

The Skipped Beat: A Study of Sociopragmatic Understanding in LLMs for 64 Languages

Chiyu Zhang^ξ Khai Duy Doan^{λ,*} Qisheng Liao^{λ,*} Muhammad Abdul-Mageed^{ξ,λ}

^ξDeep Learning & Natural Language Processing Group, The University of British Columbia

^λDepartment of Natural Language Processing & Department of Machine Learning, MBZUAI

{chiyuzh@mail, muhammad.mageed@}.ubc.ca,

{duy.doan, qisheng.liao}@mbzuai.ac.ae

Abstract

Instruction tuned large language models (LLMs), such as ChatGPT, demonstrate remarkable performance in a wide range of tasks. Despite numerous recent studies that examine the performance of instruction-tuned LLMs on various NLP benchmarks, there remains a lack of comprehensive investigation into their ability to understand cross-lingual sociopragmatic meaning (SM), i.e., meaning embedded within social and interactive contexts. This deficiency arises partly from SM not being adequately represented in any of the existing benchmarks. To address this gap, we present SPARROW, an extensive multilingual benchmark specifically designed for SM understanding. SPARROW comprises 169 datasets covering 13 task types across six primary categories (e.g., anti-social language detection, emotion recognition). SPARROW datasets encompass 64 different languages originating from 12 language families representing 16 writing scripts. We evaluate the performance of various multilingual pretrained language models (e.g., mT5) and instruction-tuned LLMs (e.g., BLOOMZ, ChatGPT) on SPARROW through fine-tuning, zero-shot, and/or few-shot learning. Our comprehensive analysis reveals that existing open-source instruction tuned LLMs still struggle to understand SM across various languages, performing close to a random baseline in some cases. We also find that although ChatGPT outperforms many LLMs, it still falls behind task-specific finetuned models with a gap of 12.19 SPARROW score. Our benchmark is available at: <https://github.com/UBC-NLP/SPARROW>

1 Introduction

Multilingual LLMs have recently transformed NLP, due to their powerful capabilities on a

* Equal contribution

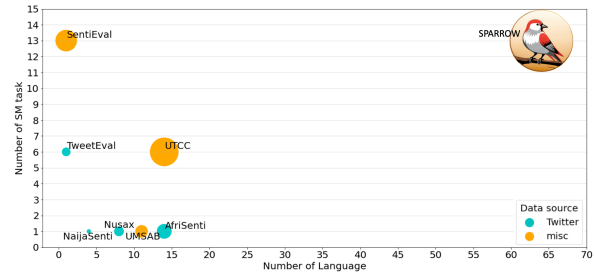


Figure 1: Comparison of SM benchmarks with leaderboards. The bubble size indicates the number of datasets. Previous works: TweetEval (Barbieri et al., 2020), UMSAB (Barbieri et al., 2022), Nusax (Winata et al., 2022), UTCC (Risch et al., 2021), NaijaSenti (Muhammad et al., 2022), AfriSenti (Muhammad et al., 2023a), SentiEval (Zhang et al., 2023b).

wide range of tasks (Xue et al., 2021; Scao et al., 2022). Methods such instruction tuning and reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022) have further enhanced the zero-shot generalizability of these models. Notably, ChatGPT exhibits impressive capabilities in this regard. Human language, however, is intrinsically ambiguous and far from solved. In fact, some forms of meaning are deeply embedded in social interactions. We collectively refer to this type of meaning as *sociopragmatic meaning* (SM). To appreciate SM, consider how the meaning of an utterance in social interaction (e.g., on social media) can be highly subtle and how it incorporates both the social variation related to language users (from a sociolinguistics perspective) (Tagliamonte, 2015) and their communicative intentions (from a pragmatics perspective) (Boxer and Cortés-Conde, 2021). Although SM is quite established within linguistics, NLP systems still struggle with this type of meaning that is intertwined in social and interactive contexts (Zhang and Abdul-Mageed, 2022). The extent to which instruction tuned models such as ChatGPT can capture SM across languages re-

mains largely unclear as these models are yet to be evaluated on appropriate datasets under standardized conditions easy to replicate.

To facilitate evaluation of LLMs and enhance fairness of model comparisons and reproducibility, early work introduces evaluation benchmarks. Most existing benchmarks, however, focus on the monolingual setting. These include GLUE (Wang et al., 2019), SentEval (Conneau and Kiela, 2018), and TweetEval (Barbieri et al., 2020) for English, ARLUE (Abdul-Mageed et al., 2021) for Arabic, CLUE (Xu et al., 2020a) for Chinese, and IndoNLU (Wilie et al., 2020) for Indonesian. Although XTREME (Hu et al., 2020) and XGLUE (Liang et al., 2020) introduce multilingual benchmarks, they only include a few SM tasks for a limited number of languages. They are also limited to standard language use (e.g., Wikipedia). Barbieri et al. (2022) propose a multilingual sentiment analysis benchmark (UMSAB), but it solely contains tweet sentiment analysis datasets in only eight languages. As such, absence of a unified, diverse, and comprehensive benchmark and a fragmented evaluation landscape hamper NLP work for cross-lingual SM.

Another challenge for SM research is the issue of *data inaccessibility* (Assenmacher et al., 2022). Although many studies release the IDs of posts (e.g., tweets), substantial amounts of these social posts become inaccessible over time due to deletion, etc. (Zhang et al., 2022). In our benchmark, we attempt to re-collect text contents of 25 datasets by using their IDs but can only retrieve 58% samples on average (see Table 8 in Appendix). This data decay also hinders fair comparisons in NLP for SM research. This issue has already become worse as corporations such as Twitter tighten access to their API, making it even harder to collect historical data. To address this bottleneck, we introduce a massively multilingual SM evaluation benchmark, dubbed *SPARROW*, that comprises 169 datasets covering 64 languages from 12 language families, 16 types of scripts, across diverse online platforms (e.g., Twitter, YouTube, and Weibo). We then perform an extensive evaluation of ChatGPT, comparing it to 13 other models ranging in size between 110M-7B parameters. Our evaluations allow us to answer multiple questions related to how it is that LLMs fare across languages on SM. To facilitate future comparisons, we also design a modular, interac-

Studies	Lang.	Tasks	SM Tasks	Dataset	Models	LeaderBrd
Zhong et al. (2023)	en	5	1	8	5	✗
Qin et al. (2023)	en	7	1	20	29	✗
Ahuja et al. (2023)	70	10	3	16	11	✗
Laskar et al. (2023)	12	12	2	140	27	✗
Bang et al. (2023)	8	8	1	23	3	✗
Lai et al. (2023)	37	7	0	8	7	✗
Das et al. (2023)	11	2	2	2	1	✗
Wang et al. (2023)	en	5	5	18	3	✗
Zhang et al. (2023b)	en	13	13	26	5	✓
Ziems et al. (2023)	en	24	18	24	13	✗
Ours	64	13	13	169	14	✓

Table 1: Our work in comparison.

tive leaderboard on top of our new benchmark.

To summarize, the contributions of this paper are as follows: **(1)** We collect, uniformize, and responsibly release massively multilingual SM datasets in a *new benchmark*; **(2)** Our SPARROW benchmark is essentially an archive of SM datasets that alleviates the serious issue of *data decay*; **(3)** We evaluate a wide range of models on our SPARROW benchmark via fine-tuning SoTA encoder-only pretrained language models and zero-shot learning of a number of generative models, including instruction tuned models (e.g., BLOOMZ) as well as ChatGPT; and **(4)** We establish standard settings for future research in this area across a large number of languages and tasks, through a *public leaderboard*.

2 Related Work

Evaluation of LLMs. There have been many attempts to evaluate ChatGPT and instruction tuned LLMs. Qin et al. (2023); Laskar et al. (2023); Zhong et al. (2023); Wu et al. (2023) utilize existing English evaluation benchmarks, such as GLUE (Wang et al., 2019) and BigBench (Srivastava et al., 2022), to evaluate LLMs’ capacities on various NLP tasks. These studies find that although ChatGPT performs less effectively than the models finetuned specifically for each task, it demonstrates superior capabilities compared to other instruction tuned LLMs (e.g., FLAN (Chung et al., 2022)). Ahuja et al. (2023); Bang et al. (2023); Laskar et al. (2023); Lai et al. (2023); Huang et al. (2023) evaluate LLMs on more diverse languages using existing multilingual benchmarks (e.g., XNLI, PAWS-X, XLSum) involving monolingual NLP tasks and crosslingual tasks (e.g., machine translation). Their findings point to a large gap in performance of instruction tuned LLMs and ChatGPT, especially on low-resource languages and those with non-Latin scripts.

SM is still not adequately represented in existing benchmarks, hindering comprehensive evaluations on more languages. As we summarize in Table 1, benchmarks used for listed evaluations only include a few SM tasks focusing on sentiment analysis. Wang et al. (2023); Zhang et al. (2023b) investigate LLMs on a number of SM tasks (e.g., offensive language detection), but only on English. Ziems et al. (2023) investigate ChatGPT performance on a range of computational social science tasks covering subjects such as sociology, psychology, and linguistics, but they again focus only on English. Das et al. (2023) extend evaluation of ChatGPT on hate speech detection to 11 languages. Compared to these works, our objective is to investigate more diverse SM tasks on a massively multilingual setting.

Sociopragmatic Meaning Benchmarks. Many previous works introduce unified benchmarks such as GLUE (Wang et al., 2019), SentEval (Conneau and Kiela, 2018), XTREME (Hu et al., 2020), and XGLUE (Liang et al., 2020). These benchmarks include a wide range of NLP tasks, but comprise a sole SM task (i.e., sentiment analysis). Some recent studies started to construct benchmarks focusing on SM: Barbieri et al. (2020) introduce TweetEval benchmark that contains seven English datasets of six SM tasks; Zhang et al. (2023b) develop SentiEval that involves 26 English datasets of 13 sentiment-related tasks. Beyond English, NusaX (Winata et al., 2022), NaijaSenti (Muhammad et al., 2022), and AfriSenti (Muhammad et al., 2023a) propose benchmarks for sentiment analysis with eight Indonesian languages, four African languages, and 14 African languages, respectively. UMSAB introduced by Barbieri et al. (2022) contains 11 sentiment analysis datasets in 11 languages. For detecting antisocial online comments, Risch et al. (2021) introduces a toxic comment collection that contains 43 datasets of six antisocial detection tasks in 14 languages. Compared to these works, our SPARROW benchmark includes significantly more SM tasks and languages, from more diverse sources (refer to Figure 1 for a comparison).

3 SPARROW Benchmark

In this section, we describe clusters of tasks in our benchmark as well as our preprocessing and unification. SPARROW consists of 13 types of tasks in six main categories. It contains 169 datasets

	Tasks	Dataset	Lang.	LF	Ser
Antisocial	Aggressive	1	1	1	1
	Dangerous	1	1	1	1
	Hate	16	11	6	5
	Offense	7	6	3	3
	H/O-Group	3	3	2	3
	H/O-Target	8	8	4	7
	Antisocial	36	20	7	10
	Emotion	26	17	7	5
	Humor	4	4	1	2
Irn & Sarc	Irony	9	7	3	3
	Sarcasm	10	4	3	3
	Irony-Type	1	1	1	1
	Irony&Sarcasm	20	8	3	3
	Sentiment	77	58	10	15
	Subjectivity	6	5	2	2
	SPARROW	169	64	12	16

Table 2: Summary of datasets in SPARROW. **Lang**: number of languages, **LF**: number of language families, **Ser**: number of scripts.

from diverse online platforms and covers a wide period of time (1986-2022). We group different tasks in our benchmark by what we perceive to be an affinity between these tasks. For example, we group tasks of hate speech, offensive language, and dangerous language detection as anti-social language detection. Meanwhile, we keep particular tasks (such as sentiment analysis and emotion recognition) distinct due to the popularity of these tasks and since there are multiple datasets representing each of them. Table 2 summarizes statistics of SPARROW. We now briefly introduce our task clusters. We provide more information about languages in SPARROW in Table 7 of the Appendix. We also provide detailed descriptions with full citations of all our datasets in Tables 9, 10, 11, 12, 13, and 14 in Appendix.

3.1 Task Clusters

Antisocial Language Detection. The proliferation of antisocial language (e.g., hate speech) toxifies public discourse, incites violence, and undermines civil society (Sap et al., 2019; Vidgen and Derczynski, 2020). Antisocial language detection is thus a useful task. We include under the umbrella of antisocial language the following: (1) aggressive language (Kumar et al., 2018), (2) dangerous language (Alshehri et al., 2020), (3) hate speech (e.g., Waseem and Hovy (2016); Deng et al. (2022)), (4) offensive language (e.g., Mubarak et al. (2020); Kralj Novak et al. (2021)), (5) offense type identification (e.g., Zampieri et al. (2019)), and (6) offense target identification (e.g., Ousidhoum et al. (2019); Jeong et al. (2022)).

Emotion Recognition. Emotion affects our decision-making as well as mental and physical health (Abdul-Mageed and Ungar, 2017). SPARROW includes 26 emotion datasets in 17 languages (e.g., Kajava (2018); Bianchi et al. (2021)).

Humor Detection. Humor is a type of figurative language which induces amusing effects, such as laughter or well-being sensations. We include four humor detection datasets in four languages (e.g., Blinov et al. (2019); Meaney et al. (2021)).

Irony & Sarcasm Detection. Irony and sarcasm also involve figurative language. An ironic/sarcastic expression intentionally uses diametric language to signify implied meaning. We include (1) nine irony detection datasets in seven languages (e.g., Xiang et al. (2020)), (2) ten sarcasm detection datasets in four languages (e.g., Walker et al. (2012)), and (3) an irony type identification dataset (Van Hee et al., 2018).

Subjectivity and Sentiment Analysis. Subjectivity analysis aims to understand the opinions, feelings, judgments, and speculations expressed via language (Abdul-Mageed et al., 2014). Our benchmark includes six subjectivity analysis datasets in five different languages (e.g., Pang and Lee (2004); Pribán and Steinberger (2022)). Subjectivity incorporates sentiment. Sentiment analysis (Poria et al., 2020) is one of the most popular tasks in SM understanding where the objective is to identify the underlying sentiment of a given text. Our benchmark contains 77 sentiment analysis datasets in 58 languages (e.g., Pang and Lee (2005); Marreddy et al. (2022)).

3.2 Preprocessing, Splits, and Metrics

We apply light normalization on all the samples by converting user mentions and hyperlinks to ‘USER’ and ‘URL’, respectively. We standardize label names for consistency across datasets without reassigning nor aggregating the original labels of the datasets. For instance, in certain sentiment analysis datasets, we map ‘0’ and ‘1’ to ‘Negative’ and ‘Positive’ respectively. Regarding data splits, if the dataset already has Train, Dev, and Test sets, we maintain the same splits. If the original dataset does not include a Dev set, we randomly select 10% of training data to be a Dev set. In cases without pre-defined splits, we use an 80% Train, 10% Dev, and 10% Test random split. For computing constraints, we also prepare a smaller Test set for

each dataset by randomly sampling 500 samples from Test (if its size exceeds 500). We refer to this smaller test set as Test-S.

We evaluate on each dataset using its original metric as Tables 9, 10, 11, 12, 13, and 14 in Appendix summarize.¹ We report the performance on individual datasets, aggregate datasets into 13 tasks, and report an average score over each task. Moreover, we introduce a metric for each main category, calculated as the average of dataset-specific metrics within that category. Inspired by previous evaluation benchmarks like GLUE (Wang et al., 2019), we define a global metric called *SPARROW score*, which represents the unweighted average of all dataset-specific metrics. The SPARROW score provides an overall indication of performance on SM tasks.

4 Evaluation Methods

4.1 Finetuning on Encoder-Only Models

We evaluate the following Transformer-encoder-based multilingual models on SPARROW: (1) **Multilingual-BERT** (mBERT) (Devlin et al., 2019) (Base, 110M parameters), (2) **XLM-RoBERTa_{Base}** (XLM-R) (Conneau et al., 2020) (270M parameters), (3) **Bernice** (DeLucia et al., 2022), a 270M-parameter model trained with 2.5B tweets in 66 languages, and (4) **InfoDCL** (Zhang et al., 2023a), a SoTA for SM understanding, which further trains XLM-R with 100M tweets in 66 languages with contrastive learning. More details about all models are in Appendix B.

4.2 Zero- and Few-Shot on LLMs

We investigate zero-shot performance on a wide range of generative models, including *pre-trained generative models*: (1) **BLOOM** (Scao et al., 2022), (2) **mT5** (Xue et al., 2021), (3) **LLaMA** (Touvron et al., 2023), *instruction tuned models*: (4) **BLOOMZ** (Scao et al., 2022), a BLOOM-initialized model tuned with multilingual xP3 corpus, (5) **BLOOMZ-P3** (Muenighoff et al., 2022), a BLOOM-initialized model tuned with English-only P3 corpus, (6) **BLOOM-Bactrian** (Li et al., 2023), a BLOOM-initialized model tuned with 3.4M instruction-following samples in 52 languages, (7) **mT0** (Muenighoff et al., 2022), an mT5 model tuned with xP3 corpus, (8) **Alpaca** (Taori et al., 2023), a

¹We select the macro-average F_1 score as the main metric if the original paper utilizes more than one metric.

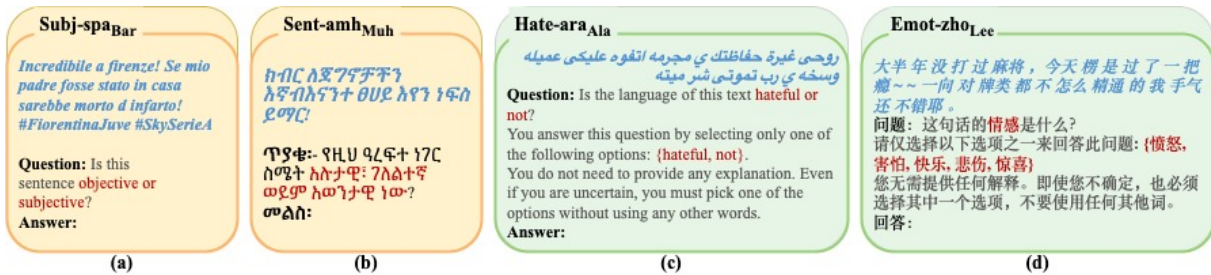


Figure 2: Examples of prompts used for zero-shot evaluation with `lm-evaluation-harness` (yellow) and ChatGPT (green). We use an English prompt (Figures a, c) and machine translated the prompt in the corresponding language (Figures b, d), respectively. The prompts construct each task as question-and-answer tasks. The actual input sample is in blue, and the label options are in red.

LLaMA-initialized model tuned with 52K English instruction-following samples, (9) **Vicuna** (Chiang et al., 2023), a LLaMA-initialized model on 70K conversational data, and (10) **ChatGPT**, for which we use the `gpt-3.5-turbo-0301` version via OpenAI API.² We use 7B-size version of BLOOM- and LLaMA-based models and 4B-size version of mT5-based models. We also evaluate six open-source LLMs (i.e., BLOOM, BLOOMZ-P3, mT5, mT0, LLaMA, and Vicuna) via few-shot in-context learning.

5 Experiments

5.1 Implementation

Finetuning. To keep computation reasonable, we randomly select 45 datasets for hyperparameter tuning and only tune the learning rate of each model.³ For all experiments, we finetune a pretrained model with an arbitrary batch size of 32 and sequence length of 128 tokens. Each model is finetuned on the full Train set of each dataset for 20 epochs (with patience = 5) based on performance on Dev set. We run each experiment *three* times with different random seeds and identify the best model on Dev in each run. We report the average performance on Test-S over three runs.⁴

Zero-shot Evaluation. We perform a zero-shot evaluation on SPARROW for BLOOM-, mT5-, and LLaMA-based models using language model evaluation harness (`lm-evaluation-harness` Gao et al. (2021)).⁵ While we do not tailor prompts

specifically for each model, customized prompts are employed for each set of tasks. These prompts follow the structure of question-and-answer tasks, where we present sample content alongside a task-specific question, as shown in Figure 2. The prompts are summarized in Appendix Table 15. We then instruct the model to generate an answer based on the provided option labels. Each option label represents a potential answer, and we calculate the log-likelihood for each candidate. The prediction with the highest log-likelihood is chosen as the model’s final prediction. For the evaluation of ChatGPT, we draw inspiration from previous practices for prompt design (Ziems et al., 2023), and incorporate additional instructions to guide its generation of the desired labels. As shown in Figure 2, we provide an instruction that compels ChatGPT to select a single label for the given input text without providing any additional explanation. We set the temperature to 0 to generate *deterministic and reproducible results* from ChatGPT. For a few instances, we observe that ChatGPT is unable to provide a direct answer. In these cases, we randomly assign a false label to each sample. In addition, we also use machine translation to translate English prompts and label names to the corresponding language of each dataset.⁶

Few-shot Evaluation. We utilize `lm-evaluation-harness` tool with the same prompts employed in zero-shot evaluation to explore the few-shot in-context learning abilities of open-source LLMs. Before the actual test examples, we prepend m examples from the Train set. Each example consists

²In rest of this paper, ChatGPT refers to `gpt-3.5-turbo-0301`.

³For more information, refer to Section C.1 in Appendix.

⁴We also report the performance of Dev and standard deviations in Appendix Table 16.

⁵<https://github.com/EleutherAI/lm-evaluation-harness>

⁶We use Google Translate for most languages. NLLB model is used to translate the languages of `ace`, `ban`, `bjn`, `bug`, and `min` because Google Translate does not cover these. The prompts of `pcm` datasets are translated by a native speaker.

of an input text, task-specific instruction, and the corresponding answer. We set m to either 3 or 5.

5.2 Results

We present the aggregated performance of Test-S on each task and main category, respectively, in Table 3. We also present test results on all datasets and compare to dataset-specific SoTA performance in Tables 17, 18, 19, 20, 21, and 22 in Appendix.

(1) How is the overall performance over different models? *All the fully finetuned models surpass the zero-shot generative models as well as ChatGPT, as shown in Table 3.* The most superior among the finetuned models is InfoDCL, which achieves a SPARROW score of 71.60 and outperforms ChatGPT with 11.56 points SPARROW score. On the other hand, the open-source models (i.e., BLOOM, mT5 and LLaMA) still significantly lag behind on multilingual SM understanding with performance close to a random baseline. Meanwhile, the instruction tuned multilingual LLMs (BLOOMZ and mT0) only slightly perform better than the random baseline.

(2) Can instruction tuning enhance LLMs' ability on SM understanding? *Yes, but it depends on the instruction training data.* Following instruction tuning on the English-only P3 dataset, BLOOMZ-P3 demonstrates an improvement of 7.76 SPARROW score compared to BLOOM. Also, BLOOMZ improves 5.85 points over BLOOM (but falls short of BLOOMZ-P3). mT0 also outperforms mT5. However, there remains a substantial gap between all instruction tuned models and finetuned models. BLOOM-Bactrian performs worse than BLOOMZ and BLOOMZ-P3, which are instruction tuned with NLP tasks. This indicates that the general purpose instruction-response dataset is not very useful for SM understanding.

To further probe how instruction tuning improves BLOOM-based models, we compare BLOOM with BLOOMZ-P3 and BLOOMZ in terms of individual tasks, finding sentiment analysis to exhibit the most significant improvement. BLOOMZ-P3 and BLOOMZ achieve a sentiment score improvement of 16.37 and 12.36, respectively, based on average calculation across 77 sentiment analysis datasets. However, BLOOM-Bactrian obtains an improvement of only 1.79 sentiment score, perhaps implying that the Bactrian

instruction-response data is not all that useful for some SM tasks. After tuning mT5 on xP3 dataset, mT0 also experiences a 13.88 improvement in the sentiment score. These may be stemming from inclusion of five English sentiment analysis datasets in both P3 and xP3 during the training phase. For example, we observe that BLOOM, BLOOMZ, BLOOMZ-P3, mT5, and mT0 obtain an accuracy of 56.4, 92.2, 93.00, 49.00, and 76.8 on SentengSoc (not included in either xP3 or P3), respectively and that BLOOM-Bactrian still performs poorly (accuracy= 53.60) after instruction tuning. Again, these results indicate that it is still important to include task-related datasets in the instruction tuning stage.

(3) How do LLMs perform across different SM tasks? *They are inferior at humor and antisocial language detection while being relatively better at sentiment and emotion recognition tasks.* BLOOMZ-P3, BLOOMZ, and mT0 exhibit considerable enhancements (> 5 points) on sentiment and emotion when compared to their respective initialization models. On the other hand, we find that instruction tuned models perform significantly worse on aggressive language detection and humor detection tasks. BLOOMZ-P3, BLOOMZ, BLOOM-Bactrian, and mT0 all incur a loss of more than 5 points on these two tasks. Upon investigating the predictions of instruction tuned models, we find that they tend to assign negative labels (i.e., non-aggressive or non-humor) which results in many false negative predictions. For a concrete example, we show that BLOOMZ-P3 predict most samples as non-humor in Figure 3a shows.

ChatGPT outperforms the open-source LLMs on all tasks except dangerous language detection. Comparing ChatGPT to InfoDCL, we find gaps favoring InfoDCL in subjectivity analysis (a difference of 9.47), emotion recognition (a difference of 10.68), and irony & sarcasm detection (a difference of 10.70). ChatGPT also largely lags behind InfoDCL in humor detection (a difference of 15.40) and antisocial language detection (a difference of 14.06). As the example shows in Figure 3b, ChatGPT makes more false positive errors (classifies non-hateful as hateful).

(4) How do LLMs perform across different languages? We now examine the impact of instruction finetuning on the model's language-wise performance. We categorize the performance of each

Tasks	Rand.		Finetuning				Zero-shot													
	—	mB.	X-R	Ber.	InfoD	BM	BMZ	BMZ (MT)	BMZ P3	BM Bac.	mT5	mT0	mT0 (MT)	LLa.	Alp.	Vic.	CG	CG (MT)		
	110M	270M	270M	270M	270M	7B	7B	7B	7B	7B	4B	4B	4B	7B	7B	7B	175B	175B		
Antisocial	Aggressive	43.14	72.71	74.64	75.45	73.96	51.06	15.82	15.82	18.72	16.37	53.67	15.82	22.00	18.31	49.29	25.07	63.53	54.36	
	Dangerous	42.06	62.36	63.57	67.13	65.23	46.87	46.87	50.84	46.87	46.87	49.31	46.87	46.87	46.87	46.87	46.87	37.93	33.68	58.74
	Hate	43.62	72.97	74.37	76.76	75.85	39.83	39.44	37.76	38.52	42.23	23.29	37.33	39.05	37.80	44.31	41.59	66.06	58.74	
	Offense	39.48	77.53	75.88	78.45	78.88	41.06	40.42	20.28	38.59	40.43	24.99	39.90	21.11	39.85	16.82	48.70	67.31	52.70	
	H/O-Group	14.82	46.18	42.39	51.15	50.24	13.63	17.26	14.23	21.23	14.81	7.02	16.25	17.01	12.35	14.13	9.26	39.66	26.74	
	H/O-Target	20.39	53.16	57.67	60.96	60.79	18.73	19.03	18.74	16.89	18.77	6.69	20.58	17.99	19.32	16.83	17.01	35.89	28.67	
	AS	35.20	66.92	67.99	71.14	70.61	33.70	32.80	27.93	31.97	33.79	20.14	32.02	28.79	31.68	30.55	34.50	56.55	47.40	
Emotion	15.86	61.42	66.87	68.13	69.27	9.71	17.18	13.85	15.07	15.19	7.75	27.87	24.21	15.14	31.80	18.12	59.58	50.85		
Humor	49.65	84.35	85.19	86.75	87.05	41.78	33.12	33.82	33.17	33.04	35.91	43.60	33.12	39.78	41.72	46.19	71.65	72.70		
I&S	Irony	42.39	64.24	65.53	66.88	68.38	36.63	35.15	38.69	44.46	36.18	36.52	34.69	33.99	40.78	27.49	47.48	58.23	56.24	
	Sarcasm	45.48	72.41	73.40	74.78	74.94	43.00	41.62	32.23	32.22	41.68	46.34	36.09	41.62	41.17	32.48	47.67	65.34	65.55	
	Irony-Type	22.36	47.35	46.43	56.04	57.58	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	30.81	30.81	
I&S	42.93	67.48	68.51	70.29	71.12	38.92	37.57	34.46	41.79	40.39	35.42	37.36	32.35	39.87	29.56	46.14	60.41	59.63		
Sentiment	34.68	66.34	69.58	70.44	71.64	26.67	39.03	28.61	43.03	28.46	20.77	34.65	32.76	27.55	25.84	25.02	60.34	54.94		
Subjectivity	41.41	72.54	74.45	74.80	75.73	44.12	29.45	30.69	30.73	39.65	37.35	41.64	36.16	42.30	30.44	38.73	66.26	59.33		
SPARROW	33.47	66.60	69.38	70.85	71.60	27.94	33.79	27.17	35.70	29.45	21.45	33.63	30.85	28.75	28.79	29.36	60.04	53.90		

Table 3: SPARROW benchmark Test-S results. We report the average of dataset-specific metrics in a task and a category, respectively. **Rand.:** random baseline, **mB.:** mBERT, **X-R:** XLM-R, **Ber.:** Bernice, **InfoD:** InfoDCL, **BM:** BLOOM, **LLa.:** LLaMA, **Alp.:** Alpaca, **Vic.:** Vicuna, **CG:** ChatGPT, **MT:** using machine translated prompts. The best performance in each setting is **Bold**. The **red font** denotes a performance lower than the random baseline.

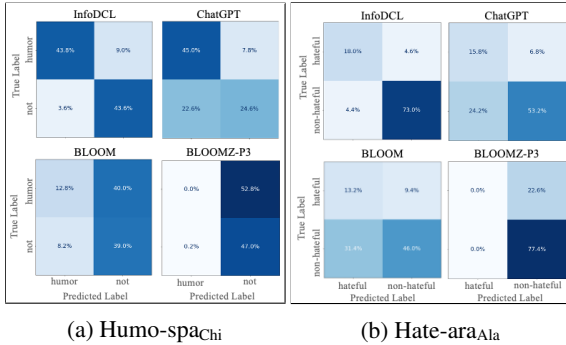


Figure 3: Confusion matrices of two datasets.

dataset based on language and calculate the average language scores across all datasets within a language. Since each language contains different tasks and datasets, a direct comparison across languages is not feasible. Therefore, we compare the relative performance between different models for each language. By comparing the instruction tuned models to their initial models, we observe that *most languages experience improvement*. However, we also observe a significant decline in performance for the Amharic (amh) dataset among these models. Specifically, BLOOMZ-P3, BLOOMZ, and mT0 experience a deterioration of 36.07, 24.99, and 26.12 points, respectively, compared to their respective initial models. We hypothesize that this deterioration can be attributed to catastrophic forgetting after instruction tuning, where Amharic was not included in the training set and does not share the writing scripts with the other included languages.

Lang	Random	InfoDCL	BMZ-P3	mT0	Vicuna	CG	CG-MT
amh	37.95	65.68	16.05	22.49	2.99	20.62	46.82
bug	30.77	71.55	34.60	18.27	12.90	34.63	30.86
ell	41.24	79.13	46.71	45.47	48.21	60.94	34.98
eng	37.90	75.48	43.32	39.23	39.75	66.51	—
fil	52.37	79.01	34.47	34.47	34.47	69.13	66.67
heb	47.60	95.80	71.20	76.60	40.80	84.20	57.40
hin	35.24	67.55	28.92	26.20	29.06	52.63	48.30
mal	31.68	82.70	43.84	41.65	24.85	44.03	31.44

Table 4: Language-wise model performance for sample languages. The complete results are in Table 23 in Appendix. Best performance in each language is **bold**, and the second best is in **green highlight**. The **red font** denotes a performance lower than the random baseline.

Similarly, the Filipino (fil) tasks exhibit an average decline of approximately 11 points on both BLOOMZ-P3 and BLOOMZ, as Filipino is not included in the xP3 dataset. Although Hindi is included in the xP3 dataset, the three instruction tuned models still show a decline in performance. Upon examining the individual performance of Hindi datasets, we find that the major deteriorations occur in the aggressive language detection and humor detection tasks, while the emotion recognition and sentiment analysis tasks show improvement. The instruction-response data for training Alpaca and Vicuna consist solely of English language. Therefore, we compare the performance of Alpaca and Vicuna to that of LLaMA using both English and non-English datasets. We observe that Alpaca and Vicuna outperform LLaMA when evaluated on English datasets, achieving

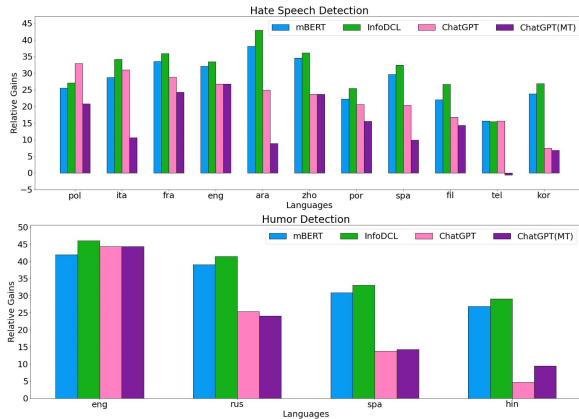


Figure 4: A comparison of different models on multiple tasks across various languages. We show the relative gain of each model compared to the random baseline.

average scores of 8.30 and 5.51, respectively. However, their performance declines when tested on non-English datasets, resulting in average decreases of 1.53 and 0.33, respectively. Compared to task-specific InfoDCL, ChatGPT performs poorly in 63 out of 64 languages, sometimes with a large gap (e.g., 45.06 lower on Amharic, 38.67 lower on Malayalam, and 36.91 lower on Buginese), as Table 4 shows.

We also investigate how different models perform on SM tasks across various languages. Results for two tasks, hate speech detection (top) and humor detection (bottom), are presented in Figure 4. The dataset for each task is grouped according to language, and the average score of each language is obtained. The relative gain of each model against the random baseline is shown, allowing us to compare across these languages.⁷ We observe that InfoDCL is the best model across various tasks and languages, with the exception of hate speech in Polish where ChatGPT outperforms it. As Figure 4 shows, ChatGPT performs better for Western languages on hate speech detection. We can also observe wider gaps in hate speech detection between ChatGPT and InfoDCL on Arabic and Korean. Similarly, while ChatGPT demonstrates satisfactory performance in English humor, it remains at significant distance behind InfoDCL in Hindi humor.

(5) Do machine translated prompts help LLMs? *Not in general, but they do help in a few cases.* We find, in Table 3, that the SPARROW score of ChatGPT with machine translated

⁷We note that different annotation artifacts across the different languages still make direct comparisons challenging.

Task	#Sample	CG (sample MT)	GPT-4
Hate-engWas	96	—	28.82
Hate-engDav	168	—	76.30
Hate-araAla	153	38.82	35.76
Hate-itaBos	104	29.75	38.44
Hate-filCab	145	15.76	23.96
Hate-aramMul	127	29.49	36.33
Hate-engBas	178	—	19.23
Hate-spaBas	174	21.06	17.35
Hate-porFor	142	35.37	37.34
Hate-polPta	54	35.62	49.57
Hate-korMoo	250	14.20	15.85
Hate-aramMub	89	20.54	16.74
Hate-zhoDen	125	29.24	39.24
Hate-korJeo	169	27.75	44.84
Hate-telMar	2	100.00	100.00
Sexi-fraChi	113	34.10	43.69
Humo-hinAgg	220	22.73	30.91
Humo-rusBii	136	26.89	47.31
Humo-spaChi	152	30.68	36.80
Humo-engMea	48	—	50.34

Table 5: Case study on using machine translated input and GPT-4 on samples mispredicted by ChatGPT.

prompts is 6.14 points lower than ChatGPT with English prompts. Meanwhile, a few tasks such as humor and sarcasm acquire improvements. We also observe a similar pattern for BLOOMZ and mT0, as Table 3 shows. The low-resource languages with non-Latin scripts experience more performance drops in general, which is in line with findings by Lai et al. (2023). Hebrew (heb) and Greek (ell) get the largest performance drops (over 25 points in each case), as shown in Table 4.

(6) Does GPT-4 outperform ChatGPT? *Yes, it does.* We provide a study on probing GPT-4’s capacities. We exploit 20 datasets from two tasks (i.e., hate speech and humor detection) in 12 languages, only choosing samples whose labels ChatGPT predicted incorrectly. We refer to this test set as GPTHard and provide samples from it to GPT-4 in their original language, employing the same English prompts as those used by ChatGPT. As Table 5 shows, GPT-4 significantly outperforms ChatGPT (McNemar’s test with $\alpha < 0.01$) on 19 datasets.⁸

(7) Can translating input samples into English help improve ChatGPT’s predictions? *Yes, it can.* Here, we use the non-English part of GPTHard (16 datasets). We translate these test samples into English using ChatGPT and subsequently employ the translated text and English prompt for classification. As Table 5 shows, we acquire a noteworthy enhancement in ChatGPT’s performance (McNemar’s test with $\alpha < 0.01$)

⁸An exception is one dataset where a significance test is not possible due to small sample size ($n = 2$).

Tasks	Zero-shot						Three-shot						Five-shot						
	BM	BMZ P3	mT5	mT0	LLa.	Vic.	BM	BMZ P3	mT5	mT0	LLa.	Vic.	BM	BMZ P3	mT5	mT0	LLa.	Vic.	
Antisocial	Aggressive	51.06	18.72	53.67	15.82	18.31	25.07	46.43	43.93	40.73	41.88	43.70	44.53	47.92	47.63	33.80	37.52	50.66	55.27
	Dangerous	46.87	46.87	49.31	46.87	46.87	46.87	46.87	46.87	45.68	46.87	46.87	46.87	46.87	46.87	48.91	46.87	46.87	46.87
	Hate	39.83	38.52	23.29	37.33	37.80	41.59	38.83	38.30	39.43	37.82	43.51	49.17	37.95	37.14	39.53	37.70	41.87	48.37
	Offense	41.06	38.59	24.99	39.90	39.85	48.70	43.94	40.25	21.99	41.53	46.36	54.49	42.42	40.59	34.10	41.67	43.83	51.72
	H/O-Group	13.63	21.23	7.02	16.25	12.35	9.26	11.43	13.04	7.92	15.98	11.81	14.19	9.68	11.76	7.23	14.50	12.58	16.27
	H/O-Target	18.73	16.89	6.69	20.58	19.32	17.01	17.09	18.41	10.48	16.55	20.56	24.84	17.44	17.45	9.32	16.56	20.09	23.60
	AS	33.70	31.97	20.14	32.02	31.68	34.50	33.14	32.55	27.19	32.36	36.42	41.69	32.43	31.88	29.17	32.09	35.35	40.99
	Emotion	9.71	15.07	7.75	27.87	15.14	18.12	17.08	12.17	10.66	23.12	32.12	40.28	18.48	12.35	10.07	25.57	34.20	41.79
	Humor	41.78	33.04	43.60	33.12	39.78	46.19	33.67	33.12	44.70	38.19	55.20	57.15	34.06	33.12	40.20	37.08	53.86	58.75
I&S	Irony	36.63	44.46	36.52	34.69	40.78	47.48	42.21	41.58	42.61	35.18	36.76	39.78	44.34	44.14	39.67	34.82	38.40	41.61
	Sarcasm	43.00	41.68	36.09	41.62	41.17	47.67	46.14	42.91	46.72	48.42	49.75	52.55	45.43	43.05	45.75	39.88	49.03	52.51
	Irony-Type	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83
	I&S	38.92	41.79	35.42	37.36	39.87	46.14	43.01	41.11	43.48	40.98	42.36	45.12	43.61	42.33	41.67	36.55	42.74	45.92
	Sentiment	26.67	43.03	20.77	34.65	27.55	25.02	34.81	35.79	24.76	31.37	37.73	34.53	33.15	37.71	23.17	29.25	40.88	39.37
Subjectivity	44.12	30.73	37.35	41.64	42.30	38.73	36.50	30.66	37.11	34.36	44.20	54.77	33.70	30.72	39.36	31.42	46.53	56.15	
SPARROW	27.94	35.70	21.45	33.63	28.75	29.36	32.76	31.91	26.12	31.71	37.82	39.44	32.03	32.84	25.47	30.44	39.48	41.97	

Table 6: Evaluating open-source LLMs on SPARROW with few-shot in-context learning. The best performance in each setting is **Bold**. The **red font** denotes a performance lower than the random baseline.

when using the translated input. We also observe that when fed with these English-translated samples, ChatGPT is able to surpass GPT4 with the original inputs in three datasets (i.e., Hate-ara_{Ala}, Hate-spa_{Bas}, Hate-ara_{Mub}). These results suggest that although ChatGPT has inferior ability on several languages in terms of detecting SM, a translate-then-detect approach may be possible.

(8) How do open-source LLMs perform with few-shot in-context learning? As Table 6 shows, we compare three-shot and five-shot results with zero-shot results. Based on SPARROW score, we observe that few-shot learning does enhance the performance of BLOOM, mT5, LLaMA, and Vicuna. With the increasing number of shots, the performance of LLaMA and Vicuna increases. Vicuna obtains SPARROW scores of 29.36, 39.44, and 41.97 with zero, three, and five shots, respectively. However, BLOOMZ-P3 and mT0 do not improve with few-shot learning. We suspect this is because the instruction finetuning of these two models only uses a zero-shot template that hurts their few-shot learning capacities. BLOOMZ-P3 and mT0 are also different from BLOOM and LLaMA in that they are finetuned on several NLP tasks only one of which is an SM task (i.e., sentiment analysis). This probably biases the behavior of these two models.

(9) Are the open-source LLMs sensitive to prompts used? We carry out a study to probe the open-source LLMs’ sensitivity to prompts. We curate 55 datasets across four tasks from SPARROW and evaluate six models with the prompts we used for evaluating ChatGPT. As Ta-

ble 24 in Appendix shows, we find that BLOOM, LLaMA, and Vicuna incur sizable performance drops (> 6 points decrease across 55 datasets), while BLOOMZ-P3, mT5, and mT0 demonstrate performance levels akin to those observed in previous experiments (< 2 points different). We leave a more comprehensive evaluation of prompt sensitivity as future work.

6 Public Leaderboard

To facilitate future work, we design a public leaderboard for scoring models on SPARROW. Our leaderboard is *interactive* and offers *rich metadata* about the various datasets in our benchmark. It also encourages users to submit information about their models (e.g., number of parameters, time to convergence, pretraining datasets). We also distribute a new *modular toolkit* for finetuning or evaluating models on SPARROW.

7 Conclusion

In order to understand the abilities of ChatGPT and other instruction tuned LLMs on capturing sociopragmatic meaning, we introduced a massively multilingual evaluation benchmark, dubbed SPARROW. The benchmark involves 169 datasets covering 64 languages from 12 language families and 16 scripts. Evaluating ChatGPT on SPARROW, we find it struggles with different languages. We also reveal that task-specific models finetuned on SM (much smaller than ChatGPT) consistently outperform larger models by a significant margin even on English.

8 Limitations

Benchmark Construction. Our SPARROW benchmark only includes text classification tasks related to SM. Despite our best efforts, we acknowledge that our benchmark has not covered existing SM datasets exhaustively. We will continue expanding this benchmark and welcome future datasets or metric contributions to it. We also plan to extend SPARROW to more types of tasks related to SM, such as span-based sentiment analysis (Xu et al., 2020b), affective language generation (Goswamy et al., 2020), and conversational sentiment analysis (Ojamaa et al., 2015). We only include text-based SM tasks. Another improvement direction is to extend this benchmark to more tasks that involve more modalities, such as affective image captioning (Mohamed et al., 2022) and multi-modal emotion recognition (Firdaus et al., 2020).

Model Selection. Due to computation constraints, we cannot evaluate on model sizes $> 7B$. However, we hope SPARROW will be used in the future to evaluate larger-sized models. Again, due to budget constraints, we only conduct a relatively small case study on GPT-4 and do not evaluate more diverse commercial instruction tuned models that are more expensive (e.g., `text-davinci-003` by OpenAI).

Experiments. While we customize prompts employed for each task, we do not tailor prompts specifically for each model. We acknowledge that the performance of models may be influenced by different prompt variants. In future work, we will test diverse prompt variations for more robust results. We only experiment with machine translated prompts in our analyses and acknowledge that the performance drop may stem from the poor quality of machine translation. We will investigate the utility of human translated prompts in a future study. In this paper, we only evaluate LLMs on zero-shot learning. The adoption of few-shot in-context learning may enhance performance, which we also leave to future work.

Ethics Statement and Broad Impacts

Data Collection and Releasing. All the 169 datasets are produced by previous research. Since there are large numbers of datasets and languages in SPARROW, it is hard to manually verify the

quality of all the datasets. As a quality assurance measure, we only include in SPARROW datasets that are introduced in peer-reviewed published research. To facilitate access to information about each dataset, we link to each published paper describing each of these datasets in Tables 9, 10, 11, 12, 13, and 14.

Following privacy protection policies, we anonymize all SPARROW data as described in Section 3.2. With reference to accessibility of the original individual dataset, SPARROW data can be categorized into three releasing strategies: (1) In the case of datasets requiring approval by the original authors, we require future researchers to obtain approval first and will share our splits once approval has been obtained. We indicate these nine datasets in our data description tables. (2) For the 25 datasets (see Table 8 in Appendix) that are shared via tweet IDs, we share our obtained data for research use. By doing so, we expect to mitigate the issue of data decay and allow fair comparisons. (3) We will share the other 135 publicly accessible datasets upon request. We will also require a justification for responsible use of the datasets. Each dataset will be shared in our Train, Dev, and Test splits along with a dataset card to indicate the original publication of the dataset.

Intended Use. The intended use of SPARROW benchmark is to construct a scoring board to facilitate model comparisons as well as enhance fairness and reproducibility across different languages and tasks. We also aim to mitigate data decay issues in social media research. SPARROW could help researchers investigate model’s capacity on SM tasks across languages. SPARROW may also be used to investigate model transferability across a wide range of tasks and diverse languages in different settings (such as zero- or few-shot settings and prompting).

Potential Misuse and Bias. We notice that some annotations in the datasets of SPARROW (e.g., for hate speech task (Waseem and Hovy, 2016)) can carry annotation and temporal biases. We recommend that any dataset in SPARROW not be used for research or in applications without careful consideration of internal biases of the datasets and potential biases of the resulting systems. We also suggest that users of SPARROW not only focus on the overall SPARROW score but also a model performance on each task and

dataset. The SPARROW score is an unweighted average score over all the dataset-specific metrics, which may lose the fine-grained information and be dominated by the largest task cluster (i.e., sentiment analysis) or languages (e.g., languages from Indo-European language family).

Acknowledgements

We acknowledge support from Canada Research Chairs (CRC), the Natural Sciences and Engineering Research Council of Canada (NSERC; RGPIN-2018-04267), the Social Sciences and Humanities Research Council of Canada (SSHRC; 435-2018-0576; 895-2020-1004; 895-2021-1008), Canadian Foundation for Innovation (CFI; 37771), Digital Research Alliance of Canada,⁹ and UBC ARC-Sockeye.¹⁰

References

- Muhammad Abdul-Mageed, Mona T. Diab, and Sandra Kübler. 2014. **SAMAR: subjectivity and sentiment analysis for arabic social media**. *Comput. Speech Lang.*, 28(1):20–37.
- Muhammad Abdul-Mageed, AbdelRahim A. Elmadany, and El Moatez Billah Nagoudi. 2021. **ARBERT & MARBERT: deep bidirectional transformers for arabic**. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 7088–7105. Association for Computational Linguistics.
- Muhammad Abdul-Mageed and Lyle Ungar. 2017. **EmoNet: Fine-grained emotion detection with gated recurrent neural networks**. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 718–728, Vancouver, Canada. Association for Computational Linguistics.
- Muhammad Abdul-Mageed, Chiyu Zhang, AbdelRahim Elmadany, Houda Bouamor, and Nizar Habash. 2022. **NADI 2022: The third nuanced Arabic dialect identification shared task**. In *Proceedings of the The Seventh Arabic Natural Language Processing Workshop (WANLP)*, pages 85–97, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Muhammad Abdul-Mageed, Chiyu Zhang, Azadeh Hashemi, and El Moatez Billah Nagoudi. 2020. **AraNet: A deep learning toolkit for Arabic social media**. In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 16–23, Marseille, France. European Language Resource Association.
- Ibrahim Abu Farha and Walid Magdy. 2020. **From Arabic sentiment analysis to sarcasm detection: The Ar-Sarcasm dataset**. In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 32–39, Marseille, France. European Language Resource Association.
- Akshita Aggarwal, Anshul Wadhawan, Anshima Chaudhary, and Kavita Maurya. 2020. **“did you really mean what you said?” : Sarcasm detection in Hindi-English code-mixed data using bilingual word embeddings**. In *Proceedings of the Sixth Workshop on Noisy User-generated Text (W-NUT 2020)*, pages 7–15, Online. Association for Computational Linguistics.
- Kabir Ahuja, Rishav Hada, Millicent Ochieng, Prachi Jain, Harshita Diddee, Samuel Maina, Tanuja Ganu, Sameer Segal, Maxamed Axmed, Kalika Bali, and Sunayana Sitaram. 2023. **MEGA: multilingual evaluation of generative AI**. *CoRR*, abs/2303.12528.
- Azalden Alakrot, Liam Murray, and Nikola S. Nikolov. 2018. **Dataset construction for the detection of anti-social behaviour in online communication in arabic**. In *Fourth International Conference On Arabic Computational Linguistics, ACLING 2018, November 17-19, 2018, Dubai, United Arab Emirates*, volume 142 of *Procedia Computer Science*, pages 174–181. Elsevier.
- Ali Alshehri, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2020. **Understanding and detecting dangerous speech in social media**. In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 40–47, Marseille, France. European Language Resource Association.
- Adam Amram, Anat Ben David, and Reut Tsarfaty. 2018. **Representations and architectures in neural sentiment analysis for morphologically rich languages: A case study from Modern Hebrew**. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2242–2252, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Seyed Arad Ashrafi Asli, Behnam Sabeti, Zahra Majdabadi, Prenti Golazizian, reza fahmi, and Omid Moemenzadeh. 2020. **Optimizing annotation effort using active learning strategies: A sentiment analysis case study in Persian**. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 2855–2861, Marseille, France. European Language Resources Association.
- Dennis Assenmacher, Derek Weber, Mike Preuss, André Calero Valdez, Alison Bradshaw, Björn Ross,

⁹<https://alliancecan.ca>

¹⁰<https://arc.ubc.ca/ubc-arc-sockeye>

- Stefano Cresci, Heike Trautmann, Frank Neumann, and Christian Grimme. 2022. [Benchmarking crisis in social media analytics: A solution for the data-sharing problem](#). *Social Science Computer Review*, 40(6):1496–1522.
- David Bamman and Noah A. Smith. 2015. [Contextualized sarcasm detection on twitter](#). In *Proceedings of the Ninth International Conference on Web and Social Media, ICWSM 2015, University of Oxford, Oxford, UK, May 26-29, 2015*, pages 574–577. AAAI Press.
- 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. 2023. [A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity](#). *CoRR*, abs/2302.04023.
- Francesco Barbieri, Valerio Basile, Danilo Croce, Malvina Nissim, Nicole Novielli, and Viviana Patti. 2016. [Overview of the evalita 2016 sentiment polarity classification task](#). In *Proceedings of Third Italian Conference on Computational Linguistics (CLiC-it 2016) & Fifth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2016), Napoli, Italy, December 5-7, 2016*, volume 1749 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. [TweetEval: Unified benchmark and comparative evaluation for tweet classification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1644–1650, Online. Association for Computational Linguistics.
- Francesco Barbieri, Luis Espinosa Anke, and Jose Camacho-Collados. 2022. [XLM-T: Multilingual language models in Twitter for sentiment analysis and beyond](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 258–266, Marseille, France. European Language Resources Association.
- Valerio Basile, Andrea Bolioli, Viviana Patti, Paolo Rosso, and Malvina Nissim. 2014. Overview of the evalita 2014 sentiment polarity classification task. *Overview of the Evalita 2014 SENTiment POLarity Classification Task*, pages 50–57.
- Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. [SemEval-2019 task 5: Multilingual detection of hate speech against immigrants and women in Twitter](#). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 54–63, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Federico Bianchi, Debora Nozza, and Dirk Hovy. 2021. [FEEL-IT: Emotion and sentiment classification for the Italian language](#). In *Proceedings of the Eleventh Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 76–83, Online. Association for Computational Linguistics.
- Federico Bianchi, Debora Nozza, and Dirk Hovy. 2022. [XLM-EMO: Multilingual emotion prediction in social media text](#). In *Proceedings of the 12th Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media Analysis*, pages 195–203, Dublin, Ireland. Association for Computational Linguistics.
- Vladislav Blinov, Valeria Bolotova-Baranova, and Pavel Braslavski. 2019. [Large dataset and language model fun-tuning for humor recognition](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 4027–4032. Association for Computational Linguistics.
- Cristina Bosco, Felice Dell’Orletta, Fabio Poletto, Manuela Sanguinetti, and Maurizio Tesconi. 2018. [Overview of the EVALITA 2018 hate speech detection task](#). In *Proceedings of the Sixth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2018) co-located with the Fifth Italian Conference on Computational Linguistics (CLiC-it 2018), Turin, Italy, December 12-13, 2018*, volume 2263 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Diana Boxer and Florencia Cortés-Conde. 2021. *Social Groups and Relational Networks*, Cambridge Handbooks in Language and Linguistics, page 227–246. Cambridge University Press.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Henrico Bertini Brum. 2018. *Expansão de recursos para análise de sentimentos usando aprendizado semi-supervisionado*. Ph.D. thesis, Universidade de São Paulo.
- Henrico Bertini Brum and Maria das Graças Volpe Nunes. 2018. [Building a sentiment corpus of tweets in Brazilian Portuguese](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation, LREC 2018, Miyazaki,*

- Japan, May 7-12, 2018. European Language Resources Association (ELRA).
- Neil Vicente Cabasag, Vicente Raphael Chan, Sean Christian Lim, Mark Edward Gonzales, and Charibeth Cheng. 2019. Hate speech in philippine election-related tweets: Automatic detection and classification using natural language processing. *Philippine Computing Journal*, XIV No, 1.
- Bharathi Raja Chakravarthi, Ruba Priyadarshini, Vigneshwaran Muralidaran, Navya Jose, Shardul Suryawanshi, Elizabeth Sherly, and John P. McCrae. 2022. Dravidiancodemix: sentiment analysis and offensive language identification dataset for dravidian languages in code-mixed text. *Lang. Resour. Evaluation*, 56(3):765–806.
- Qianben Chen, Richong Zhang, Yaowei Zheng, and Yongyi Mao. 2022. Dual contrastive learning: Text classification via label-aware data augmentation. *CoRR*, abs/2201.08702.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.
- Patricia Chiril, Véronique Moriceau, Farah Benamara, Alda Mari, Gloria Origi, and Marlène Coulomb-Gully. 2020. An annotated corpus for sexism detection in french tweets. In *Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020*, pages 1397–1403. European Language Resources Association.
- Luis Chiruzzo, Santiago Castro, Santiago Góngora, Aiala Rosá, J. A. Meaney, and Rada Mihalcea. 2021. Overview of HAHA at IberLEF 2021: Detecting, rating and analyzing humor in Spanish. *Proces. del Leng. Natural*, 67:257–268.
- Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4299–4307.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. *CoRR*, abs/2210.11416.
- Alessandra Teresa Cignarella, Simona Frenda, Valerio Basile, Cristina Bosco, Viviana Patti, and Paolo Rosso. 2018. Overview of the EVALITA 2018 task on irony detection in italian tweets (ironita). In *Proceedings of the Sixth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2018) co-located with the Fifth Italian Conference on Computational Linguistics (CLiC-it 2018), Turin, Italy, December 12-13, 2018*, volume 2263 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Alexandra Ciobotaru and Liviu P. Dinu. 2021. Red: A novel dataset for romanian emotion detection from tweets. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 296–305, Varna, Bulgaria. INCOMA Ltd.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Alexis Conneau and Douwe Kiela. 2018. SentEval: An evaluation toolkit for universal sentence representations. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Diogo Cortiz, Jefferson O. Silva, Newton Calegari, Ana Luísa Freitas, Ana Angélica Soares, Carolina Botelho, Gabriel Gaudencio Rêgo, Waldir Sampaio, and Paulo Sergio Boggio. 2021. A weak supervised dataset of fine-grained emotions in portuguese. *CoRR*, abs/2108.07638.
- Mithun Das, Saurabh Kumar Pandey, and Animesh Mukherjee. 2023. Evaluating chatgpt’s performance for multilingual and emoji-based hate speech detection. *CoRR*, abs/2305.13276.
- Thomas Davidson, Dana Warmusley, Michael W. Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the Eleventh International Conference on Web and Social Media, ICWSM 2017, Montréal, Québec, Canada, May 15-18, 2017*, pages 512–515. AAAI Press.
- Alexandra DeLucia, Shijie Wu, Aaron Mueller, Carlos Aguirre, Mark Dredze, and Philip Resnik. 2022. Bernice: A multilingual pre-trained encoder for Twitter. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December, 2022*, pages 6191–6205. Association for Computational Linguistics.

- Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. [GoEmotions: A dataset of fine-grained emotions](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4040–4054, Online. Association for Computational Linguistics.
- Jiawen Deng, Jingyan Zhou, Hao Sun, Fei Mi, and Minlie Huang. 2022. [COLD: A benchmark for chinese offensive language detection](#). *CoRR*, abs/2201.06025.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Mountaga Diallo, Chayma Fourati, and Hatem Haddad. 2021. [Bambara language dataset for sentiment analysis](#).
- Alexiei Dingli and Nicole Sant. 2016. [Sentiment analysis on maltese using machine learning](#). In *Proceedings of The Tenth International Conference on Advances in Semantic Processing (SEMAPRO 2016)*, pages 21–25.
- Stefan Daniel Dumitrescu, Andrei-Marius Avram, and Sampo Pyysalo. 2020. [The birth of romanian BERT](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 4324–4328. Association for Computational Linguistics.
- AbdelRahim A. Elmadany, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2022. [ORCA: A challenging benchmark for arabic language understanding](#). *CoRR*, abs/2212.10758.
- Ibrahim Abu Farha, Wajdi Zaghouni, and Walid Magdy. 2021. [Overview of the WANLP 2021 shared task on sarcasm and sentiment detection in arabic](#). In *Proceedings of the Sixth Arabic Natural Language Processing Workshop, WANLP 2021, Kyiv, Ukraine (Virtual), April 9, 2021*, pages 296–305. Association for Computational Linguistics.
- Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. 2017. [Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1615–1625, Copenhagen, Denmark. Association for Computational Linguistics.
- Mauajama Firdaus, Hardik Chauhan, Asif Ekbal, and Pushpak Bhattacharyya. 2020. [MEISD: A multi-modal multi-label emotion, intensity and sentiment dialogue dataset for emotion recognition and sentiment analysis in conversations](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4441–4453, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Paula Fortuna, João Rocha da Silva, Juan Soler-Company, Leo Wanner, and Sérgio Nunes. 2019. [A hierarchically-labeled Portuguese hate speech dataset](#). In *Proceedings of the Third Workshop on Abusive Language Online*, pages 94–104, Florence, Italy. Association for Computational Linguistics.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2021. [A framework for few-shot language model evaluation](#).
- Bilal Ghanem, Jihen Karoui, Farah Benamara, Véronique Moriceau, and Paolo Rosso. 2019. [IDAT at FIRE2019: overview of the track on irony detection in arabic tweets](#). In *FIRE '19: Forum for Information Retrieval Evaluation, Kolkata, India, December, 2019*, pages 10–13. ACM.
- Prenti Golazizian, Behnam Sabeti, Seyed Arad Ashrafi Asli, Zahra Majdabadi, Omid Momenzadeh, and Reza Fahmi. 2020. [Irony detection in Persian language: A transfer learning approach using emoji prediction](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 2839–2845, Marseille, France. European Language Resources Association.
- Xiaochang Gong, Qin Zhao, Jun Zhang, Ruibin Mao, and Ruifeng Xu. 2020. [The design and construction of a chinese sarcasm dataset](#). In *Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020*, pages 5034–5039. European Language Resources Association.
- Raymond G Gordon Jr. 2005. *Ethnologue, languages of the world*. <http://www.ethnologue.com/>.
- Tushar Goswamy, Ishika Singh, Ahsan Barkati, and Ashutosh Modi. 2020. [Adapting a language model for controlled affective text generation](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2787–2801, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Zekeriya Güven, Banu Diri, and Tolgahan Çakaloğlu. 2020. Comparison of n-stage latent dirichlet allocation versus other topic modeling methods for emotion analysis. *Journal of the Faculty of Engineering and Architecture of Gazi University*, 35(4).
- Vong Anh Ho, Duong Huynh-Cong Nguyen, Danh Hoang Nguyen, Linh Thi-Van Pham,

- Duc-Vu Nguyen, Kiet Van Nguyen, and Ngan Luu-Thuy Nguyen. 2019. [Emotion recognition for Vietnamese social media text](#). In *Computational Linguistics - 16th International Conference of the Pacific Association for Computational Linguistics, PACLING 2019, Hanoi, Vietnam, October 11-13, 2019, Revised Selected Papers*, volume 1215 of *Communications in Computer and Information Science*, pages 319–333. Springer.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. [XTREME: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 4411–4421. PMLR.
- Haoyang Huang, Tianyi Tang, Dongdong Zhang, Wayne Xin Zhao, Ting Song, Yan Xia, and Furu Wei. 2023. [Not all languages are created equal in llms: Improving multilingual capability by cross-lingual-thought prompting](#). *CoRR*, abs/2305.07004.
- Md. Asif Iqbal, Avishek Das, Omar Sharif, Mohammed Moshikul Hoque, and Iqbal H. Sarker. 2022. [Bemoc: A corpus for identifying emotion in bengali texts](#). *SN Comput. Sci.*, 3(2):135.
- Khondoker Ittehadul Islam, Sudipta Kar, Md Saiful Islam, and Mohammad Ruhul Amin. 2021. [SentNoB: A dataset for analysing sentiment on noisy Bangla texts](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3265–3271, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hayeon Jang, Munhyong Kim, and Hyopil Shin. 2013. [KOSAC: A full-fledged korean sentiment analysis corpus](#). In *Proceedings of the 27th Pacific Asia Conference on Language, Information and Computation, PACLIC 27, Taipei, Taiwan, November 21-24, 2013*. National Chengchi University, Taiwan.
- Younghoon Jeong, Juhyun Oh, Jaimeen Ahn, Jongwon Lee, Jihyung Moon, Sungjoon Park, and Alice Oh. 2022. [KOLD: korean offensive language dataset](#). *CoRR*, abs/2205.11315.
- Kaisla Kajava. 2018. [Cross-lingual sentiment preservation and transfer learning in binary and multi-class classification](#). Master’s thesis, University of Helsinki.
- Md. Rezaul Karim, Sumon Kanti Dey, Tanhim Islam, Sagor Sarker, Mehadi Hasan Menon, Kabir Hosain, Md. Azam Hossain, and Stefan Decker. 2021. [DeepHateExplainer: Explainable hate speech detection in under-resourced bengali language](#). In *8th IEEE International Conference on Data Science and Advanced Analytics, DSAA 2021, Porto, Portugal, October 6-9, 2021*, pages 1–10. IEEE.
- Pei Ke, Haozhe Ji, Siyang Liu, Xiaoyan Zhu, and Minlie Huang. 2020. [SentiLARE: Sentiment-aware language representation learning with linguistic knowledge](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6975–6988, Online. Association for Computational Linguistics.
- Jan Kocoń, Piotr Miłkowski, and Monika Zaśko-Zielińska. 2019. [Multi-level sentiment analysis of PolEmo 2.0: Extended corpus of multi-domain consumer reviews](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 980–991, Hong Kong, China. Association for Computational Linguistics.
- Petra Kralj Novak, Igor Mozetič, and Nikola Ljubešić. 2021. [Slovenian twitter hate speech dataset IMSyPP-sl](#). Slovenian language resource repository CLARIN.SI.
- Atharva Kulkarni, Meet Mandhane, Manali Likhitkar, Gayatri Kshirsagar, and Raviraj Joshi. 2021. [L3cubemahasant: A marathi tweet-based sentiment analysis dataset](#). *CoRR*, abs/2103.11408.
- Ritesh Kumar, Aishwarya N. Reganti, Akshit Bhatia, and Tushar Maheshwari. 2018. [Aggression-annotated corpus of hindi-english code-mixed data](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation, LREC 2018, Miyazaki, Japan, May 7-12, 2018*. European Language Resources Association (ELRA).
- Viet Dac Lai, Nghia Trung Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. 2023. [Chatgpt beyond english: Towards a comprehensive evaluation of large language models in multilingual learning](#). *CoRR*, abs/2304.05613.
- Md. Tahmid Rahman Laskar, M. Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Xiangji Huang. 2023. [A systematic study and comprehensive evaluation of chatgpt on benchmark datasets](#). *CoRR*, abs/2305.18486.
- Sophia Lee and Zhongqing Wang. 2015. [Emotion in code-switching texts: Corpus construction and analysis](#). In *Proceedings of the Eighth SIGHAN Workshop on Chinese Language Processing*, pages 91–99, Beijing, China. Association for Computational Linguistics.
- Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. 2023. [Bactrian-x : A multilingual replicable instruction-following model with low-rank adaptation](#). *CoRR*, abs/2305.15011.
- Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, Xiaodong Fan, Ruofei Zhang, Rahul Agrawal, Edward Cui, Sining Wei, Taroon Bharti, Ying Qiao, Jiun-Hung Chen, Winnie Wu, Shuguang Liu, Fan Yang, Daniel Campos, Rangan Majumder, and Ming Zhou. 2020.

- XGLUE: A new benchmark dataset for cross-lingual pre-training, understanding and generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 6008–6018. Association for Computational Linguistics.
- Mounika Marreddy, Subba Reddy Oota, Lakshmi Sireesha Vakada, Venkata Charan Chinni, and Radhika Mamidi. 2022. [Am i a resource-poor language? data sets, embeddings, models and analysis for four different nlp tasks in telugu language](#). *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 22(1).
- J. A. Meaney, Steven Wilson, Luis Chiruzzo, Adam Lopez, and Walid Magdy. 2021. [SemEval 2021 task 7: HaHackathon, detecting and rating humor and offense](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 105–119, Online. Association for Computational Linguistics.
- Youssef Mohamed, Faizan Farooq Khan, Kilichbek Haydarov, and Mohamed Elhoseiny. 2022. [It is okay to not be okay: Overcoming emotional bias in affective image captioning by contrastive data collection](#). In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 21231–21240. IEEE.
- Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. [SemEval-2018 task 1: Affect in tweets](#). In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 1–17, New Orleans, Louisiana. Association for Computational Linguistics.
- Jihyung Moon, Won Ik Cho, and Junbum Lee. 2020. [BEEP! Korean corpus of online news comments for toxic speech detection](#). In *Proceedings of the Eighth International Workshop on Natural Language Processing for Social Media*, pages 25–31, Online. Association for Computational Linguistics.
- Igor Mozetič, Miha Grčar, and Jasmina Smailović. 2016. Multilingual twitter sentiment classification: The role of human annotators. *PloS one*, 11(5):e0155036.
- Hamdy Mubarak, Kareem Darwish, Walid Magdy, Tamer Elsayed, and Hend Al-Khalifa. 2020. [Overview of OSACT4 Arabic offensive language detection shared task](#). In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 48–52, Marseille, France. European Language Resource Association.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M. Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2022. [Crosslingual generalization through multitask finetuning](#). *CoRR*, abs/2211.01786.
- Shamsuddeen Hassan Muhammad, Idris Abdulmumin, Abinew Ali Ayele, Nedjma Ousidhoum, David Ifeoluwa Adelani, Seid Muhie Yimam, Ibrahim Sa'id Ahmad, Meriem Beloucif, Saif M. Mohammad, Sebastian Ruder, Oumaima Hourrane, Pavel Brazdil, Felermimo Dário Mário António Ali, Davis David, Salomey Osei, Bello Shehu Bello, Falalu Ibrahim, Tajuddeen Gwadabe, Samuel Rutunda, Tadesse Destaw Belay, Wendimu Baye Mestelle, Hailu Beshada Balcha, Sisay Adugna Chala, Hagos Tesfahun Gebremichael, Bernard Opoku, and Steven Arthur. 2023a. [Afrisenti: A twitter sentiment analysis benchmark for african languages](#). *CoRR*, abs/2302.08956.
- Shamsuddeen Hassan Muhammad, Idris Abdulmumin, Seid Muhie Yimam, David Ifeoluwa Adelani, Ibrahim Sa'id Ahmad, Nedjma Ousidhoum, Abinew Ali Ayele, Saif M. Mohammad, Meriem Beloucif, and Sebastian Ruder. 2023b. [SemEval-2023 Task 12: Sentiment Analysis for African Languages \(AfriSenti-SemEval\)](#). In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*. Association for Computational Linguistics.
- Shamsuddeen Hassan Muhammad, David Adelani, Sebastian Ruder, Ibrahim Sa'id Ahmad, Idris Abdulmumin, Shehu Bello Bello, Monojit Choudhury, Chris Chinenye Emezue, Saheed Salahuddeen Abdullahi, Anuoluwapo Aremu, Alipio George, and Pavel Brazdil. 2022. [Naijasenti: A nigerian twitter sentiment corpus for multilingual sentiment analysis](#). In *Proceedings of the 13th Language Resources and Evaluation Conference*, pages 590–602, Marseille, France. European Language Resources Association.
- Hala Mulki, Hatem Haddad, Chedi Bechikh Ali, and Halima Alshabani. 2019. [L-HSAB: A Levantine Twitter dataset for hate speech and abusive language](#). In *Proceedings of the Third Workshop on Abusive Language Online*, pages 111–118, Florence, Italy. Association for Computational Linguistics.
- Sebastian Nordhoff and Harald Hammarström. 2011. [Glottolog/langdoc: Defining dialects, languages, and language families as collections of resources](#). In *Proceedings of the First International Workshop on Linked Science 2011, Bonn, Germany, October 24, 2011*, volume 783 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Emily Öhman, Marc Pàmies, Kaisla Kajava, and Jörg Tiedemann. 2020. [XED: A multilingual dataset for sentiment analysis and emotion detection](#). In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 6542–6552. International Committee on Computational Linguistics.

- Birgitta Ojamaa, Päivi Kristiina Jokinen, and Kadri Muischenk. 2015. [Sentiment analysis on conversational texts](#). In *Proceedings of the 20th Nordic Conference of Computational Linguistics, NODALIDA 2015, Institute of the Lithuanian Language, Vilnius, Lithuania, May 11-13, 2015*, volume 109 of *Linköping Electronic Conference Proceedings*, pages 233–237. Linköping University Electronic Press / Association for Computational Linguistics.
- Shereen Oraby, Vrindavan Harrison, Lena Reed, Ernesto Hernandez, Ellen Riloff, and Marilyn Walker. 2016. [Creating and characterizing a diverse corpus of sarcasm in dialogue](#). In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 31–41, Los Angeles. Association for Computational Linguistics.
- Reynier Ortega-Bueno, Francisco Rangel, D Hernández Farias, Paolo Rosso, Manuel Montes-y Gómez, and José E Medina Pagola. 2019. [Overview of the task on irony detection in spanish variants](#). In *Proceedings of the Iberian Languages Evaluation Forum co-located with 35th Conference of the Spanish Society for Natural Language Processing, IberLEF@SEPLN 2019, Bilbao, Spain, September 24th, 2019*, volume 2421 of *CEUR Workshop Proceedings*, pages 229–256. CEUR-WS.org.
- Nedjma Ousidhoum, Zizheng Lin, Hongming Zhang, Yangqiu Song, and Dit-Yan Yeung. 2019. [Multilingual and multi-aspect hate speech analysis](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 4674–4683. Association for Computational Linguistics.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *NeurIPS*.
- Wuraola Fisayo Oyewusi, Olubayo Adekanmbi, and Olalekan Akinsande. 2020. [Semantic enrichment of nigerian pidgin english for contextual sentiment classification](#).
- Bo Pang and Lillian Lee. 2004. [A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts](#). In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics, 21-26 July, 2004, Barcelona, Spain*, pages 271–278. ACL.
- Bo Pang and Lillian Lee. 2005. [Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales](#). In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)*, pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.
- Braja Gopal Patra, Dipankar Das, and Amitava Das. 2018. [Sentiment analysis of code-mixed indian languages: An overview of sail_code-mixed shared task @icon-2017](#). *CoRR*, abs/1803.06745.
- Flor Miriam Plaza del Arco, Carlo Strapparava, L. Alfonso Urena Lopez, and Maite Martin. 2020. [Emo-Event: A multilingual emotion corpus based on different events](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1492–1498, Marseille, France. European Language Resources Association.
- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, and Rada Mihalcea. 2020. [Beneath the tip of the iceberg: Current challenges and new directions in sentiment analysis research](#). *arXiv preprint arXiv:2005.00357*.
- Pavel Pribán and Josef Steinberger. 2022. [Czech dataset for cross-lingual subjectivity classification](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference, LREC 2022, Marseille, France, 20-25 June 2022*, pages 1381–1391. European Language Resources Association.
- Tomáš Ptáček, Ivan Habernal, and Jun Hong. 2014. [Sarcasm detection on Czech and English Twitter](#). In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 213–223, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.
- Michał Ptaszynski, Agata Pieciukiewicz, and Paweł Dybała. 2019. Results of the PolEval 2019 shared task 6: First dataset and open shared task for automatic cyberbullying detection in Polish twitter. *Proceedings of the PolEval 2019 Workshop*, page 89.
- Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. [Is chatgpt a general-purpose natural language processing task solver?](#) *CoRR*, abs/2302.06476.
- Ashwin Rajadesingan, Reza Zafarani, and Huan Liu. 2015. [Sarcasm detection on twitter: A behavioral modeling approach](#). In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, WSDM 2015, Shanghai, China, February 2-6, 2015*, pages 97–106. ACM.
- Luis Rei, Dunja Mladenic, and Simon Krek. 2016. [A multilingual social media linguistic corpus](#). In *Proceedings of the 4th Conference on CMC and Social Media Corpora for the Humanities, Ljubljana, Slovenia*.
- Ellen Riloff, Ashequl Qadir, Prafulla Surve, Lalindra De Silva, Nathan Gilbert, and Ruihong Huang. 2013. [Sarcasm as contrast between a positive sentiment and negative situation](#). In *Proceedings of*

- the 2013 Conference on Empirical Methods in Natural Language Processing, pages 704–714, Seattle, Washington, USA. Association for Computational Linguistics.
- Julian Risch, Philipp Schmidt, and Ralf Krestel. 2021. [Data integration for toxic comment classification: Making more than 40 datasets easily accessible in one unified format](#). In *Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021)*, pages 157–163, Online. Association for Computational Linguistics.
- Sara Rosenthal, Noura Farra, and Preslav Nakov. 2017. [SemEval-2017 task 4: Sentiment analysis in Twitter](#). In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 502–518, Vancouver, Canada. Association for Computational Linguistics.
- Piotr Rybak, Robert Mroczkowski, Janusz Tracz, and Ireneusz Gawlik. 2020. [KLEJ: Comprehensive benchmark for polish language understanding](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1191–1201, Online. Association for Computational Linguistics.
- Nazanin Sabri, Reyhane Akhavan, and Behnam Bahrak. 2021. [EmoPars: A collection of 30k emotion-annotated persian social media texts](#). pages 167–173.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. [Multitask prompted training enables zero-shot task generalization](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. [The risk of racial bias in hate speech detection](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678, Florence, Italy. Association for Computational Linguistics.
- Mei Silviana Saputri, Rahmad Mahendra, and Mirna Adriani. 2018. [Emotion classification on indonesian twitter dataset](#). In *2018 International Conference on Asian Language Processing, IALP 2018, Bandung, Indonesia, November 15-17, 2018*, pages 90–95. IEEE.
- Alexander G. Sboev, Aleksandr Naumov, and Roman B. Rybka. 2020. [Data-driven model for emotion detection in russian texts](#). In *Proceedings of the 2020 Annual International Conference on Brain-Inspired Cognitive Architectures for Artificial Intelligence, BICA 2020, Eleventh Annual Meeting of the BICA Society, November 10-15, 2020, Virtual Event / Natal, Rio Grande do Norte, Brazil*, volume 190 of *Procedia Computer Science*, pages 637–642. Elsevier.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilic, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelan, and et al. 2022. [BLOOM: A 176b-parameter open-access multilingual language model](#). *CoRR*, abs/2211.05100.
- Iyanuoluwa Shode, David Ifeoluwa Adelan, and Anna Feldman. 2022. [YOSM: A new Yorùbá Sentiment Corpus for Movie Reviews](#). *AfricaNLP 2022 @ICLR*.
- Debaditya Shome. 2021. [Emohind: Fine-grained multilabel emotion recognition from hindi texts with deep learning](#). In *12th International Conference on Computing Communication and Networking Technologies, ICCCNT 2021, Kharagpur, India, July 6-8, 2021*, pages 1–5. IEEE.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Santilli, Andreas

- Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K. Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Ghohlamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakas, and et al. 2022. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#). *CoRR*, abs/2206.04615.
- Yu Sun, Shuohuan Wang, Yu-Kun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. 2020. [ERNIE 2.0: A continual pre-training framework for language understanding](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8968–8975. AAAI Press.
- Varsha Suresh and Desmond C. Ong. 2021. [Not all negatives are equal: Label-aware contrastive loss for fine-grained text classification](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 4381–4394. Association for Computational Linguistics.
- Arthit Suriyawongkul, Ekapol Chuangsuwanich, Pattarawat Chormai, and Charin Polpanumas. 2019. [Pythainlp/wisesight-sentiment: First release](#).
- Haruya Suzuki, Yuto Miyauchi, Kazuki Akiyama, Tomoyuki Kajiwara, Takashi Ninomiya, Noriko Takemura, Yuta Nakashima, and Hajime Nagahara. 2022. [A Japanese dataset for subjective and objective sentiment polarity classification in micro blog domain](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 7022–7028, Marseille, France. European Language Resources Association.
- Sali A Tagliamonte. 2015. *Making waves: The story of variationist sociolinguistics*. John Wiley & Sons.
- Songbo Tan and Jin Zhang. 2008. [An empirical study of sentiment analysis for chinese documents](#). *Expert Syst. Appl.*, 34(4):2622–2629.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- Mike Thelwall, Kevan Buckley, and Georgios Paltoglou. 2012. [Sentiment strength detection for the social web](#). *J. Assoc. Inf. Sci. Technol.*, 63(1):163–173.
- Hao Tian, Can Gao, Xinyan Xiao, Hao Liu, Bolei He, Hua Wu, Haifeng Wang, and Feng Wu. 2020. [SKEP: Sentiment knowledge enhanced pre-training for sentiment analysis](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4067–4076, Online. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#). *CoRR*, abs/2302.13971.
- Cynthia Van Hee, Els Lefever, and Véronique Hoste. 2018. [SemEval-2018 task 3: Irony detection in English tweets](#). In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 39–50, New Orleans, Louisiana. Association for Computational Linguistics.
- Erik Velldal, Lilja Øvrelid, Eivind Alexander Bergem, Cathrine Stadsnes, Samia Touileb, and Fredrik Jørgensen. 2018. [NoReC: The Norwegian review corpus](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Bertie Vidgen and Leon Derczynski. 2020. Directions in abusive language training data, a systematic review: Garbage in, garbage out. *Plos one*, 15(12):e0243300.
- Deepanshu Vijay, Aditya Bohra, Vinay Singh, Syed Sarfaraz Akhtar, and Manish Shrivastava. 2018. [A dataset for detecting irony in hindi-english code-mixed social media text](#). In *Proceedings of 4th Workshop on Sentiment Computing, Sentiment Analysis, Opinion Mining, and Emotion Detection (EMSASW 2018) Co-located with the 15th Extended Semantic Web Conference 2018 (ESWC 2018), Heraklion, Greece, June 4, 2018*, volume 2111 of *CEUR Workshop Proceedings*, pages 38–46. CEUR-WS.org.
- Marilyn Walker, Jean Fox Tree, Pranav Anand, Rob Abbott, and Joseph King. 2012. [A corpus for research on deliberation and debate](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 812–817, Istanbul, Turkey. European Language Resources Association (ELRA).
- Harald G Wallbott and Klaus R Scherer. 1986. How universal and specific is emotional experience? evidence from 27 countries on five continents. *Social science information*, 25(4):763–795.
- Shuo Wan, Bohan Li, Anman Zhang, Wenhuan Wang, and Donghai Guan. 2020. [S²ap: Sequential sentiweibo analysis platform](#). In *Database Systems for Advanced Applications - 25th International Conference, DASFAA 2020, Jeju, South Korea, September*

- 24-27, 2020, *Proceedings, Part III*, volume 12114 of *Lecture Notes in Computer Science*, pages 745–749. Springer.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Zengzhi Wang, Qiming Xie, Zixiang Ding, Yi Feng, and Rui Xia. 2023. [Is chatgpt a good sentiment analyzer? A preliminary study](#). *CoRR*, abs/2304.04339.
- Zeeraq Waseem and Dirk Hovy. 2016. [Hateful symbols or hateful people? predictive features for hate speech detection on Twitter](#). In *Proceedings of the NAACL Student Research Workshop*, pages 88–93, San Diego, California. Association for Computational Linguistics.
- Bryan Wilie, Karissa Vincentio, Genta Indra Winata, Samuel Cahyawijaya, Xiaohong Li, Zhi Yuan Lim, Sidik Soleman, Rahmad Mahendra, Pascale Fung, Syafri Bahar, and Ayu Purwarianti. 2020. [Indonlu: Benchmark and resources for evaluating indonesian natural language understanding](#). In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2020, Suzhou, China, December 4-7, 2020*, pages 843–857. Association for Computational Linguistics.
- Genta Indra Winata, Alham Fikri Aji, Samuel Cahyawijaya, Rahmad Mahendra, Fajri Koto, Ade Romadhony, Kemal Kurniawan, David Moeljadi, Radityo Eko Prasajo, Pascale Fung, Timothy Baldwin, Jey Han Lau, Rico Sennrich, and Sebastian Ruder. 2022. [Nusax: Multilingual parallel sentiment dataset for 10 indonesian local languages](#). *CoRR*, abs/2205.15960.
- Minghao Wu, Abdul Waheed, Chiyu Zhang, Muhammad Abdul-Mageed, and Alham Fikri Aji. 2023. [Lamini-lm: A diverse herd of distilled models from large-scale instructions](#). *CoRR*, abs/2304.14402.
- Rong Xiang, Xuefeng Gao, Yunfei Long, Anran Li, Emmanuele Chersoni, Qin Lu, and Chu-Ren Huang. 2020. [Ciron: a new benchmark dataset for chinese irony detection](#). In *Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020*, pages 5714–5720. European Language Resources Association.
- Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi, Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang, Weijian Xie, Yanting Li, Yina Patterson, Zuoyu Tian, Yiwen Zhang, He Zhou, Shaowei-hua Liu, Zhe Zhao, Qipeng Zhao, Cong Yue, Xinrui Zhang, Zhengliang Yang, Kyle Richardson, and Zhenzhong Lan. 2020a. [CLUE: A chinese language understanding evaluation benchmark](#). In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 4762–4772. International Committee on Computational Linguistics.
- Lu Xu, Lidong Bing, Wei Lu, and Fei Huang. 2020b. [Aspect sentiment classification with aspect-specific opinion spans](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3561–3567, Online. Association for Computational Linguistics.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [mt5: A massively multilingual pre-trained text-to-text transformer](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 483–498. Association for Computational Linguistics.
- Seid Muhie Yimam, Hizkiel Mitiku Alemayehu, Abinew Ayele, and Chris Biemann. 2020. [Exploring Amharic sentiment analysis from social media texts: Building annotation tools and classification models](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1048–1060, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. [Predicting the type and target of offensive posts in social media](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1415–1420, Minneapolis, Minnesota. Association for Computational Linguistics.
- Marcos Zampieri, Preslav Nakov, Sara Rosenthal, Pepa Atanasova, Georgi Karadzhov, Hamdy Mubarak, Leon Derczynski, Zeses Pitenis, and Çağrı Çöltekin. 2020. [Semeval-2020 task 12: Multilingual offensive language identification in social media \(offenseval 2020\)](#). In *Proceedings of the Fourteenth Workshop on Semantic Evaluation, SemEval@COLING 2020, Barcelona (online), December 12-13, 2020*, pages 1425–1447. International Committee for Computational Linguistics.
- Chiyu Zhang and Muhammad Abdul-Mageed. 2022. [Improving social meaning detection with pragmatic masking and surrogate fine-tuning](#). In *Proceedings of the 12th Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media Analysis*, pages 141–156, Dublin, Ireland. Association for Computational Linguistics.

- Chiyu Zhang, Muhammad Abdul-Mageed, and Ganesh Jawahar. 2023a. [Contrastive learning of sociopragmatic meaning in social media](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2405–2439, Toronto, Canada. Association for Computational Linguistics.
- Chiyu Zhang, Muhammad Abdul-Mageed, and El Moatez Billah Nagoudi. 2022. [Decay no more: A persistent twitter dataset for learning social meaning](#). In *Workshop Proceedings of the 16th International AAAI Conference on Web and Social Media, ICWSM 2022 Workshops, Atlanta, Georgia, USA [hybrid], June 6, 2022*.
- Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. 2023b. [Sentiment analysis in the era of large language models: A reality check](#). *CoRR*, abs/2305.15005.
- Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. 2023. [Can chatgpt understand too? A comparative study on chatgpt and fine-tuned BERT](#). *CoRR*, abs/2302.10198.
- Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2023. [Can large language models transform computational social science?](#) *CoRR*, abs/2305.03514.

Appendices

A Benchmark

Table 7 summarizes language distribution of datasets in SPARROW and taxonomy of these language according to Ethnologue (Gordon Jr, 2005) and Glottolog (Nordhoff and Hammarström, 2011). Tables 9, 10, 11, 12, 13, and 14 describe the datasets in tasks of antisocial language detection, emotion recognition, humor detection, irony and sarcasm detection, sentiment analysis, and subjectivity analysis, respectively.

We empirically characterize the issue of data inaccessibility by re-collecting tweets content via tweet IDs. Table 8 shows the data decay issue of 25 datasets.

B Models

B.1 Finetuning on Encoder-only LLMs

We evaluate the following Transformer-encoder-based multilingual PLMs on SPARROW. We finetune each PLMs on the full training set and update all the parameters of the model during the training.

(1) **Multilingual-BERT** (mBERT) (Devlin et al., 2019) is trained on a Wikipedia corpus including 104 languages with masked language modelling (MLM) and next sentence prediction objectives. It contains 110M parameters. mBERT tokenizes text by using WordPiece with a vocabulary size of 172K.

(2) **XLM-RoBERTa_{Base}** (XLM-R) (Conneau et al., 2020) is trained on CommonCrawl data involving 100 languages with MLM objective. It uses a SentencePiece tokenizer with a vocabulary size of 250K and contains 270M parameters.

(3) **Bernice** (DeLucia et al., 2022) is trained with 2.5B tweets in 66 languages and MLM objective. Bernice consists of 270M parameters and a tweet-specific SentencePiece tokenizer including a vocabulary size of 250K.

(4) **InfoDCL** (Zhang et al., 2023a) further trains XLM-R with 100M tweets in 66 languages with two contrastive learning, MLM, and distant label prediction objectives. InfoDCL shows that it effectively learns language representations for understanding SM.

B.2 Zero-shot Setting on LLMs

We also investigate the zero-shot performance on a wide range of LLMs:

(1) **BLOOM** (Scao et al., 2022) is a Transformer decoder-only model trained on the ROOTS corpus consisting of 46 natural and 13 programming languages. BLOOM uses a multilingual vocabulary with 250K tokens and is trained with auto-regressive language modelling objectives.

(2) **Multilingual T5 (mT5)** (Xue et al., 2021) is Transformer encoder-decoder model trained on CommonCrawl data involving 101 languages and contains a vocabulary with 250K tokens. It trained with sequence-to-sequence MLM objective.

(3) **LLaMA** (Touvron et al., 2023) is a Transformer decoder-only model pretrained on 1.4T tokens where the majority are English and a small amount of data in 20 other languages. We utilize LLaMA with 7B parameters and a vocabulary with 30K tokens.

(4) **BLOOMZ** (Muennighoff et al., 2022) is also an instruction finetune model. It further finetunes BLOOM on xP3 corpus that contains 13 type of tasks in 46 languages with English prompt. We benchmark SPARROW on the BLOOM-based models with a size of 7.1B parameters.

(5) **BLOOMZ-P3** (Muennighoff et al., 2022) is an instruction finetuned model. It is initialized by BLOOM and further finetunes on English-only P3 corpus (Sanh et al., 2022) containing 2,073 natural language prompts for eight types of NLP tasks.

(6) **BLOOM-Bactrian** (Li et al., 2023) tune BLOOM on a 3.4M instruction-following dataset in 52 languages with low-rank adaptation modules. Li et al. (2023) translate the English 67K instructions from Alpaca and Dolly datasets into 51 languages and utilize ChatGPT API to generate responses in the corresponding language.

(7) **mT0** (Muennighoff et al., 2022) is instruction fine-tuned mT5 model with xP3 corpus. We evaluate the mT5-based models with XL size (with 3.7B parameters).

(8) **Alpaca** (Taori et al., 2023) further tune LLaMA on a 52K instruction-following dataset that is generated by gpt-3.5-turbo of OpenAI API. The dataset includes diverse English instruction-following tasks, e.g., question answering and programming.

(9) **Vicuna** (Chiang et al., 2023) further tune LLaMA on 70K diverse user-shared conversations with ChatGPT in English.

(10) **ChatGPT** is a conversation-based LLM trained GTP-3 (Brown et al., 2020) through reinforcement learning with human feedback (Ouyang

et al., 2022; Christiano et al., 2017). We exploit gpt-3.5-turbo-0301 via OpenAI API.¹¹

C Experiments

C.1 Hyperparameters

To be computation friendly, we only tune the peak learning rate of each model in a set of $\{1e-4, 5e-5, 3e-5, 1e-5\}$ and randomly select 45 datasets for hyper-parameter tuning. We fine-tune a PLM with an arbitrary batch size of 32, sequence length of 128 tokens, and 20 epochs with patience of five epochs based on the model performance on Dev set. We fine-tune each dataset three time with different seeds and identify the best model based on Dev set performance. The best learning rate for each model is identified based on the average score of Dev set of the 45 datasets. The best peak learning rate is $3e-5$ for mBERT, XLM-T, and Bernice and $1e-5$ for other models.

C.2 Prompts

The prompts we use in our experiments are summarized in Table 15.

C.3 Results

Table 16 shows aggregated performance of fine-tuned models on Dev and Test-S. We report the average of dataset-specific metrics and standard deviation in a task and a category. We also report the Test-S performance of tasks of antisocial language detection, emotion recognition, humor detection, irony and sarcasm detection, sentiment analysis, and subjectivity analysis in Tables 17, 18, 19, 20, 21, and 22, respectively.

We provide a concise study to probe the sensitivity of open-source LLMs to prompts and present the results in Table 24.

¹¹<https://openai.com/>

Lang. Family	Lang.	Code	# dataset	Script
Afro-Asiatic	Amharic	amh	1	Ethiopic
	Arabic	ara	15	Arabic
	Darija	ary	1	Arabic
	Dziriya	arq	1	Arabic
	Hausa	hau	1	Latin
	Hebrew	heb	1	Hebrew
	Maltese	mlt	1	Latin
Atlantic-Congo	Bambara	bam	1	Latin
	Igbo	ibo	1	Latin
	Kinyarwanda	kin	1	Latin
	Swahili	swh	1	Latin
	Twi	twi	1	Latin
	Tsonga	tso	1	Latin
	Yoruba	yor	2	Latin
Austroasiatic	Vietnamese	vie	1	Latin
Austronesian	Acehnese	ace	1	Latin
	Balinese	ban	1	Latin
	Banjarese	bjn	1	Latin
	Buginese	bug	1	Latin
	Filipino	fil	1	Latin
	Indonesian	ind	3	Latin
	Javanese	jav	1	Latin
	Madurese	mad	1	Latin
	Minangkabau	min	1	Latin
	Ngaju	nij	1	Latin
Dravidian	Sundanese	sun	1	Latin
	Toba batak	bbc	1	Latin
	Kannada	kan	2	Kannada, Latin
	Malayalam	mal	2	Malayalam, Latin
	Tamil	tam	2	Tamil, Latin
	Telugu	tel	2	Telugu
Indo-European	Albanian	sqi	1	Latin
	Bosnian	bos	1	Latin
	Bulgarian	bul	1	Cyrillic
	Bengali	ben	4	Bengali, Latin
	Croatian	hrv	1	Latin
	Czech	ces	2	Latin
	Danish	dan	1	Latin
	English	eng	27	Latin
	French	fra	6	Latin
	German	deu	3	Latin
	Greek	ell	1	Greek
	Hindi	hin	5	Devanagari, Latin
	Italian	ita	11	Latin
	Marathi	mar	1	Devanagari
	Nigerian Pidgin	pcm	2	Latin
	Norwegian	nor	1	Latin
	Persian	fas	3	Arabic
	Portuguese	por	4	Latin
	Polish	pol	4	Latin
	Romanian	ron	2	Latin
Russian	rus	3	Cyrillic	
Spanish	spa	9	Latin	
Serbian	srp	1	Cyrillic	
Slovak	slk	1	Latin	
Slovenian	slv	2	Latin	
Swedish	swe	1	Latin	
Japonic	Japanese	jpn	1	Han, Hir., Kat.
Koreanic	Korean	kor	5	Hangul
Sino-Tibetan	Chinese	zho	6	Han
Tai-Kadai	Thai	tha	1	Thai
Turkic	Turkish	tur	2	Latin
Uralic	Finnish	fin	3	Latin
	Hungarian	hun	1	Latin

Table 7: Summary of languages covered in SPARROW. **Lang.:** Language. Language code is marked by ISO 639-3 code. Language information is retrieved from Ethnologue (Gordon Jr, 2005) and Glottolog (Nordhoff and Hammarström, 2011). The column **# dataset** shows the number of datasets covered by SPARROW per language. **Hir.:** Hiragana, **Kat.:** Katakana

Dataset	Study	Year	Original	Retrieval	Decay %
Sarc-engRil	Riloff et al. (2013)	2013	3K	1K	0.41
Sarc-cesPta	Ptáček et al. (2014)	2013	7K	4K	0.29
Sarc-engPta	Ptáček et al. (2014)	2013	100K	89K	0.11
Sarc-engBam	Bamman and Smith (2015)	2015	19K	14K	0.24
Sent-bulMoz	Mozetič et al. (2016)	2016	67K	27K	0.59
Sent-bosMoz	Mozetič et al. (2016)	2016	44K	20K	0.54
Sent-deuMoz	Mozetič et al. (2016)	2016	109K	52K	0.52
Sent-engMoz	Mozetič et al. (2016)	2016	103K	43K	0.58
Sent-spaMoz	Mozetič et al. (2016)	2016	275K	153K	0.44
Sent-hrvMoz	Mozetič et al. (2016)	2016	97K	66K	0.32
Sent-hunMoz	Mozetič et al. (2016)	2016	109K	40K	0.63
Sent-polMoz	Mozetič et al. (2016)	2016	223K	109K	0.51
Sent-porMoz	Mozetič et al. (2016)	2016	157K	49K	0.69
Sent-rusMoz	Mozetič et al. (2016)	2016	107K	41K	0.62
Sent-slkMoz	Mozetič et al. (2016)	2016	70K	38K	0.46
Sent-slvMoz	Mozetič et al. (2016)	2016	133K	74K	0.44
Sent-sqiMoz	Mozetič et al. (2016)	2016	53K	36K	0.31
Sent-srpMoz	Mozetič et al. (2016)	2016	73K	27K	0.63
Sent-sweMoz	Mozetič et al. (2016)	2016	58K	34K	0.42
Hate-engWas	Waseem and Hovy (2016)	2016	16K	10K	0.36
Sent-porBru	Brum (2018)	2017	157K	56K	0.64
Sent-engRos	Rosenthal et al. (2017)	2017	50K	42K	0.15
Iron-hinVij	Vijay et al. (2018)	2018	3K	2K	0.10
Sexi-freChi	Chiril et al. (2020)	2018	12K	9K	0.22
Humo-hinAgg	Aggarwal et al. (2020)	2018	7K	5K	0.30

Table 8: Data decay issue in social media data. These 25 datasets are distribute by tweet IDs. We retrieve these tweet on Nov. 2020 - Jan. 2022 and find that 42% samples are inaccessible.

Dataset	Study	Lang.	Source / Domain	Year	#Lb	Labels	Data Slipt	Metric
Aggr-hinKum	Kumar et al. (2018)	hin	Twitter, Facebook	2018	2	{Aggressive, Not}	9,306/1,163/1,164	M-F1
Dang-araAIs	Alshehri et al. (2020)	ara	Twitter	2020	2	{Dangerous, Not}	3,474/615/663	M-F1
Hate-engWas	Waseem and Hovy (2016)	eng	Twitter	2016	3	{Not, Racism, Sexism}	8,683/1,086/1,085	W-F1
Hate-engDav	Davidson et al. (2017)	eng	Twitter	2017	3	{Hate, Not, Offensive}	19,826/2,478/2,479	W-F1
Hate-araAla	Alakrot et al. (2018)	ara	YouTube comment	2017	2	{Hate, Not}	9,014/1,127/1,127	M-F1
Hate-itaBos*	Bosco et al. (2018)	ita	Twitter	2018	2	{Hate, Not}	2,700/300/1,000	M-F1
Hate-flCab	Cabasag et al. (2019)	fil	Twitter	2016	2	{Hate, Not}	10,000/4,232/4,232	M-F1
Hate-araMul	Mulki et al. (2019)	ara	Twitter	2019	3	{Abusive, Hate, Not}	4,208/468/1,170	M-F1
Hate-engBas	Basile et al. (2019)	eng	Twitter	2019	2	{Hate, Not}	9,000/1,000/3,000	M-F1
Hate-spaBas	Basile et al. (2019)	spa	Twitter	2019	2	{Hate, Not}	4,500/500/1,600	M-F1
Hate-porFor	Fortuna et al. (2019)	por	Twitter	2019	2	{Hate, Not}	4,536/567/567	M-F1
Hate-polPla	Ptaszynski et al. (2019)	pol	Twitter	2019	2	{Hate, Not}	9,037/1,004/1,000	M-F1
Hate-korMoo	Moon et al. (2020)	kor	News comment	2020	3	{Hate, Not, Offensive}	7,106/790/471	M-F1
Hate-araMub	Mubarak et al. (2020)	ara	Twitter	2020	2	{Hate, Not}	6,839/1,000/2,000	M-F1
Hate-zhoDen	Deng et al. (2022)	zho	Weibo	2022	2	{Hate, Not}	25,726/6,431/5,323	M-F1
Hate-korJeo	Jeong et al. (2022)	kor	News and YouTube comment	2022	2	{Hate, Not}	32,343/4,043/4,043	M-F1
Hate-telMar	Marreddy et al. (2022)	tel	Misc	2022	2	{Hate, Not}	24,599/3,510/7,033	M-F1
Sexi-fraChi	Chiril et al. (2020)	fra	Twitter	2018	2	{Not, Sexism}	7,670/959/959	M-F1
Offe-engZam	Zampieri et al. (2019)	eng	Twitter	2019	2	{Not, Offensive}	11,916/1,324/860	M-F1
Offe-araZam	Zampieri et al. (2020)	ara	Twitter	2019	2	{Not, Offensive}	7,055/784/1,827	M-F1
Offe-danZam	Zampieri et al. (2020)	dan	Misc	2019	2	{Not, Offensive}	2,664/296/329	M-F1
Offe-ellZam	Zampieri et al. (2020)	ell	Twitter	2019	2	{Not, Offensive}	7,869/874/1,544	M-F1
Offe-turZam	Zampieri et al. (2020)	tur	Twitter	2019	2	{Not, Offensive}	28,149/3,128/3,515	M-F1
Offe-araMub	Mubarak et al. (2020)	ara	Twitter	2020	2	{Not, Offensive}	6,839/1,000/2,000	M-F1
Offe-slvNov	Kralj Novak et al. (2021)	slv	Twitter	2020	4	{Appropriate, Inappropriate, Not, Offensive}	65,021/8,127/8,128	M-F1
Offe-G-engZam	Zampieri et al. (2019)	eng	Twitter	2019	3	{Group, Individual, Others}	3,485/391/213	M-F1
Hate-G-araOus	Ousidhoum et al. (2019)	ara	Twitter	2019	13	{African_descent, Arabs, Asians, Christian, Gay, Immigrants, Indian/hindu, Individual, Jews, Muslims, Others, Refugees, Women}	2,682/334/335	M-F1
Hate-G-fraOus	Ousidhoum et al. (2019)	fra	Twitter	2019	16	{African_descent, Arabs, Asians, Christian, Gay, Gspanics, Immigrants, Indian/hindu, Individual, Jews, Left_wing_people, Muslims, Others, Refugees, Special_needs, Women}	3,211/401/402	M-F1
Offe-T-engZam	Zampieri et al. (2019)	eng	Twitter	2019	2	{Targeted, Untargeted}	3,963/437/240	M-F1
Hate-T-araOus	Ousidhoum et al. (2019)	ara	Twitter	2019	4	{Gender, Origin, Others, Religion}	2,682/334/336	M-F1
Hate-T-fraOus	Ousidhoum et al. (2019)	fra	Twitter	2019	6	{Disability, Gender, Origin, Others, Religion, Sexual_Orientation}	3,211/401/402	M-F1
Hate-T-benKar	Karim et al. (2021)	ben	Misc	2020	4	{Geopolitical, Personal, Political, Religion}	4,558/570/570	M-F1
Offe-T-kanCha	Chakravarthi et al. (2022)	kan	YouTube comment	2019	5	{Group, Individual, Not, Others, Untargeted}	4,694/586/593	M-F1
Offe-T-malCha	Chakravarthi et al. (2022)	mal	YouTube comment	2019	4	{Group, Individual, Not, Untargeted}	14,723/1,836/1,844	M-F1
Offe-T-tamCha	Chakravarthi et al. (2022)	tam	YouTube comment	2019	5	{Group, Individual, Not, Others, Untargeted}	33,685/4,216/4,232	M-F1
Hate-T-korJeo	Jeong et al. (2022)	kor	News and YouTube comment	2022	4	{Group, Individual, Other, Untargeted}	16,239/2,049/2,022	M-F1

Table 9: Description of 36 antisocial language detection datasets. **Lang.:** Language is marked by ISO 639-3, **#Lb:** the label size of a dataset. **M-F1:** Macro-F1, **W-F1:** Weighted-F1. * indicates that data sharing needs approval from the original authors.

Dataset	Study	Lang.	Source / Domain	Year	#Lb	Labels	Data Slipt	Metric
Emot-engWal	Wallbott and Scherer (1986)	eng	Questionnaire	1986	7	{Anger, Disgust, Fear, Guilt, Joy, Sadness, Shame}	6,132/767/767	M-F1
Emot-zhoLee	Lee and Wang (2015)	zho	Weibo	2015	5	{Anger, Fear, Happy, Sadness, Surprise}	3,122/347/418	Accuracy
Emot-finKaj	Kajava (2018)	fin	Subtitle	2016	8	{Anger, Anticipation, Disgust, Fear, Joy, Sadness, Surprise, Trust}	5,197/577/653	M-F1
Emot-fraKaj	Kajava (2018)	fra	Subtitle	2016	8	{Anger, Anticipation, Disgust, Fear, Joy, Sadness, Surprise, Trust}	5,198/577/653	M-F1
Emot-itaKaj	Kajava (2018)	ita	Subtitle	2016	8	{Anger, Anticipation, Disgust, Fear, Joy, Sadness, Surprise, Trust}	5,197/577/653	M-F1
Emot-araAbd	Abdul-Mageed et al. (2020)	ara	Twitter	2016	8	{Anger, Anticipation, Disgust, Fear, Joy, Sadness, Surprise, Trust}	50,000/910/941	M-F1
Emot-engMoh	Mohammad et al. (2018)	eng	Twitter	2018	4	{Anger, Joy, Optimism, Sadness}	3,257/374/1,421	M-F1
Emot-araMoh	Mohammad et al. (2018)	ara	Twitter	2018	4	{Anger, Joy, Fear, Sadness}	2,284/490/1,188	M-F1
Emot-spaMoh	Mohammad et al. (2018)	spa	Twitter	2018	4	{Anger, Joy, Fear, Sadness}	2,708/479/1,696	M-F1
Emot-indSap	Saputri et al. (2018)	ind	Twitter	2019	5	{Anger, Fear, Happy, Love, Sadness}	3,520/440/441	M-F1
Emot-turGuv	Güven et al. (2020)	tur	Twitter	2020	5	{Anger, Fear, Happy, Sadness, Surprise}	3,200/400/400	Accuracy
Emot-indWil	Wilie et al. (2020)	ind	Twitter	2018	5	{Anger, Fear, Happy, Love, Sadness}	3,169/352/440	M-F1
Emot-vieHo	Ho et al. (2019)	vie	Facebook	2019	7	{Anger, Disgust, Fear, Joy, Others, Sadness, Surprise}	5,548/686/693	W-F1
Emot-engPla	Plaza del Arco et al. (2020)	eng	Twitter	2020	7	{Anger, Disgust, Fear, Joy, Others, Sadness, Surprise}	5,842/730/731	M-F1
Emot-spaPla	Plaza del Arco et al. (2020)	spa	Twitter	2020	7	{Anger, Disgust, Fear, Joy, Others, Sadness, Surprise}	6,727/841/841	M-F1
Emot-finOhm	Öhman et al. (2020)	fin	Subtitle	2020	8	{Anger, Anticipation, Disgust, Fear, Joy, Sadness, Surprise, Trust}	8,864/1,118/1,086	M-F1
Emot-engDem	Demszky et al. (2020)	eng	Reddit	2020	27	{Admiration, Amusement, Anger, Annoyance, Approval, Caring, Confusion, Curiosity, Desire, Disappointment, Disapproval, Disgust, Embarrassment, Excitement, Fear, Gratitude, Grief, Joy, Love, Nervousness, Optimism, Pride, Realization, Relief, Remorse, Sadness, Surprise6}	23,485/2,956/2,984	M-F1
Emot-itaBia	Bianchi et al. (2021)	ita	Twitter	2021	4	{Anger, Fear, Joy, Sadness}	1,629/204/204	M-F1
Emot-ronCio	Ciobotaru and Dinu (2021)	ron	Twitter	2020	4	{Anger, Fear, Joy, Sadness}	2,600/318/324	M-F1
Emot-hinDeb	Shome (2021)	hin	Machine Translation	2021	27	{Admiration, Amusement, Anger, Annoyance, Approval, Caring, Confusion, Curiosity, Desire, Disappointment, Disapproval, Disgust, Embarrassment, Excitement, Fear, Gratitude, Grief, Joy, Love, Nervousness, Optimism, Pride, Realization, Relief, Remorse, Sadness, Surprise}	23,485/2,956/2,984	M-F1
Emot-porCor	Cortiz et al. (2021)	por	Twitter	2021	28	{Admiration, Amusement, Anger, Annoyance, Approval, Compassion, Confusion, Curiosity, Desire, Disappointment, Disapproval, Disgust, Embarrassment, Envy, Excitement, Fear, Gratitude, Grief, Joy, Longing, Love, Nervousness, Optimism, Pride, Relief, Remorse, Sadness, Surprise}	24,919/2,769/12,966	M-F1
Emot-fasSab	Sabri et al. (2021)	fas	Twitter	2021	6	{Anger, Fear, Happy, Hatred, Sadness, Wonder}	4,180/523/523	M-F1
Emot-rusSbo	Sboev et al. (2020)	rus	Misc	2021	5	{Anger, Fear, Joy, Sadness, Surprise}	3,951/427/1,128	M-F1
Emot-benIqb	Iqbal et al. (2022)	ben	Misc	2022	6	{Anger, Disgust, Fear, Joy, Sadness, Surprise}	5,600/700/700	M-F1
Emot-fraBia*	Bianchi et al. (2022)	fra	Machine Translation	2018	4	{Anger, Fear, Joy, Sadness}	3,798/476/476	M-F1
Emot-deuBia*	Bianchi et al. (2022)	deu	Machine Translation	2018	4	{Anger, Fear, Joy, Sadness}	3,798/476/476	M-F1

Table 10: Description of 26 emotion recognition datasets. **Lang.:** Language is marked by ISO 639-3, **#Lb:** the label size of a dataset. **M-F1:** Macro-F1, **W-F1:** Weighted-F1.

Dataset	Study	Lang	Source / Domain	Year	#Lb	Labels	Data Slipt	Metric
Humo-hinAgg	Aggarwal et al. (2020)	hin	Twitter	2018	2	{Humor, Not}	4,187/524/523	Accuracy
Humo-rusBli	Blinov et al. (2019)	rus	Misc	2018	2	{Humor, Not}	251,416/61,794/1,877	M-F1
Humo-spaChi	Chiruzzo et al. (2021)	spa	Twitter	2019	2	{Humor, Not}	24,000/6,000/6,000	M-F1
Humo-engMea	Meaney et al. (2021)	eng	Twitter	2021	2	{Humor, Not}	8,000/1,000/1,000	M-F1

Table 11: Description of four humor detection datasets. **Lang:** Language is marked by ISO 639-3, **#Lb:** the label size of a dataset. **M-F1:** Macro-F1.

Dataset	Study	Lang.	Source / Domain	Year	#Lb	Labels	Data Slipt	Metric
Iron-ita _{Bas}	Basile et al. (2014)	ita	Twitter	2014	2	{Irony, Not}	4,062/453/1,936	M-F1
Iron-spa _{Bar}	Barbieri et al. (2016)	spa	Twitter	2014	2	{Irony, Not}	6,669/741/1,997	M-F1
Iron-eng _{Hee}	Van Hee et al. (2018)	eng	Twitter	2018	2	{Irony, Not}	3,450/384/784	F1-irony
Iron-ita _{Cig}	Cignarella et al. (2018)	ita	Twitter	2018	2	{Irony, Not}	3,579/398/872	M-F1
Iron-hin _{Vij}	Vijay et al. (2018)	hin	Twitter	2018	2	{Irony, Not}	2,217/277/277	M-F1
Iron-ara _{Gha}	Ghanem et al. (2019)	ara	Twitter	2019	2	{Irony, Not}	3,622/402/1,006	M-F1
Iron-spa _{Ort}	Ortega-Bueno et al. (2019)	spa	Twitter	2019	2	{Irony, Not}	2,160/240/600	M-F1
Iron-fas _{Gol} *	Golazizian et al. (2020)	fas	Twitter	2019	2	{Irony, Not}	2,352/295/294	Accuracy
Iron-zho _{Xia} *	Xiang et al. (2020)	zho	Weibo	2020	5	{Insufficient_Evidence, Irony, Not, Unlikely_Ironic, Weakly_Irony}	7,014/876/876	M-F1
Sarc-eng _{Wal}	Walker et al. (2012)	eng	Debate Forum	2012	2	{Not, Sarcasm}	900/100/995	M-F1
Sarc-eng _{Ril}	Riloff et al. (2013)	eng	Twitter	2013	2	{Not, Sarcasm}	1,413/177/177	F1-sarcasm
Sarc-ces _{Pta}	Ptáček et al. (2014)	ces	Twitter	2013	2	{Not, Sarcasm}	3,977/497/497	M-F1
Sarc-eng _{Pta}	Ptáček et al. (2014)	eng	Twitter	2013	2	{Not, Sarcasm}	71,433/8,929/8,930	M-F1
Sarc-eng _{Bam}	Bamman and Smith (2015)	eng	Twitter	2015	2	{Not, Sarcasm}	11,864/1,483/1,484	Accuracy
Sarc-eng _{Raj}	Rajadesingan et al. (2015)	eng	Twitter	2015	2	{Not, Sarcasm}	41,261/5,158/5,158	Accuracy
Sarc-eng _{Ora}	Oraby et al. (2016)	eng	Debate Forum	2016	2	{Not, Sarcasm}	900/100/2,260	M-F1
Sarc-zho _{Gon} *	Gong et al. (2020)	zho	News comment	2019	2	{Not, Sarcasm}	3,978/497/497	M-F1
Sarc-ara _{Abu}	Abu Farha and Magdy (2020)	ara	Twitter	2020	2	{Not, Sarcasm}	7,593/844/2,110	M-F1
Sarc-ara _{Far}	Farha et al. (2021)	ara	Twitter	2020	2	{Not, Sarcasm}	11,293/1,255/3,000	M-F1
Iron-T-eng _{Hee}	Van Hee et al. (2018)	eng	Twitter	2018	4	{Ironic_by_clash, Not, Other_irony, Situational_irony}	3,450/384/784	M-F1

Table 12: Description of 20 irony and sarcasm detection datasets. **Lang.:** Language is marked by ISO 639-3, **#Lb:** the label size of a dataset. **M-F1:** Macro-F1. * indicates that data sharing needs an approval from the original authors.

Dataset	Study	Lang.	Source / Domain	Year	#Lb	Labels	Data Splt	Metric
Sent-engPan	Pang and Lee (2005)	eng	Movie review	2005	2	{Negative, Positive}	8,529/1,066/1,067	Accuracy
Sent-zhoTan	Tan and Zhang (2008)	zho	Misc	2008	2	{Negative, Positive}	9,600/1,200/1,200	M-F1
Sent-T-engThe	Thelwall et al. (2012)	eng	Twitter	2012	2	{Negative, Positive}	900/100/1,113	Accuracy
Sent-Y-engThe	Thelwall et al. (2012)	eng	YouTube comment	2012	2	{Negative, Positive}	900/100/1,142	Accuracy
Sent-5-engSoc	Socher et al. (2013)	eng	Movie review	2013	5	{Negative, Neutral, Positive, Very_Negative, Very_Positive}	8,544/1,101/2,210	Accuracy
Sent-korJan*	Jang et al. (2013)	kor	News article	2013	4	{Complex, Negative, Neutral, Positive}	4,187/523/524	M-F1
Sent-engSoc	Socher et al. (2013)	eng	Movie review	2013	2	{Negative, Positive}	6,920/872/1,821	Accuracy
Sent-itBas	Basile et al. (2014)	ita	Twitter	2014	2	{Negative, Positive}	2,376/265/1,207	M-F1
Sent-itBas	Barbieri et al. (2016)	ita	Twitter	2016	2	{Negative, Positive}	3,738/416/1,018	M-F1
Sent-mltJin	Dingli and Sant (2016)	mlt	Movie review	2016	2	{Negative, Positive}	596/85/171	M-F1
Sent-bulMoz	Mozetić et al. (2016)	bul	Twitter	2016	3	{Negative, Neutral, Positive}	22,184/2,773/2,773	M-F1
Sent-bosMoz	Mozetić et al. (2016)	bos	Twitter	2016	3	{Negative, Neutral, Positive}	16,335/2,042/2,042	M-F1
Sent-deuMoz	Mozetić et al. (2016)	deu	Twitter	2016	3	{Negative, Neutral, Positive}	42,010/5,251/5,252	M-F1
Sent-engMoz	Mozetić et al. (2016)	eng	Twitter	2016	3	{Negative, Neutral, Positive}	34,538/4,317/4,318	M-F1
Sent-spaMoz	Mozetić et al. (2016)	spa	Twitter	2016	3	{Negative, Neutral, Positive}	122,410/15,301/15,302	M-F1
Sent-hrvMoz	Mozetić et al. (2016)	hrv	Twitter	2016	3	{Negative, Neutral, Positive}	52,971/6,621/6,622	M-F1
Sent-hunMoz	Mozetić et al. (2016)	hun	Twitter	2016	3	{Negative, Neutral, Positive}	32,717/4,089/4,090	M-F1
Sent-polMoz	Mozetić et al. (2016)	pol	Twitter	2016	3	{Negative, Neutral, Positive}	87,941/10,993/10,992	M-F1
Sent-porMoz	Mozetić et al. (2016)	por	Twitter	2016	3	{Negative, Neutral, Positive}	39,525/4,941/4,940	M-F1
Sent-rusMoz	Mozetić et al. (2016)	rus	Twitter	2016	3	{Negative, Neutral, Positive}	32,941/4,117/4,118	M-F1
Sent-slkMoz	Mozetić et al. (2016)	slk	Twitter	2016	3	{Negative, Neutral, Positive}	30,694/3,837/3,837	M-F1
Sent-slvMoz	Mozetić et al. (2016)	slv	Twitter	2016	3	{Negative, Neutral, Positive}	59,924/7,491/7,490	M-F1
Sent-sqiMoz	Mozetić et al. (2016)	sqi	Twitter	2016	3	{Negative, Neutral, Positive}	29,375/3,672/3,672	M-F1
Sent-srpMoz	Mozetić et al. (2016)	srp	Twitter	2016	3	{Negative, Neutral, Positive}	22,124/2,765/2,766	M-F1
Sent-sweMoz	Mozetić et al. (2016)	swe	Twitter	2016	3	{Negative, Neutral, Positive}	27,277/3,409/3,410	M-F1
Sent-deuRei	Rei et al. (2016)	deu	Twitter	2016	3	{Negative, Neutral, Positive}	2,701/337/338	M-F1
Sent-spaRei	Rei et al. (2016)	spa	Twitter	2016	3	{Negative, Neutral, Positive}	6,099/767/762	M-F1
Sent-itaRei	Rei et al. (2016)	ita	Twitter	2016	3	{Negative, Neutral, Positive}	6,818/853/852	M-F1
Sent-engRos	Rosenthal et al. (2017)	eng	Twitter	2017	3	{Negative, Neutral, Positive}	42,756/4,751/12,284	M-Recall
Sent-benPat*	Patra et al. (2018)	ben	Twitter	2015	3	{Negative, Neutral, Positive}	2,250/250/3,038	M-F1
Sent-hinPat*	Patra et al. (2018)	hin	Twitter	2015	3	{Negative, Neutral, Positive}	11,642/1,293/5,525	M-F1
Sent-hebAmr	Amram et al. (2018)	heb	Facebook	2018	2	{Negative, Positive}	8,951/995/2,488	Accuracy
Sent-porBru	Brum and das Graças Volpe Nunes (2018)	por	Twitter	2017	3	{Negative, Neutral, Positive}	45,127/5,585/5,637	M-F1
Sent-finKaj	Kajava (2018)	fin	Subtitle	2016	2	{Negative, Positive}	5,197/577/653	M-F1
Sent-fraKaj	Kajava (2018)	fra	Subtitle	2016	2	{Negative, Positive}	5,198/577/653	M-F1
Sent-itaKaj	Kajava (2018)	ita	Subtitle	2016	2	{Negative, Positive}	5,197/577/653	M-F1
Sent-norVel	Veldal et al. (2018)	nor	Online review	2018	6	{Negative1, Negative2, Negative3, Positive4, Positive5, Positive6}	34,903/4,360/4,351	M-F1
Sent-polKoc	Kocoo et al. (2019)	pol	Customer review	2019	4	{Complex, Negative, Neutral, Positive}	5,170/574/1,217	M-F1
Sent-thaSur	Suriyawongkul et al. (2019)	tha	Facebook	2019	3	{Negative, Neutral, Positive}	21,152/2,362/2,614	M-F1
Sent-zhoWan*	Wan et al. (2020)	zho	Weibo	2019	2	{Negative, Positive}	95,990/11,999/11,999	M-F1
Sent-fasAsh	Ashrafi Asli et al. (2020)	fas	Customer review	2020	3	{Negative, Neutral, Positive}	75,094/9,387/9,387	M-F1
Sent-ronDum	Dumitrescu et al. (2020)	ron	Customer review	2020	2	{Negative, Positive}	16,146/1,795/1,005	M-F1
Sent-pcmOye	Oyewusi et al. (2020)	pcm	Twitter	2020	3	{Negative, Neutral, Positive}	11,200/1,400/1,400	M-F1
Sent-polRyb	Rybak et al. (2020)	pol	Customer review	2020	5	{Negative, Neutral, Positive, Very_Negative, Very_Positive}	8,619/958/1,002	M-F1
Sent-indWil	Wilie et al. (2020)	ind	Misc	2019	3	{Negative, Neutral, Positive}	9,900/1,100/1,260	M-F1
Sent-araAbd	Abdul-Mageed et al. (2021)	ara	Twitter	2021	3	{Negative, Neutral, Positive}	49,301/4,443/4,933	M-F1
Sent-bamDia	Diallo et al. (2021)	bam	Misc	2021	3	{Negative, Neutral, Positive}	2,436/305/305	M-F1
Sent-benIsl	Islam et al. (2021)	ben	Twitter	2021	3	{Negative, Neutral, Positive}	12,575/1,567/1,586	M-F1
Sent-marKul	Kulkarni et al. (2021)	mar	Twitter	2020	3	{Negative, Neutral, Positive}	12,114/1,500/2,250	Accuracy
Sent-kanCha	Chakravarthi et al. (2022)	kan	YouTube comment	2019	2	{Negative, Positive}	3,995/505/502	M-F1
Sent-malCha	Chakravarthi et al. (2022)	mal	YouTube comment	2019	2	{Negative, Positive}	8,410/1,044/1,039	M-F1
Sent-tamCha	Chakravarthi et al. (2022)	tam	YouTube comment	2019	2	{Negative, Positive}	24,063/2,966/3,047	M-F1
Sent-araMuh	Abdul-Mageed et al. (2022)	ara	Twitter	2021	3	{Negative, Neutral, Positive}	1,500/500/3,000	Accuracy
Sent-amhMuh	Yimam et al. (2020)	amh	Twitter	2020	3	{Negative, Neutral, Positive}	5,984/1,497/1,999	W-F1
Sent-aryMuh	Muhammad et al. (2023b)	ary	Twitter	2021	3	{Negative, Neutral, Positive}	5,583/494/2,961	W-F1
Sent-ardMuh	Muhammad et al. (2023b)	arq	Twitter	2021	3	{Negative, Neutral, Positive}	1,651/414/958	W-F1
Sent-hauMuh	Muhammad et al. (2022)	hau	Twitter	2021	3	{Negative, Neutral, Positive}	14,172/2,677/5,303	W-F1
Sent-iboMuh	Muhammad et al. (2022)	ibo	Twitter	2021	3	{Negative, Neutral, Positive}	10,192/1,841/3,682	W-F1
Sent-pcmMuh	Muhammad et al. (2022)	pcm	Twitter	2021	3	{Negative, Neutral, Positive}	5,121/1,281/4,154	W-F1
Sent-kinMuh	Muhammad et al. (2022)	kin	Twitter	2021	3	{Negative, Neutral, Positive}	3,302/827/1,026	W-F1
Sent-swhMuh	Muhammad et al. (2022)	swh	Twitter	2021	3	{Negative, Neutral, Positive}	1,810/453/748	W-F1
Sent-tsoMuh	Muhammad et al. (2022)	tso	Twitter	2021	3	{Negative, Neutral, Positive}	804/203/254	W-F1
Sent-twiMuh	Muhammad et al. (2022)	twi	Twitter	2021	3	{Negative, Neutral, Positive}	3,481/388/949	W-F1
Sent-yorMuh	Muhammad et al. (2022)	yor	Twitter	2021	3	{Negative, Neutral, Positive}	8,522/2,090/4,515	W-F1
Sent-yorSho	Shode et al. (2022)	yor	Misc	2021	2	{Negative, Positive}	800/200/500	M-F1
Sent-jpnSuz	Suzuki et al. (2022)	jpn	SNS post	2021	5	{Negative, Neutral, Positive, Very_Negative, Very_Positive}	30,000/2,500/2,500	Accuracy
Sent-aceWin	Winata et al. (2022)	ace	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-banWin	Winata et al. (2022)	ban	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-bbcWin	Winata et al. (2022)	bbc	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-bjnWin	Winata et al. (2022)	bjn	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-bugWin	Winata et al. (2022)	bug	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-javWin	Winata et al. (2022)	jav	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-madWin	Winata et al. (2022)	mad	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-minWin	Winata et al. (2022)	min	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-nijWin	Winata et al. (2022)	nij	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-sunWin	Winata et al. (2022)	sun	Trans. of online com.	2022	3	{Negative, Neutral, Positive}	500/100/400	M-F1
Sent-telMar	Marreddy et al. (2022)	tel	Misc	2022	3	{Negative, Neutral, Positive}	24,599/3,510/7,033	M-F1

Table 13: Description of 77 sentiment analysis datasets. **Lang.:** Language is marked by ISO 639-3, **#Lb:** the label size of a dataset. **M-F1:** Macro-F1, **M-Recall:** Macro-Recall, **Trans. of online com.:** Trans. of online com. * indicates that data sharing needs approval from the original authors.

Dataset	Study	Lang.	Source / Domain	Year	#Lb	Labels	Data Splt	Metric
Subj-engPan	Pang and Lee (2004)	eng	Moview review	2004	2	{Objective, Subjective}	8,100/900/1,000	Accuracy
Subj-korJan*	Jang et al. (2013)	kor	News article	2013	7	{Agreement, Argument, Emotion, Intention, Judgment, Others, Speculation}	4,284/535/536	M-F1
Subj-itaBas	Basile et al. (2014)	ita	Twitter	2014	2	{Objective, Subjective}	4,061/452/1,935	M-F1
Subj-itaBas	Barbieri et al. (2016)	ita	Twitter	2016	2	{Objective, Subjective}	6,669/741/1,943	M-F1
Subj-spaBar	Barbieri et al. (2016)	spa	Twitter	2014	2	{Objective, Subjective}	6,669/741/1,998	M-F1
Subj-cesPri	Pribán and Steinberger (2022)	ces	Moview review	2021	2	{Objective, Subjective}	7,500/500/2,000	Accuracy

Table 14: Description of four subjectivity analysis datasets. **Lang.:** Language is marked by ISO 639-3, **#Lb:** the label size of a dataset. **M-F1:** Macro-F1. * indicates that data sharing needs an approval from the original authors.

Dataset	Prompt
Dang-araAls	{Content} Question: Is the language of this sentence {labels}? Answer:
Aggr-hinKum, Hate-engDav, Hate-engWas, Hate-araAla, Hate-itaBos, Hate-filCab Hate-araMul, Hate-engBas, Hate-spaBas, Hate-porFor, Hate-polPta, Hate-korMoo Hate-araMub, Hate-zhoDen, Hate-korJeo, Hate-telMar, Sexi-fraChi, Ofie-engZam Ofie-araZam, Ofie-danZam, Ofie-ellZam, Ofie-turZam, Ofie-araMub, Ofie-slvNov	{Content} Question: Is the language of this text {labels}? Answer:
Offe-T-kanCha, Offe-G-engZam, Hate-T-korJeo	{Content} Question: Does this offensive text target {labels}? Answers:
Hate-G-araOus, Hate-G-fraOus	{Content} Question: Does this hate speech target {labels}? Answer:
Offe-T-engZam	{Content} Question: Is this offensive text {labels} insult? Answer:
Hate-T-araOus, Hate-T-fraOus	{Content} Question: Does this hate speech text insult against people based on their attribute of {labels}? Answer:
Hate-T-benKar	{Content} Question: Does this text express {labels}? Answer:
Offe-T-malCha, Offe-T-tamCha	{Content} Question: Is this sentence hate speech or not? If yes, does this sentence target individual, group or not? Answer:
Emot-engWal, Emot-zhoLee, Emot-finKaj, Emot-fraKaj, Emot-itaKaj, Emot-msaHus, Emot-araAbd Emot-engMoh, Emot-araMoh, Emot-spaMoh, Emot-indSap, Emot-turGuv, Emot-spaMoh, Emot-indWil Emot-vieHo, Emot-engPla, Emot-spaPla, Emot-finOhm, Emot-engDem, Emot-itaBia, Emot-ronCio Emot-hinDeb, Emot-porCor, Emot-fasSab, Emot-rusSbo, Emot-benIqb, Emot-fraBia, Emot-deuBia	{Content} Question: Is the emotion of this sentence {labels}? Answer:
Humo-hinAgg, Humo-rusBij, Humo-spaChi, Humo-engMea	{Content} Question: Is this sentence {labels}? Answer:
Iron-itaBas, Iron-spaBar, Iron-engHee, Iron-itaCig, Iron-hinVij, Iron-araGha, Iron-spaOrt Iron-fasGol, Iron-zhoXia, Sarc-engWal, Sarc-engRil, Sarc-cesPta, Sarc-engPta, Sarc-engBam Sarc-engRaj, Sarc-engOra, Sarc-zhoGon, Sarc-araAbu, Sarc-araFar	{Content} Question: Is this sentence {labels}? Answer:
Iron-T-engHee	{Content} Question: Is the type of this text {labels}? Answer:
Sent-engPan, Sent-zhoTan, Sent-korJan, Sent-engSoc, Sent-itaBas, Sent-benPat, Sent-hinPat Sent-itaBas, Sent-mltDin, Sent-bulMoz, Sent-bosMoz, Sent-deuMoz, Sent-engMoz, Sent-spaMoz Sent-hrvMoz, Sent-hunMoz, Sent-polMoz, Sent-porMoz, Sent-rusMoz, Sent-slkMoz, Sent-slvMoz Sent-sqIMoz, Sent-srpMoz, Sent-sweMoz, Sent-deuRej, Sent-spaRej, Sent-itaRej, Sent-engRos Sent-hebAmr, Sent-porBru, Sent-finKaj, Sent-fraKaj, Sent-itaKaj, Sent-polKoe, Sent-thaSur Sent-zhoWan, Sent-fasAsh, Sent-ronDum, Sent-pcmOye, Sent-polRyb, Sent-indWij, Sent-araAbd Sent-bamDia, Sent-benIsl, Sent-marKul, Sent-kanCha, Sent-malCha, Sent-tamCha, Sent-araMuh Sent-amhMuh, Sent-aryMuh, Sent-arqMuh, Sent-hauMuh, Sent-iboMuh, Sent-pcmMuh, Sent-kinMuh Sent-swhMuh, Sent-tsoMuh, Sent-twiMuh, Sent-yorMuh, Sent-yorSho, Sent-jpnSuz, Sent-aceWin Sent-banWin, Sent-bbcWin, Sent-bjwWin, Sent-bugWin, Sent-javWin, Sent-madWin, Sent-minWin Sent-nijWin, Sent-sunWin, Sent-telMar, Sent-T-engThe, Sent-Y-engThe, Sent-5-engSoc	{Content} Question: Is the sentiment of this sentence {labels}? Answer:
Sent-norVel	{Content} Question: Is this text rated as {labels}? Higher is better. Answer:
Subj-korJan	{Content} Question: Does this sentence express {label}? Answer:
Subj-engPan, Subj-itaBas, Subj-itaBas, Subj-spaBar, Subj-cesPri	{Content} Question: Is this sentence {labels}? Answer:

Table 15: Prompts use for zero-shot evaluation with lm-evaluation-harness.

Dataset	Metric	Random	mBERT	XLM-R	Bernice	InfoDCL	BLOOM	BLOOMZ	BLOOMZ- (MT)	BLOOMZ- P3	BLOOMZ- Bactrian	mT5	mT0	mT0 (MT)	LLaMA	Alpaca	Vicuna	ChatGPT	ChatGPT (MT)	SoTA	SoTA study
Humo-hin _{Agg}	Acc.	51.40	78.27	78.87	78.27	80.40	58.60	39.60	39.80	39.60	40.80	52.20	39.60	42.40	56.80	60.40	37.20	56.00	61.00	69.30	Aggarwal et al. (2020)
Humo-rus _{Hi}	M-F1	47.38	86.47	87.60	88.13	88.80	34.12	34.12	34.12	34.12	36.38	45.09	34.12	34.12	35.36	32.52	54.60	72.75	71.44	89.00	Blinov et al. (2019)
Humo-spa _{Chi}	M-F1	54.58	85.46	84.45	87.20	87.66	48.25	32.07	32.07	31.97	39.39	134.28	32.07	32.07	38.56	34.55	50.10	68.28	68.80	88.50	Chiruzzo et al. (2021)
Humo-eng _{Mea}	M-F1	45.25	87.19	89.86	93.41	91.34	26.16	26.69	26.69	26.47	27.06	42.83	26.69	26.69	28.39	39.39	42.86	89.56	89.56	98.54	Meaney et al. (2021)
Average	—	49.65	84.35	85.19	86.75	87.05	41.78	33.12	33.17	33.04	35.91	43.60	33.12	33.82	39.78	41.72	46.19	71.65	72.70	—	—

Table 19: Full Test-S results on humor detection. **SoTA**: Previous SoTA performance on each respective dataset. **Underscore** indicates that we have different data splits to the SoTA model.

Dataset	Metric	Random	mBERT	XLM-R	Bernice	InfoDCL	BLOOM	BLOOMZ	BLOOMZ- (MT)	BLOOMZ- P3	BLOOMZ- Bactrian	mT5	mT0	mT0 (MT)	LLaMA	Alpaca	Vicuna	ChatGPT	ChatGPT (MT)	SoTA	SoTA study
Iron-hin _{Agg}	M-F1	41.87	63.38	61.54	63.16	65.67	46.92	46.92	50.38	50.97	46.92	49.59	46.92	46.92	49.02	13.27	55.49	59.91	55.16	59.59	Basile et al. (2014)
Iron-spa _{Par}	M-F1	40.70	57.90	58.01	61.66	67.52	46.47	46.47	48.65	56.77	46.47	48.89	46.47	46.47	49.34	13.88	52.78	66.10	63.97	54.12	Barbieri et al. (2016)
Iron-eng _{Hee}	F1-iro.	45.00	59.99	63.43	67.43	68.25	0.00	0.00	0.00	25.41	0.00	1.03	0.00	0.00	5.29	55.06	41.19	59.00	59.00	70.50	Van Hee et al. (2018)
Iron-ita _{Cig}	M-F1	51.80	70.37	72.66	77.44	75.92	34.67	34.21	45.82	48.21	34.21	138.78	34.21	34.21	48.22	33.24	56.22	73.32	74.20	73.10	Cignarella et al. (2018)
Iron-hin _{Vij}	M-F1	46.88	70.90	73.25	72.90	71.16	56.66	43.99	57.75	51.30	55.80	46.28	42.53	46.63	52.02	23.61	57.73	52.89	57.92	77.00	Vijay et al. (2018)
Iron-ara _{Cha}	M-F1	48.39	82.19	83.95	82.45	83.08	40.10	34.24	32.61	51.43	32.61	33.80	32.61	32.61	39.17	34.40	35.94	68.78	67.40	84.40	Abdul-Mageed et al. (2020)
Iron-spa _{ort}	M-F1	46.84	67.42	71.18	72.79	73.88	39.47	39.47	53.40	54.46	39.47	40.22	39.47	39.47	60.64	29.47	56.77	62.92	61.69	71.67	Ortega-Bueno et al. (2019)
Iron-fas _{Pol}	Acc.	48.30	74.04	74.60	73.58	76.53	51.70	56.46	56.46	51.70	56.46	56.46	56.46	56.46	49.66	43.88	57.82	62.59	56.80	83.10	Golzarizian et al. (2020)
Iron-zho _{Cia}	M-F1	11.71	31.96	31.19	30.52	33.36	13.65	14.56	3.09	9.89	13.65	13.65	13.55	3.17	13.61	0.63	13.40	18.58	10.06	57.20	Xiang et al. (2020)
Sarc-eng _{Wal}	M-F1	53.33	63.32	65.36	67.45	63.65	45.21	32.80	32.80	33.23	41.54	145.01	32.80	32.80	34.55	50.47	33.38	70.79	70.79	69.00	Felbo et al. (2017)
Sarc-eng _{Al}	F1-sar.	27.00	46.76	52.50	54.89	57.46	10.53	0.00	0.00	0.00	0.00	36.36	0.00	0.00	0.00	38.97	14.81	50.00	50.00	51.00	Riloff et al. (2013)
Sarc-ces _{Pa}	M-F1	35.92	66.12	60.17	65.97	67.89	46.90	49.03	3.68	49.03	49.75	33.94	49.03	3.68	34.46	6.22	55.00	51.20	52.48	58.20	Piaček et al. (2014)
Sarc-eng _{Pa}	M-F1	49.85	94.28	95.76	94.99	95.56	41.77	38.42	38.42	39.49	37.05	45.25	38.42	38.42	40.45	33.01	48.17	74.30	74.30	92.37	Piaček et al. (2014)
Sarc-eng _{Bam}	Acc.	52.20	79.73	80.40	82.40	82.27	48.60	52.00	52.00	51.80	50.60	49.00	52.00	52.00	49.20	52.00	52.20	64.60	64.60	85.10	Bamman and Smith (2015)
Sarc-eng _{Raj}	Acc.	48.80	94.20	95.33	96.27	95.67	77.60	91.20	91.20	90.80	87.80	119.60	91.20	91.20	87.00	15.00	83.60	74.60	74.60	92.94	Rajadesingan et al. (2015)
Sarc-eng _{Ora}	M-F1	49.00	72.73	75.87	74.69	75.64	41.48	32.89	32.89	32.80	42.11	41.83	32.89	32.89	33.71	47.08	36.08	73.78	73.78	75.00	Felbo et al. (2017)
Sarc-zho _{Gon}	M-F1	48.25	71.01	70.63	72.75	70.53	49.04	33.29	33.20	33.20	56.10	138.89	33.29	33.29	36.25	41.41	49.36	53.50	51.05	73.68	Gong et al. (2020)
Sarc-ara _{abu}	M-F1	44.26	69.06	69.57	69.57	71.87	27.16	44.57	16.39	44.57	45.15	20.66	44.57	16.39	53.38	17.21	53.74	74.38	75.47	76.30	Abdul-Mageed et al. (2021)
Sarc-ara _{Far}	M-F1	46.22	66.87	68.45	68.80	68.81	41.67	42.00	21.63	41.93	53.31	30.34	42.00	21.63	42.65	23.47	50.37	66.24	68.43	73.10	Farha et al. (2021)
Iron-T-eng _{Hee}	M-F1	22.36	47.35	46.43	56.04	57.58	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	18.83	30.81	30.81	50.70	Van Hee et al. (2018)
Average	—	42.93	67.48	68.51	70.29	71.12	38.92	37.57	34.46	41.79	40.39	35.42	37.36	32.35	39.87	29.56	46.14	60.41	59.63	—	—

Table 20: Full Test-S results on irony & sarcasm detection. **SoTA**: Previous SoTA performance on each respective dataset. **Underscore** indicates that we have different data splits to the SoTA model. Best model of each dataset is in **bold**.

Lang Fam.	Lang	Random	InfoDCL	BMZ-P3	mT0	Vicuna	CG	CG-MT
Afro-Asiatic	ara	34.05	73.53	36.78	35.61	33.61	60.81	52.53
	amh	37.95	65.68	16.05	22.49	2.99	20.62	46.82
	arq	34.23	71.25	52.02	18.05	5.33	63.89	67.58
	ary	35.94	53.44	37.40	23.41	16.64	52.19	51.66
	hau	35.22	72.18	30.14	20.93	17.39	55.52	34.13
	heb	47.60	95.80	71.20	76.60	40.80	84.20	57.40
	mlt	47.51	68.01	47.70	39.25	48.98	78.47	77.45
Atlantic-C.	bam	27.61	65.57	36.70	24.31	14.40	40.27	37.11
	ibo	32.90	76.75	23.21	11.52	28.37	57.55	33.46
	kin	35.83	56.69	28.06	14.84	22.23	53.78	29.03
	swh	33.72	61.29	17.60	14.02	45.55	54.39	53.84
	twi	34.04	64.51	46.56	29.34	5.28	51.12	32.06
	tso	34.03	52.55	45.54	30.74	6.50	42.58	35.70
	yor	41.68	74.88	43.10	35.83	26.38	64.77	42.31
Austroasi.	vie	16.12	64.58	12.22	27.81	9.52	54.69	32.96
Austrones.	ace	34.78	77.36	37.89	24.73	12.79	52.63	58.05
	ban	30.14	79.49	41.90	29.91	13.82	60.91	42.28
	bjn	30.77	84.50	50.51	27.78	14.89	69.34	75.43
	bug	30.77	71.55	34.60	18.27	12.90	34.63	30.86
	fil	52.37	79.01	34.47	34.47	34.47	69.13	66.67
	ind	22.85	83.05	42.86	36.77	20.69	75.29	64.14
	jav	31.62	84.79	48.06	37.65	15.21	73.03	78.56
	mad	28.64	78.36	45.44	21.85	13.14	61.07	61.14
	min	34.41	84.07	49.93	32.41	14.95	69.80	62.91
	nij	34.86	77.22	42.86	22.89	15.21	57.64	57.07
	sun	32.18	81.71	44.83	37.65	12.93	64.97	68.76
	bbc	30.60	73.58	36.42	19.20	13.86	38.43	40.65
Dravidian	kan	32.14	61.79	39.23	27.80	19.12	44.81	30.39
	mal	31.68	82.70	43.84	41.65	24.85	44.03	31.44
	tam	29.45	57.81	38.53	34.27	17.60	45.85	33.20
	tel	32.29	59.61	40.99	35.15	36.61	53.33	38.68
Indo-Euro.	sqi	29.39	47.04	31.78	26.59	16.11	46.82	33.84
	bos	34.37	68.31	31.80	21.16	16.72	63.22	48.15
	bul	32.48	65.09	27.35	19.37	23.69	59.32	54.39
	ben	24.53	69.49	24.71	25.18	15.62	53.25	37.33
	hrv	30.23	67.83	35.18	27.47	13.67	62.02	55.66
	ces	41.06	80.15	50.91	46.61	51.00	65.30	63.24
	dan	41.14	82.09	46.68	46.68	51.84	66.91	66.79
	eng	37.90	75.48	43.32	39.23	39.75	66.51	—
	fra	24.65	66.58	23.29	32.20	30.46	62.08	56.84
	deu	24.40	67.55	20.00	26.51	29.31	52.24	35.41
	ell	41.24	79.13	46.71	45.47	48.21	60.94	34.98
	hin	35.24	67.55	28.92	26.20	29.06	52.63	48.30
	ita	39.83	74.78	38.76	41.45	45.05	70.29	64.37
	mar	29.60	86.47	43.80	35.40	34.40	68.60	50.20
	pcm	33.36	66.65	47.16	24.25	13.92	61.13	42.75
	nor	16.24	42.36	19.84	35.73	9.80	42.11	41.41
	fas	33.03	62.43	35.01	38.73	26.64	55.02	51.17
	por	29.17	66.10	24.45	20.14	19.16	41.10	39.37
	pol	30.43	68.04	31.82	31.96	19.15	58.25	49.43
	ron	38.16	89.83	48.97	67.14	30.54	89.11	77.34
rus	31.85	84.31	25.49	28.33	28.51	71.34	69.08	
spa	35.98	70.20	30.80	32.17	37.71	57.16	56.64	
srp	34.09	56.85	27.45	20.55	19.05	55.84	52.54	
slk	28.57	75.71	31.89	27.85	10.46	57.46	56.18	
slv	26.32	58.62	14.37	16.84	16.64	46.25	36.97	
	swe	33.35	70.98	29.31	20.77	16.05	60.10	62.58
Japonic	jpn	18.00	61.47	44.20	47.60	26.80	55.80	44.40
Koreanic	kor	28.55	57.46	19.48	21.24	16.26	42.90	37.06
Sino-Tib.	zho	36.77	76.10	48.19	46.94	38.19	63.43	61.39
Tai-Kadai	tha	31.26	75.17	20.16	23.94	24.40	52.81	36.32
Turkic	tur	29.56	87.71	32.32	48.07	35.44	82.27	64.84
Uralic	fin	25.68	63.85	15.84	30.45	23.08	63.17	57.48
	hun	27.80	68.37	30.31	27.19	16.39	53.22	43.89

Table 23: Language-wise model performance. The best performance in each language is **bold**, and the second best is in **green highlight**. The **red font** denotes a performance lower than the random baseline. **BMZ**: BLOOMZ, **CG**: ChatGPT, **MT**: using machine translated prompts.

Task	lm-evaluation-harness Prompts						ChatGPT Prompts					
	BLOOM	BLOOMZ	mT5	mT0	LLaMA	Vicuna	BLOOM	BLOOMZ	mT5	mT0	LLaMA	Vicuna
Hate	39.83	38.52	23.29	37.33	37.80	41.59	18.39	31.34	28.96	38.19	18.37	37.33
Emotion	9.71	15.07	7.75	27.87	15.14	18.12	8.61	20.07	7.57	29.63	17.61	8.65
Humor	41.78	33.04	43.60	33.12	39.78	46.19	41.59	44.99	41.34	34.13	41.59	33.12
Irony	36.63	44.46	36.52	34.69	40.78	47.48	27.33	26.02	36.08	34.31	26.02	34.70
Average	25.21	28.01	19.58	32.12	27.72	31.80	16.92	26.14	20.91	33.21	20.95	23.03

Table 24: Study on model sensitivity to prompts used for zero-shot evaluation.