

# A Hybrid Approach for Detecting AI-Generated Text in English

Yanyang Gong, Yunfan Zhou

Department of Informatics, University of Zurich

{yanyang.gong, yunfan.zhou}@uzh.ch

## Abstract

Detecting AI-generated text is becoming harder as generators improve and writing styles vary across domains. To address the COLING 2025 GenAI Content Detection Task 1 (English), we propose a hybrid detector that combines a train-free statistical signal with a train-based classifier, allowing the two to complement each other. Experiments show that the hybrid model consistently outperforms both individual baselines on validation and test datasets. We further analyze the model’s performance across different domains and sources to assess its robustness.

## 1 Introduction

Rapid advances in large language models (LLMs) have made AI-generated text widespread, raising concerns about academic integrity and content attribution. At the same time, they blur the boundary between human- and machine-written text, making detection increasingly important and challenging.

This project therefore addresses the GenAI Content Detection Task 1 from COLING 2025 (Wang et al., 2025). The task consists of two subtasks: Monolingual subtask A and Multilingual subtask B. In this work, we focus on subtask A, which requires binary classification of English text as either AI-generated or human-written.

To address this task, we review two representative approaches for AI-generated text detection: train-based discriminative models (e.g., DeBERTa classifiers (Gritsai et al., 2025)) and train-free statistical methods (e.g., perplexity-based approaches (Zhu et al., 2025)). Train-based models have been the mainstream approach in this task. However, prior studies point out that these models often overfit to in-distribution features, resulting in poor generalization to out-of-distribution texts (Zhu et al., 2025). On the other hand, train-free methods are also vulnerable: their performance can be undermined when AI-generated text closely mimics hu-

man writing styles (Wang et al., 2025) or when applied across different domains.

Motivated by the limitations, our project aims to investigate: *Can we combine complementary signals from train-free and train-based methods to improve robustness and overall detection performance?* To answer this question, we propose a hybrid model that systematically combines these two approaches.

The main contribution of this paper can be summarized as follows: (1) We propose a hybrid classifier framework that combines train-based and train-free methods, leveraging their complementary strengths. (2) We show that the hybrid model consistently outperforms each individual method on both the validation and test sets, including the top-ranked train-based model in the COLING 2025 Task 1 and a train-free baseline. (3) We further examine classifier performance across domains, providing additional insights and analyzing model robustness accordingly.

## 2 Related work

Existing approaches for detecting machine-generated text can be broadly divided into train-based and train-free methods.

Train-based methods build classifiers using neural networks or fine-tune pre-trained models, such as RoBERTa (Solaiman et al., 2019), Llama (Doan and Inui, 2025), or DeBERTa (Gritsai et al., 2025). However, these methods often suffer from overfitting and poor generalization to out-of-distribution texts.

Train-free methods rely on statistical features without requiring model training. Representative approaches include DetectGPT (Mitchell et al., 2023), which leverages probability curvature, and DNA-DetectLLM (Zhu et al., 2025), which calculates cumulative repair effort to transform input into an “ideal” AI-generated sequence. A key ad-

vantage is their low deployment cost, though performance can vary substantially across domains and generation models.

Some methods incorporate statistical features in addition to classifiers. For instance, a recent study extracts linguistic and statistical markers and feeds them into a classifier such as Random Forest (Joseph et al., 2025). Our work follows a similar idea: starting from a train-based DeBERTa classifier (Gritsai et al., 2025), we incorporate a train-free statistical score (computed by DNA-DetectLLM (Zhu et al., 2025)) and fuse it with the model representations for final classification.

### 3 Methodology

#### 3.1 Baseline Methods

**DNA-DetectLLM** (Zhu et al., 2025) is a train-free method designed to detect AI-generated text. The core idea is based on the concept of "mutation repair." For any given input text, an "ideal" AI-generated sequence is first constructed by selecting the most probable token at each position. The difference between the input text and the ideal sequence represents the "mutated positions." The cumulative effort required to "repair" these mutations is quantified as the repair score, which serves as the signal for classification. If the repair score exceeds a calibrated threshold, the text is classified as human-written; otherwise, it is classified as AI-generated.

The repair score is mathematically defined as:

$$R(s) = \frac{1}{T+1} \sum_{t=0}^T \sigma(s_t | s) \quad (1)$$

where  $s$  is the original sequence,  $s_t$  is the sequence after  $t$  repair steps,  $T$  is the total number of mutated tokens to be corrected, and  $\sigma(s_t | s)$  represents the conditional score (calculated by log-perplexity and cross-perplexity) at each step. Alternatively, the score can be simplified as:

$$R(s) = \frac{1}{2} (\sigma(s) + \sigma(\hat{s} | s)), \quad (2)$$

where  $\hat{s}$  is the ideal sequence. The proof of this simplification can be found in the original paper (Zhu et al., 2025).

This method demonstrates strong performance and robustness compared to state-of-the-art detection techniques. It is computationally efficient, as it does not require training, but the choice of the

large language model (LLM) used for inference and perplexity calculation is critical. Additionally, the calibrated threshold significantly impacts performance and is sensitive to the dataset distribution

**DeBERTa-Based Classifier** (Gritsai et al., 2025) is proposed by Advacheck team, who participated in the GenAI Detection Task 1 competition and achieved the best performance in the English sub-task. The model uses DeBERTa as the backbone and incorporates additional classification heads. These heads are not only used for binary classification (human-written vs. AI-generated text) but also for domain-specific classification within the dataset. This architecture is efficient, adaptable, and flexible, as the classification heads can be easily adjusted to suit specific datasets and tasks. However, the model's robustness is limited, as performance drops were observed on the test set. Additionally, the final classification also relies on a threshold.

Figures for the baseline models architectures can be found in Appendix B.

#### 3.2 Proposed Hybrid Model

Inspired by the strengths of both DNA-DetectLLM and the DeBERTa-based classifier, we propose a hybrid model that combines their complementary advantages. Our goal is to design a model that is flexible (capable of leveraging pre-trained models), lightweight (able to use smaller LLMs for inference to reduce computational costs), and robust (enhanced by incorporating statistical signals from train-free methods).

The proposed hybrid model integrates the repair score from DNA-DetectLLM as an additional feature into the DeBERTa-based classifier. The architecture consists of the following steps:

- Repair Score Calculation: a set of lightweight LLM is used to calculate a simplified repair score for each input text.
- Feature-Level Integration: the repair score is concatenated with the DeBERTa pooled output embeddings as an additional feature.
- Modified DeBERTa Model: the DeBERTa-v3 model is modified to accept the repair score as an input feature.
- Pre-trained weights are loaded for the DeBERTa backbone, while the classification head is reinitialized to accommodate the additional repair score feature.

- Fine-Tuning: the classification layer is fine-tuned on a small dataset. Instead of relying on a threshold, the model uses the argmax function to output the most probable class.

Mathematically, the hybrid model can be expressed as:

$$y = \text{argmax} (W \cdot [f_{\text{DeBERTa}}(x), R(s)] + b) \quad (3)$$

where  $f_{\text{DeBERTa}}(x)$  represents the pooled output embeddings from the pretrained DeBERTa model,  $R(s)$  is the repair score.

## 4 Experiments

### 4.1 Dataset

We conduct experiments on the subsets of COLING 2025 MGT Detection dataset (English sub-task) (Wang et al., 2025), which consists of training, development, and test sets collected from multiple domains and sources. Notably, there exists a distribution shift between the training/development sets and the test set, as the test data are collected from a different source.

We apply source-stratified proportional sampling to construct balanced subsets that preserve the original distribution within computational constraints. Detailed dataset breakdowns are provided in Appendix A.

### 4.2 Implementation Details

**DNA-DetectLLM (Baseline 1):** We use two sets of lightweight language models (gpt2-medium and gpt2-large, EleutherAI/pythia-160m and EleutherAI/pythia-410m) to compute repair scores, which are then standardized using training set statistics. The classification threshold is optimized on the training set by maximizing F1-score.

**DeBERTa-Based Classifier (Baseline 2):** We employ the pretrained DeBERTa-v3-base model from Hugging Face Transformers (OU-Advacheck, 2024), optimizing its classification threshold on the validation set. We report the best result over multiple trials<sup>1</sup>.

**Proposed Hybrid Model:** Our hybrid model integrates the DNA-DetectLLM repair score as an additional feature into the DeBERTa-v3-base architecture. The model’s classification head was modified to accept concatenated inputs: the 768-dimensional pooled output from DeBERTa and the

<sup>1</sup>We note that the pretrained model’s label mapping required verification, we tested both configurations and used the one yielding better performance

1-dimensional normalized repair score, resulting in a 769-dimensional feature vector. The classification head consists of three linear layers (768 + 1, 512, 256, 2) with GELU activation and dropout ( $p=0.5$ ).

During training, the DeBERTa backbone and pooler are frozen, and only the classification head is updated. We employ the AdamW optimizer with a learning rate of  $3 \times 10^{-4}$ , coupled with a linear learning rate schedule and 50 warm-up steps. Training is performed with a batch size of 32 for a single epoch. Unlike baseline methods, our hybrid model relies on the argmax operation for classification, eliminating the need for threshold.

All experiments were implemented in PyTorch using the Hugging Face Transformers library and conducted on Google Colab with T4 GPUs.

**Evaluation Metrics:** All models were evaluated using Accuracy and F1 Score.

## 5 Results

### 5.1 Overall Results

Table 1 compares the performance of baseline and hybrid models on the Dev and Test sets. Both Hybrid-GPT and Hybrid-Pythia consistently outperform Baseline1 and Baseline2 across all metrics on both splits, demonstrating the effectiveness of the proposed hybrid approach.

A general performance degradation from Dev to Test is observed for most models, reflecting the distribution shift between the two sets and aligning with prior findings. Although the hybrid models maintain clear advantages over the baselines on the Test set, the performance gap between Dev and Test remains comparable to that of Baseline2, suggesting limited improvement in robustness to distribution shift.

Hybrid-GPT performs better in baseline 1 than Hybrid-Pythia, indicates the performance of DNA-Detect LLM relies on the LLM used as we analyzed before, and they two achieve very similar results in hybrid model, indicating that lightweight LLMs with comparable capacities yield similar gains in the hybrid framework. Finally, we observed unstable performance of Baseline2, thus we report the best result obtained across multiple independent runs.

### 5.2 Finer Analysis by Sub Source

We further evaluate the performance by sub\_source (domain for test set). Figure 1 presents the F1

Model	Dev		Test	
	F1	Accuracy	F1	Accuracy
Baseline1 (GPT)	74.61	60.34	73.81	63.06
Baseline1 (Pythia)	74.52	60.20	70.14	55.72
Baseline2	84.07	81.26	73.93	67.46
Hybrid-GPT	<b>96.50</b>	<b>95.60</b>	80.10	76.52
Hybrid-Pythia	96.43	95.52	<b>80.42</b>	<b>77.02</b>

Table 1: Performance comparison of baseline and hybrid models on the Dev and Test sets. For Baseline-2, we report the best performance over multiple trials.

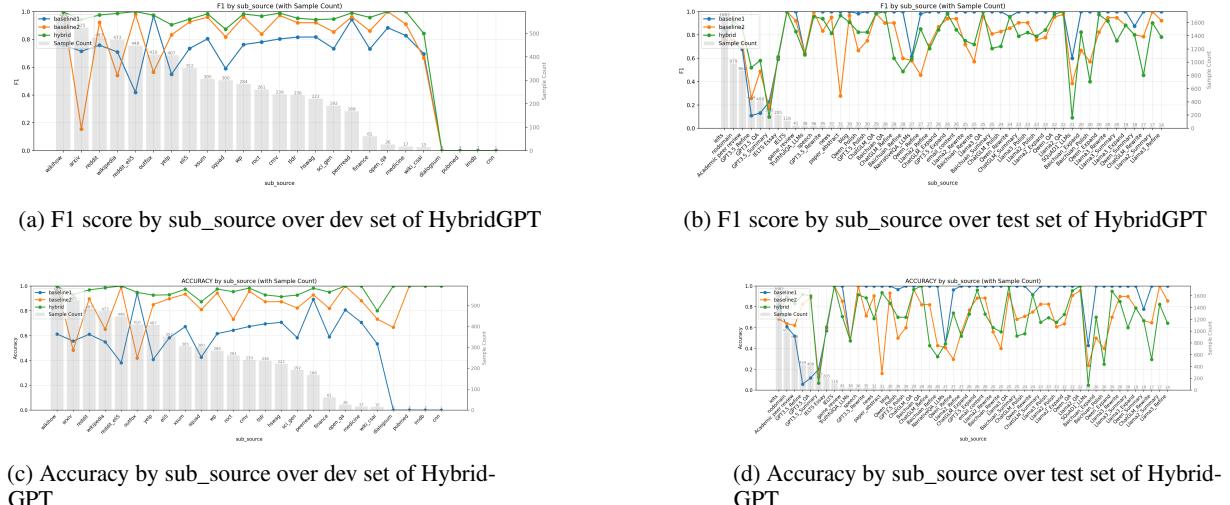


Figure 1: Performance by sub\_source of HybridGPT

scores of the baseline models and the hybrid model (GPT) across different sub\_sources, with bars indicating the sample counts for the validation and test sets (the train set exhibits a similar distribution to the validation set, see Appendix A).

On the dev set, the hybrid models consistently outperform the baselines across all sub\_sources. On the Test set, however, baseline1 achieves the best performance in most sub\_sources, except for the major sub\_sources with more samples, suggesting that baseline1 is relatively robust to distribution shifts, whereas the hybrid model still relies on training data for optimal performance.

We also evaluated model performance across multiple metrics, both by sub\_source and by model. Additional results and detailed breakdowns are provided in the Appendix C.

## 6 Conclusion and Limitations

In this work, we proposed a hybrid framework for AI-generated text detection that combines train-free statistical scores with train-based discriminative models. Across both the Dev and Test sets, the

hybrid models consistently outperform both baselines, showing that signals from these two types of methods can be effectively integrated.

A key advantage of our approach is that it achieves performance improvements using relatively lightweight models without requiring complex architectural modifications or extensive re-training. However, our results also show that while the hybrid approach achieves stable improvements on domains with sufficient training data, it shows limited gains or even underperforms baselines under low-resource or out-of-domain conditions, highlighting remaining sensitivity to distribution shift.

Finally, this work has several limitations. First, due to computational resource constraints, we trained and evaluated our models on a limited subset of the data. As a result, our results are not directly comparable to the official competition submissions, which use the full dataset. In addition, the detector models used in our experiments are relatively small compared to state-of-the-art LLMs, which may limit the performance of our approach. Future work could explore scaling the hybrid frame-

work to larger detector architectures and training on larger and more diverse datasets.

## References

Nhi H. Doan and Kentaro Inui. 2025. Grape at genai detection task 1: Leveraging compact models and linguistic features for robust machine-generated text detection. In *Proceedings of the 1st Workshop on GenAI Content Detection (GenAIDetect)*, pages 209–217.

G. Gritsai, A. Voznuyk, I. Khabutdinov, and A. Grabovoy. 2025. Advacheck at genai detection task 1: AI detection powered by domain-aware multi-tasking. In *Proceedings of the 1st Workshop on GenAI Content Detection (GenAIDetect)*, pages 236–243.

Emmanuel Joseph, Micheal Bennet, and Thomas Kingsley. 2025. Feature-based detection of ai-generated text: An analysis of stylometric and perplexity markers in contemporary large language models. ResearchGate. Preprint.

Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. Detectgpt: Zero-shot machine-generated text detection using probability curvature. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, pages 24950–24962. PMLR.

OU-Advacheck. 2024. Deberta-v3-base fine-tuned for ai-generated text detection. <https://huggingface.co/OU-Advacheck/deberta-v3-base-draigenc-mgt1a>. Hugging Face model.

Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, and 1 others. 2019. Release strategies and the social impacts of language models. *arXiv preprint arXiv:1908.09203*.

Yuxia Wang, Artem Shelmanov, Jonibek Mansurov, Akim Tsvigun, Vladislav Mikhailov, Rui Xing, Zhubhan Xie, Jiahui Geng, Giovanni Puccetti, Ekaterina Artemova, Jinyan Su, Minh Ngoc Ta, Mervat Abassy, Kareem Ashraf Elozeiri, Saad El Dine Ahmed El Eter, Maiya Goloburda, Tarek Mahmoud, Raj Vardhan Tomar, Nurkhan Laiyk, and 7 others. 2025. Genai content detection task 1: English and multilingual machine-generated text detection: Ai vs. human. In *Proceedings of the 1st Workshop on GenAI Content Detection (GenAIDetect)*, pages 244–261, Abu Dhabi, UAE. International Conference on Computational Linguistics.

Xiaowei Zhu, Yubing Ren, Fang Fang, Qingfeng Tan, Shi Wang, and Yanan Cao. 2025. DNA-detectLLM: Unveiling AI-generated text via a DNA-inspired mutation-repair paradigm. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.

## A Dataset Breakdown Analysis

This section provides additional statistics of the dataset.

Sub_source	Train	Val
<b>*** Train + Val Sub_source Statistics ***</b>		
arxiv	1791 (8.76%)	525 (8.75%)
cmv	752 (3.76%)	239 (3.98%)
cnn	100 (0.48%)	31 (0.53%)
dialogsum	12 (0.06%)	3 (0.05%)
fill	100 (5.97%)	30 (5.25%)
Finance	139 (0.76%)	61 (1.02%)
hawkg	851 (4.26%)	222 (3.70%)
lamb	2 (0.01%)	2 (0.01%)
medicine	56 (0.28%)	17 (0.28%)
open_qa	120 (0.60%)	26 (0.43%)
pubmed	665 (3.53%)	158 (2.80%)
reddit_all	481 (2.28%)	448 (7.47%)
reddit_slangs	1488 (7.44%)	242 (4.25%)
news	100 (0.52%)	32 (0.52%)
scl_gen	784 (3.52%)	192 (3.20%)
squad	100 (0.52%)	30 (0.50%)
tiny	100 (0.52%)	23 (0.37%)
wikit_cse	36 (0.18%)	15 (0.25%)
wikipedia	1757 (8.79%)	565 (9.42%)
wikispecies	100 (0.52%)	43 (0.75%)
wp	954 (4.77%)	254 (4.73%)
xsum	969 (4.44%)	305 (5.09%)
yelp	1383 (6.92%)	407 (6.79%)
TOTAL	19999	5998
	TOTAL	25997

Figure 2: Sample distribution across sub\_sources and models for the train and validation sets.

Sub_source	Test
<b>*** Test Sub_source Statistics ***</b>	
Academic_peer_review	866 (44.48%)
Baichuan_Expand	20 (0.33%)
Baichuan_Polish	20 (0.33%)
Baichuan_QA	28 (0.47%)
Baichuan_Refine	28 (0.47%)
Baichuan_Rewrite	25 (0.42%)
Baichuan_Summary	25 (0.42%)
ChatGPT_Expand	26 (0.43%)
ChatGPT_Polish	24 (0.40%)
ChatGPT_Refine	25 (0.43%)
ChatGPT_Rewrite	28 (0.47%)
ChatGPT_Summary	17 (0.28%)
IELTS	116 (1.93%)
IELTS_Essay	205 (3.42%)
Llama2_Expand	22 (0.37%)
Llama2_Polish	22 (0.37%)
Llama2_QA	22 (0.37%)
Llama2_Refine	27 (0.45%)
Llama2_Rewrite	25 (0.42%)
Llama2_Summary	17 (0.29%)
Llama3_Expand	19 (0.32%)
Llama3_Polish	23 (0.38%)
Llama3_QA	25 (0.42%)
Llama3_Refine	14 (0.23%)
Llama3_Rewrite	20 (0.33%)
Llama3_Summary	20 (0.33%)
NarrativeQa_Expand	27 (0.45%)
Qwen_Expand	26 (0.33%)
Qwen_Polish	26 (0.33%)
Qwen_QA	22 (0.37%)
Qwen_Refine	27 (0.45%)
Qwen_Rewrite	24 (0.40%)
Qwen_Summary	18 (0.30%)
SQuAD_LLMs	21 (0.39%)
TruthfulQA_LLMs	38 (0.63%)
blip	38 (0.50%)
email_content	26 (0.45%)
game_review	41 (0.65%)
ilets	1682 (28.85%)
news	32 (0.53%)
nodomain	979 (16.33%)
paper_abstract	31 (0.52%)
speech	36 (0.60%)
TOTAL	5996

Figure 3: Sample distribution across sub\_sources and models for the test set.

## B Baseline Model Architecture

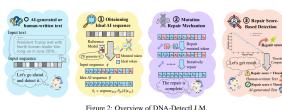


Figure 4: Architecture of the train-free baseline mode. Adapted from (Zhu et al., 2025)

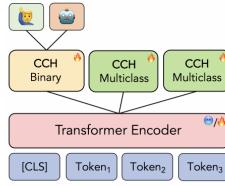


Figure 1: Overview of the proposed multi-task architecture. Modules marked only with  $\heartsuit$  are trainable at all stages. The weights of the Transformer Encoder are frozen  $\heartsuit$  at the first stage of training and trainable  $\heartsuit$  at the second one. The Custom Classification Head (CCH) described in Appendix A is used for predictions.

Figure 5: Architecture of the train-based baseline mode.  
Adapted from ([OU-Advacheck](#), 2024)

## C Additional Performance Metrics

This section presents additional evaluation metrics for the hybrid model, including performance breakdowns by generation model and detailed precision/recall analysis across sub\_sources.

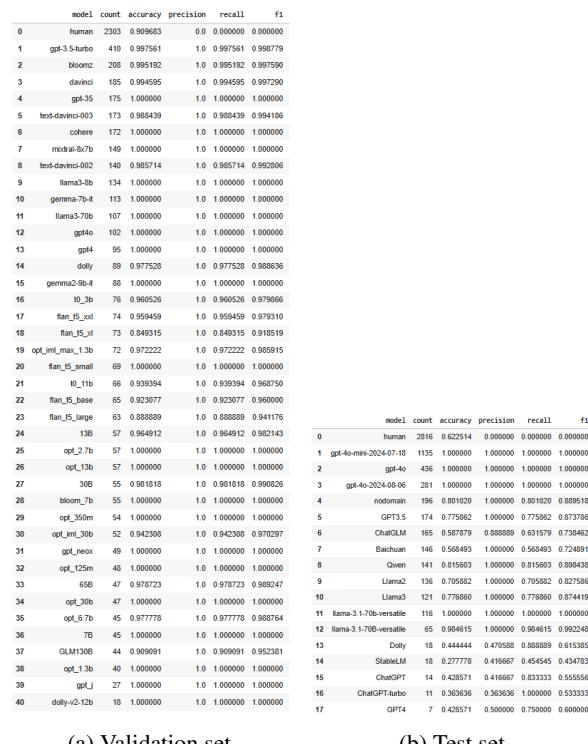


Figure 6: Hybrid-GPT performance by generation model on validation and test sets

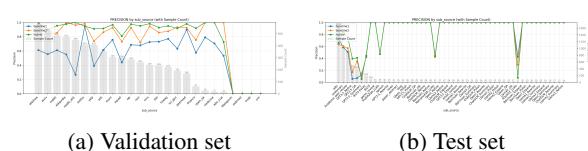


Figure 7: Precision comparison by sub\_source for Hybrid-GPT across validation and test sets.

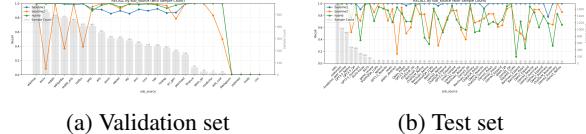


Figure 8: Recall comparison by sub\_source for Hybrid-GPT across validation and test sets.

	model	count	accuracy	precision	recall	f1
0	human	2353	0.90249	0.0	0.00000	0.00000
1	gpt-3.5-turbo	410	0.987561	1.0	0.97956	0.987735
2	bloomz	208	0.986912	1.0	0.986912	0.987920
3	davinci	185	0.984593	1.0	0.984593	0.985493
4	gpt-35	175	0.100000	1.0	1.00000	1.00000
5	text-davinci-003	173	0.984849	1.0	0.984849	0.984849
6	cohere	172	0.100000	1.0	1.00000	1.00000
7	mintral-8x7a	149	0.100000	1.0	1.00000	1.00000
8	text-davinci-002	140	0.986714	1.0	0.986714	0.988006
9	llama3-8b	134	0.100000	1.0	1.00000	1.00000
10	germea-7b-1R	130	0.100000	1.0	1.00000	1.00000
11	llama3-7b0	107	0.100000	1.0	1.00000	1.00000
12	gpt40	102	0.100000	1.0	1.00000	1.00000
13	gpt4	95	0.100000	1.0	1.00000	1.00000
14	dolly	89	0.986262	1.0	0.986262	0.985285
15	germea-2b-8b	88	0.100000	1.0	1.00000	1.00000
16	10_3b	76	0.980526	1.0	0.980526	0.980526
17	flan_t5_xxl	75	0.984593	1.0	0.984593	0.979303
18	flan_t5_xl	73	0.835616	1.0	0.835616	0.914446
19	opti_imax_max_13b	72	0.972222	1.0	0.972222	0.986915
20	flan_t5_small	69	0.100000	1.0	1.00000	1.00000
21	to_1B	65	0.939384	1.0	0.939384	0.939384
22	flan_t5_base	65	0.923077	1.0	0.923077	0.960000
23	flan_t5_large	63	0.888899	1.0	0.888899	0.911790
24	13B	57	0.984692	1.0	0.984692	0.982142
25	opt_2.7b	57	0.100000	1.0	1.00000	1.00000
26	opt_15b	57	0.100000	1.0	1.00000	1.00000
27	30B	55	0.983636	1.0	0.983636	0.984887
28	bloomz_7B	55	0.100000	1.0	1.00000	1.00000
29	opt_350m	54	0.100000	1.0	1.00000	1.00000
30	opt_imax_300	52	0.942308	1.0	0.942308	0.970293
31	gpt_neox	49	0.100000	1.0	1.00000	1.00000
32	opt_125m	48	0.100000	1.0	1.00000	1.00000
33	65B	47	0.978723	1.0	0.978723	0.982333
34	opt_30B	47	0.100000	1.0	1.00000	1.00000
35	opt_6.7b	45	0.977778	1.0	0.977778	0.988764
36	7B	45	0.977778	1.0	0.977778	0.988764
37	GLM1303	44	0.909091	1.0	0.909091	0.923838
38	opt_1.3b	40	0.100000	1.0	1.00000	1.00000
39	gpt_L	27	0.100000	1.0	1.00000	1.00000

(a) Validation set (b) Test set

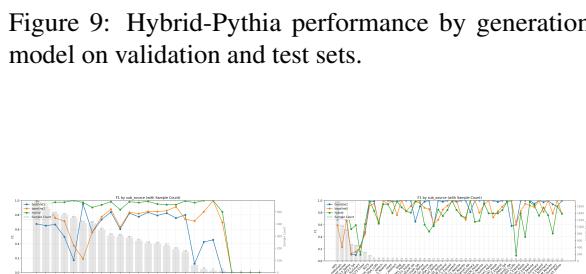


Figure 9: Hybrid-Pythia performance by generation model on validation and test sets.

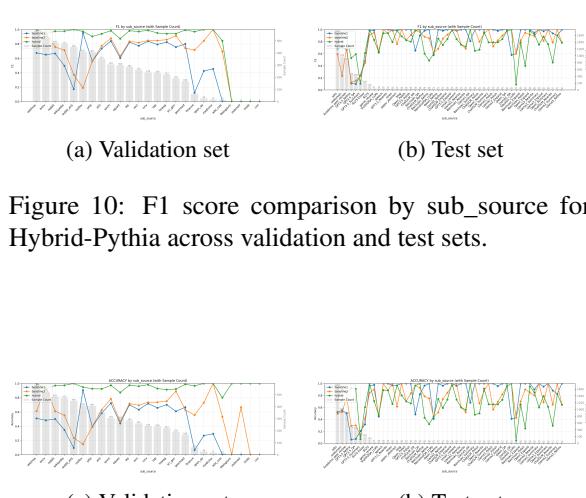


Figure 11: Accuracy comparison by sub\_source for Hybrid-Pythia across validation and test sets.

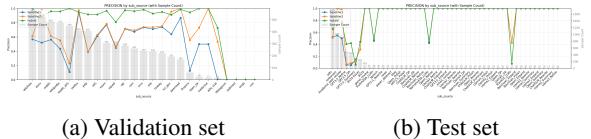


Figure 12: Precision comparison by sub\_source for Hybrid-Pythia across validation and test sets.

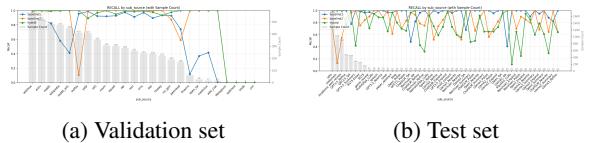


Figure 13: Recall comparison by sub\_source for Hybrid-Pythia across validation and test sets.