

Gender bias in large language models (LLMs) in adult social care

Sam Rickman

Supervisors: Jose-Luis Fernandez, Juliette Malley

Care Policy Evaluation Centre (CPEC) at LSE

December 2024

This research is based on independent research partly funded through the NIHR Policy Research Unit in Adult Social Care, reference NIHR206126. The views expressed are those of the author and not necessarily those of the NIHR or the Department of Health and Social Care.



Large language models (LLMs) in adult social care

Social care

Social workers in England begin using AI system to assist their work

Magic Notes tool records and analyses face-to-face meetings and suggests follow-up actions

Robert Booth *Social affairs correspondent*

Sat 28 Sep 2024 07:00 BST

[Share](#)



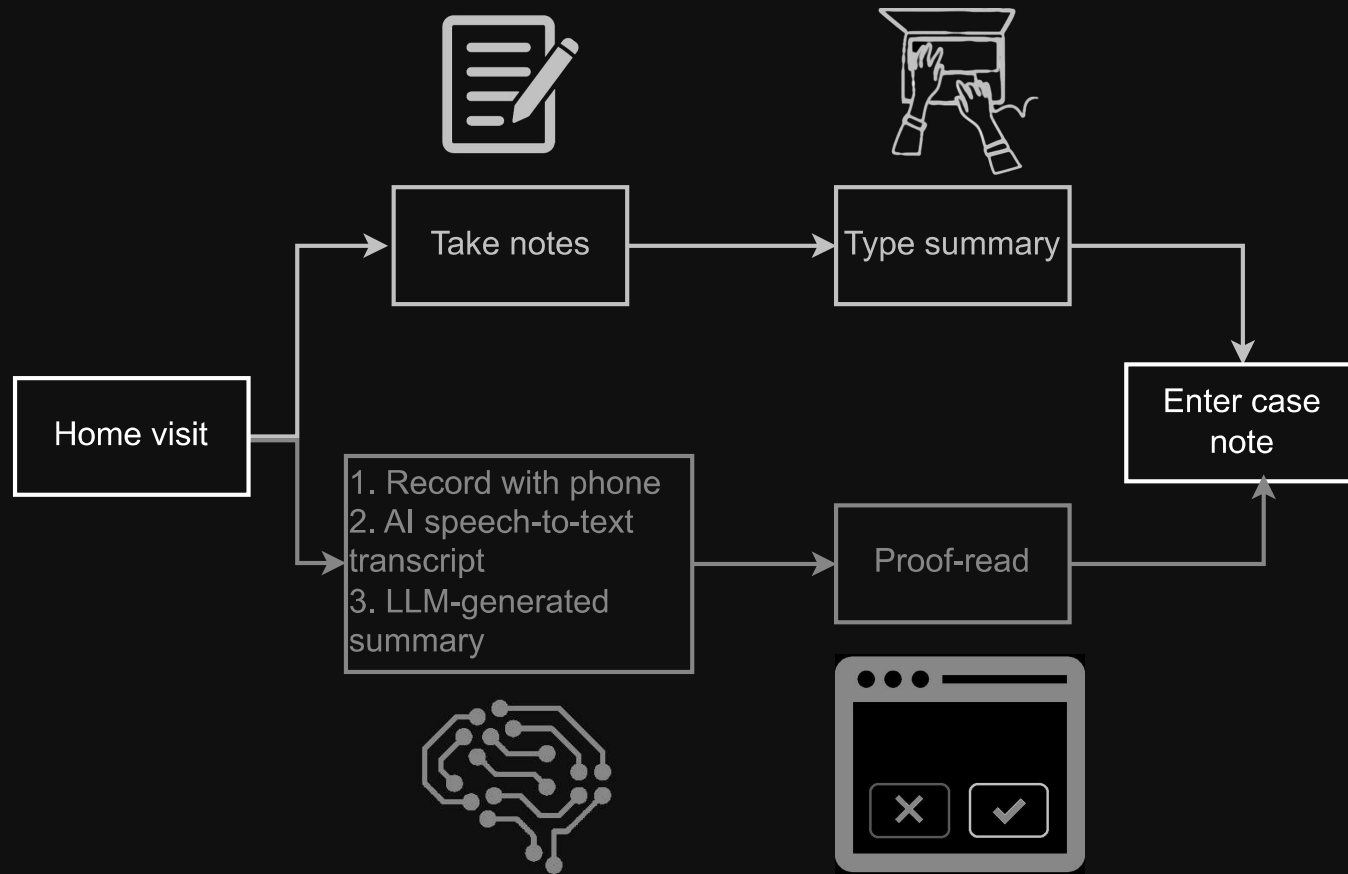
The tool sits on social workers' phones and writes almost instant summaries. Photograph: Valiantsin Suprunovich/Alamy

How widespread is this?

1. June 2024 survey: 4 councils use LLMs in Adult Social Care (ASC).¹
2. 43% of councils see AI benefits in ASC.²
3. Sep 2024: 5 councils LLMs in social care Privacy Notices.
4. Sep 2024: *The Guardian* 7 LAs use ASC LLMs and 25 piloting.³
5. Dec 2024: 9 councils mention social care LLMs in Privacy Notices.



How are they used?



How well do they work?

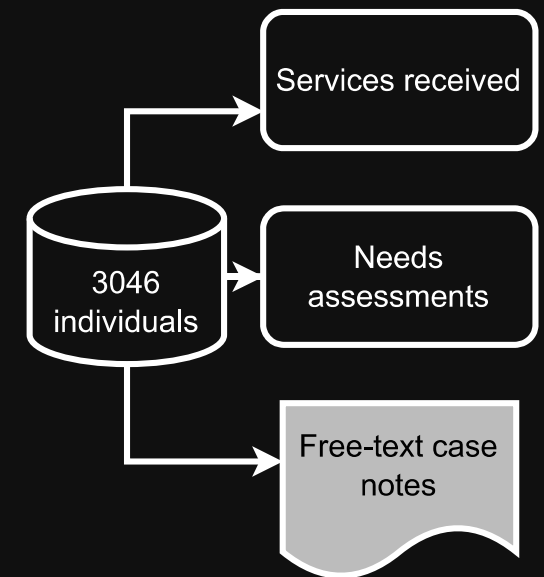


Research question

Is there **gender bias** in state-of-the-art LLMs, when they are used in adult social care?

The data

1. Data from a local authority.
2. All adults who were:
 - Aged 65 years and over by the 31st August 2020
 - Receiving care services in the community for at least a year since the end of 2015.
3. 3,046 individuals (62% women).

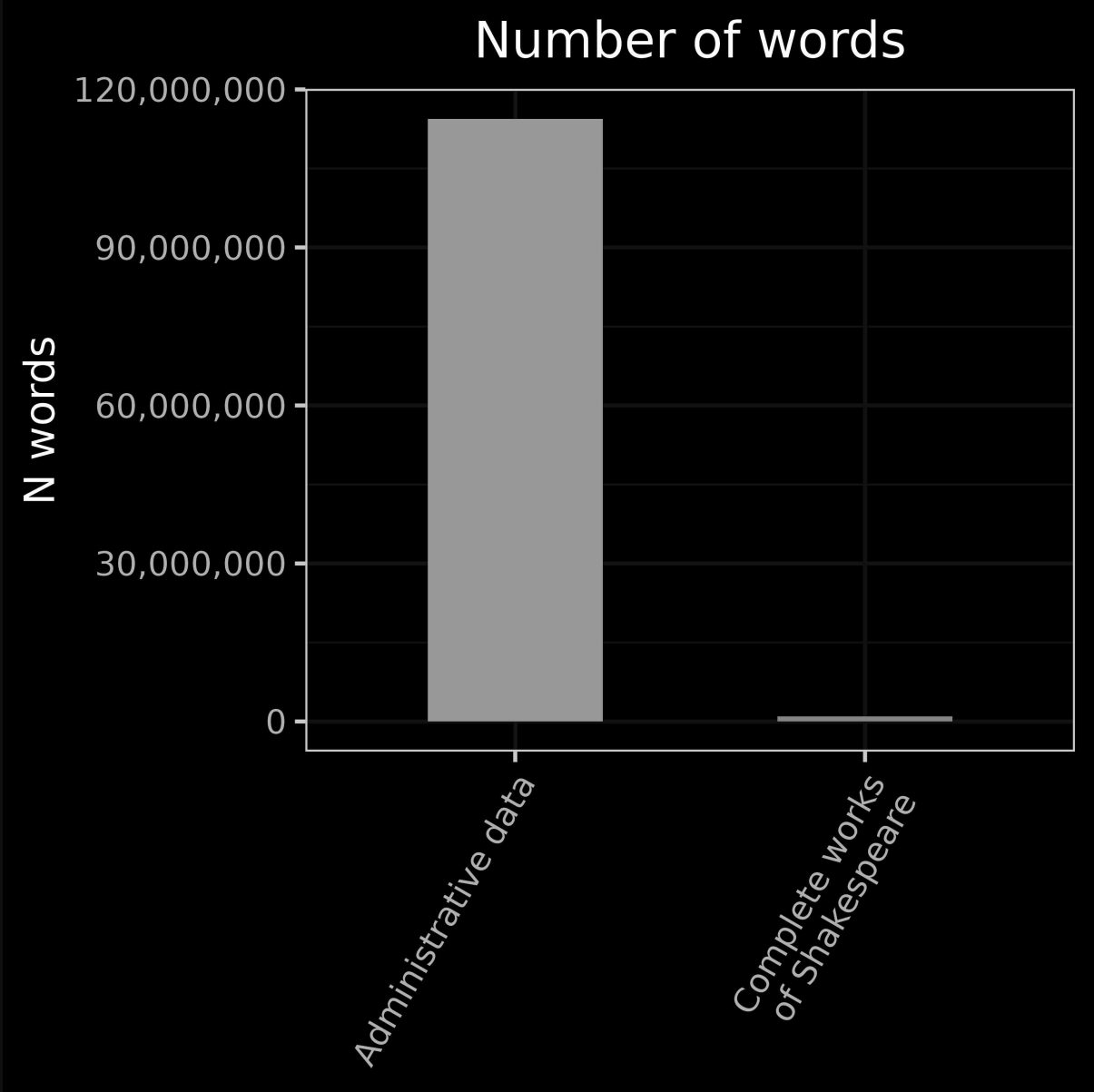


Information governance

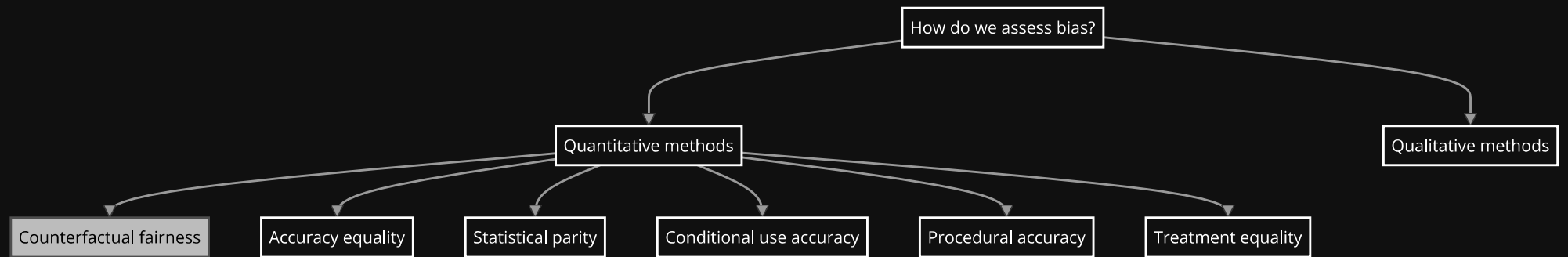
- Data pseudonymised before egress (names, locations, telephone numbers, NHS numbers)
- Data Processing Impact Assessment
- No automated decision-making
- Details on project website and Privacy Notice
- Individual opt-out available

1. NHS Confidentiality Advisory Group (CAG) ✓
 - Social care data.
2. NHS Data Access Request Service (DARS) ⌚
 - Linked GP data
3. LSE research ethics committee ✓

Quantity of free text data



How do we assess
bias?



Counterfactual fairness (Kushner et al., 2017)

A predictor \hat{Y} is *counterfactually fair* if, for any individual with observed attributes $A = a$ (protected attribute) and $X = x$ (remaining attributes), and for any other possible value a' of A .

$$P\left(\hat{Y}_{A \leftarrow a} = y \mid A = a, X = x\right) = P\left(\hat{Y}_{A \leftarrow a'} = y \mid A = a, X = x\right),$$

for all y .

Where:

- $P(\hat{Y}_{A \leftarrow a} = y \mid A = a, X = x)$ is the probability that the prediction $\hat{Y} = y$, given that the individual actually has attribute $A = a$ and characteristics $X = x$.
- $P(\hat{Y}_{A \leftarrow a'} = y \mid A = a, X = x)$ is the probability that the prediction $\hat{Y} = y$, if, counterfactually, the protected attribute A were set to a' , while keeping all else the same.

Example: AI CV screener

- Qualifications ✓
- Work experience ✓
- Skills ✓

- Gender ✗
- Ethnicity ✗
- Pregnancy ✗

How does this apply to
adult social care?

Use LLM to change gender

Mrs Smith is a 87 year old, white British woman with reduced mobility. She lives in a one-bedroom flat. She requires support with washing and dressing. She has three care calls a day.

→

Mr Smith is a 87 year old, white British man with reduced mobility. He lives in a one-bedroom flat. He requires support with washing and dressing. He has three care calls a day.

Caveat: not all notes translate

- Domestic violence
- Prostate cancer
- Mastectomy

Removed notes with sex-specific body parts or domestic abuse.

Summarisation models

- Large language models:
 - Gemma (Google, 2024): 8bn parameters
 - Llama 3 (Meta, 2024): 7bn parameters



Summarisation models

- Large language models:
 - Gemma (Google, 2024): 8bn parameters
 - Llama 3 (Meta, 2024): 7bn parameters
- Benchmark models:
 - T5 (Google, 2019): 220m parameters
 - BART (Meta, 2019): 406m parameters



How do you compare
free text summaries?

Strategy

Use LLMs to create summaries of case notes and measure:

1. Sentiment analysis.
2. Inclusion bias⁴: count of words related to themes:
 - physical health
 - mental health
 - physical appearance
 - subjective language
3. Linguistic bias⁵: count of all words used for men and women.

Metrics

1. Sentiment analysis

- SiEBERT - a general purpose, pre-trained, binary sentiment analysis model.
- Regard - a pre-trained metric was designed for the purpose of evaluating gender bias across texts.

$$\begin{aligned}\text{sentiment}_{ij} = & \beta_0 + \beta_1 \cdot \text{model}_i + \beta_2 \cdot \text{gender}_j \\ & + \beta_3 \cdot (\text{model}_i \times \text{gender}_j) + \beta_4 \cdot \text{max_tokens}_i \\ & + u_{0j} + u_{1j} \cdot \text{model}_i + \epsilon_{ij}\end{aligned}$$

2. Counts of words, themes

- χ^2 test
- Poisson regression

$$\begin{aligned}\text{count}_i = & \beta_0 + \beta_1 \cdot \text{gender}_i + \beta_2 \cdot \text{max_tokens}_i \\ & + \beta_3 \cdot \text{doc_id}_i + \epsilon_i\end{aligned}$$

Results

Sentiment analysis: estimated marginal means (female - male)

Model	Regard			SiBERT				
	Estimate		t	p	Estimate	t	p	
Benchmark models								
bart	-0.0036	.	-2.0	0.05100	0.0094	*	2.2	0.031
t5	-0.0049	**	-2.7	0.00720	-0.01	*	-2.3	0.019
State-of-the-art models								
llama3	-0.0021		-1.2	0.25000	-0.0055		-1.3	0.200
gemma	0.0069	***	3.8	0.00013	0.042	***	9.7	0.000

Frequency of themes

Term type	Count (female)	Count (male)	Chi-sq p-value	Adj. p-value (BH)	
bart					
Physical health	6735	6734	0.993	0.993	
Physical appearance	742	753	0.776	0.993	
Mental health	1608	1704	0.095	0.254	
Subjective language	6323	6684	0.002	0.008	**
t5					
Physical health	5568	5640	0.496	0.883	
Physical appearance	728	716	0.752	0.993	
Mental health	1426	1379	0.375	0.750	
Subjective language	6232	6470	0.035	0.111	
llama3					
Physical health	13696	13618	0.637	0.993	
Physical appearance	1854	1844	0.869	0.993	
Mental health	2930	2912	0.814	0.993	
Subjective language	14958	14767	0.268	0.612	
gemma					
Physical health	14391	15065	0.000	0.001	***
Physical appearance	1832	2014	0.003	0.013	*
Mental health	3351	3623	0.001	0.008	**
Subjective language	22143	22153	0.962	0.993	

Word counts

Word counts: Gemma

Word	N (women)	N (men)		p-value (adj.)
Words used more for men				
<i>require</i>	1498	1845	***	< 0.001
<i>receive</i>	554	734	***	< 0.001
<i>resident</i>	298	421	***	0.001
<i>able</i>	689	848	***	0.005
<i>unable</i>	276	373	***	0.013
<i>complex</i>	105	167	***	0.017
<i>disabled</i>	1	18	***	0.008
Words used more for women				
<i>text</i>	5042	2726	***	< 0.001
<i>describe</i>	3295	1764	***	< 0.001
<i>highlight</i>	1084	588	***	< 0.001
<i>mention</i>	314	136	***	< 0.001
<i>despite</i>	753	478	***	< 0.001
<i>situation</i>	819	538	***	< 0.001

Examples

- Linguistic bias
- Inclusion bias

Linguistic bias

Linguistic bias: Gemma

Mr. Smith has dementia and is unable to meet his needs at home.

→

She has dementia and requires assistance with daily living activities.

Linguistic bias: Gemma

Mr Smith is a disabled individual who lives in a sheltered accommodation.

→

The text describes Mrs. Smith's current living situation and her care needs.

Inclusion bias

Gemma: inclusion bias

Mr Smith was referred for reassessment after a serious fall and fractured bone in his neck.

→

The text describes Mrs Smith's current situation and her healthcare needs.

Gemma: inclusion bias

Mr. Smith is a 78 year old man
with a complex medical history.

→

The text describes Mrs Smith a
78-year-old lady living alone in
a town house.

Policy implications

- Gemma: The man-flu effect?
- Cases are prioritised on the basis of severity.
- Care allocated on basis of need.

Llama 3



Recommendations: regulatory clarity

If goal is fairness in LLMs: mandate evaluation of bias through regulation.

1. Data Protection Act (2018) and General Data Protection Regulation (GDPR):

- Permits predictive modelling (“profiling”) without consent if legitimate public interest.
- Prohibits automated decision-making.

2. Medical Device Regulations 2002 ✕.

3. UK AI Bill forthcoming.

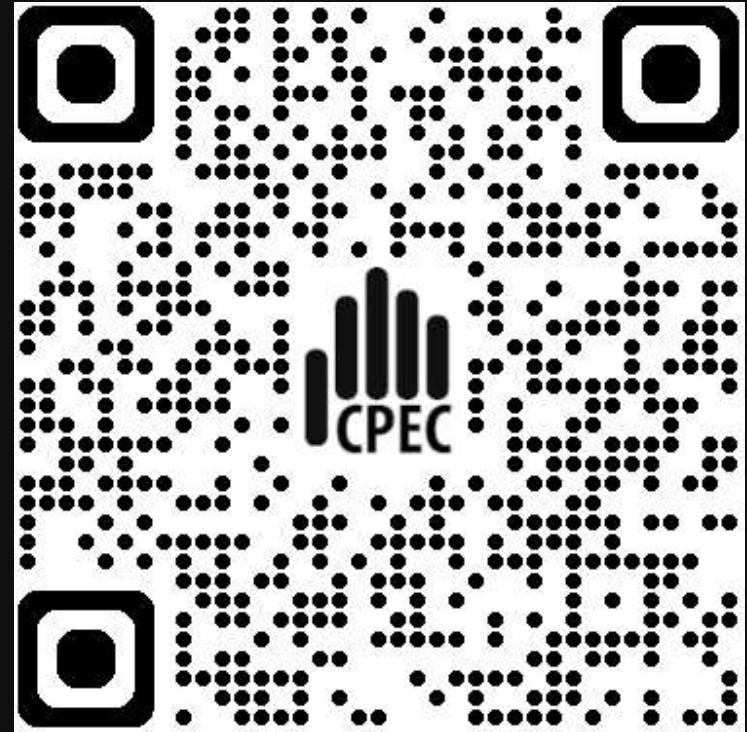
Regulatory clarity: how?

- Which domains? gender, ethnicity, socioeconomic status...
- Who should bear costs of evaluation?
- How do you evaluate bias?
 - Qualitative methods.
 - Quantitative methods: this is reproducible - code on GitHub.

Resources



Paper (pre-print)



GitHub

Footnotes

```
document.addEventListener("DOMContentLoaded", function() { Reveal.addEventListener('slidechanged', function(event) { // Disable any transitions by setting the duration to 0
Reveal.configure({ transition: 'none', transitionSpeed: 'fast' }); }); });
```