# Beyond speculation: Measuring the growing presence of LLM-generated texts in multilingual disinformation

*Increased sophistication of large language models (LLMs) and the consequent quality of generated multilingual text raises concerns about potential disinformation misuse. While humans struggle to distinguish LLM-generated content from human-written texts, the scholarly debate about their impact remains divided. Some argue that heightened fears are overblown due to natural ecosystem limitations, while others contend that specific "longtail" contexts face overlooked risks. Our study bridges this debate by providing the first empirical evidence of LLM presence in the latest real-world disinformation datasets, documenting the increase of machine-generated content following ChatGPT's release, and revealing crucial patterns across languages, platforms, and time periods.*

Authors: Dominik Macko (1), Aashish Anantha Ramakrishnan (2), Jason Samuel Lucas (2), Robert Moro (1), Ivan Srba (1), Adaku Uchendu (3), Dongwon Lee (2)
Affiliations: (1) Kempelen Institute of Intelligent Technologies, Slovakia, (2) The Pennsylvania State University, PA, USA, (3) MIT Lincoln Laboratory, USA[1]

## Research questions

- To what extent are machine-generated texts present in existing disinformation datasets?
- Within labeled disinformation datasets, what is the distribution of machine-generated text in "false" and "true" content?
- How does the prevalence of machine-generated texts vary across languages, information sources, and time periods?

## Essay summary

Our study makes several key contributions to understanding LLM-generated disinformation:

- By validation on broader datasets, our detection methods establish a robust analytical framework for examining real-world disinformation content, confirming both the increasing presence and prevalence of machine-generated texts in disinformation datasets over time.
- The distribution of LLM-generated content varies significantly across languages and platforms, revealing targeted patterns of misuse rather than uniform effects. This provides empirical validation for previously speculated concerns and unsupported fears about increased LLM deployment in disinformation campaigns.
- Most importantly, our findings underscore the urgent need for continued investigation and improved countermeasures, including enhanced detection methods and credibility assessment systems to preserve information integrity in our evolving digital landscape.

## Implications

Since the widespread availability of LLMs through user-friendly interfaces like ChatGPT, concerns about their misuse have grown significantly (Zellers et al., 2019; Crothers et al., 2023; Zhuo et al., 2023). These concerns are supported by evidence of LLM deployment across scholarly papers (Haider et al., 2024), scientific peer reviews (Liang et al., 2024), Wikipedia articles (Brooks et al., 2024), and content-driven social media (Sun et al., 2024). While research demonstrates LLMs' capability to generate misinformation (Lucas et al., 2023; Williams et al., 2024; Vykopal et al., 2024; Heppell et al., 2024; Zugecova et al., 2024) and the computing field is exploring regulations and policies (Nahar et al., 2025), the actual extent of real-world misuse remains unknown, limiting our understanding of the true impact. Simon et al. (2023) argue that heightened fears are overblown due to natural ecosystem limitations, while Lucas et al. (2024) contend that specific "longtail" contexts face overlooked risks. Our study addresses this gap by providing the first comprehensive assessment of LLM-generated content prevalence in real-world disinformation contexts.

The capabilities of LLMs can be used for various legitimate use cases (e.g., polishing or translation) to support different content production pipelines. However, due to LLMs' increased scalability, lower costs, and ease of use, there are various reasons and ways they can be misused for promoting disinformation. These range from disinformation text creation in multiple languages (Vykopal et al., 2024), through text translation and polishing (especially for non-native speakers; Shafayat et al., 2024), modification (producing the same disinformation narrative in various forms; Dash et al., 2024), or targeting specific groups of audiences (e.g., personalization or micro-targeting; Zugecova et al., 2024). Recent findings highlight concerning trends of LLM misuse. Multi-modal generative models, when used in news contexts can lead to both the unintentional (Anantha Ramakrishnan et al., 2024) and intentional (Weikmann et al., 2023) production of misinformation. Barman et al. (2024) demonstrate how multi-modal content generation enables sophisticated disinformation at scale, while Coeckelbergh (2025) identifies potential threats to electoral processes. The spread of LLM-generated misinformation can also lead to data corpora pollution, inducing factual errors and hallucinations in newer models trained on this data (Pan et al., 2023; Guo et al., 2024). Our analysis reveals increasing prevalence of LLM-generated content (1.85% in 2023) in the MultiClaim dataset (containing multilingual social-media texts matched with fact-checked claims), confirming these are not merely speculated concerns. The societal implications are significant such that LLM-generated misinformation can erode democratic processes and societal trust (Chen et al., 2024) and threaten online safety (Chen and Shu, 2024).

Our findings reveal three key implications for understanding and addressing AI's role in disinformation. First, **machine-generated content manifests unevenly across languages and platforms**, with some contexts showing significantly higher prevalence. While broad ecosystem protections matter, targeted interventions (e.g., posts in a specific language of a specific platform) are needed in vulnerable contexts, particularly in languages such as Polish and French, where prevalence rates reach 4.7%.

Second, **the increasing temporal trend signals the importance of proactive monitoring**. The documented rise from 0.93% in 2021 to 1.85% in 2023 in the MultiClaim dataset (Pikuliak et al., 2023) demonstrates a clear trajectory requiring attention, particularly in high-prevalence contexts that may indicate broader trends. This suggests focusing early warning systems on languages and platforms showing higher concentrations of machine-generated content.

Although the detection methods can inherently have detection errors, such as false positives and false negatives, we have shown their capabilities on existing labeled multilingual benchmark datasets. The increasing prevalence in time corresponds with a common hypothesis that after ChatGPT's release (November 2022), the prevalence of machine-generated texts has increased. Indeed, the most significant increase has been observed in 2023, when the proliferation of many modern LLMs with conversational interfaces enabled simplified usage for a wider community.

Third, **our empirical approach emphasizes the value of evidence-based assessment**. Rather than relying on speculation, our study provides concrete measurements of machine-generated content's prevalence and distribution. It enables more targeted and effective responses to AI-generated disinformation.

Our results have demonstrated that LLM-generated texts are already present in the existing disinformation datasets coming from the real world, such as MultiClaim. Although the results do not indicate that LLMs are used for disinformation in significantly higher volume than for legitimate use, the fears of LLM misuse are justified. The MultiClaim dataset consists of fact-checked claims, but fact-checkers can only verify a small sample of the most viral or suspicious claims due to the vast daily influx of information. This creates a significant imbalance between true and false claims, complicating systematic comparisons.

The evidence of the presence of LLM-generated texts in datasets like USC-X (Balasubramanian et al., 2024) or FIGNEWS (Zaghouani et al., 2024) may also imply that LLMs are already being misused for propaganda for elections or warfare. To enhance trust in the online information space, this further underscores the importance of indicating whether a given text has been generated or altered by an LLM. It will enable the readers to assess its credibility as such, considering that disinformation as well as unintentional misinformation can be generated by LLMs.

Future research should examine factors driving varying prevalence across contexts and investigate correlations with information ecosystem vulnerabilities. Maintaining and improving detection capabilities will be crucial as AI technology evolves, especially for vulnerable languages and platforms. Understanding these patterns and developing effective detection methods remains essential for preserving the integrity of information ecosystems.

# Findings

*Finding 1: There are machine-generated texts in the existing disinformation dataset MultiClaim — a dataset of multilingual social-media texts containing fact-checked claims — and their prevalence is increasing over time.*

One of the key objectives of this study is to evaluate whether the LLMs are already being misused for mis-/disinformation in the real world. Two multilingual detectors, fine-tuned on different data for machine-generated text detection (a binary classification task), show an increase of mean probability (Mean Score) of the texts grouped per year as being generated by the LLMs (Figure 1). This provides evidence that in recent years (after the release of most-advanced popular LLMs, like ChatGPT), LLMs are increasingly being misused for disinformation.

Although the years 2021-2023 contain approximately the same number of text samples (around 25,000), the Mean Score is increased from 0.04 and 0.27 in 2021 to 0.05 and 0.36 in 2023 for the Gemma_MultiDomain and Gemma_GenAI detectors, respectively. Such a 30% increase of Mean Score roughly represents adding or changing of 1.3% (for Gemma_MultiDomain) to 13% (for Gemma_GenAI) of scores to 1.0 (i.e., maximal detector confidence). Since the detectors do not always provide a probability of 1.0 for each machine-generated text, the prevalence of LLM-generated texts between 2021 and 2023 has most probably increased even more. Based on this increase, we estimate that at least 1.5% to 15% of texts in 2023 of the MultiClaim dataset have been generated by LLMs.
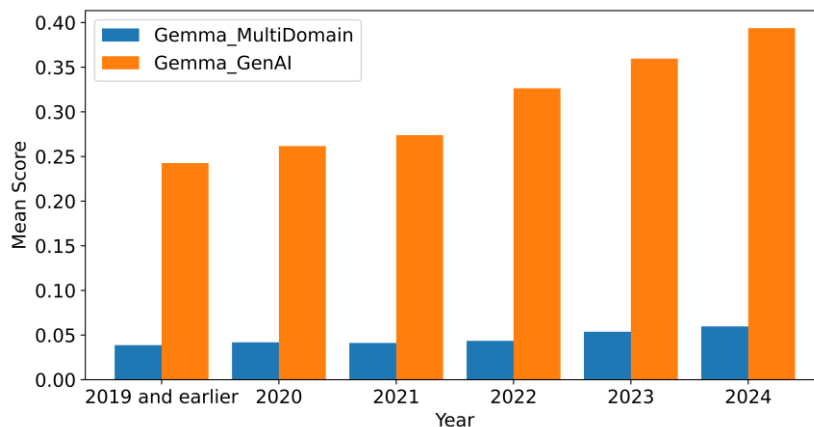
***Figure 1. Per-year Mean Score of the two fine-tuned detectors for MultiClaim texts. Both detectors independently show increasing Mean Scores in recent years.***

To further analyze the proportion of the LLM-generated texts in the dataset, we combined their highest-confidence predictions (as described in the Methods section). Figure 2 illustrates the proportion of samples predicted (by highest certainty) to be generated by LLMs. The proportion has increased from 0.93% in 2021 to 1.85% in 2023 (i.e., a relative increase by 99%). Given the worst-case calculated precision of the positive class (i.e., machine-generated) of such combined prediction on the existing multilingual labeled datasets (0.93), we can be highly confident that at least 1.7% of texts in 2023 have been generated or modified by LLMs. The samples detected as "generated" in earlier years (especially 2019 and earlier) can form a baseline, representing either false positives or texts generated/modified by pre-LLM machines (e.g., translators).
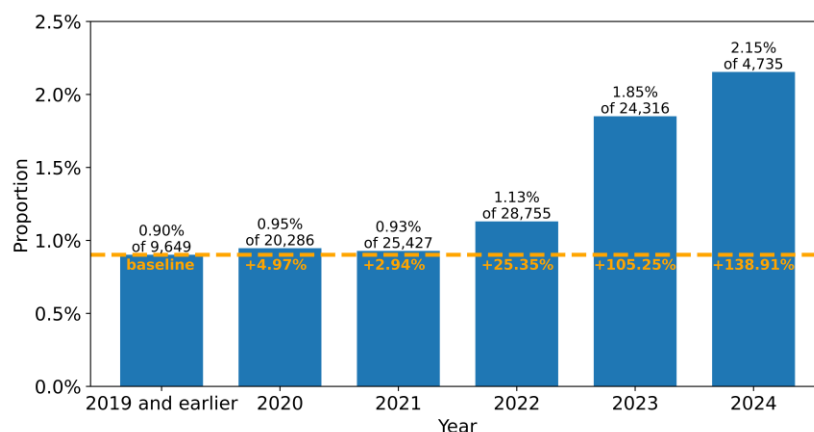


***Figure 2. Per-year proportion of the MultiClaim texts detected to be machine-generated. Proportion is increasing in recent years, with the highest increase in 2023 (after ChatGPT release).***

*Finding 2: There are differences in the prevalence of machine-generated texts in the existing disinformation dataset MultiClaim among languages as well as sources.*

Analogously to the previous combined prediction, we have aggregated the results based on the languages, for which at least 1,000 samples are included in the dataset. The results are illustrated in Figure 3, where

although English and Spanish have the highest absolute number of texts predicted to be positive, the relative proportion shows to be the highest in the case of Polish (4.7%) and French (4.2%).
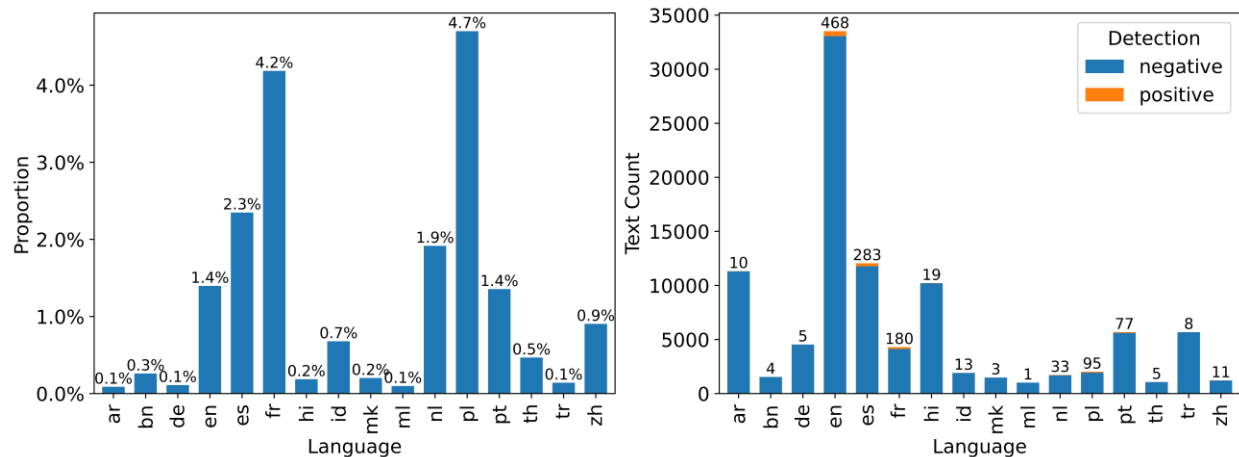


***Figure 3. Left: Per-language proportion of the texts detected to be machine-generated (for languages including enough samples). Proportion differs among languages and is the highest in case of Polish and French. Right: Per-language number of samples showing also the distribution of samples classified as human-written (negative) and machine-generated (positive).***

Regarding the contained social-media platforms, the differences are less severe, ranging from 0.64% in Twitter (X), 1.29% in Facebook, to 1.5% in Instagram. Although Telegram contains over 1.7% of LLM-generated texts, it covers less than 1,000 samples in the dataset and, thus, is not objectively comparable to the others. Based on fact-checked claim verdict rating categories, the differences are not clearly observable, since "False", "Misleading", and "Not categorized" covered over 90% of all samples. Among these three categories, the first one contains the highest proportion (1.36%) of LLM-generated texts. Although the "True" category covers much less samples (786), out of them, 0.76% of texts are identified as generated.

*Finding 3: There are machine-generated texts also present in other existing datasets and their prevalence differs among them.*

We have similarly analyzed the presence of LLM-generated texts in other existing datasets, potentially including disinformation or propaganda. The results illustrated in Figure 4 indicate that LLMs have generated or modified some of the texts also in these datasets. Although the proportions are not directly comparable to each other (due to containing a different number of samples of different kinds), FIGNEWS seems to contain the highest portion of machine-generated texts (3.16%).
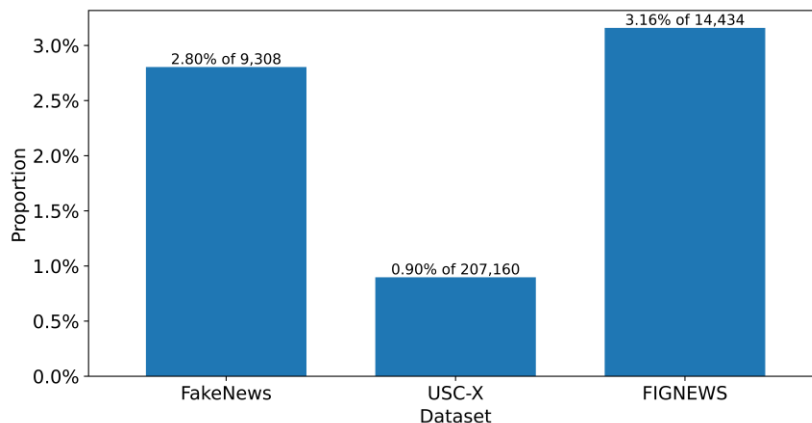
*Figure 4. Proportion of the texts detected to be machine-generated in different datasets.*

The FakeNews dataset contains the labels representing the text to be real or fake (labels verified by humans), including 2.61% and 3.27% of LLM-generated texts, respectively. It further provides evidence that LLMs are already being used for disinformation. Based on FIGNEWS dataset analysis, we found over 10% of the French texts were being generated by LLMs, which is significantly more than around 4.7% for English. On the other hand, the Hebrew, Hindi, and Arabic texts contained under 0.5% of LLM-generated texts. We further found a higher prevalence of LLM-generated texts in the texts labeled as "Biased against Israel" (5.11%) compared to the texts labeled as "Biased against Palestine" (3.13%).

Although there is a rather low portion of LLM-generated texts in the USC-X dataset (0.9%), it significantly differs across the months in the election year of 2024 (illustrated in Figure 5). It has increased from 0.44% in January to 2.39% in November (i.e., an increase of 4.5 times). There is approximately the same number of samples in each month (8,700 to 10,000); thus, the proportions are reliably comparable.
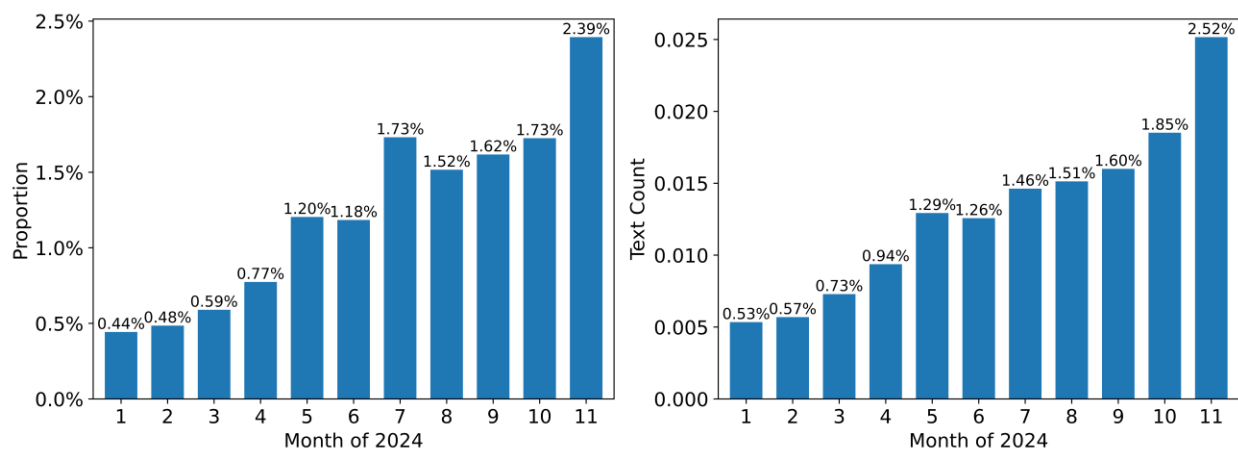


*Figure 5. The per-month proportion of the texts detected to be machine-generated in the USC-X dataset for the election year of 2024 (all languages on the left, English texts only on the right). Proportion is increasing towards the election date in November 2024.*

## Methods

In this study, two multilingual detectors, fine-tuned on different data for machine-generated text detection are used. Specifically, the Gemma-2-9b-it (Gemma Team, 2024) model from Google is fine-tuned for the binary classification task on train split of the GenAI dataset (the detector called

Gemma_GenAI) and on train splits of MULTITuDE combined with MultiSocial datasets (the detector called Gemma_MultiDomain). This model has been selected for its best performance in the study by Macko et al. (2025). For the efficient fine-tuning process, we utilized the QLoRA technique (Dettmers et al., 2024). Since the detectors are trained on different data, their detection capabilities complement each other.

The detectors are firstly evaluated on the existing multilingual datasets (test splits, not included in the detectors training), which are labeled (the ground truth of machine-generated vs human-written texts) for benchmarking of detection methods. After showing the detectors are robust enough to be used in the wild (on out-of-distribution data), we use them on four existing datasets, presumably containing disinformation or propaganda. The used datasets are briefly described in Table 1.

**Table 1.** *Overview of the used datasets. The upper part of the table includes the datasets that we have used for estimating the prevalence of the machine-generated texts, and the lower part includes the specialized benchmark datasets for evaluating the detection performance of detectors.*

| Dataset | Description |
|---|---|
| MultiClaim (Pikuliak et al., 2023) | A dataset of previously fact-checked claims by professional fact-checkers along with social-media posts gathered from the wild matching the claims. Besides the publicly available data, we have requested from the authors the more recent data (up to March 2024). The final data covers a small portion of texts from prior to 2019 and thousands of samples for the years up to 2024. |
| FakeNews (Raza et al., 2024) | A dataset of English news articles about US elections, collected from April to October 2023, containing GPT-4 annotations (reviewed by humans) considering the texts to be real or fake. |
| USC_X (Balasubramanian et al., 2024) | A sampled portion of a multilingual US election Twitter/X dataset (usc-x-24-us-election), collected from May to November 2024 (the last tweets from November 6th, 2024). The pseudo-random sampling covers up to 10,000 tweets per each month. |
| FIGNEWS (Zaghouani et al., 2024) | A dataset of the multilingual shared task on news media narratives (covering bias and propaganda), including Facebook posts in 5 languages on Israel War on Gaza that were publicly posted between October 7, 2023, and January 31, 2024. |
| MULTITuDE (Macko et al., 2023) | An official test split of a dataset of real news articles of various media sources in 11 languages, containing also the generated counterparts (based on their headlines) by 8 LLMs. |
| MultiSocial (Macko et al., 2024) | An official test split of a dataset of real social-media texts in 22 languages from 5 platforms, containing also the generated counterparts (by using 3 iterations of paraphrasing) by 7 LLMs. |
| SemEval (Wang et al., 2024) | An official test split of the multilingual track of the shared task 8 of SemEval-2024 workshop, containing student essays and news articles in 4 languages written by humans and generated by 7 LLMs. |
| GenAI (Wang et al., 2025) | An official test split of the multilingual subtask of the COLING workshop of GenAI content detection task 1, containing texts of various domains of 29 sources in 15 languages written by humans and generated by 19 LLMs. |
| MIX (Macko et al., 2025) | A mixture of existing datasets on a binary machine-generated text detection task, covering various domains, 7 languages and over 75 generators. |

Fine-tuned detectors are known to not generalize well to out-of-distribution data, requiring their classification thresholds to be carefully calibrated on expected data distributions. For measuring the prevalence of LLM-generated texts independently from distribution-specific calibrations, we use the Mean Score metric, which represents the probability (values from 0.0 to 1.0) of the text being generated by a machine (a higher value represents a higher probability).

Furthermore, we use a combined confident detection utilizing two fine-tuned detectors for the actual prediction of machine-generated texts. The prediction is positive if one of the detectors predicts the positive class with a probability of 1.0 and the other detector is not fully confident with the negative prediction (i.e., a positive class probability is > 0.0). Moreover, the Gemma_GenAI confident positive prediction is disqualified (not taken into account) for languages that exceed a 0.1 false positive rate on existing labeled datasets, specifically for Arabic, German, Italian, and Russian. The other detector has not outreached such a threshold of false positive rate for any language.

Table 2 shows the performance of such detection on existing multilingual benchmark datasets, containing human-written (the number of samples is provided in the Human column) as well as machine-generated (the number of samples is provided in the Machine column) samples. The table provides three performance metrics: False Positive Rate (FPR) is a ratio of how many human-written texts have been incorrectly predicted as positive, True Positive Rate (TPR), a measure equivalent to Recall of the positive class, is a ratio of how many of machine-generated samples have been correctly predicted as positive, and Precision reflects a ratio of how many of the predicted positive samples are actually machine-generated.

In regard to the evaluation using unlabeled data (where the truth is unknown), the key indicator is Precision, as it indicates the reliability of the identified positive samples. As indicated in the table, we can assume that only 93% (the worst case) of the predicted positives are texts generated by an LLM.

**Table 2.** *Performance of the combined confident detection of machine-generated texts on the existing multilingual labeled datasets. The worst-case performance is boldfaced.*

| Dataset | Precision | FPR | TPR | |Human| | |Machine| |
|---|---|---|---|---|---|
| MULTITuDE | 0.9985 | 0.0093 | 0.7442 | 3236 (11%) | 26059 (89%) |
| MultiSocial | 0.9969 | 0.0108 | **0.4988** | 17367 (13%) | 121460 (87%) |
| SemEval | 0.9490 | **0.0579** | 0.9852 | 20238 (48%) | 22140 (52%) |
| GenAI | 0.9955 | 0.0025 | 0.5194 | 73634 (49%) | 77791 (51%) |
| MIX | **0.9358** | 0.0476 | 0.6916 | 99759 (50%) | 100000 (50%) |

We aggregate the analyzed texts in the datasets for various aspects (e.g., language, time period) to see differences. For the crucial era distinction in regard to massive usage of LLMs for text generation, we consider the release of ChatGPT to the general public in November 2022. For the sake of the pre-/post-ChatGPT release era, we count the year 2022 as a part of the pre-ChatGPT period.

# Bibliography

Anantha Ramakrishnan, A., Huang, S. X., & Lee, D. (2024). ANNA: Abstractive Text-to-Image Synthesis with Filtered News Captions. Proceedings of the Third Workshop on Advances in Language and Vision Research. https://doi.org/10.48550/arXiv.2301.02160

Balasubramanian, A., Zou, V., Narayana, H., You, C., Luceri, L., & Ferrara, E. (2024). A public dataset tracking social media discourse about the 2024 US presidential election on Twitter/X. arXiv preprint arXiv:2411.00376. https://doi.org/10.48550/arXiv.2411.00376

Barman, D., Guo, Z. and Conlan, O., 2024. The dark side of language models: Exploring the potential of llms in multimedia disinformation generation and dissemination. Machine Learning with Applications, p.100545. https://doi.org/10.1016/j.mlwa.2024.100545

Brooks, C., Eggert, S., & Peskoff, D. (2024). The Rise of AI-Generated Content in Wikipedia. In Proceedings of the First Workshop on Advancing Natural Language Processing for Wikipedia (pp. 67–79). https://doi.org/10.18653/v1/2024.wikinlp-1.12

Chen, C., & Shu, K. (2024). Combating misinformation in the age of LLMs: Opportunities and challenges. AI Magazine, 45(3), 354–368. https://doi.org/10.1002/aaai.12188

Coeckelbergh, M. (2025). LLMs, Truth, and Democracy: An Overview of Risks. Science and Engineering Ethics, 31(1), 1–13. https://doi.org/10.1007/s11948-025-00529-0

Crothers, E. N., Japkowicz, N., & Viktor, H. L. (2023). Machine-generated text: A comprehensive survey of threat models and detection methods. IEEE Access, 11, 70977–71002. https://doi.org/10.1109/ACCESS.2023.3294090

Dash, S., Xu, Y., & Spiro, E. (2024). AI-Paraphrasing Increases Perceptions of Social Consensus & Belief in False Information. https://doi.org/10.31219/osf.io/8zb5p

Dettmers, T., Pagnoni, A., Holtzman, A., & Zettlemoyer, L. (2024). QLoRA: Efficient finetuning of quantized LLMs. Advances in Neural Information Processing Systems, 36.

Gemma Team (2024). Gemma. https://doi.org/10.34740/KAGGLE/M/3301

Guo, Y., Shang, G., Vazirgiannis, M., & Clavel, C. (2024, June). The Curious Decline of Linguistic Diversity: Training Language Models on Synthetic Text. In NAACL 2024 Findings-Annual Conference of the North American Chapter of the Association for Computational Linguistics. https://doi.org/10.18653/v1/2024.findings-naacl.228

Haider, J., Söderström, K. R., Ekström, B., & Rödl, M. (2024). GPT-fabricated scientific papers on Google Scholar: Key features, spread, and implications for preempting evidence manipulation. Harvard Kennedy School (HKS) Misinformation Review. https://doi.org/10.37016/mr-2020-156

Liang, W., Izzo, Z., Zhang, Y., Lepp, H., Cao, H., Zhao, X., ... & Zou, J. Y. (2024). Monitoring AI-Modified Content at Scale: A Case Study on the Impact of ChatGPT on AI Conference Peer Reviews. In Proceedings of the 41st International Conference on Machine Learning (pp. 29575–29620). https://dl.acm.org/doi/10.5555/3692070.3693262

Lucas, J. S., Maung, B.M., Tabar, M., McBride, K. and Lee, D. (2024). The Longtail Impact of Generative AI on Disinformation: Harmonizing Dichotomous Perspectives. IEEE Intelligent Systems, 39(5), pp.12–19. https://doi.org/10.1109/MIS.2024.3439109

Lucas, J., Uchendu, A., Yamashita, M., Lee, J., Rohatgi, S. & Lee, D. (2023). Fighting Fire with Fire: The Dual Role of LLMs in Crafting and Detecting Elusive Disinformation. In 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023 (pp. 14279–14305). Association for Computational Linguistics (ACL). https://doi.org/10.18653/v1/2023.emnlp-main.883

Macko, D., Kopal, J., Moro, R., & Srba, I. (2024). MultiSocial: Multilingual Benchmark of Machine-Generated Text Detection of Social-Media Texts. arXiv preprint arXiv:2406.12549. https://doi.org/10.48550/arXiv.2406.12549

Macko, D., Moro, R., & Srba, I. (2025). Increasing the Robustness of the Fine-tuned Multilingual Machine-Generated Text Detectors. arXiv preprint arXiv: 2503.15128. https://doi.org/10.48550/arXiv.2503.15128

Macko, D., Moro, R., Uchendu, A., Lucas, J., Yamashita, M., Pikuliak, M., ... & Bieliková, M. (2023). MULTITuDE: Large-Scale Multilingual Machine-Generated Text Detection Benchmark. In

Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (pp. 9960–9987). https://doi.org/10.18653/v1/2023.emnlp-main.616

Nahar, M., Lee, S., Guillen, R., & Lee, D. (2025). Generative AI Policies under the Microscope: How CS Conferences Are Navigating the New Frontier in Scholarly Writing. ACM Comm. of the ACM (CACM). https://doi.org/10.48550/arXiv.2410.11977

Nimmo, B. (2024). AI and covert influence operations: Latest trends. OpenAI: Technical report. https://downloads.ctfassets.net/kftzwdyauwt9/5IMxzTmUclSOAcWUXbkVrK/3cfab518e6b10789ab8843bcca18b633/Threat_Intel_Report.pdf

Pan, Y., Pan, L., Chen, W., Nakov, P., Kan, M.-Y., & Wang, W. (2023). On the risk of misinformation pollution with large language models. Findings of the Association for Computational Linguistics: EMNLP 2023, 1389–1403. https://doi.org/10.18653/v1/2023.findings-emnlp.97

Pikuliak, M., Srba, I., Moro, R., Hromadka, T., Smoleň, T., Melišek, M., Vykopal, I., Šimko, J., Podroužek, J., & Bieliková, M. (2023). MultiClaim: Multilingual Previously Fact-Checked Claim Retrieval (1.0) [Data set]. Zenodo. https://doi.org/10.5281/zenodo.7737983

Raza, S., Khan, T., Chatrath, V., Paulen-Patterson, D., Rahman, M., & Bamgbose, O. (2024). FakeWatch: a framework for detecting fake news to ensure credible elections. Social Network Analysis and Mining, 14(1), 142. https://doi.org/10.1007/s13278-024-01290-1

Shafayat, S., Kim, E., Oh, J., & Oh, A. (2024). Multi-FAct: Assessing Factuality of Multilingual LLMs using FActScore. In First Conference on Language Modeling. https://openreview.net/pdf?id=lkrH6ovzsj

Simon, F. M., Altay, S., & Mercier, H. (2023). Misinformation reloaded? Fears about the impact of generative AI on misinformation are overblown. Harvard Kennedy School (HKS) Misinformation Review. https://doi.org/10.37016/mr-2020-127

Sun, Z., Zhang, Z., Shen, X., Zhang, Z., Liu, Y., Backes, M., ... & He, X. (2024). Are We in the AI-Generated Text World Already? Quantifying and Monitoring AIGT on Social Media. arXiv preprint arXiv:2412.18148. https://doi.org/10.48550/arXiv.2412.18148

Wang, Y., Mansurov, J., Ivanov, P., Su, J., Shelmanov, A., Tsvigun, A., ... & Arnold, T. (2024). SemEval-2024 Task 8: Multidomain, Multimodel and Multilingual Machine-Generated Text Detection. In Proceedings of the 18th International Workshop on Semantic Evaluation (pp. 2057–2079). https://doi.org/10.18653/v1/2024.semeval-1.279

Wang, Y., Shelmanov, A., Mansurov, J., Tsvigun, A., Mikhailov, V., Xing, R., ... & Nakov, P. (2025). GenAI Content Detection Task 1: English and Multilingual Machine-Generated Text Detection: AI vs. Human. In Proceedings of the 1stWorkshop on GenAI Content Detection (GenAIDetect), 244–261. https://aclanthology.org/2025.genaidetect-1.27/

Weikmann, T., & Lecheler, S. (2023). Visual disinformation in a digital age: A literature synthesis and research agenda. New Media & Society, 25(12), 3696–3713. https://doi.org/10.1177/14614448221141648

Zaghouani, W., Jarrar, M., Habash, N., Bouamor, H., Zitouni, I., Diab, M., ... & AbuOdeh, M. (2024). The FIGNEWS Shared Task on News Media Narratives. In Proceedings of The Second Arabic Natural Language Processing Conference (pp. 530–547). https://doi.org/10.18653/v1/2024.arabicnlp-1.56

Zellers, R., Holtzman, A., Rashkin, H., Bisk, Y., Farhadi, A., Roesner, F., & Choi, Y. (2019). Defending against neural fake news. In Proceedings of the 33rd International Conference on Neural Information Processing Systems (pp. 9054–9065). https://dl.acm.org/doi/10.5555/3454287.3455099

Zhuo, T. Y., Huang, Y., Chen, C., & Xing, Z. (2023). Red teaming ChatGPT via jailbreaking: Bias, robustness, reliability and toxicity. arXiv preprint arXiv:2301.12867. https://doi.org/10.48550/arXiv.2301.12867

Zugecova, A., Macko, D., Srba, I., Moro, R., Kopal, J., Marcincinova, K., Mesarcik, M. (2024). Evaluation of LLM Vulnerabilities to Being Misused for Personalized Disinformation Generation. arXiv preprint arXiv:2412.13666. https://doi.org/10.48550/arXiv.2412.13666

**Competing interests**
The authors declare no competing interests.

**Ethics**
No human or animal experiments have been conducted in this research. No human subjects were involved, requiring informed consent or ethical review. FakeNews is published under MIT license, USC_X and FIGNEWS have CC BY-NC-SA 4.0 license, SemEval and GenAI datasets have Apache-2.0 license. A part of these datasets may, however, consist of copyrighted material. The remaining datasets have restricted access, can be obtained upon request for research purpose only and do not allow their resharing. Considering potential inclusion of copyrighted materials as well as restrictions under which datasets are published, we decided to follow data minimization principle when sharing the data necessary to replicate our work. More specifically, we publish the IDs of original dataset records (without resharing the human-written/machine-generated texts) together with all metadata required for replication (identification of languages, social media platforms, etc.) and our predicted labels from both machine-generated text detectors (including the prediction probabilities). All datasets and models have been used in a way following their terms of use and licensing (used for research purpose only). Replication materials produced as part of this research do not disclose personally identifiable information. When working with the models' predictions, we prioritized lowering the number of false positives; thus, the results presented in the paper may be biased toward underestimating the actual prevalence of machine-generated texts rather than to inflating it.

# Appendix: Supplementary data

Figure A1 illustrates the proportion of machine-generated texts (as detected by the combined confident prediction) in the individual social-media platforms available in the MultiClaim dataset. Figure A2 illustrates the proportion of machine-generated texts in the MultiClaim dataset across the assigned fact-checking rating categories. Similarly, Figure A3, Figure A4, and Figure A5 illustrate the proportion of machine-generated texts in the FakeNews and FIGNEWS datasets based on available labels and language identification.



*Figure A1. Per-platform proportion of the texts detected to be machine-generated (Telegram contains less than 1,000 samples). Proportion differs among platforms and is the lowest in Twitter (X).*
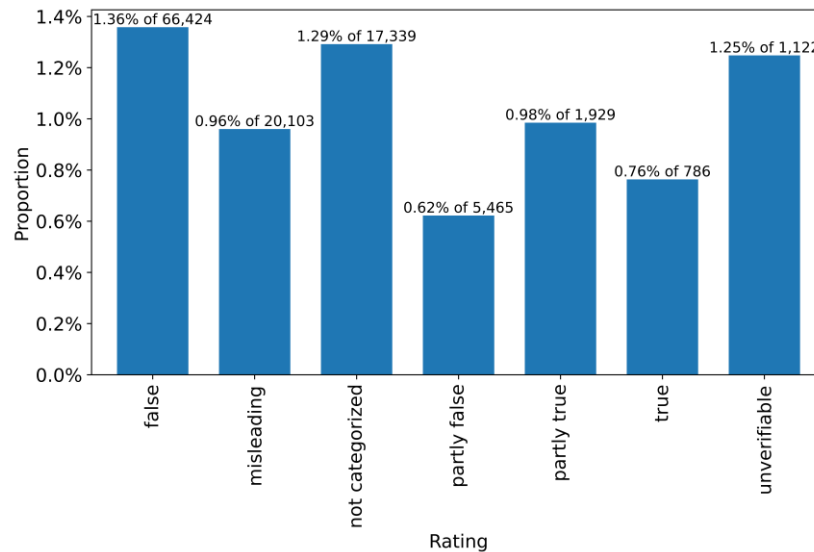


*Figure A2. Per-rating proportion of the texts detected to be machine-generated. Proportion is the highest in case of "false" and "not categorized" texts. The "true" rating group covers less than 1,000 samples.*
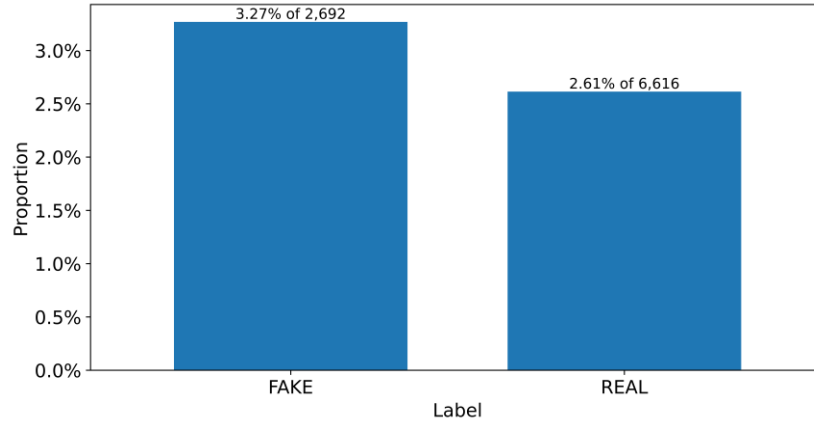
**Figure A3. Per-label proportion of the texts detected to be machine-generated in the FakeNews dataset. Proportion is slightly higher in the "FAKE" texts.**
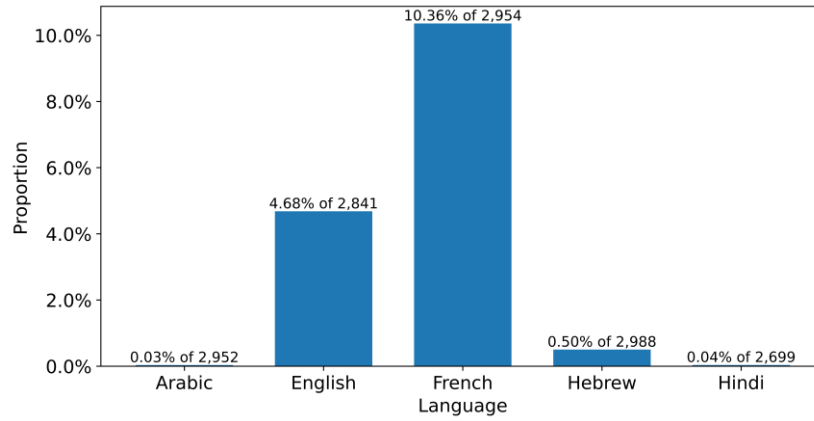
**Figure A4. Per-language proportion of the texts detected to be machine-generated in the FIGNEWS dataset. Proportion is the highest in case of French and English.**
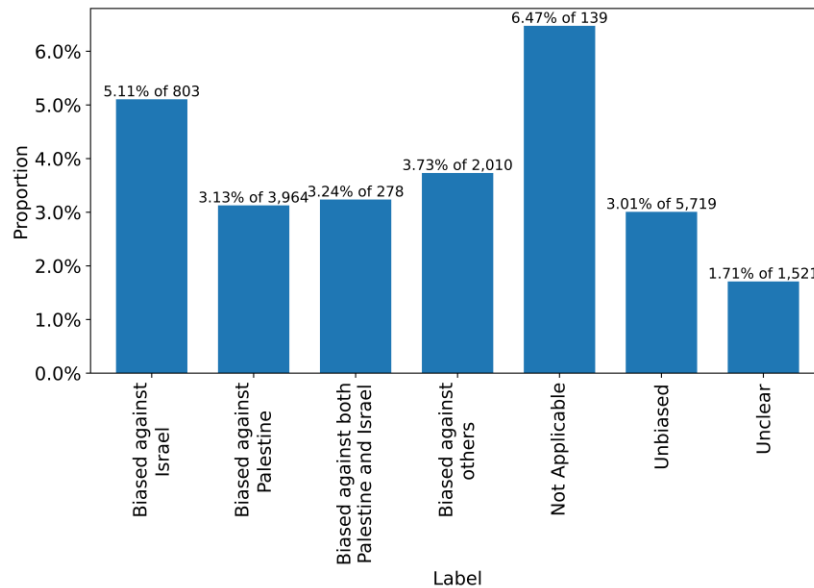
**Figure A5. Per-label proportion of the texts detected to be machine-generated in the FIGNEWS dataset. Proportion is the highest for "Not Applicable" texts and for "Biased against Israel", both covering less than 1,000 samples.**