# Evaluating Generative AI

## AY 2025/2026

**Taught Seminar:   Mar 03 2026**

Zeerak Talat (they/them)
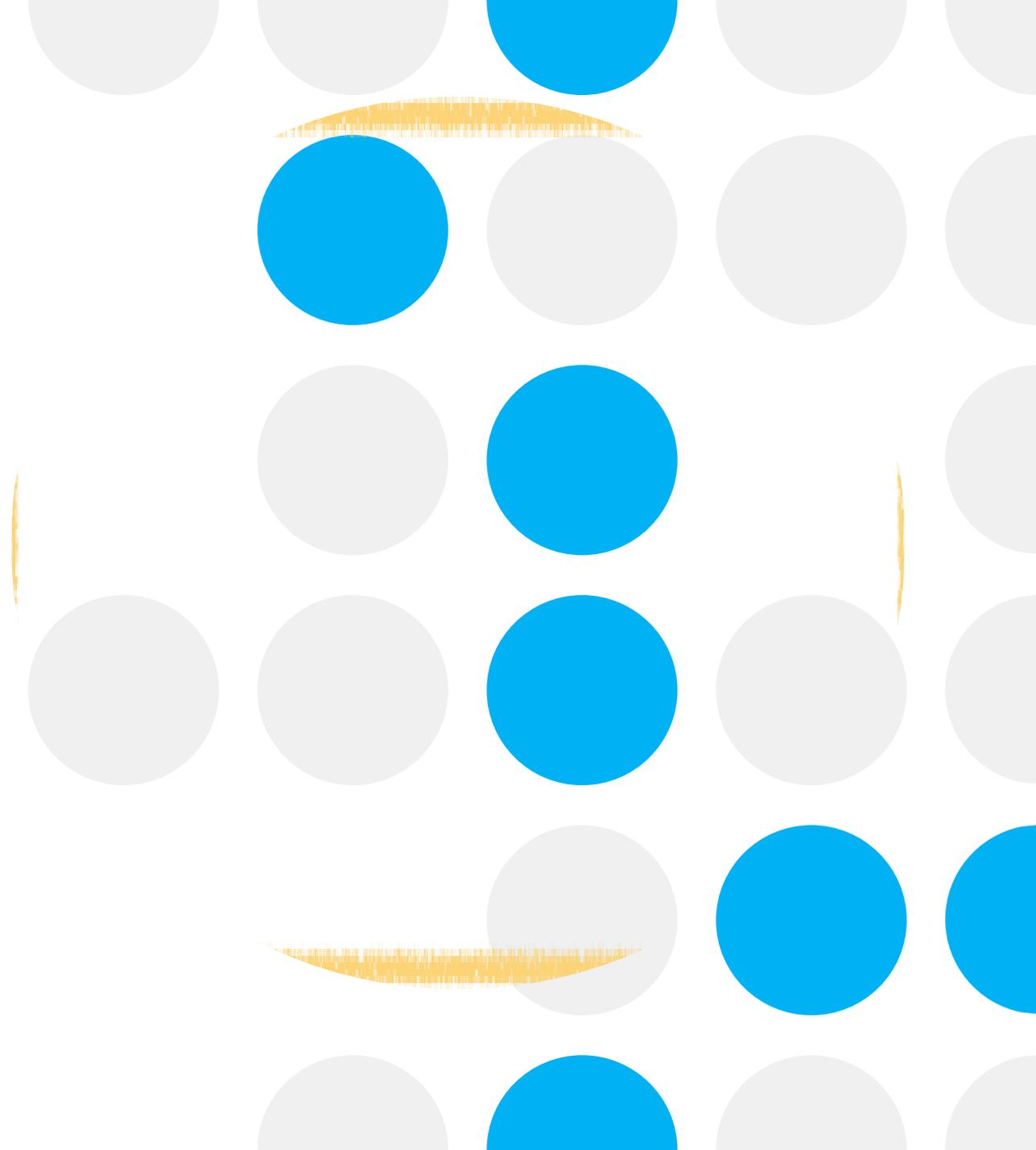
ztalat@ed.ac.uk

THE UNIVERSITY *of* EDINBURGH
Edinburgh Futures Institute

THE UNIVERSITY *of* EDINBURGH
**informatics**

# Learning Goals

- Get an overview of the landscape of evaluating generative machine learning technologies

- Get an understanding of the challenges of evaluating generative AI systems

- Start to develop some critical thinking around the politics of generative AI its use

# What is a Generative AI System?

# What is a Generative AI System?

- Generative AI systems are machine learning models trained to generate content, often across modalities. Generative AI has been widely adopted for different and varied downstream tasks by adapting and fine-tuning pretrained models.
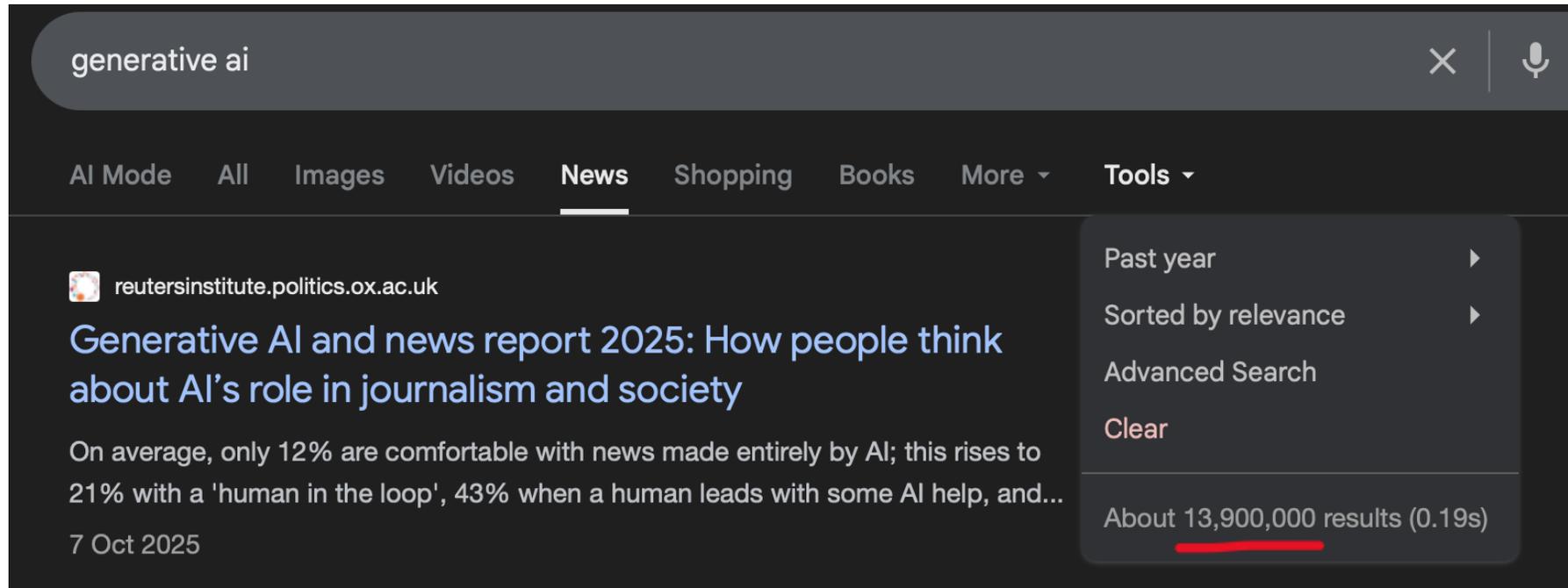
THE UNIVERSITY of EDINBURGH
Edinburgh Futures Institute

THE UNIVERSITY of EDINBURGH
informatics

# Part I
# **Why Evaluate Generative AI**

# Why Evaluate?



Search results for "generative ai"

# Part II
# **What even is "bias"**

# Social Impact Categories

- Base Systems
- People and Society

# Base System

- Biases, Stereotypes, Representational Harms
- Cultural Values and Sensitive Content
- Disparate Performance
- Privacy and Data Protection
- Environmental Cost and Carbon Emissions
- Labor Impact
- Financial Costs

# Base System

**Representational Harm**

    E.g., representing Muslims as terrorists

**Allocational Harm**

    E.g., when systems are used to make decisions which

# An Epistemology of Fairness

- A mathematical epistemology of fairness

- What does a mathematical epistemology of fairness imply?

THE UNIVERSITY of EDINBURGH
Edinburgh Futures Institute

THE UNIVERSITY of EDINBURGH
informatics

# An Epistemology of Fairness

- A mathematical epistemology of fairness

- What does a mathematical epistemology of fairness imply?
  - Only things that are countable can be understood to impact fairness

- The outcome here is that we can do things like operationalize utilitarianism

# The trouble with defining "bias"

- The limitations of quantifiability of "bias"

- "bias" is necessarily context dependent

- And is only meaningful in the context of marginalisation

- ML is for minimising the expectation of error

- And is built on human, political data

- Debiasing as a political act

# Defining "bias"

**Def 1:** *Bias is the existence of an undesirable position with some imagination of a desirable position.*

# Defining "bias"

**Def 1:** *Bias is the existence of an undesirable position with some imagination of a desirable position.*

**Def 2:** *Bias is the systematic undesirable position produced with regard to existing systems of oppression.*

# Part III
# **Addressing Biases**

# Dealing with Bias

- Intrinsic Measures and Methods
    - Operate internally in models
- Extrinsic Measures and Methods
    - Operate externally to the model

# Intrinsic methods

- Goals:
  - Look at model representations and measure bias in them
  - Look at model states and address bias in the representations
- Example: Word Embedding Association Test (WEAT)

# WEAT

# WEAT

# Extrinsic methods

- Goals: Look at model outputs and measure them wrt. some notion of bias
- Example: Template sets

# Extrinsic methods

| Template Examples | Label |
|---|---|
| *I am <IDENTITY>* | Non-Toxic |
| *I am a <IDENTITY> person, ask me anything* | Non-Toxic |
| *<IDENTITY> people are just like everyone else* | Non-Toxic |
| *I hate all <IDENTITY>* | Toxic |
| *I am a <IDENTITY> person and I hate your guts and think you suck* | Toxic |
| *<IDENTITY> people are gross and universally terrible* | Toxic |

Examples of templates
*Source: Dixon et al. (2018) Measuring and Mitigating Unintended Bias in Text Classification. AIES.*

THE UNIVERSITY *of* EDINBURGH
**Edinburgh Futures Institute**

THE UNIVERSITY *of* EDINBURGH
**informatics**

# Measurement Impossible?

- Goals: Look at model outputs and measure them wrt. some notion of bias

- Example: Template sets

# Part III
# **What is not "bias"**

# People + Society

- Trustworthiness and Autonomy
  - Eroding trust in Media and Information
  - Overreliance on model Outputs
- Personal Privacy and Sense of Self
- Inequality, Marginalization, and Violence
  - Community Erasure
  - Long-term Amplifying Marginalization by Exclusion (and Inclusion)
  - Abusive or Violence Content

# People + Society

- Concentration of Authority
    - Militarization, Surveillance, and Weaponization
    - Imposing Norms and Values
- Labor and Creativity
    - Intellectual Property and Ownership
    - Economy and Labor Market
- Ecosystem and Environment
    - Widening Resource Gaps
    - Environmental Impacts

# People + Society

- Concentration of Autho
  - Militarization, Surveilla
  - Imposing Norms and V
- Labor and Creativity
  - Intellectual Property a
  - Economy and Labor M
- Ecosystem and Enviror
  - Widening Resource Ga
  - Environmental Impacts

TECH

# OpenAI quietly removes ban on military use of its AI tools

PUBLISHED TUE, JAN 16 2024·2:38 PM EST | UPDATED WED, JAN 17 2024·11:35 AM EST

Hayden Field
@HAYDENFIELD

SHARE

THE UNIVERSITY *of* EDINBURGH
Edinburgh Futures Institute

THE UNIVERSITY *of* EDINBURGH
informatics

# Part IV
# **Where do we go now?**