



Universität





Motivation

of all times

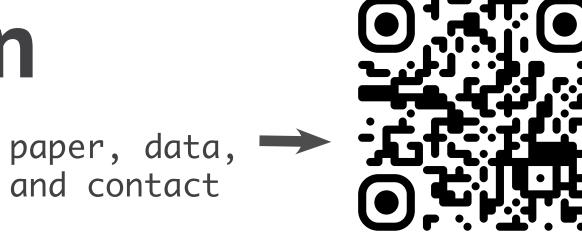
## AustroTox

# Findings **ACL 2024**

## A Dataset for Target-Based

## Austrian German Offensive Language Detection

Pia Pachinger, Janis Goldzycher, Anna Maria Planitzer, Wojciech Kusa, Allan Hanbury, Julia Neidhardt



www.pia.wien

Offensive / Toxic: derogatory remarks,

# or incites to hate or violence

#### "Bei Vielen ist der Schädel gut mit Gehirn gefüllt... Nur der BIMAZ, der hat noch viel Platz" target: individual best secretary of the interior

"Many people's skulls are well filled with brains... Only the BIMAZ still has plenty of room"

blow up in our faces."

Need for country-specific toxicity detection

"27-year-old [Nationality]. Stopped reading there" target: group

Need for target-aware toxicity detection



Need for vulgarity-aware toxicity



Need for toxicity detection aware of country-specific vulgarities

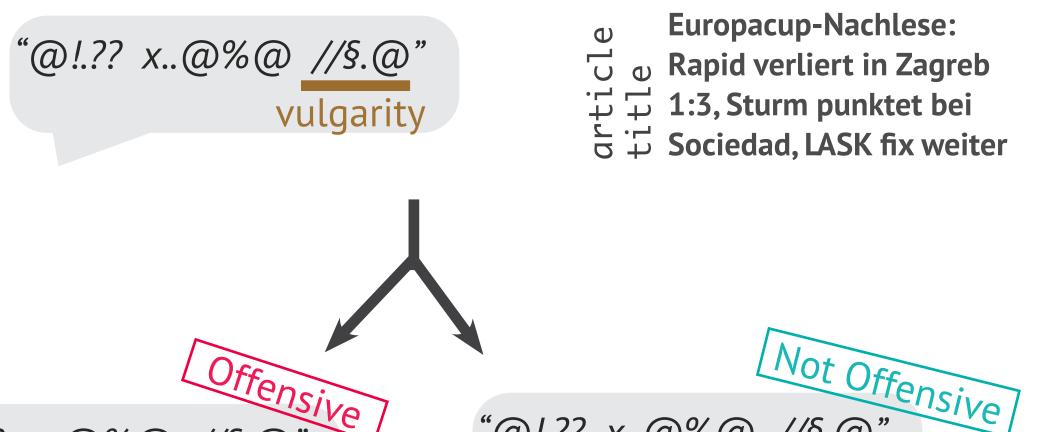
## **Dataset Creation**



#### derStandard.at

newspaper online forum

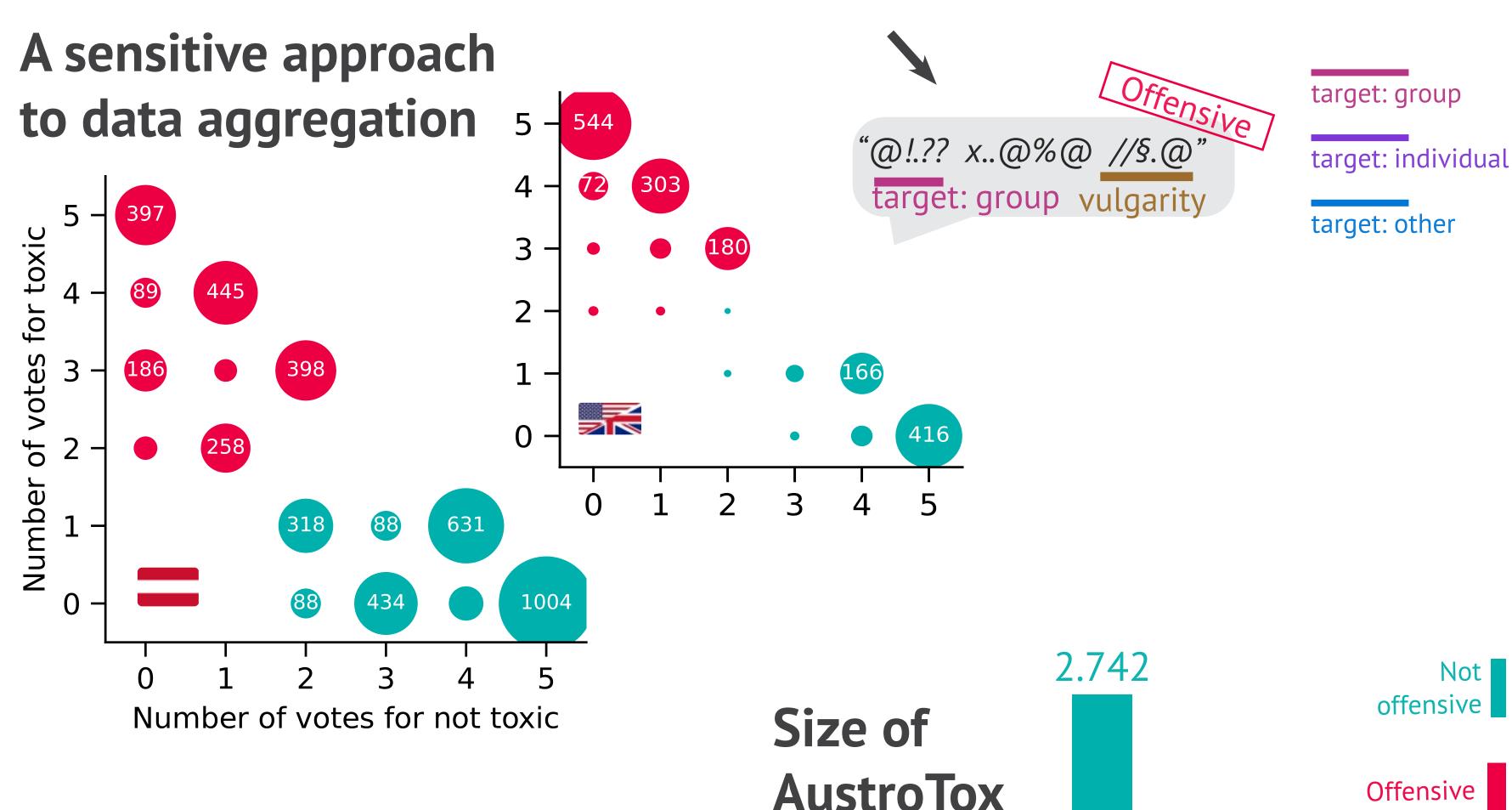
Jigsaw Toxic Comment Classification Challenge



"@!.?? x..@%@ //§.@'

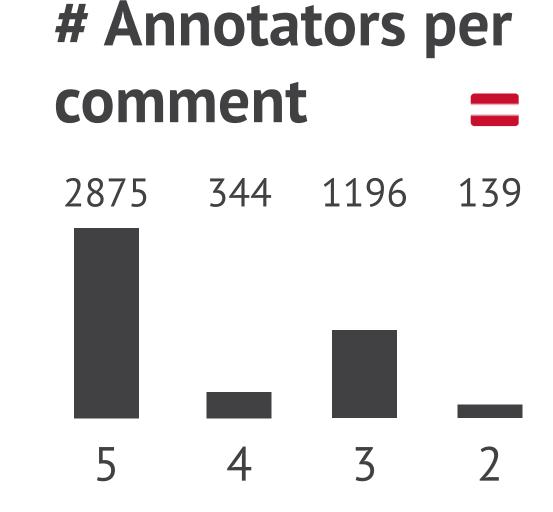
vulgarity

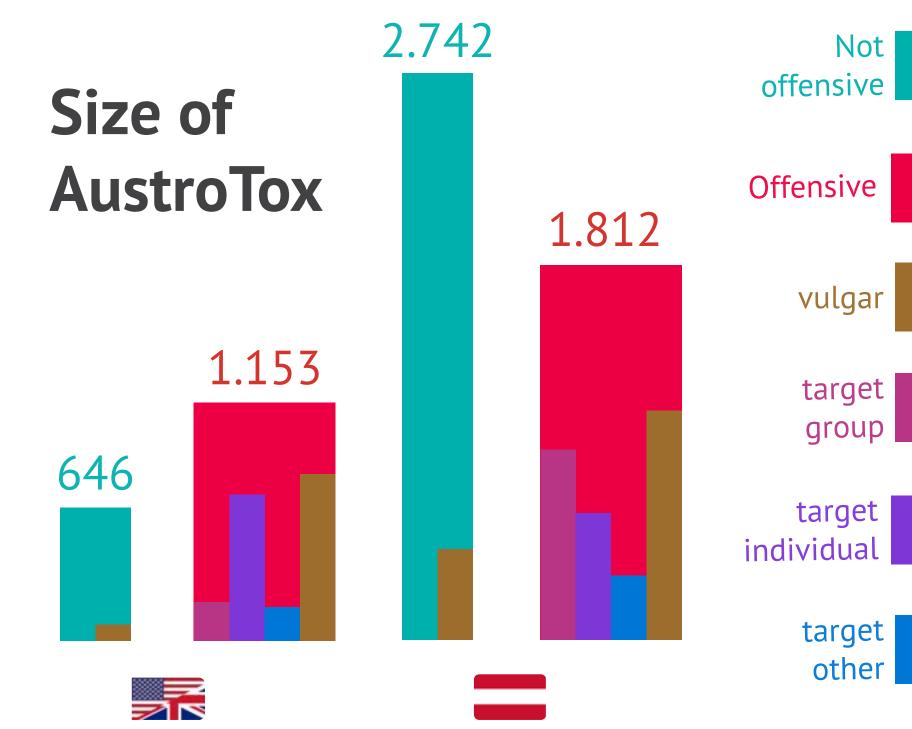
**Annotation strategy** 



"@!.?? x..@%@ //§.@'

vulgarity





### **Experiments Using AustroTox**

#### System prompt for generative models

You receive a user comment from a newspaper article. Your task is to answer the followig questions about the comment:

- 1. Is the comment toxic? (1 = toxic, 0 = non-toxic)Definition: ...
- 2. Who or what is the target of toxicity?...
- 3. Mark vulgarities with "Vulgarity"....

Respond in JSON format with the following fields:

F1 scores		# Params	Offensive		Vulgarity		Target	
fine-tuned	BERT en BERT de BERT db Gelectra GBERT Gelectra GBERT Electra Roberta # Shots	110M 110M 110M 110M 335M 335M 335M 335M	.64 .69 .71 .50 .72 .52 .73	- - - - - .88 .90	.69 .69 .69 .69 .56 .71	- - - .64 .77	.16 .22 .23 .24 .23 .19 .21	- - .08 .27
not fine-tuned	LeoLM 0	7B	.61 .52	-	-	-	-	-
	Mistral 0	7.24B	.30	-	-	-	-	-
	Llama 3 0	8B	.67 .59	.78 .82	-	-	-	-
	GPT 3.5 0	-	.68 .72	.89	.40 .43	.46 .47	.17 .20	.16 .18
	GPT 4 0 5	_	.70 . <b>76</b>	.87 .89	.36	.41 .43	.20	.15