ZEERAK TALAT | UNIVERSITY OF SHEFFIELD | 03.09.2019

# "IT AIN'T ALL GOOD"
# MARGINALISATION IN THE NAME OF PROTECTION

"THANGS AIN'T GOIN' LIKE YOU THINK THEY SHOULD, IT'S ALL ON YOU"
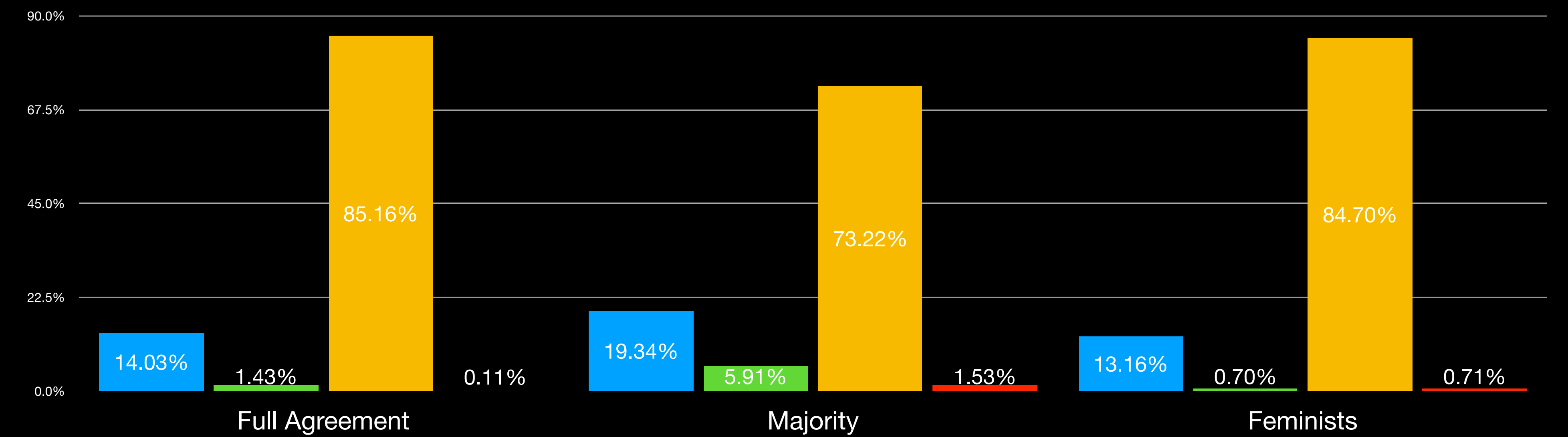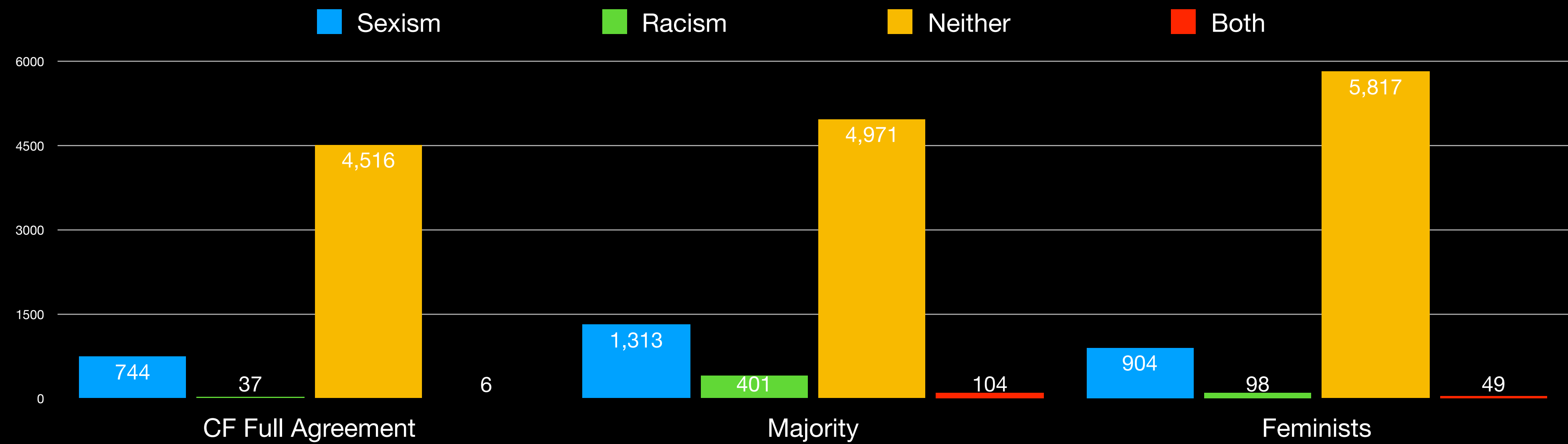
"All Good" – De La Soul featuring Chaka Khan

# LAW

# ANNOTATION

1. uses a sexist or racial slur.
2. attacks a minority.
3. seeks to silence a minority.
4. criticizes a minority (without a well founded argument).
5. promotes, but does not directly use, hate speech or violent crime.
6. criticizes a minority and uses a straw man argument.
7. blatantly misrepresents truth or seeks to distort views on a minority with unfounded claims.
8. shows support of problematic hash tags. E.g. "#BanIslam", "#whoriental", "#whitegenocide"
9. negatively stereotypes a minority.
10. defends xenophobia or sexism.
11. contains a screen name that is offensive, as per the previous criteria, the tweet is ambiguous (at best), and the tweet is on a topic that satisfies any of the above criteria.

| | Full Agreement v. Feminists | Majority v. Feminists | All CF Annotators |
|---|---|---|---|
| Cohen's kappa | 0.70 | 0.34 | 0.57 |
| Krippendorf's alpha | 0.70 | 0.32 | |

# F1-SCORES

| | Amateur (Majority Vote) | Expert/Feminists |
|---|---|---|
| *Character* n-*gram* | 86.41 | 91.24 |
| *Token* n-*gram* | 86.37 | 91.55 |
| **Binary Gender** | 76.64 | 77.77 |
| **GenderProbability** | 86.37 | 81.30 |
| ***Brown Clusters*** | 84.50 | 87.74 |
| AHST | 71.71 | 55.40 |

# ETHICS AND BIAS

▸ The Risk of Racial Bias in Hate Speech Detection. Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi & Noah A Smith. ACL (2019)

▸ Racial Bias in Hate Speech and Abusive Language Detection Datasets. **Thomas Davidson, Debasmita Bhattacharya and Ingmar Weber**. Third Workshop on Abusive Language Online (2019)

▸ Lipstick on a Pig: Debiasing Methods Cover up Systematic Gender Biases in Word Embeddings But do not Remove Them. **Hila Gonen and Yoav Goldberg**. NAACL (2019)

▸ Hateful symbols or hateful people? predictive features for hate speech detection on twitter. **Zeerak Waseem and Dirk Hovy**. NAACL SRW (2016)

▸ Are you a racist or am I seeing things? Annotator influence on hate speech detection on Twitter. **Zeerak Waseem**. The First Workshop on NLP and Computational Social Science (2016)

▸ Understanding Abuse: A Typology of Abusive Language Detection Subtasks. **Zeerak Waseem, Thomas Davidson, Dana Warmsley, and Ingmar Weber**. The First Workshop on Abusive Language Online (2017).

▸ Bridging the Gaps: Multi Task Learning for Domain Transfer of Hate Speech Detection. **Zeerak Waseem, James Thorne, and Joachim Bingel**. In Jennifer Golbeck (editor), Online Harassment (2018)

▸ A Reductions Approach to Fair Classification. **Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, Hanna Wallach**. ICML (2018)

▸ What's in a Name? Reducing Bias in Bios Without Access to Protected Attributes. A**lexey Romanov, Maria De-Arteaga, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, Anna Rumshisky and Adam Kalai**. NAACL (2019)