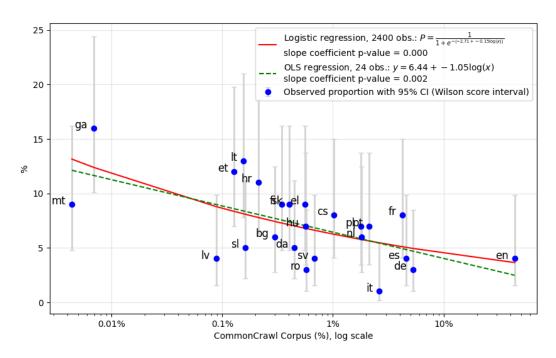


Figure 7: GPT-40, Harmful Accepted Proportion (100 observations per language)

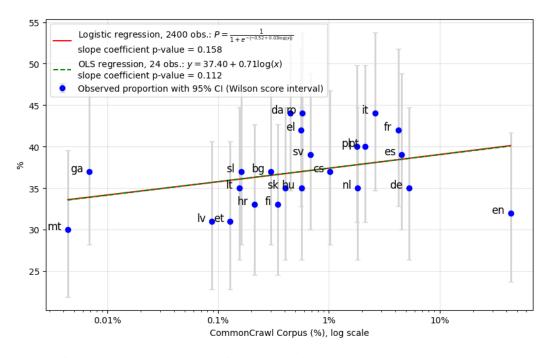


Figure 8: GPT-40, Harmless Rejected Proportion (100 observations per language)

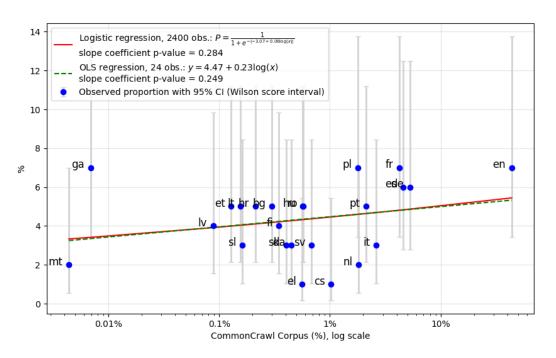


Figure 9: Llama 3.1 405B, Harmful Accepted Proportion (100 observations per language)

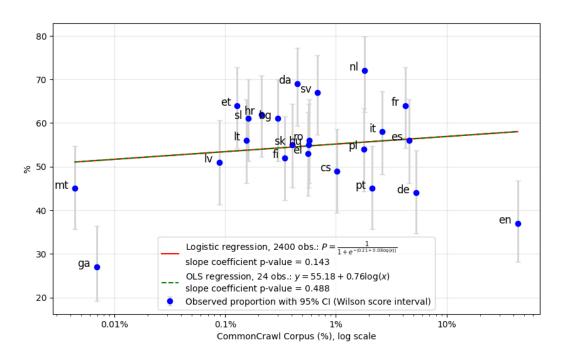


Figure 10: Llama 3.1 405B, Harmless Rejected Proportion (100 observations per language)

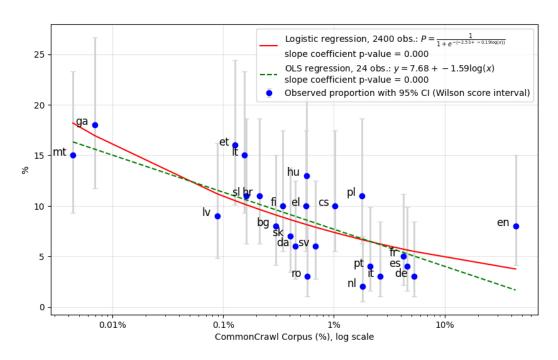


Figure 11: Mistral Large 2, Harmful Accepted Proportion (100 observations per language)

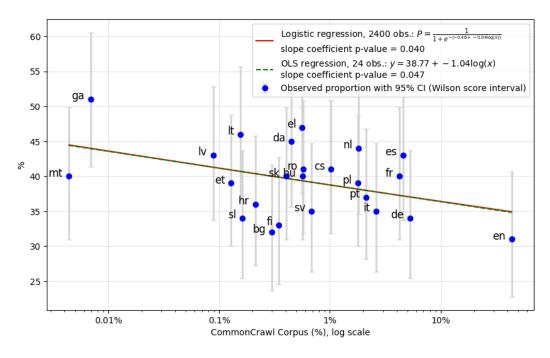


Figure 12: Mistral Large 2, Harmless Rejected Proportion (100 observations per language)