# "Yeah well, that's just your opinion man"

## The fallouts of abusive language detection

27 10 2020

Zeerak Talat

University of Sheffield | Digital Democracies Institute

Twitter: @zeerakw

## Clean up the Internet

Clean up the internet is an independent, UK-based organisation concerned about the degradation in online discourse and its implications for democracy. We campaign for evidence-based action to increase civility and respect online, and to reduce online bullying, trolling, intimidation, and misinformation.

## Unilever warns social media to clean up "toxic" content

Toxic Content on Digital Platfor...
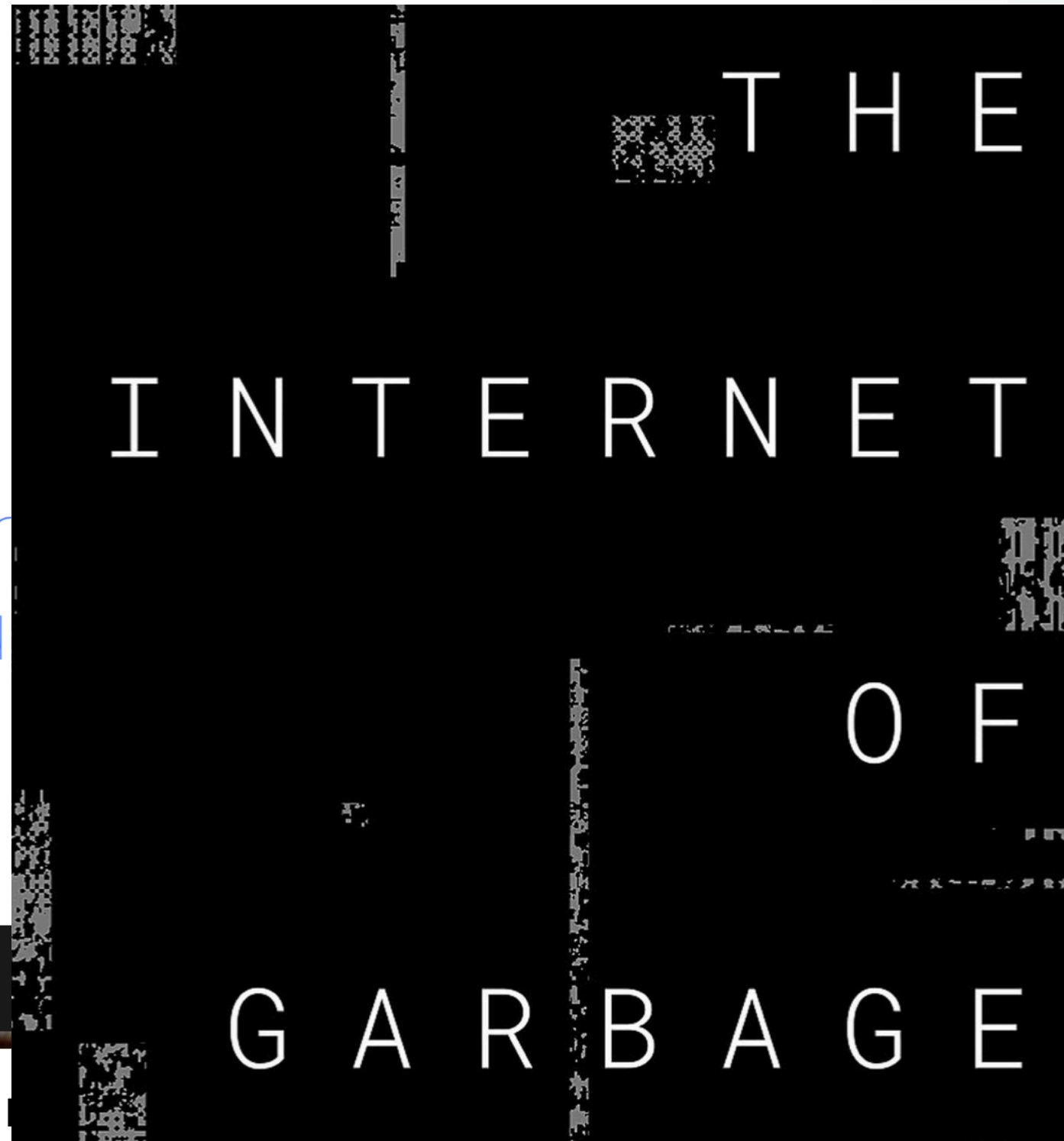the rise. How Can Brands Rebu...
Consumer Trust?

**Two ways social networks could control toxic content**

: The Internet's Great Spring-

AMNESTY
INTERNATIONAL

WHO WE ARE

THE

INTERNET

OF

GARBAGE

THE

NERS

TOXIC TWITTER

PLACE FOR WO...

Abuse, racism and hate speech - help us clean up the comments online

# Computational Modelling

# Goals of abuse detection systems

### Imaginaries

- Protecting marginalized communities and people from abuse

- Encode ideological positions into labelled data

- Obtain an understanding of the subjective nature of abuse.

### Realities

- Minority language (e.g. sociolects and dialects) use is penalized

- Use of crowd-sourcing labels as a source of ground truth

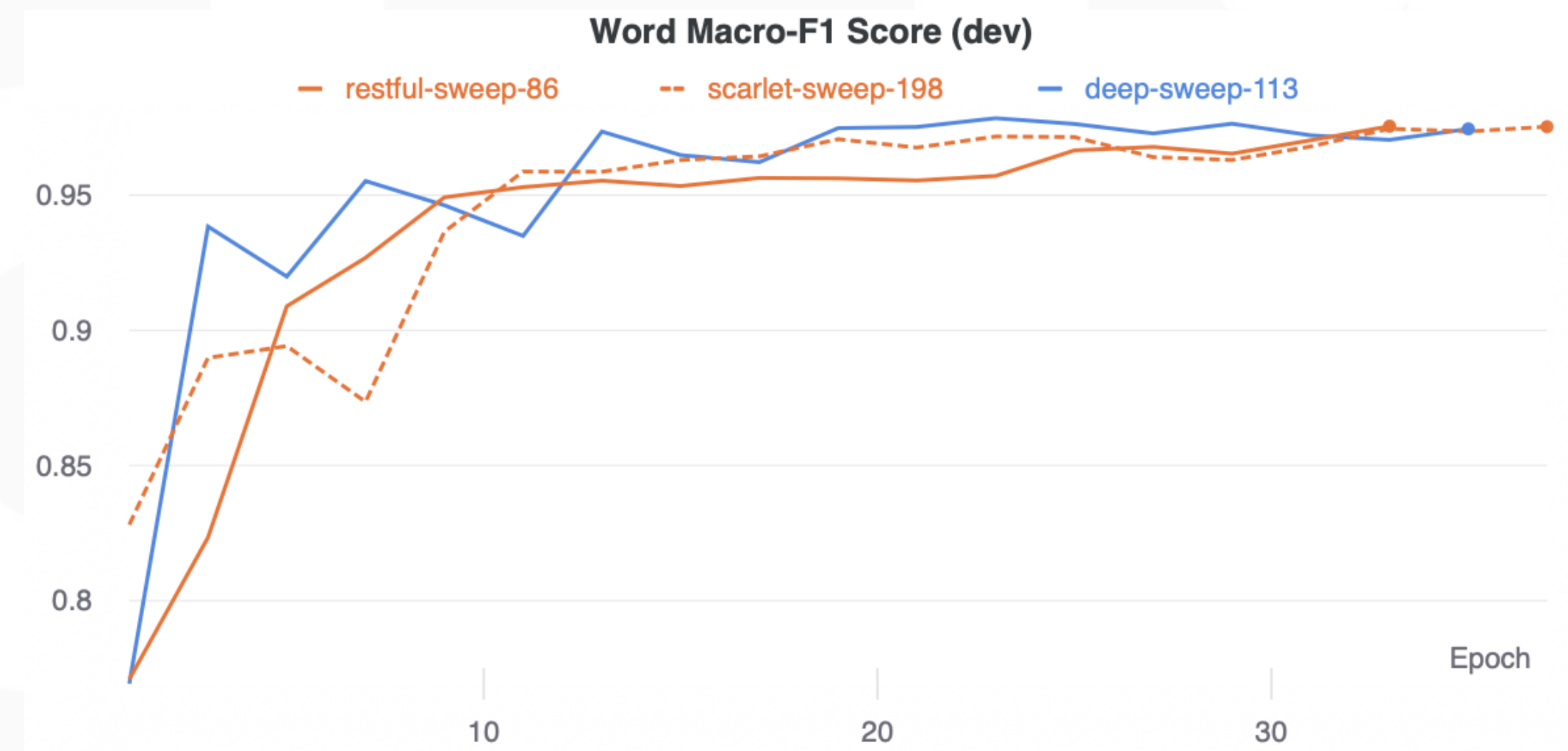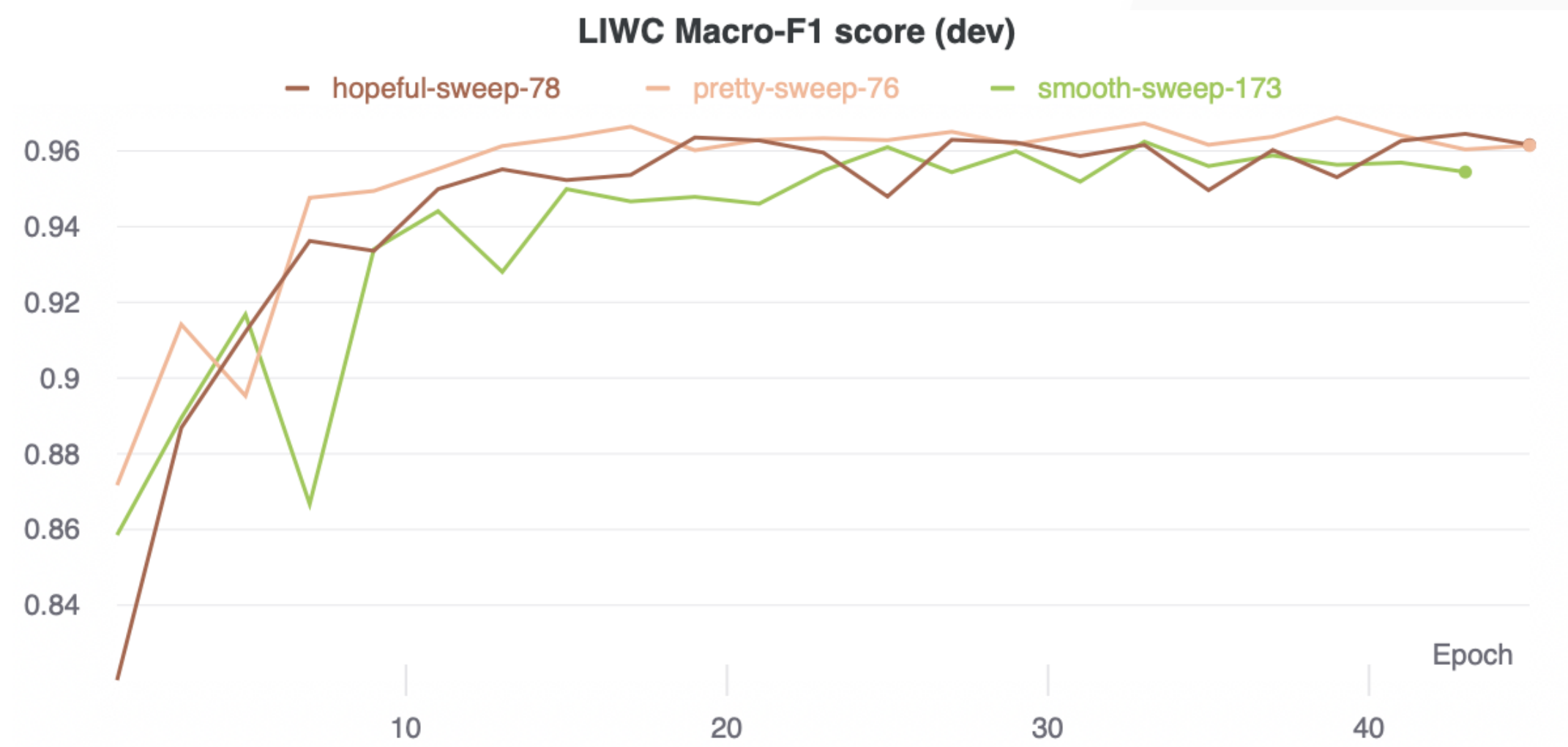- Hardwiring (simple) correlations into the model

Modelling Trends

# **The increasing complexity of abuse detection models**

- From feature engineering to inductive

 biases

- Pre-trained language models

# Very preliminary findings from my dissertation

Models and inputs

**LIWC Macro-F1 score (dev)**

— hopeful-sweep-78    — pretty-sweep-76    — smooth-sweep-173

Epoch

**Word Macro-F1 Score (dev)**

— restful-sweep-86    - - scarlet-sweep-198    — deep-sweep-113

Epoch

# "Don't you know that you're toxic?"

ddi | digital democracies institute
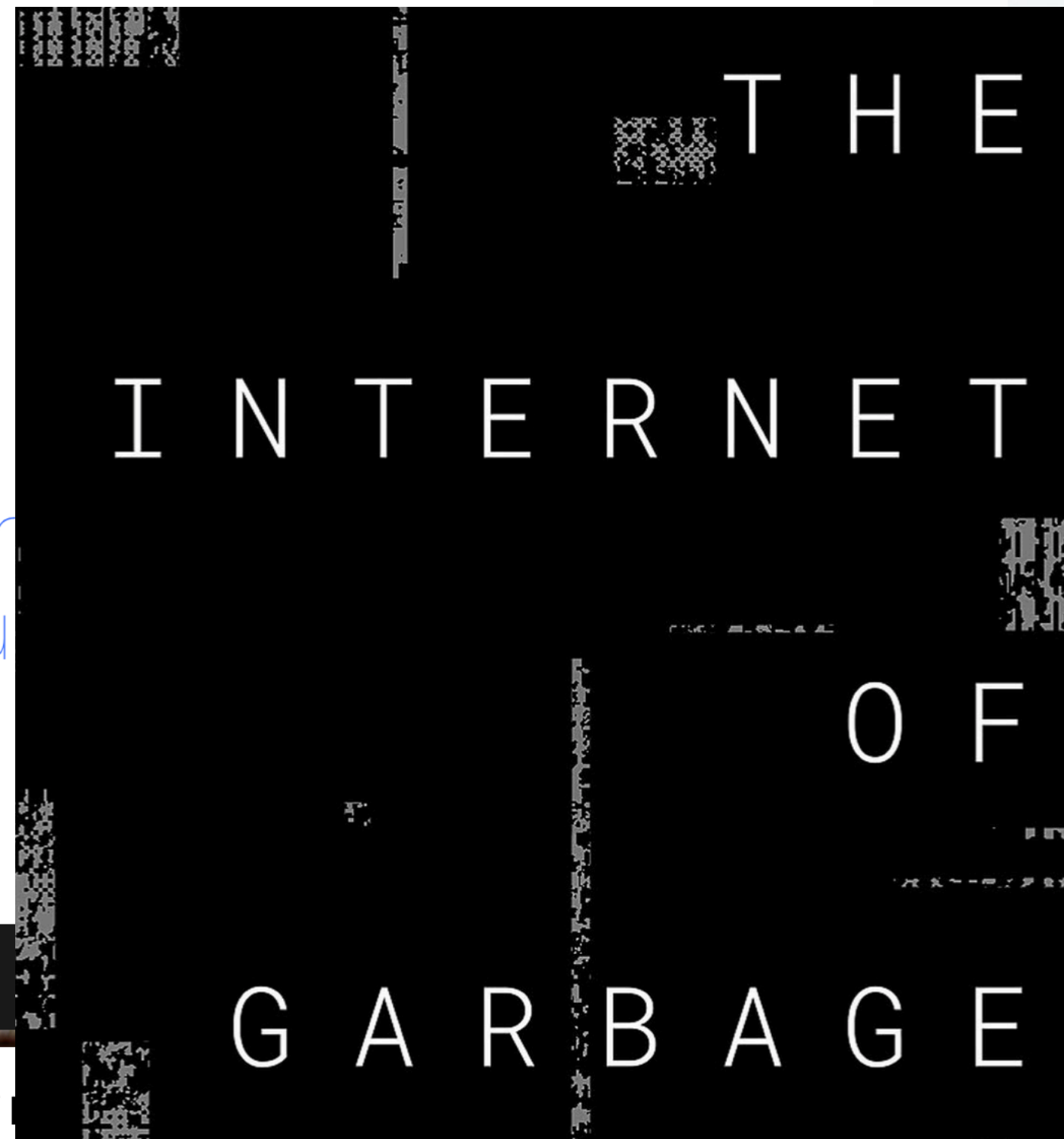
# Clean up the Internet

Clean up the internet is an independent, UK-based organisation concerned about the degradation in online discourse and its implications for democracy. We campaign for evidence-based action to increase civility and respect online, and to reduce online bullying, trolling, intimidation, and misinformation.

# Unilever warns social media to clean up "toxic" content

Toxic Content on Digital Platfor... the rise. How Can Brands Rebu... Consumer Trust?

AMNESTY INTERNATIONAL

WHO WE ARE

# Two ways social networks could control toxic content

: The Internet's Great Spring-

THE INTERNET OF GARBAGE

THE NERS

TOXIC TWITTER A TO...

PLACE FOR WO...

Abuse, racism and hate speech - help us clean up the comments online

# The case of content moderation

The politics of toxicity

- The "toxicity" of words

- Definitional power of "toxic"

- Independent third parties



IF YOU DROP SOAP ON THE GROUND, IS THE SOAP DIRTY?

OR IS THE DIRT CLEAN?

# Content moderation technology as vehicles of oppression

# "No one ever accused the God of monotheism of objectivity, only of indifference"

Donna Harraway (1988),
*Situated Knowledges: The Science Question in Feminism and the Privilege of Partial Perspective.*

# Embodiment in machine learning pipelines

## Designer

- Disassociated from data, model, and outcomes
- "Lack of diversity  cited as reason for biased models

## Data

- Collected from embodied humans, adjudicated on and disentangled from those contexts
- Subjectivity always finds its way
- Annotation == truth (to the model)
- Encode dominant discourses on objects

# Future directions of content moderation technologies