# Hate Speech Detection

Group 15

Carsten Gieshoff, Zhaolin Li, David Pomerenke

# Problem

# Problem

## Hate Speech

@JuanYeez shut yo beaner ass up spic and hop your faggot ass back across the border little nigga

@MDreyfus @NatFascist88 Shit your ass your moms pussy u Jew bastard. Ur times coming. Heil Hitler!

My advice of the day: If your a tranny… go fuck your self!

## Offensive Language

When you realize how curiosity is a bitch #CuriosityKilledMe

Why you worried bout that other hoe? Cuz that other hoe aint worried bout another hoe

I knew Kendrick Lamar was onto something when he said "I call a bitch a bitch, a hoe a hoe, a woman a woman"

# Problem

## Hate Speech

@JuanYeez shut yo beaner ass up spic and hop your faggot ass back across the border little nigga

@MDreyfus @NatFascist88 Shit your ass your moms pussy u Jew bastard. Ur times coming. Heil Hitler!

My advice of the day: If your a tranny... go fuck your self!

## Offensive Language

When you realize how curiosity is a bitch #CuriosityKilledMe

Why you worried bout that other hoe? Cuz that other hoe aint worried bout another hoe

I knew Kendrick Lamar was onto something when he said "I call a bitch a bitch, a hoe a hoe, a woman a woman"
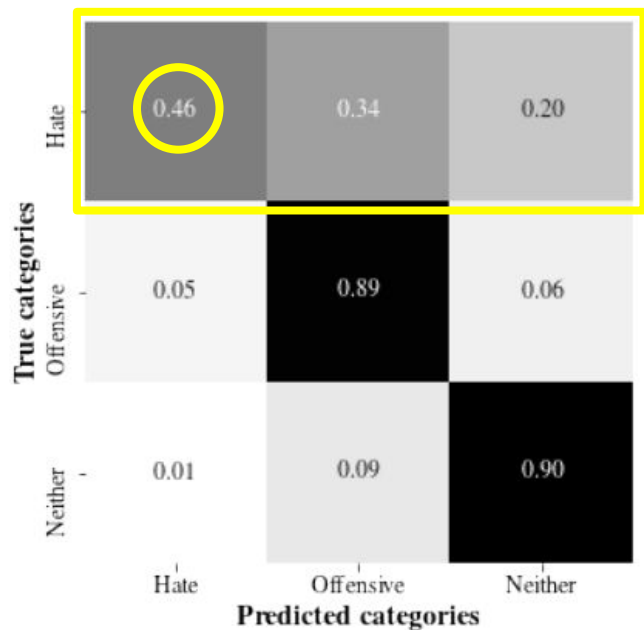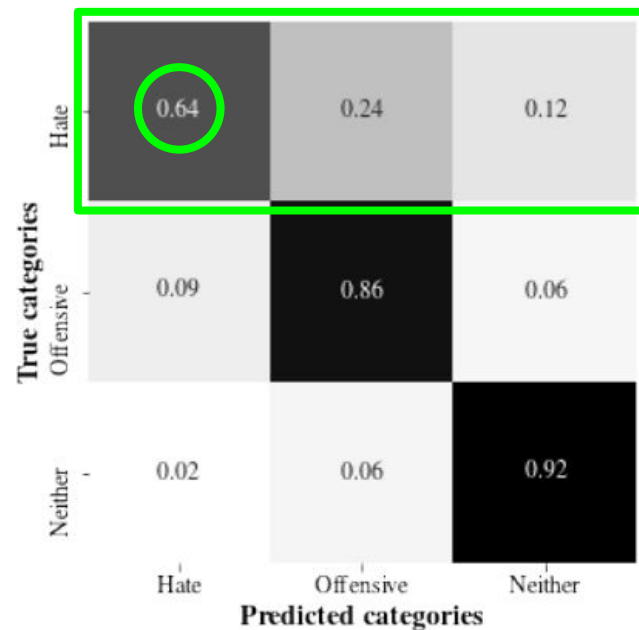
HateBase Data Set

# Approach

1. Improve existing approach (Davidson & al 2017)
   - Word-based vs character-based tokenization
   - Different machine learning models
   - Different weighting of categories

# Approach

1. Improve existing approach (Davidson & al 2017)
   - Word-based vs character-based tokenization
   - Different machine learning models
   - Different weighting of categories
2. Build classifier with transformer language model
   - BERT
   - BERTweet (RoBERTa)

# Approach

1. Improve existing approach (Davidson & al 2017)
   - Word-based vs character-based tokenization
   - Different machine learning models
   - Different weighting of categories
2. Build classifier with transformer language model
   - BERT
   - BERTweet (RoBERTa)
3. Transfer classification from English to other languages
   - mBART-50 for translation from Spanish, Italian, Turkish
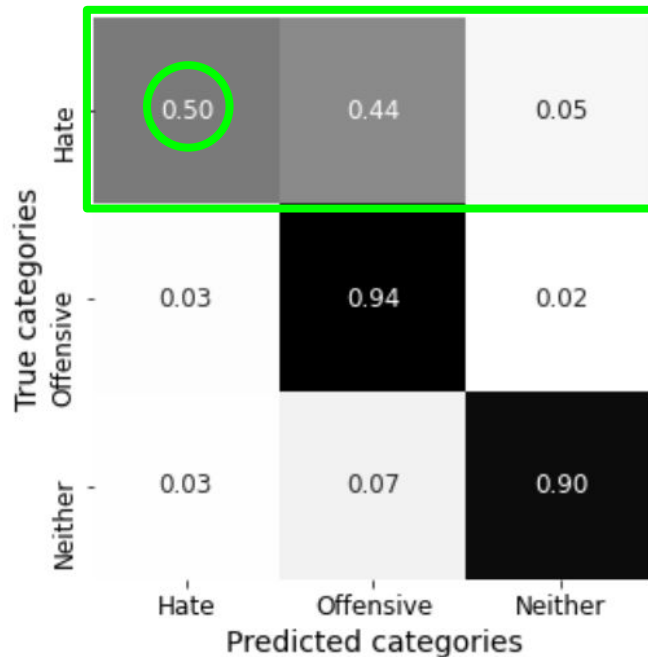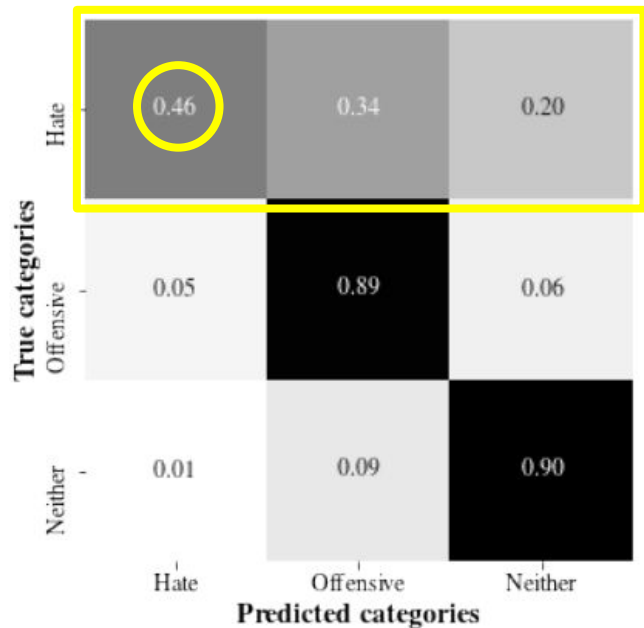   - Evaluation on detoxify data set

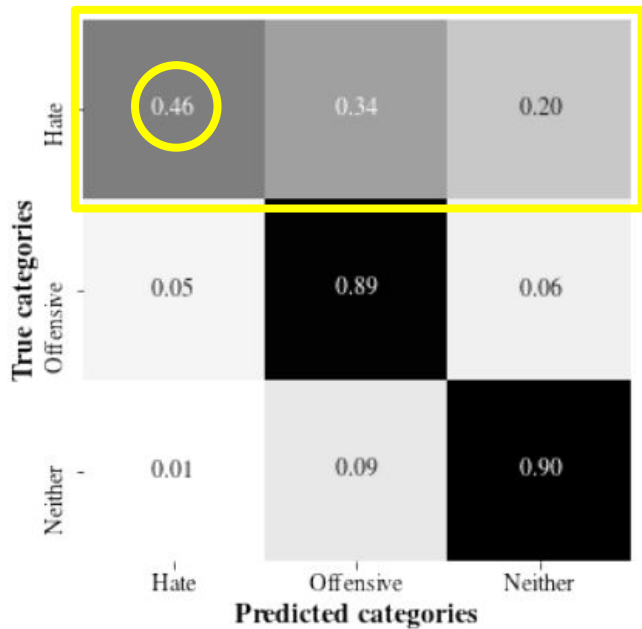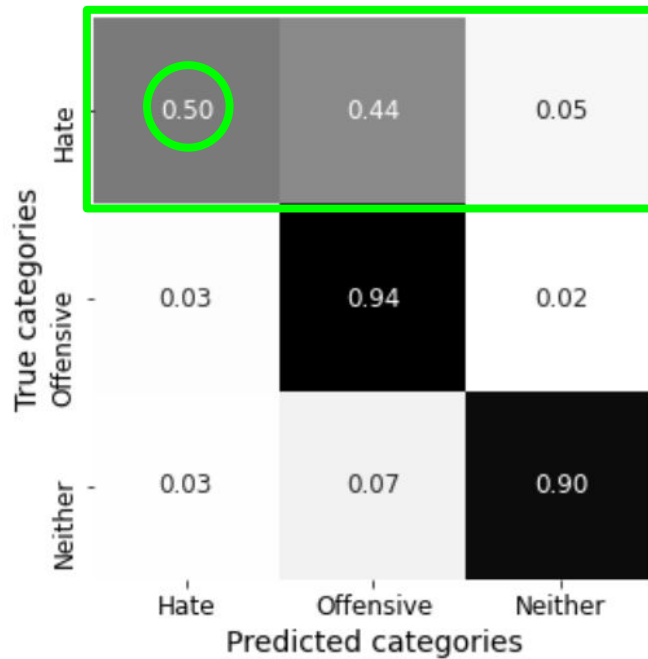# Results: Different Weighting



(a) $E = 1$

(c) $E = 1.5$
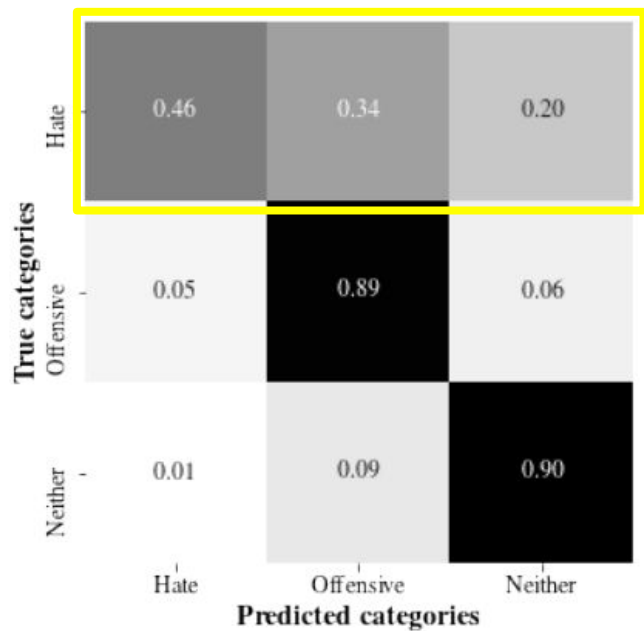
# Results: Original vs BERT

# Results: Original vs BERT
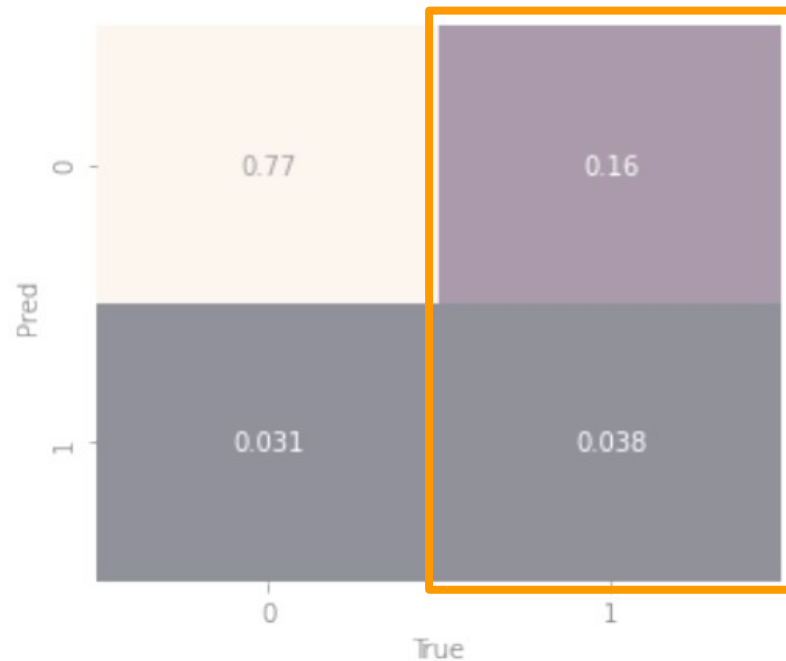


84% accuracy



91% accuracy

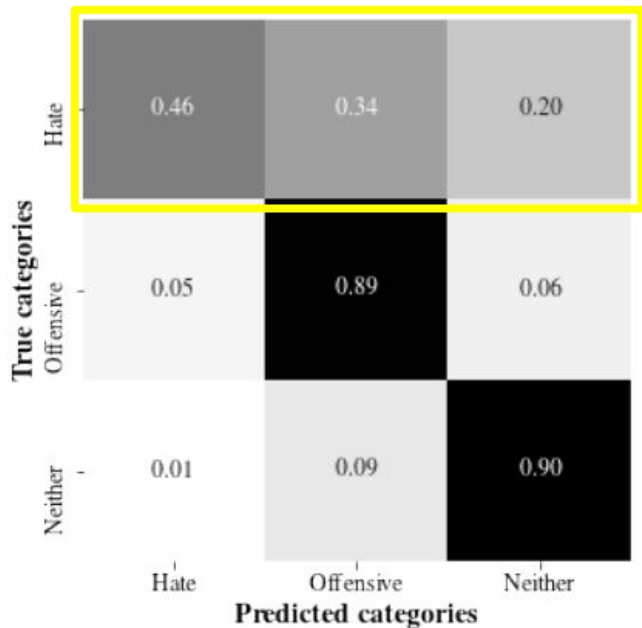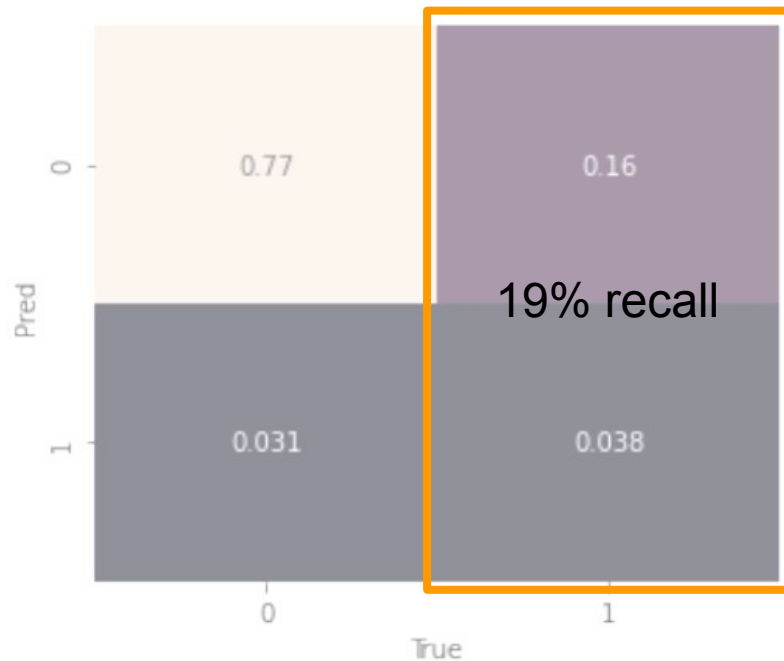# Results: Original vs transfer to Italian



84% accuracy

81% accuracy

# Results: Original vs transfer to Italian



84% accuracy

81% accuracy

# Hate Speech Detection

Carsten Gieshoff, Zhaolin Li, David Pomerenke

carstengieshoff/HS-Detection-Project

- We improve hate speech recall by aggressive weighting
- Our BERT-based model beats the baseline from 2017
- We perform multilingual detection based on only English training data
- Future work: Improving recall using data augmentation