

SHARE: A LEXICON OF HARMFUL EXPRESSIONS BY SPANISH SPEAKERS

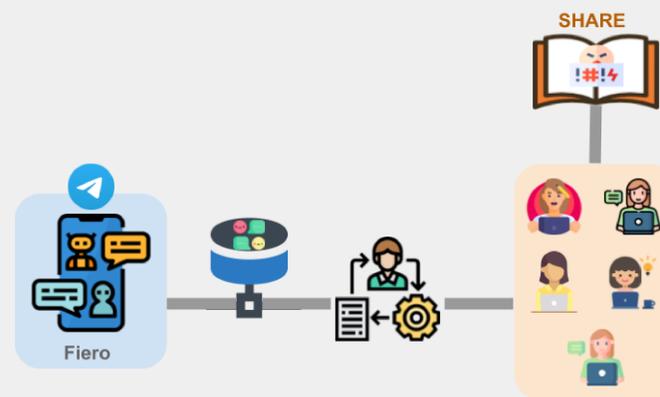
Flor Miriam Plaza-del-Arco, Ana Belén Parras-Portillo, Pilar López-Úbeda, Beatriz Botella Gil, María-Teresa Martín-Valdivia

Universidad de Jaén, Campus Las Lagunillas, 23071, Jaén (Spain)
 R+D+I department. HT medica. Carmelo Torres nº2, 23007, Jaén, Spain
 Department of Software and Computing System, University of Alicante, Alicante, Spain
 International Conference on Language Resources and Evaluation (LREC 2022)

INTRODUCTION

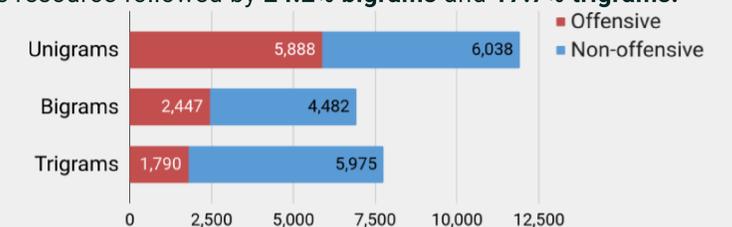
- **SHARE (Spanish HARMful Expressions)** is a new lexical resource composed of **insults and offensive expressions** collected using the **Fiero chatbot** and then **manually labeled by 5 annotators**.
- We used SHARE to release **OffenES_spans** which is the **OffendES** corpus **automatically annotated** with offensive **entities** relying on **SHARE**.
- We explore the usefulness of **SHARE** for the **interpretability** of offensive comments by comparing it with a BERT-based fine-tuning model.

DATA COLLECTION AND ANNOTATION



STATISTICS

SHARE is composed of **10,125 offensive unigrams and expressions**. The number of offensive **unigrams** represents **58.2%** of the resource followed by **24.2% bigrams** and **17.7% trigrams**.



OFFENSIVE ENTITY RECOGNITION

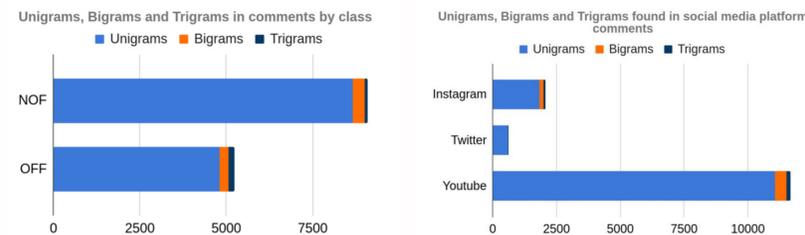
OffendES_spans

We **automatically annotated** the existing **OffendES** corpus with the terms included in **SHARE**. Two types of **entities** are found in the ANN files: Two types of entities are found within the ANN files: **OFFENSIVE_EXPRESSION** and **OFFENSIVE_TERM**.

Comment: *Das puto asco escoria* (You fucking disgusting scum)

	Start character offset	End character offset	Mention string
T1	OFFENSIVE_TERM	14 21	escoria (scum)
T2	OFFENSIVE_TERM	4 8	puto (fucking)
T3	OFFENSIVE_EXPRESSION	4 13	puto asco (fucking disgusting)
T4	OFFENSIVE_TERM	9 13	asco (disgust)

ANALYSIS



The 12 most frequent entries in OffendES_spans

Term	Freq. ↓	Term	Freq. ↓
mierda (shit)	1480	asco (disgust)	385
puto (whore)	804	loca (crazy)	341
puta (bitch)	706	gorda (fat)	336
mala (bad)	510	coño (pussy)	331
malo (bad)	442	basura (trash)	254
pringada (sucker)	440	falsa (false)	239

TOXIC SPANS DETECTION

Model	P (%)	R (%)	F ₁ (%)
BERT	91.01	91.11	91.07

ID	BERT-LIME	SHARE
818	Das puta pena dalas lo de siempre You're a fucking pity dalas as usual.	puta, das puta pena bitch, you're fucking pitiful
1227	Maldito enano rikillo Damn dwarf rikillo	maldito, enano Damn, dwarf
1545	presa es donde debes estar, pendeja loca. prison is where you belong, you crazy asshole.	pendeja, loca asshole, crazy

CONCLUSION

We release SHARE, a new lexical resource composed of offensive words and expressions for Spanish. The annotation process by five annotators obtained an agreement of 78.8%. We leverage SHARE to release OffenES_spans by automatically labeled with the terms and expressions found in SHARE. We believe that these new resources will contribute to the offensive language research community, particularly in Spanish, where there is a great scarcity of resources compared to English.

ACKNOWLEDGEMENTS

This work has been partially supported by Big Hug project (P20_00956, PAIDI 2020) and WeLee project (1380939, FEDER Andalucía 2014-2020) funded by the Andalusian Regional Government, LIVING-LANG project (RTI2018-094653-B-C21) funded by MCIN/AEI/10.13039/501100011033 and by ERDF A way of making Europe, and the scholarship (FPI-PRE2019-089310) from the Ministry of Science, Innovation, and Universities of the Spanish Government.