

# When Does Hate Transfer? A Large-Scale Study of Cross-Lingual Hate Speech Detection

Anonymous ACL submission

## Abstract

Effective detection of hate content is essential for maintaining healthy online communities, yet performance can vary substantially across languages. This work investigates a core practical question: Can training data from one language support hate-speech detection in another language with comparable performance? To answer this question, we conduct a large-scale study across 14 languages. Our analysis shows that transfer effectiveness could depend on several factors—including, but not limited to, dataset size, resources and temporal overlap. At the same time, we observe consistent patterns in which certain language pairs exhibit stronger transfer performance. Upon further analysis, we find that these patterns can be attributed to cultural cues or shared societal characteristics. Our analysis includes low-resource languages, offering practical insights into when cross-lingual transfer can support them and where its limitations emerge. These findings provide guidance for designing more reliable and generalizable hate detection systems.

**Warning:** This paper contains examples of hate language that may disturb some readers. These examples are included for research on hate speech detection.

## 1 Introduction

Hate language detection aims to identify harmful expressions in online text, such as hate speech, harassment, and threats (Hoang et al., 2024). As digital communication expands globally, exposure to such content has increased, motivating efforts to build automated systems capable of supporting safer online communities. However, what counts as hate speech often depends on several criteria and even human annotators vary in how they perceive hate across contexts (Sap et al., 2022a). These differences pose a challenge for developing models that generalize reliably across languages.

Large Language Models (LLMs) offer strong multilingual representations and are increasingly

used in hate detection (Khondaker et al., 2023; Kumar et al., 2024; Abaskohi et al., 2024). However, the cross-lingual performance is not uniform, and the factors shaping transfer effectiveness should be studied (Meyer and Buys, 2024). This raises a practical and underexplored question: Can training data from one language support hate-speech detection in another language with comparable performance? Several factors could influence transfer success, including but not limited to dataset size, label definitions, temporal alignment, and linguistic similarity. Preliminary work suggests that broader social or regional contexts may also shape how hate expressions are formed and recognized (Bokaei et al., 2025; Pawar et al., 2024).

In this study, we conduct a large-scale analysis involving 14 languages from diverse regions and domains. We first build a comprehensive cross-lingual transfer matrix by training hate-detection classifiers in each language and evaluating them across all others to find how effective cross-lingual transfer learning is for hate-speech detection across languages. We then perform two complementary experiments: (1) integrating all languages into a single multilingual training set, (2) constructing a balanced multilingual dataset with equal representation per language. These setups allow us to examine how data volume and representational balance influence cross-lingual generalization.

The current study addresses the following research questions (RQs):

- RQ1. How does the performance of cross-lingual training compare to in-language training for hate-speech detection?
- RQ2. Would using training data from multiple languages improve performance compared to relying on a single source language?
- RQ3. What could be the potential reasons for the

081	success or failure of transfer learning in hate-	et al., 2019; Mandl et al., 2019; Mubarak et al.,	131
082	speech detection?	2022; Wiegand et al., 2018).	132
083	These questions aim to clarify whether cross-	<b>Cross-lingual Transfer Learning for Hate</b>	133
084	lingual transfer is effective for hate-speech detec-	<b>Language.</b> Transfer learning has been applied to	134
085	tion, and to identify the conditions under which it	a range of NLP tasks, including sentiment analy-	135
086	succeeds or fails, providing insights for developing	sis (Fsih et al., 2022; Husain et al., 2022), irony	136
087	more reliable and generalizable hate detection sys-	detection (Golazizian et al., 2020), machine trans-	137
088	tems. We use LLaMA 3 as our model, based on pre-	lation (Zoph et al., 2016; Kim et al., 2019; Ade-	138
089	liminary experiments indicating that it achieves the	bara and Abdul-Mageed, 2021), fake news detec-	139
090	best average performance across our 14 languages:	tion (Cruz et al., 2020), and offensive language	140
091	Turkish, Urdu, Persian, Arabic, Hindi, Bengali,	classification (Zhou et al., 2023). For multilingual	141
092	Indonesian, Malay, English, Spanish, French, Ger-	hate detection, prior work evaluates cross-lingual	142
093	man, Korean, and Chinese.	transfer from English to multiple low-resource lan-	143
094	<b>2 Related Work</b>	guages using contextual embeddings (Ranasinghe	144
095	<b>Hate Language Detection.</b> Early work on hate	and Zampieri, 2021), compares multilingual ver-	145
096	language detection focused primarily on English,	sus translation-based pipelines for Arabic–English	146
097	using Machine Learning and Deep Learning ap-	(El-Alami et al., 2022), and improves low-resource	147
098	proaches to classify hate speech (Asogwa et al.,	performance with nearest-neighbor retrieval (Ghor-	148
099	2022; Davidson et al., 2017; Mullah and Zainon,	banpour et al., 2025). A recent survey provides a	149
100	2021; Malik et al., 2024; Zimmerman et al., 2018;	comprehensive overview of cross-lingual transfer	150
101	Zhou et al., 2020; Roy et al., 2020; Zhang et al.,	approaches and datasets over different languages	151
102	2018; Bade et al., 2024; Aiyanyo et al., 2020; Cao	by analyzing 67 papers on cross-lingual offensive	152
103	et al., 2020; Risch et al., 2020; Wang et al., 2020;	language detection (Jiang and Zubiaga, 2024).	153
104	Pamungkas and Patti, 2019; Van Hee et al., 2015;	Research explicitly examining cultural or social	154
105	Guo and Gauch, 2024; Cano Basave et al., 2013).	influences in transfer is more limited. Zero-shot	155
106	Subsequent research expanded to a wider range	analyses show that language-specific taboo or col-	156
107	of languages, including Indonesian (Ibrohim and	loquial expressions often fail to transfer across	157
108	Budi, 2019), Danish (Sigurbergsson and Derczyn-	languages (Nozza, 2021). Studies comparing English,	158
109	ski, 2020), Arabic (Mubarak et al., 2021; Bensalem	Chinese, and Korean demonstrate culture-specific	159
110	et al., 2023), Korean (Jeong et al., 2022), Chinese	biases that affect transfer, with few-shot adapta-	160
111	(Deng et al., 2022), Greek (Pitenis et al., 2020),	tion offering improvements (Zhou et al., 2023). Multi-	161
112	Persian (Delbari et al., 2024), and Hindi (Gupta	modal work further highlights cross-cultural vari-	162
113	et al., 2022; Kapoor et al., 2019).	ation in hate-speech judgments (Bui et al., 2025).	163
114	Recent work highlights key challenges in multi-	Zhou et al. (2023) showed that Korean-to-Chinese	164
115	lingual hate detection, including comparisons of	transfer can outperform English-to-Chinese sug-	165
116	translation-based and multilingual pipelines (Bell	gest that contextual or social alignment may some-	166
117	et al., 2025), limitations in zero-shot transfer for	times support better generalization. Bokaei et al.	167
118	taboo expressions (Nozza, 2021), and the influ-	(2025) examines cross-lingual transfer between	168
119	ence of annotator identity and beliefs on hate judg-	Arabic, Indonesian, Persian, and English, report-	169
120	ments (Sap et al., 2022b). Multimodal datasets	ing that transfer into Persian tends to be stronger	170
121	such as Multi3Hate reveal substantial cross-cultural	from Arabic and Indonesian than from English, high-	171
122	variation in hate-speech perception (Bui et al.,	lighting that cross-lingual effectiveness can vary	172
123	2025), while surveys document strong geographic	substantially across source languages.	173
124	and cultural imbalances in existing hate-speech re-	Our study complements the literature by analy-	174
125	sources (Tonneau et al., 2024). Additional stud-	zing transfer behavior across 14 languages and	175
126	ies explore the feasibility of parallel multilingual	examining correlations between transfer effective-	176
127	hate-speech data via machine translation (Korre	ness and different factors.	177
128	et al., 2024). Shared workshops such as OffensE-	<b>3 Dataset</b>	178
129	val, HASOC, OSACT5, and GermEval continue to	We use publicly available datasets from peer-	179
130	support benchmarking across languages (Zampieri	reviewed publications and shared evaluation tasks	180
		to ensure data quality and comparability across	181

Language	Data Size	Train	Test	Train <sub>Hate</sub>	Test <sub>Hate</sub>	Source	Period	Ref.
Bengali (BN)	30000	9007	1002	4504	501	YT, FB	2017–2020	(Romim et al., 2020)
Urdu (UR)	10041	9000	1000	4500	300	X	–	(Rizwan et al., 2020)
Hindi (HI)	8192	4171	467	2222	241	X, FB, WA	–	(Bhardwaj et al., 2020)
Persian (FA)	7000	2060	230	1002	117	X	2020–2022	(Delbari et al., 2024)
Indonesian (IN)	13169	9005	1003	4502	502	X	≤ 2018	(Ibrohim and Budi, 2019)
Malay (MS)	5014	3403	378	1701	189	X	2022–2023	(Maity et al., 2023)
Arabic (AR)	15014	9005	1002	4503	502	X	2016–2018	(Albadi et al., 2018)
Turkish (TR)	52672	9010	1003	4505	502	X, IG	2016–2023	(Toraman et al., 2022)
Korean (KD)	9341	3065	300	1520	142	X	2018–2020	(Moon et al., 2020; Lee et al., 2022)
Chinese (ZH)	8065	4950	318	2470	151	Web	–	(Moosa and Najiba, 2022)
Spanish (ES)	29856	9000	1000	4731	625	X	2016–2023	(Tonneau et al., 2024)
German (DE)	58025	9008	1003	4503	502	X, FB	2016–2023	(Tonneau et al., 2024)
French (FR)	18072	7483	829	3739	413	X	2016–2023	(Tonneau et al., 2024)
English (EN)	360494	9000	1038	4500	502	FB, YT, RD, X	2016–2023	(Tonneau et al., 2024)

Table 1: Hate language detection datasets. Abbreviations of sources: FB = Facebook, YT = YouTube, X = Twitter, WA: WhatsApp, IG = Instagram, RD = Reddit.

languages. Our study includes fourteen languages drawn from diverse linguistic families and geographical regions. Most datasets originate from Twitter, though some languages include additional sources such as Facebook and YouTube. Their collection periods span roughly 2016–2023. Table 1 summarizes the dataset sizes, the training set sizes, the number of hate samples used for fine-tuning, and their sources and collection periods.

All datasets exhibit class imbalance, with varying proportions of hate versus non-hate content. To conduct controlled comparisons across languages, we standardize the amount of training data wherever possible. Specifically, among the languages with sufficiently large datasets, we identify the smallest available number of hate examples (approximately 4,500) and sample an equal number of non-hate instances, producing balanced datasets of roughly 10,000 examples per language. We use an 80/10/10 split for train, validation, and test sets. Balancing all training sets to a similar size helps isolate the effect of language rather than dataset volume, and avoids situations where high-resource languages (such as English or German) would appear to transfer better simply because they have much larger datasets. It also keeps the computational requirements reasonable for training fourteen models. All balancing steps were applied only when the underlying dataset contained enough annotated hate examples to support this procedure. However, this was not the case for all languages. For Hindi, Malay, Korean, Chinese, and Persian, the available datasets contain fewer hate instances. To enable consistent evaluation across datasets with different annotation schemes (e.g., hostile, hate, offensive, abusive), we focus on hate-speech samples. Most datasets primarily contain explicit hate

speech; however, the Hindi and Korean datasets include a proportion of implicit instances.

## 4 Experimental Setup

We investigate cross-lingual transfer in hate detection using the datasets summarized in Table 1. Our experimental design includes (1) single-language fine-tuning with cross-lingual evaluation, (2) zero-shot baselines, and (3) two complementary multilingual fine tuning setups to investigate alternative explanations for transfer behavior.

**Model Selection and Baselines:** Before selecting a primary model to be fine-tuned, we initially experimented multiple baselines: 1) Zero-shot baselines, including GPT-4o, Gemma-3, and LLaMA-3-Instruct, and 2) Fine-tuning baselines, including Gemma-3 and LLaMA-3-8B. Fine-tuned baselines using both Gemma and LLaMA outperformed all zero-shot baselines across all the 14 languages. Further, fine-tuning LLaMA-3-8B yields higher average macro-F1 compared to Gemma 3, which aligns with recent work (Guo and Sarker, 2025; Bokaei et al., 2025) highlighting LLaMA-3’s multilingual strength. Thus, we adopt LLaMA-3-8B for all our upcoming fine-tuning experiments. Table 3 in the Appendix reports the baseline results.

**Fine-tuning Configuration:** All fine-tuning uses the official LLaMA-3-8B pretrained checkpoint. We apply AdamW with a learning rate of  $1e-5$ , batch size 8, and a maximum sequence length of 512 tokens. Training proceeds for up to 10 epochs with early stopping (patience = 3) based on validation F1 score. A fixed random seed (42) ensures reproducibility. All experiments are conducted on a single NVIDIA Quadro RTX 8000 (48 GB). Performance is reported using macro-F1.

**Single-Language Fine-Tuning:** For each lan-

guage, we fine-tune a separate LLaMA-3-8B model on its own training set and evaluate it on all 14 test sets. For example, a model fine-tuned on English is evaluated on English and on the remaining 13 languages. This procedure produces a  $14 \times 14$  cross-lingual transfer matrix, allowing us to examine empirical transfer patterns. Because LLaMA is pretrained as a multilingual model, fine-tuning on one language does not remove its multilingual knowledge; instead, it specializes the model toward the target language’s hate characteristics.

The  $14 \times 14$  single-language fine-tuning experiment directly addresses RQ1 and reveals if any language pairs transfer well or poorly. To answer RQ2, we design two complementary experiments. The first examines whether increased data volume and multilingual diversity improve performance. We refer to this experiment as Global Integration (Imbalanced). The other evaluates whether ensuring equal representation across languages affects transfer behavior by removing dataset-size differences as a confounding variable. We refer to this experiment as Global Integration (balanced). Together, these additional setups allow us to investigate whether transfer effectiveness stems from data size, multilingual exposure, or any structural similarities observed in the main  $14 \times 14$  matrix.

**Global Integration (Imbalanced) Experiment:** We merge the full training sets from all 14 languages into one multilingual dataset and fine-tune a single LLaMA-3-8B model. The model is then evaluated separately on each language’s test set. This experiment tests whether greater data volume and diversity produce better generalization.

**Global Integration (Balanced) Experiment:** To control for dataset-size differences, we sample an equal number of hate ( $n = 1000$ ) and non-hate ( $n = 1000$ ) instances from each language and merge them into a balanced multilingual dataset. A single model is fine-tuned and evaluated on each language’s test set, isolating the effect of balanced representation by controlling for dataset size.

In addition, from the main  $14 \times 14$  transfer matrix, we identify groups of languages that show strong mutual transfer from the initial experiment. For each group, we merge the training sets and fine-tune a single model, which is then evaluated separately on each language within the group. We refer to this setting as the Transfer Groups experiment. In all conditions, we use the same train, validation, and test splits for each language to ensure comparability across models.

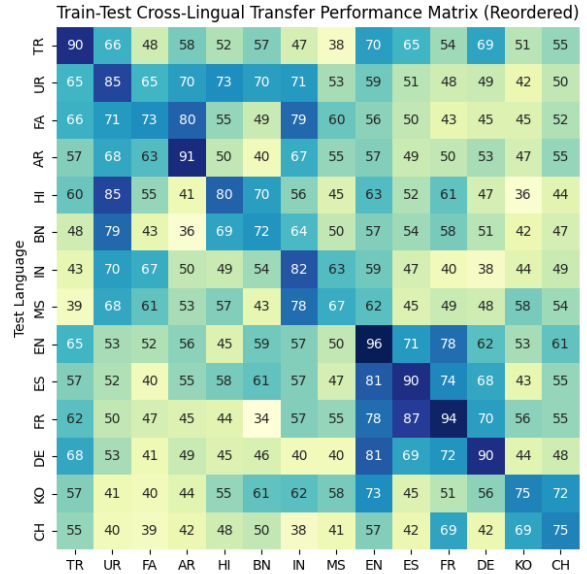


Figure 1: Cross-lingual F1 scores for hate detection across 14 languages. Language abbreviations are given in Table 1.

## 5 Results

Figure 1 presents the  $14 \times 14$  cross-lingual transfer matrix for hate content detection. Each cell reports the F1 score obtained when a model is fine-tuned on the language shown on the x-axis and evaluated on the test language on the y-axis. The matrix contains 196 train–test combinations, providing a comprehensive view of both in-language and cross-lingual performance.

Across languages, the highest F1 scores mostly occur on the diagonal, where the model is fine-tuned and tested on the same language. These in-language results are consistently strong: Turkish (90), Arabic (91), English (96), French (94), Spanish (90), and German (90) achieve F1 scores at or above 90. In a few cases, cross-lingual models outperform the in-language baseline; for example, models fine-tuned on Arabic or Indonesian yield higher scores for Persian than models fine-tuned on Persian itself.

Off-diagonal values show substantial variation in the quality of cross-lingual transfer. Based on these off-diagonal patterns, several cluster of languages emerge in which transfer performance is consistently higher within the cluster than across the cluster. These include: (1) Persian–Arabic–Urdu–Indonesian, (2) Urdu–Hindi–Bengali, (3) Indonesian–Malay, (4) Korean–Chinese, and (5) French–Spanish–English–German. As Figure 1 represents, models fine-tuned on Urdu perform

strongly on Hindi (85), and English-trained models yield relatively high scores on French (78), Spanish (74), and German (69). There are also cases of strong mutual transfer, such as *French* ↔ *English* (78) and *Bengali* ↔ *Urdu* (BN→UR: 70; UR→BN: 79). Lower scores appear for languages like Bengali, whose models perform less well on European test sets (BN→FR: 51; BN→DE: 42), and Chinese models show weaker transfer to several languages (CH→ES: 42; CH→AR: 38). Languages such as Korean, Chinese, and Malay generally show lower transfer performance both as train and test languages.

Table 2 summarizes results from the three complementary experiments. In the Global-Imbalanced setup, combining the full training data from all languages leads to performance gains for several high-resource languages, including English, French, and Spanish (average gain: +2.86), as well as improvements for languages such as Urdu. In contrast, most lower-resource languages—such as Bengali, Persian, and Malay—show declines relative to their single-language baselines (average: -5.43).

In the cluster Transfer Groups experiment, fine tuning in language-selective groups from the 14×14 matrix yields consistent performance improvements for nearly all group members (average gain: +5.62). Gains are observed in groups such as Hindi–Bengali–Urdu, Arabic–Persian–Indonesian, and English–French–Spanish–German, though the magnitude of improvement varies across languages. In the Global-Balanced setup, equal sampling from all languages reduces the effect of data imbalance. High-resource languages still improve over their single-language baselines, though gains are smaller than in the imbalanced condition (average: +1.57). For lower-resource languages, the balanced dataset mitigates some of the declines seen in the Global-Imbalanced experiment, with most showing small gains and some remaining close to baseline.

After observing the results and identifying groups of languages that show strong mutual transfer in the 14×14 transfer matrix, we merge the training sets within each group and fine-tune a single model. We then evaluate this model separately on each language in the group (Transfer Group Experiments). This setup tests whether transfer-based language grouping improves performance over single-language fine-tuning. We observe consistent gains across most group members (average +5.62), including Hindi–Bengali–Urdu, Arabic–Persian–Indonesian,

and English–French–Spanish–German, with varying magnitudes.

## 6 Analysis

The 14×14 transfer matrix in Figure 1 shows several structural patterns. We have already mentioned the five clusters of languages with relatively strong mutual transfer. While most languages achieve their highest performance through in-language learning, we also observe groups of language pairs where transfer learning produces better—or at least comparable—results. Table 4 in the Appendix presents some hate samples correctly predicted only via transfer learning. We first examine whether dataset size or temporal overlap can account for these patterns.

To understand whether dataset-related factors can account for the transfer patterns observed in Figure 1, we evaluate each cluster in terms of three properties: dataset size, resource level, and temporal overlap.

**Dataset Size:** Across the clusters, dataset size shows partial but inconsistent alignment. The French–Spanish–English–German group comprises some of the largest datasets (7k–9k hate samples), which could make within-group transfer more stable. The Urdu–Hindi–Bengali cluster also falls within a relatively similar medium range (3k–5k hate samples), which may support some internal consistency. In contrast, the Persian–Arabic–Urdu–Indonesian cluster includes languages with markedly different dataset sizes: Persian (~1k hate) is much smaller than Arabic or Indonesian (~4.5k hate), yet it clusters with them despite this imbalance. Similarly, Indonesian–Malay shows large variation, as Malay contains substantially fewer hate examples. The Korean–Chinese group also exhibits different dataset scales, with Chinese containing roughly twice as many hate examples as Korean. Taken together, dataset size sometimes aligns with cluster boundaries (notably in English- German- French - Spanish and South Asian groups), but not in clusters involving Persian, Malay, Korean, and Chinese, indicating that dataset size alone cannot explain the overall structure.

**Dataset Resource:** Many clusters, such as Urdu–Hindi–Bengali, Persian–Arabic–Urdu–Indonesian, and Indonesian–Malay, consist entirely of Twitter-based datasets. However, this pattern does not distinguish them meaningfully, because almost all languages in the study—including those that do not cluster

	TR	UR	FA	AR	HI	BN	IN	MS	EN	ES	FR	DE	KO	ZH
Solo-Source TL	90	85	80	91	85	79	82	78	96	90	94	90	<b>75</b>	75
Global-Imb	<b>92</b>	88	64	87	73	66	84	62	<b>98</b>	95	<b>96</b>	94	71	72
	[+2]	[+3]	[-16]	[-4]	[-12]	[-13]	[+2]	[-16]	[+2]	[+5]	[+2]	[+4]	[-4]	[-3]
Global-bal	91	87	67	89	81	69	83	63	96	92	94	93	73	76
	[+1]	[+2]	[-13]	[-2]	[-4]	[-10]	[+1]	[-15]	[0]	[+2]	[0]	[+3]	[-2]	[+1]
Observed Cluster	90	<b>90</b>	<b>83</b>	<b>93</b>	<b>90</b>	<b>83</b>	<b>85</b>	<b>80</b>	97	<b>96</b>	<b>96</b>	<b>95</b>	<b>75</b>	<b>78</b>
	[0]	[+5]	[+3]	[+2]	[+5]	[+4]	[+3]	[+2]	[+1]	[+6]	[+2]	[+5]	[0]	[+2]

Table 2: Performance comparison across different fine tuning strategies. Bracketed numbers indicate improvement or decline relative to Solo-Source TL.

together—are also drawn primarily from Twitter. Conversely, the Korean–Chinese cluster does not share a unified source either: Korean is drawn from Twitter, while Chinese originates from general web data. Thus, data source alignment does not systematically track the cluster structure.

**Temporal Overlap:** A similar mixed pattern appears when comparing dataset collection periods: English- German- French - Spanish- share nearly identical time spans (2016–2023), which could plausibly support more consistent transfer. The Urdu–Hindi–Bengali datasets likewise cover broadly similar years. However, the Persian–Arabic–Urdu–Indonesian cluster shows substantial variation: Persian (2020–22) does not overlap with Arabic (2016–18) or Indonesian ( $\leq 2018$ ), yet all appear in the same cluster. In addition, the Indonesian–Malay pair has almost no temporal overlap ( $\leq 2018$  vs. 2022–23), yet they form a clear transfer cluster. The analysis for Korean–Chinese is not possible since the Chinese time window is unspecified. These comparisons show that temporal alignment may support transfer in some clusters (English - German - French - Spanish and South Asian) but fails to predict clusters in others.

**Qualitative Analysis:** To better understand when transferability occurs, we performed a cross-lingual qualitative error analysis, inspecting examples that were correctly classified only after applying transfer learning and comparing them with those correctly predicted by the in-language model.

Further analysis of the cluster comprising Arabic–Indonesian–Urdu–Persian reveals that, Arabic and Urdu tend to support the detection of hate related to religious and political topics, especially sociopolitical hate that is prevalent in the Middle Eastern context. The Arabic dataset mainly provided relevant contextual cues that aligned well with Persian discourse in these domains. Urdu also

contributes effectively to the detection of political hate, often to a greater extent compared to Arabic. Indonesian, on the other hand, contributes more strongly to detecting role-based hate (e.g., targeted specific profession) and performs well on longer tweets that mix neutral and hate sentiments. Table 4 illustrates Persian examples that are correctly predicted only after applying transfer learning, alongside those captured only by the Persian- fine tuned model.

We also observe a strong transfer between Indonesian and Malay. Upon further examination we observed, both datasets include comparable expressions of political frustration, religious commentary, and identity-based accusations. Religious references are frequently used to judge or shame others, including accusations related to religious observance. Some metaphorical and idiomatic forms—such as insults referencing entities or moral faults—appear in both languages, allowing Indonesian fine-tuned models to generalize well to Malay. We also observed that both languages sometimes employ similar proverbs, insults, and metaphorical language, particularly the use of animals like "dog" and "devil" as slurs in hate instances. Social issues such as traffic caused by bazaars, public frustration over commodity scarcity, and emotional sharing on social media are likewise mirrored in both contexts.

However, the reverse direction is weaker. One possible reason is the limited amount of hate data in the Malay dataset, as discussed, which restricts the model’s ability to generalize when used as a source language for transfer learning. A similar pattern is observed with Persian, as its limited hate training data reduces its effectiveness as a target language.

Another cluster where transfer learning outperforms in-language learning includes Urdu, Hindi, and Bengali. Models fine-tuned on Urdu often

perform better on Bengali and Hindi than the respective in-language models. Further analysis reveals that transfer from Urdu is particularly effective for highly polarized, overt hate expressions involving explicit slurs, religious antagonism, or politicized attacks. In contrast, in-language models more accurately classify tweets containing idiomatic expressions or contextually grounded criticism. For instance, socially embedded sarcasm in Bengali was frequently missed by the Urdu-fine-tuned model but correctly identified by the Bengali-fine-tuned counterpart. In this setting, a key factor is that the Urdu dataset contains a high density of explicit hate markers—religious slurs, profane language, and direct references to identifiable targets. These overt expressions show lexical consistency and limited idiomatic variation, enabling the model to learn discriminative hate patterns reliably. Urdu tweets also frequently specify hate targets, making intent clearly detectable. By comparison, Bengali and Hindi datasets contain more ambiguous, sarcastic, or culturally embedded forms of expression; specifically, the Hindi dataset includes implicit hate, which is harder to generalize across languages.

We also observe that although topics such as political or religious hate appear across different datasets, the strength of transfer does not seem to depend on topical similarity alone. For example, although English and French contain religious or political hate, transfer learning from Arabic or Urdu yields larger gains when detecting religious or political hate in Persian. This pattern also appears when transferring from Arabic to English or French.

To broaden our qualitative examination, we also inspected examples from language pairs that do not fall within the same empirical cluster, focusing on how transfer behaves in settings where alignment is weaker. Our results show that even when the training and test languages are not within the same cluster (e.g., Spanish–Chinese, Persian–French), transfer learning can still detect certain types of hate speech, though to a lesser extent. Specifically, explicit, emotionally extreme, or lexically offensive expressions tend to transfer well regardless of which empirical cluster the language belongs to. For example, in our experiments, the Spanish→Chinese transfer achieved moderate performance in detecting overtly hate Chinese instances such as: “可以证明你是个没脑子的田园女权婊了”, (You’re a brainless feminist bitch) or

“微博女权婊人数还是相当多的!” (So many feminist bitches on Weibo!) These phrases involve direct insults and highly polarized language, patterns that models fine-tuned on Spanish data can still recognize due to the frequent recurrence of explicit misogynistic insults across languages. However, the Spanish-trained model struggled with more context-dependent expressions, such as: “唐明皇一世英名就毁在杨家人手上了” (Emperor Tang Minghuang’s lifetime reputation was ruined at the hands of the Yang family - A historical reference used to delegitimize women), or “我还看到一个百事256的说自己是女权的, 最近在力挺川美侵台军宅, 和粉红互动得可欢了呢。” (I also saw someone called ‘Pepsi256’ who claims to be a feminist. Lately, she’s been enthusiastically supporting the Sichuan Fine Arts Institute’s pro-invasion-of-Taiwan military otakus and happily interacting with the nationalists (a.k.a. little pinks). In contrast, training the model on Korean resulted in significantly better performance, not only on overtly hate expressions but also on more subtle, context-dependent instances like those discussed above. In other instances, we can note experiments with French and Persian, fine-tuned on Persian, were still able to correctly identify highly explicit hate content in French. For example: “Vous pouvez crever sale race” (You can die, filthy race) This phrase combines an aggressive death wish with a racial slur, making it both emotionally extreme and lexically offensive, enabling the Persian-fine-tuned model to detect it despite the different cluster. Therefore, our result shows that even across dissimilar clusters, transfer learning can succeed when the hate content is lexically explicit or emotionally extreme. However, it consistently fails on those instances that rely on background context, such as social and historical references.

As discussed, we have three complementary experiments. Global-Imb experiment demonstrates that integrating data from all languages in an imbalanced manner primarily benefits high-resource languages, taking advantage of enhanced performance for high-resource languages, it weakens specific cues for low-resource ones, as their contextually grounded may be overshadowed by dominant high-resource patterns. The Cluster experiment shows that grouping languages based on their similar cluster not only mitigates this wash-out effect but also produces the most consistent and substantial gains, especially for low-resource languages. Lastly, Global-Bal experiment, by bal-

ancing representation across languages, offers a middle ground—reducing the disadvantage faced by low-resource languages while slightly decrease the gains of high-resource ones.

## 7 Discussion

In addressing RQ1—whether cross-lingual training can match in-language hate speech detection—our results show that cross-lingual models can, in several cases, achieve performance comparable to or even exceeding in-language fine-tuning. Across the 14 languages studied, this effect is most evident within the observed clusters, where languages such as Persian, Indonesian, Urdu, Malay, Hindi, and Bengali benefit substantially from training on related languages.

Regarding RQ2—whether multilingual training improves performance beyond single-source transfer—our findings show that only specific multilingual setups yield consistent gains. Globally imbalanced training primarily benefits high-resource languages while obscuring cues needed for low-resource ones. In contrast, cluster-based training provides the most reliable improvements across languages.

RQ3 concerns potential reasons underlying successful or failed transfer. In this study we first examined two readily measurable properties—dataset size and temporal overlap—neither of which suffices to explain the structure of transfer patterns. A qualitative analysis explores whether conceptual similarities in hate discourse could shed light on transferability. Our analysis shows that, within several clusters, the transferable expressions are not random: they often reflect recurring cultural cues and communicative norms shared across languages. For instance, in the Arabic–Urdu–Indonesian–Persian setting, transferred examples frequently revolve around similar forms of political hostility, religiously framed judgments, or role-based accusations—forms of hate that carry comparable pragmatic functions across these societies. In the Indonesian–Malay pair, transferable cases similarly reflect shared idiomatic insults, moral accusations, and socially embedded metaphors. In the Urdu–Hindi–Bengali cluster, transfer succeeds when expressions rely on explicit slurs or polarized antagonism, echoing shared sociopolitical tensions in the region. This suggests that some of the empirical clusters can reflect overlapping cultural cues or communicative patterns rather than coincidental lexical overlap (Zhou et al.,

2023; Bokaei et al., 2025). Our results suggest that cross-lingual transfer can match or exceed in-language performance (RQ1), but doing so reliably requires selecting training languages that share meaningful discourse patterns (RQ2). As for why such clusters exist (RQ3), our evidence points to a combination of measurable dataset factors and deeper commonalities in social communication; however, fully isolating these influences remains an open direction for future research. Larger multilingual ethnographic corpora, culturally informed annotation studies, and controlled manipulations of social context are required to further validate these hypotheses.

## Conclusion

This study presented a comprehensive analysis of cross-lingual transfer learning for hate-speech detection across fourteen languages. While in-language fine-tuning remains the strongest baseline overall, our results show that cross-lingual models can achieve comparable—or even superior—performance in several settings, particularly within clusters of languages that exhibit recurring similarities in hate discourse. Multilingual training proves most effective when it leverages these empirically derived clusters, whereas global imbalanced integration offers limited benefits for low-resource languages. Our analysis indicates that dataset size, resource and temporal overlap alone cannot account for the structure of transfer patterns. Instead, the qualitative evidence suggests that transferability often arises from shared communicative norms, discourse practices, or sociopolitical framing that shape how hate is expressed within specific groups of languages. At the same time, explicit or highly polarized hate expressions transfer broadly even across distant languages, whereas culturally embedded or context-dependent forms remain challenging. Overall, the findings highlight both the promise and limitations of cross-lingual transfer for hate detection. They suggest that effective multilingual systems should incorporate the broader social dynamics that influence how online hate is articulated. Future work can build on these insights by developing culturally informed modeling strategies, expanding the range of low-resource languages, and exploring more direct ways of capturing the social context underlying hate language.

## 722 Limitations

723 This study has several limitations. First, the  
724 datasets used across the fourteen languages vary in  
725 source platforms, annotation practices, and topical  
726 distributions. Although we standardize dataset size  
727 where possible and map labels to a binary hate/non-  
728 hate scheme, differences in domain, explicitness,  
729 and annotation criteria may still influence transfer  
730 behavior. In particular, implicit hate appears only  
731 in a subset of languages, limiting our ability to  
732 evaluate how implicit expressions transfer across  
733 languages. Second, while our analysis considers  
734 dataset size and collection period, we do not control  
735 for all factors that may shape transfer performance,  
736 such as differences in sampling strategies or de-  
737 mographic characteristics of annotators. As such,  
738 the empirical clusters we identify reflect observed  
739 transfer patterns.

740 Third, the thematic observations drawn from our  
741 qualitative error analysis rely on manual inspec-  
742 tion of representative examples. Although these  
743 analyses provide useful insights, they may not cap-  
744 ture the full variability of hate discourse across  
745 languages. A more extensive human evalua-  
746 tion—including multiple annotators would offer  
747 a stronger foundation for interpreting the nature of  
748 transferable expressions. Finally, although we dis-  
749 cuss broad social or discourse-related similarities  
750 as one possible contributor to transfer effectiveness,  
751 we do not measure such factors directly. We rely on  
752 transfer patterns themselves as empirical signals,  
753 and therefore cannot isolate the influence of lin-  
754 guistic, social, or dataset-driven properties. Future  
755 work would be needed to validate these interpreta-  
756 tions more systematically.

## 757 Ethics Statement

758 This study analyzes publicly available hate  
759 datasets and does not involve collecting new user  
760 data. All datasets were obtained from prior peer-  
761 reviewed work or shared tasks that follow estab-  
762 lished ethical guidelines. Because hate datasets  
763 may contain harmful or offensive language, exam-  
764 ples shown in this paper are minimized and pre-  
765 sented only when necessary for scientific trans-  
766 parency. No personally identifiable information is  
767 included in our datasets or model outputs. All ex-  
768 periments were conducted using anonymized text.  
769 Models trained in this work are not intended for  
770 deployment without further evaluation, fairness au-  
771 diting, and context-specific calibration.

## References

- 772 Amirhossein Abaskohi, Sara Baruni, Mostafa Masoudi,  
773 Nesa Abbasi, Mohammad Hadi Babalou, Ali Edalat,  
774 Sepehr Kamahi, Samin Mahdizadeh Sani, Nikoo  
775 Naghavian, Danial Namazifard, Pouya Sadeghi, and  
776 Yadollah Yaghoobzadeh. 2024. [Benchmarking large  
777 language models for Persian: A preliminary study fo-  
778 cusing on ChatGPT](#). In *Proceedings of the 2024 Joint  
779 International Conference on Computational Linguis-  
780 tics, Language Resources and Evaluation (LREC-  
781 COLING 2024)*, pages 2189–2203, Torino, Italia.  
782 ELRA and ICCL. 783
- Ife Adebara and Muhammad Abdul-Mageed. 2021. [Im-  
784 proving similar language translation with transfer  
785 learning](#). In *Proceedings of the Sixth Conference on  
786 Machine Translation*, pages 273–278, Online. Asso-  
787 ciation for Computational Linguistics. 788
- Imatitukia D Aiyanyo, Hamman Samuel, and Heuseok  
789 Lim. 2020. A systematic review of defensive and of-  
790 fensive cybersecurity with machine learning. *Applied  
791 Sciences*, 10(17):5811. 792
- Nuha Albadi, Maram Kurdi, and Shivakant Mishra.  
793 2018. [Are they our brothers? analysis and detec-  
794 tion of religious hate speech in the arabic twitter-  
795 sphere](#). In *2018 IEEE/ACM International Confer-  
796 ence on Advances in Social Networks Analysis and  
797 Mining (ASONAM)*, pages 69–76. 798
- Doris Chinedu Asogwa, Chiamaka Ijeoma Chukwuneke,  
799 CC Ngene, and GN Anigbogu. 2022. Hate speech  
800 classification using SVM and Naive Bayes. *arXiv  
801 preprint arXiv:2204.07057*. 802
- Girma Bade, Olga Kolesnikova, Grigori Sidorov, and  
803 José Oropeza. 2024. [Social media hate and offen-  
804 sive speech detection using machine learning method](#).  
805 In *Proceedings of the Fourth Workshop on Speech,  
806 Vision, and Language Technologies for Dravidian  
807 Languages*, pages 240–244, St. Julian’s, Malta. Asso-  
808 ciation for Computational Linguistics. 809
- Samuel Bell, Eduardo Sánchez, David Dale, Pontus  
810 Stenertorp, Mikel Artetxe, and Marta R Costa-jussà.  
811 2025. Translate, then detect: Leveraging machine  
812 translation for cross-lingual toxicity classification.  
813 In *Proceedings of the Tenth Conference on Machine  
814 Translation*, pages 253–268. 815
- Imene Bensalem, Meryem Mout, and Paolo Rosso. 2023.  
816 Offensive language detection in Arabizi. In *Proceed-  
817 ings of ArabicNLP 2023*, pages 423–434. 818
- Mohit Bhardwaj, Md. Shad Akhtar, Asif Ekbal, Ami-  
819 tava Das, and Tanmoy Chakraborty. 2020. [Hostility  
820 detection dataset in hindi](#). *CoRR*, abs/2011.03588. 821
- Zahra Bokaei, Walid Magdy, and Bonnie Webber. 2025.  
822 [Culture matters in toxic language detection in Persian](#).  
823 In *Proceedings of the 63rd Annual Meeting of the  
824 Association for Computational Linguistics (Volume  
825 1: Long Papers)*, pages 9290–9304, Vienna, Austria.  
826 Association for Computational Linguistics. 827

828	Minh Duc Bui, Katharina Von Der Wense, and Anne	Faeze Ghorbanpour, Daryna Dementieva, and Alexan-	884
829	Lauscher. 2025. Multi <sup>3</sup> hate: Multimodal, multilin-	der Fraser. 2025. <a href="#">Data-efficient hate speech detec-</a>	885
830	gual, and multicultural hate speech detection with	<a href="#">tion via cross-lingual nearest neighbor retrieval with</a>	886
831	vision–language models. In <i>Proceedings of the 2025</i>	<a href="#">limited labeled data</a> . In <i>Proceedings of the 2025 Con-</i>	887
832	<i>Conference of the Nations of the Americas Chap-</i>	<i>ference on Empirical Methods in Natural Language</i>	888
833	<i>ter of the Association for Computational Linguistics:</i>	<i>Processing</i> , pages 29662–29680, Suzhou, China. As-	889
834	<i>Human Language Technologies (Volume 1: Long Pa-</i>	ssociation for Computational Linguistics.	890
835	<i>pers)</i> , pages 9714–9731.		
836	Amparo Elizabeth Cano Basave, Yulan He, Kang Liu,	Preni Golazizian, Behnam Sabeti, Seyed Arad	891
837	and Jun Zhao. 2013. <a href="#">A weakly supervised Bayesian</a>	Ashrafi Asli, Zahra Majdabadi, Omid Momenzadeh,	892
838	<a href="#">model for violence detection in social media</a> . In <i>Pro-</i>	and Reza Fahmi. 2020. <a href="#">Irony detection in Persian</a>	893
839	<i>ceedings of the Sixth International Joint Conference</i>	<a href="#">language: A transfer learning approach using emoji</a>	894
840	<i>on Natural Language Processing</i> , pages 109–117,	<a href="#">prediction</a> . In <i>Proceedings of the Twelfth Language</i>	895
841	Nagoya, Japan. Asian Federation of Natural Lan-	<i>Resources and Evaluation Conference</i> , pages 2839–	896
842	guage Processing.	2845, Marseille, France. European Language Re-	897
843	Rui Cao, Roy Ka-Wei Lee, and Tuan-Anh Hoang. 2020.	sources Association.	898
844	Deep hate: Hate speech detection via multi-faceted	Xiaoyu Guo and Susan Gauch. 2024. <a href="#">Using sarcasm</a>	899
845	text representations. In <i>Proceedings of the 12th ACM</i>	<a href="#">to improve cyberbullying detection</a> . In <i>Proceedings</i>	900
846	<i>Conference on Web Science</i> , pages 11–20.	<a href="#">of the Fourth Workshop on Threat, Aggression &amp;</a>	901
847	Jan Christian Blaise Cruz, Julianne Agatha Tan, and	<a href="#">Cyberbullying @ LREC-COLING-2024</a> , pages 52–	902
848	Charibeth Cheng. 2020. <a href="#">Localization of fake news</a>	59, Torino, Italia. ELRA and ICCL.	903
849	<a href="#">detection via multitask transfer learning</a> . In <i>Proce-</i>	Yuting Guo and Abeed Sarker. 2025. <a href="#">Benchmark-</a>	904
850	<i>edings of the Twelfth Language Resources and Evalua-</i>	<a href="#">ing open-source large language models on health-</a>	905
851	<i>tion Conference</i> , pages 2596–2604, Marseille, France.	<a href="#">care text classification tasks</a> . <i>arXiv preprint</i>	906
852	European Language Resources Association.	<i>arXiv:2503.15169</i> .	907
853	Thomas Davidson, Dana Warmsley, Michael Macy, and	Vikram Gupta, Sumegh Roychowdhury, Mithun Das,	908
854	Ingmar Weber. 2017. Automated hate speech de-	Somnath Banerjee, Punyajoy Saha, Binny Mathew,	909
855	tection and the problem of offensive language. In	Animesh Mukherjee, et al. 2022. Multilingual abu-	910
856	<i>Proceedings of the international AAAI conference on</i>	<a href="#">sive comment detection at scale for Indic languages</a> .	911
857	<i>web and social media</i> , volume 11, pages 512–515.	<i>Advances in Neural Information Processing Systems</i> ,	912
858	Zahra Delbari, Nafise Sadat Moosavi, and Moham-	35:26176–26191.	913
859	mad Taher Pilehvar. 2024. Spanning the spectrum of	Nhat Hoang, Xuan Long Do, Duc Anh Do, Duc Anh Vu,	914
860	hatred detection: a Persian multi-label hate speech	and Anh Tuan Luu. 2024. <a href="#">ToXCL: A unified frame-</a>	915
861	dataset with annotator rationales. In <i>Proceedings of</i>	<a href="#">work for toxic speech detection and explanation</a> . In	916
862	<i>the AAAI Conference on Artificial Intelligence</i> , vol-	<i>Proceedings of the 2024 Conference of the North</i>	917
863	ume 38, pages 17889–17897.	<i>American Chapter of the Association for Computa-</i>	918
864	Jiawen Deng, Jingyan Zhou, Hao Sun, Chujie Zheng,	<i>tional Linguistics: Human Language Technologies</i>	919
865	Fei Mi, Helen Meng, and Minlie Huang. 2022.	<i>(Volume 1: Long Papers)</i> , pages 6460–6472, Mexico	920
866	<a href="#">COLD: A benchmark for Chinese offensive language</a>	City, Mexico. Association for Computational Lin-	921
867	<a href="#">detection</a> . In <i>Proceedings of the 2022 Conference</i>	guistics.	922
868	<i>on Empirical Methods in Natural Language Process-</i>	Fatemah Husain, Hana Al-Ostad, and Halima Omar.	923
869	<i>ing</i> , pages 11580–11599, Abu Dhabi, United Arab	2022. <a href="#">A weak supervised transfer learning approach</a>	924
870	Emirates. Association for Computational Linguistics.	<a href="#">for sentiment analysis to the Kuwaiti dialect</a> . In	925
871	Fatima-zahra El-Alami, Said Ouatik El Alaoui, and	<i>Proceedings of the Seventh Arabic Natural Language</i>	926
872	Noureddine En Nahnahi. 2022. A multilingual of-	<i>Processing Workshop (WANLP)</i> , pages 161–173, Abu	927
873	fensive language detection method based on transfer	Dhabi, United Arab Emirates (Hybrid). Association	928
874	learning from transformer fine-tuning model. <i>Journal</i>	for Computational Linguistics.	929
875	<i>of King Saud University-Computer and Informa-</i>	Muhammad Okky Ibrohim and Indra Budi. 2019. <a href="#">Multi-</a>	930
876	<i>tion Sciences</i> , 34(8):6048–6056.	<a href="#">label hate speech and abusive language detection</a>	931
877	Emna Fsih, Sameh Kchaou, Rahma Boujelbane, and	<a href="#">in Indonesian Twitter</a> . In <i>Proceedings of the Third</i>	932
878	Lamia Hadrich-Belguith. 2022. <a href="#">Benchmarking trans-</a>	<i>Workshop on Abusive Language Online</i> , pages 46–	933
879	<a href="#">fer learning approaches for sentiment analysis of Ara-</a>	57, Florence, Italy. Association for Computational	934
880	<a href="#">bic dialect</a> . In <i>Proceedings of the Seventh Arabic</i>	Linguistics.	935
881	<i>Natural Language Processing Workshop (WANLP)</i> ,	Younghoon Jeong, Juhyun Oh, Jongwon Lee, Jaimeen	936
882	pages 431–435, Abu Dhabi, United Arab Emirates	Ahn, Jihyung Moon, Sungjoon Park, and Alice Oh.	937
883	(Hybrid). Association for Computational Linguistics.	2022. <a href="#">KOLD: Korean offensive language dataset</a> .	938
		In <i>Proceedings of the 2022 Conference on Empiri-</i>	939
		<i>cal Methods in Natural Language Processing</i> , pages	940

941	10818–10833, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	997
942		998
943	Aiqi Jiang and Arkaitz Zubiaga. 2024. Cross-lingual offensive language detection: A systematic review of datasets, transfer approaches and challenges. <i>arXiv preprint arXiv:2401.09244</i> .	999
944		1000
945		1001
946		1002
947	Raghav Kapoor, Yaman Kumar, Kshitij Rajput, Rajiv Ratn Shah, Ponnurangam Kumaraguru, and Roger Zimmermann. 2019. Mind your language: Abuse and offense detection for code-switched languages. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 33, pages 9951–9952.	1003
948		1004
949		1005
950		1006
951		1007
952		1008
953	Md Tawkat Islam Khondaker, Abdul Waheed, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2023. GPTAraEval: A comprehensive evaluation of ChatGPT on Arabic NLP. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 220–247, Singapore. Association for Computational Linguistics.	1009
954		1010
955		1011
956		1012
957		1013
958		1014
959		1015
960	Yunsu Kim, Petre Petrov, Pavel Petrushkov, Shahram Khadivi, and Hermann Ney. 2019. Pivot-based transfer learning for neural machine translation between non-English languages. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 866–876, Hong Kong, China. Association for Computational Linguistics.	1016
961		1017
962		1018
963		1019
964		1020
965		1021
966		1022
967		1023
968		1024
969	Katerina Korre, Arianna Muti, and Alberto Barrón-Cedeño. 2024. The challenges of creating a parallel multilingual hate speech corpus: An exploration. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 15842–15853, Torino, Italia. ELRA and ICCL.	1025
970		1026
971		1027
972		1028
973		1029
974		1030
975		1031
976	Ankit Kumar, Richa Sharma, and Punam Bedi. 2024. Towards optimal NLP solutions: Analyzing GPT and LLaMA-2 models across model scale, dataset size, and task diversity. <i>Engineering, Technology &amp; Applied Science Research</i> , 14(3):14219–14224.	1032
977		1033
978		1034
979		1035
980		1036
981	Jean Lee, Taejun Lim, Heejun Lee, Bogeun Jo, Yangsok Kim, Heegeun Yoon, and Soyeon Caren Han. 2022. K-MHaS: A multi-label hate speech detection dataset in Korean online news comment. In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 3530–3538, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	1037
982		1038
983		1039
984		1040
985		1041
986		1042
987		1043
988		1044
989	Krishanu Maity, Shaubhik Bhattacharya, Sriparna Saha, and Manjeevan Seera. 2023. A deep learning framework for the detection of malay hate speech. <i>IEEE Access</i> , 11:79542–79552.	1045
990		1046
991		1047
992		1048
993	Jitendra Singh Malik, Hezhe Qiao, Guansong Pang, and Anton van den Hengel. 2024. Deep learning for hate speech detection: a comparative study. <i>International Journal of Data Science and Analytics</i> , pages 1–16.	1049
994		1050
995		1051
996		1052
	Thomas Mandl, Sandip Modha, Prasenjit Majumder, Daksh Patel, Mohana Dave, Chintak Mandlia, and Aditya Patel. 2019. Overview of the hasoc track at fire 2019: Hate speech and offensive content identification in Indo-European languages. In <i>Proceedings of the 11th Annual Meeting of the Forum for Information Retrieval Evaluation, FIRE '19</i> , page 14–17, New York, NY, USA. Association for Computing Machinery.	997
		998
		999
		1000
		1001
		1002
		1003
		1004
		1005
	Francois Meyer and Jan Buys. 2024. A systematic analysis of subwords and cross-lingual transfer in multilingual translation. In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 2194–2200, Mexico City, Mexico. Association for Computational Linguistics.	1006
		1007
		1008
		1009
		1010
		1011
	Jihyung Moon, Won Ik Cho, and Junbum Lee. 2020. BEEP! Korean corpus of online news comments for toxic speech detection. In <i>Proceedings of the Eighth International Workshop on Natural Language Processing for Social Media</i> , pages 25–31, Online. Association for Computational Linguistics.	1012
		1013
		1014
		1015
		1016
		1017
	Wajid Hassan Moosa and Najiba. 2022. Multilingual hate speech dataset. <a href="https://www.kaggle.com/datasets/wajidhassanmoosa/multilingual-hatespeech-dataset">https://www.kaggle.com/datasets/wajidhassanmoosa/multilingual-hatespeech-dataset</a> . Open Data Commons Attribution License (ODC-By) v1.0.	1018
		1019
		1020
		1021
		1022
	Hamdy Mubarak, Hend Al-Khalifa, and Abdulmohsen Al-Thubaity. 2022. Overview of OSACT5 shared task on Arabic offensive language and hate speech detection. In <i>Proceedings of the 5th Workshop on Open-Source Arabic Corpora and Processing Tools with Shared Tasks on Qur'an QA and Fine-Grained Hate Speech Detection</i> , pages 162–166, Marseille, France. European Language Resources Association.	1023
		1024
		1025
		1026
		1027
		1028
		1029
		1030
	Hamdy Mubarak, Ammar Rashed, Kareem Darwish, Younes Samih, and Ahmed Abdelali. 2021. Arabic offensive language on Twitter: Analysis and experiments. In <i>Proceedings of the Sixth Arabic Natural Language Processing Workshop</i> , pages 126–135, Kyiv, Ukraine (Virtual). Association for Computational Linguistics.	1031
		1032
		1033
		1034
		1035
		1036
		1037
	Nanlir Sallau Mullah and Wan Mohd Nazmee Wan Zainon. 2021. Advances in machine learning algorithms for hate speech detection in social media: a review. <i>IEEE Access</i> , 9:88364–88376.	1038
		1039
		1040
		1041
	Debora Nozza. 2021. Exposing the limits of zero-shot cross-lingual hate speech detection. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 907–914.	1042
		1043
		1044
		1045
		1046
		1047
	Endang Wahyu Pamungkas and Viviana Patti. 2019. Cross-domain and cross-lingual abusive language detection: A hybrid approach with deep learning and a multilingual lexicon. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational</i>	1048
		1049
		1050
		1051
		1052



- 1165 Yanling Zhou, Yanyan Yang, Han Liu, Xiufeng Liu,  
1166 and Nick Savage. 2020. [Deep learning based fusion](#)  
1167 [approach for hate speech detection](#). *IEEE Access*,  
1168 8:128923–128929.
- 1169 Steven Zimmerman, Udo Kruschwitz, and Chris Fox.  
1170 2018. Improving hate speech detection with deep  
1171 learning ensembles. In *Proceedings of the eleventh*  
1172 *international conference on language resources and*  
1173 *evaluation (LREC 2018)*.
- 1174 Barret Zoph, Deniz Yuret, Jonathan May, and Kevin  
1175 Knight. 2016. [Transfer learning for low-resource neu-](#)  
1176 [ral machine translation](#). In *Proceedings of the 2016*  
1177 *Conference on Empirical Methods in Natural Lan-*  
1178 *guage Processing*, pages 1568–1575, Austin, Texas.  
1179 Association for Computational Linguistics.

1180 **Appendix**

method	model	En	Ar	Fa	BN	Hi	ZH	Kr	SP	Fr	De	IN	MS	Tu	UR
Zero-shot	GPT - 4o	82	64	55	53	67	67	60	75	77	75	65	56	71	45
	Llama 3	82	72	45	58	70	66	49	79	76	74	64	51	70	52
	Gemma 3	85	59	58	54	56	52	43	67	69	70	40	51	46	38
Fine-tuning	Llama 3	<b>96</b>	<b>91</b>	73	<b>72</b>	<b>80</b>	75	75	<b>90</b>	<b>94</b>	90	<b>82</b>	67	<b>90</b>	<b>85</b>
	Gemma 3	94	90	<b>75</b>	70	75	<b>82</b>	<b>76</b>	82	80	<b>94</b>	75	<b>71</b>	88	<b>85</b>

Table 3: F1 scores of different models for multilingual hate-speech detection.

Lang	TL From	Tweet	Explanation
FA	IN	<p>رندی به محضر فقیهی رسیدو حرکات رقص را جداجدا انجام می داد و می پرسید آیا حرام است؟ فقیه میگفت نه. پس رند شروع به رقصیدن کرد. فقیه گفت تجزیه اش خوب بود ولی مرده شور ترکیبش رو بردن حالا حکایت این عدالت خوار است بعضیشون عیبی بچه های خوبین ولی مرده شور ترکیبشون رو برده من برم به کارای خودم برسم خدافظ ##انتخابات مجلس</p> <p>A trickster went to a cleric and performed dance moves separately, asking if they were forbidden. The cleric said no. Then the trickster started dancing, and the cleric said, 'Breaking it down was fine, but damn the combination!' This is exactly the case with these so-called justice-seekers—some of them are actually good kids, but damn their combination! Anyway, I'll get back to my own business. Goodbye. #JusticeSeekers #ParliamentElection</p>	Long tweet containing a mixture of neutral and toxic sentiments, with shared religious content.
FA	AR	<p>وقتی یه جوان بدبخت میخاد استخدام بشه ریشه و بیخ مسائل خانوادگیش رو درمیارن اما ز بیاکلام یا دختربازیش و شراب خواریش هنوز تو دانشگاه اسلامی!!! ندریس میکنه</p> <p>When a poor young man wants to get hired, they dig up every root of his family issues, but Zibakalam, with his womanizing and wine-drinking, is still teaching at an Islamic university!</p>	Religious hypocrisy framed with moral outrage, a sentiment common to both Persian and Arabic contexts.
MS	IN	<p>Oh dan kereta mmm suka menghimpit ke motor walaupun jelas2 ade moto kt situ smpi terpaksa mengelak n kereta tu buat bodoh je. Ramai ye perangai mcm ni.</p> <p>Oh and cars really like to squeeze in on motorcycles even when it's obvious there's a bike there, to the point the rider has to avoid them, and the car just acts clueless. There are many people with this kind of attitude.</p>	Everyday frustration or complaint expressed in the form of Malay-language hate.
BN	UR	<p>বেঙ্গালবের বাচ্চা বেঙ্গালব হ্যাকার গ্রুপ বর্ণনা হ্যাকার গ্রুপ বার করে দিব লুচ্চামি কি করস শরীরটা দেখছিস নিজের ছবিটা দেখ</p> <p>Impudent child's impudent hacker group description hacker group will be kicked out, what lewdness are you doing, looking at the body, look at your own picture.</p>	Accusing a group of hackers of lewd behaviour and threatening expulsion.
HI	UR	<p>हम तेरे लिये कुछ नहीं कर सकते है दोस्त...!!! हमें अमी और मंदिर और मस्जिद बनानी है...!!!</p> <p>We can't do anything for you my friend... We still have to build more temples and mosques...!!</p>	Politicians prioritize religious projects over people's real needs.

Table 4: Hate samples correctly predicted only via transfer learning across languages.