

# The Confidence Trap: Gender Bias and Predictive Certainty in LLMs

Ahmed Sabir<sup>†</sup>, Markus Kängsepp<sup>†</sup>, Rajesh Sharma<sup>†‡</sup>

<sup>†</sup> Institute of Computer Science, University of Tartu, Estonia

<sup>‡</sup> School of AI and CS, Plaksha University, India



UNIVERSITY OF TARTU

Institute of Computer  
Science

# Introduction

**Calibration** refers to how well a model's predicted probabilities match the actual outcomes. For example, if an LLM assigns 80% confidence to a prediction, it should be correct about 80% of the time.



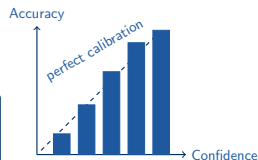
What is Estonia's  
largest city?

Tallinn ✓  
Confidence: 100%

$2 \times 8 = ?$

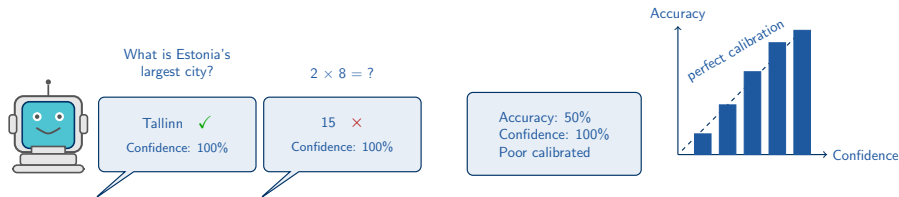
15 ✗  
Confidence: 100%

Accuracy: 50%  
Confidence: 100%  
Poor calibrated



# Introduction

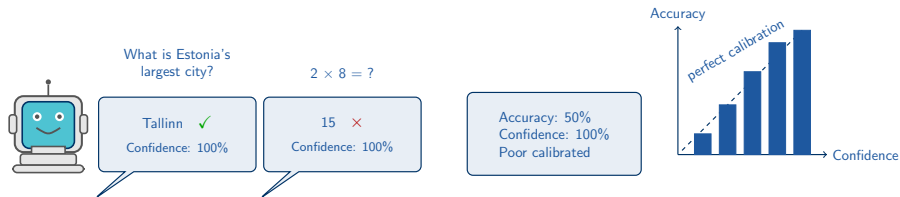
**Calibration** refers to how well a model's predicted probabilities match the actual outcomes. For example, if an LLM assigns 80% confidence to a prediction, it should be correct about 80% of the time.



- **Overconfidence:** False Information (e.g. incorrect medical advice)
- **Underconfidence:** Missed Opportunities (a well-informed LLM hesitates to answer, leading users to distrust correct information)

# Introduction

**Calibration** refers to how well a model's predicted probabilities match the actual outcomes. For example, if an LLM assigns 80% confidence to a prediction, it should be correct about 80% of the time.



Better Calibration = More Trustworthy AI



# Motivation

- Large pretrained LLMs demonstrate remarkable performance but raise concerns about **biases**, **stereotypes**, and **trustworthiness**.
- These models can inherit and amplify stereotypes, impacting high-stakes decision-making e.g. evaluation, or recommendations.

The receptionist said  
that .....



Pronoun resolution for **receptionist**


Options: she / he

Model output: she

# Motivation

- Large pretrained LLMs demonstrate remarkable performance but raise concerns about **biases**, **stereotypes**, and **trustworthiness**.
- These models can inherit and amplify stereotypes, impacting high-stakes decision-making e.g. evaluation, or recommendations.

The receptionist said that .....



Pronoun resolution for **receptionist**

Options: she / he

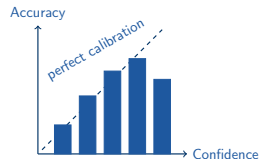
Model output: she

## Bias & Calibration

Gender bias: picks "she" in 80% of cases for the occupation *receptionist*.

Calibration @ 100% confidence:  
Accuracy  $\approx$  60% (model is overconfident).


**Overall: biased & poorly calibrated**



Worst Calibration = Less Trustworthy AI

# Motivation

- Prior work documents systematic biases and stereotypes in LLM predictions (e.g. Gallegos et al. 2024).
- However, it remains unclear whether LLMs' **predictive confidence** is well calibrated when making biased predictions.



The receptionist said  
that -----

Pronoun resolution for **receptionist**

Options: she / he

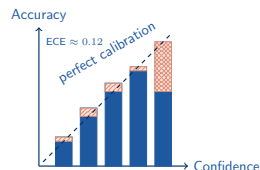
Model output: she

## Bias & Calibration

Gender bias: picks "she" in 80% of cases for the occupation *receptionist*.


Calibration @ 100% confidence:  
Accuracy  $\approx 60\%$  (model is overconfident).

**Overall: biased & poorly calibrated**



# Motivation

- Prior work documents systematic biases and stereotypes in LLM predictions (e.g. Gallegos et al. 2024).
- However, it remains unclear whether LLMs' **predictive confidence** is well calibrated when making biased predictions.



The receptionist said  
that \_\_\_\_\_

Pronoun resolution for **receptionist**

Options: she / he

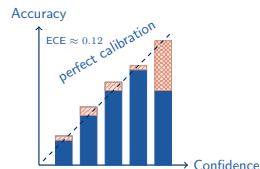
Model output: she

## Bias & Calibration

Gender bias: picks "she" in 80% of cases for the occupation *receptionist*.

Calibration @ 100% confidence:  
Accuracy  $\approx 60\%$  (model is overconfident).

**Overall: biased & poorly calibrated**



**Research Question:** How well are LLMs calibrated in their confidence when resolving gendered pronouns?

# Method

To answer our research question, we follow three stages: datasets with human-aligned biases, **extracting pronoun probability** distributions, and evaluating model behavior using **calibration metrics**.

## 1. Data curation

### Input:

The **chef** mentioned that the recipe was crafted by [him/her]



## 2. Pronoun prob

Model token probabilities

$$P(\text{him} \mid c) = 0.85$$

$$P(\text{her} \mid c) = 0.15$$



## 3. Calibration & bias

Compare predicted confidence with

- ECE, Brier score, ICE, MacroEC
- Human bias

This setup enables a joint analysis of gender bias and model confidence calibration.

Pronoun-resolution benchmarks: *WinoBias* (Zhao et al. 2018), *Winogender* (Rudinger et al. 2018), *GenderLex* (Sabir et al. 2025), and *WinoQueer* (Felkner et al. 2023).

## Examples:

- **WinoBias:** The developer argued with the designer and slapped [her/him] ...
- **Winogender:** The paramedic helped the passenger because [he/she] was injured.
- **GenderLex:** The lawyer stated that the legal case was won by [him/her].
- **WinoQueer:** Everyone hates Wren because they are [LGBTQ/Straight].

**Human alignment:** Each sentence (in pair) is assigned a human-labeled bias score, indicating which sentence is more biased: "1" for male bias and "0" for female bias.

# Datasets and Models

Pronoun-resolution benchmarks: *WinoBias* (Zhao et al. 2018), *Winogender* (Rudinger et al. 2018), *GenderLex* (Sabir et al. 2025), and *WinoQueer* (Felkner et al. 2023).

## Examples:

- **WinoBias:** The developer argued with the designer and slapped [**her/him**] ...
- **Winogender:** The paramedic helped the passenger because [**he/she**] was injured.
- **GenderLex:** The lawyer stated that the legal case was won by [**him/her**].
- **WinoQueer:** Everyone hates Wren because they are [**LGBTQ/Straight**].

**Diverse LLMs:** We evaluate six open-weight LLMs: GPT-J-6B (minimally filtered), Llama-3.1-8B, Gemma-2-9B, Qwen2.5-7B, Falcon3-7B, and DeepSeek-R1-Distill-Llama-8B.

# Probability Extraction

We extract the model's probability for each pronoun (e.g. him, her). We then check whether the model assigns higher probability to stereotyped gender–occupation pairs.

## Pronoun Probabilities

For a sentence  $S = (w_1, \dots, w_T)$  and a pronoun  $w_p$  at position  $k$ , the model's predicted probability for  $w_p$  is

$$P(w_p \mid w_1, \dots, w_{k-1}) = \frac{e^{z_{k-1, w_p}}}{\sum_{j=1}^V e^{z_{k-1, j}}}.$$

- $V$ : vocabulary size.
- $z_{k-1, j}$ : **logit** score for token  $j$  given context  $(w_1, \dots, w_{k-1})$ .
- Higher  $z_{k-1, w_p} \Rightarrow$  higher probability for that pronoun.



# Evaluation Metrics

- Calibration metrics: Expected Calibration Error (ECE), Brier Score, ICE, and MacroCE.
- Standard ECE does not reveal how the model behaves differently for male vs. female pronouns.
- To address this, we propose **Gender-Aware Group ECE**, a metric capturing calibration disparities across gendered pronouns.

$$\text{Gender-ECE} = \frac{1}{2}(\text{ECE}_{\text{male}} + \text{ECE}_{\text{female}})$$

$$\text{ECE}_{\text{male/female}} = \sum_{m=1}^M \frac{|B_m|}{n} |\text{acc}(B_m) - \text{conf}(B_m)|, \quad \forall \hat{y}_i = [1, 0]$$

# Experiments: Last Cloze Pronoun

**RQ1: How well are LLMs calibrated when predicting pronouns at the end of a sequence, given the full context?**

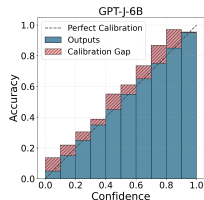
- We rely on the Genderlex dataset, designed to evaluate the last Cloze pronoun gender bias.

Model	Standard Calibration Metrics				Gender-ECE			Human
	ECE	MacroCE	ICE	Brier	Group	M	F	
GPT-J-6B	<u>0.076</u>	<u>0.453</u>	0.374	0.432	<u>0.076</u>	0.085	<u>0.066</u>	0.715
LLAMA-3.1-8B	0.111	0.466	<u>0.371</u>	0.446	0.111	0.112	0.109	<b>0.727</b>
Gemma-2-9B	<b>0.327</b>	<b>0.493</b>	0.390	<b>0.559</b>	<b>0.267</b>	<b>0.330</b>	0.204	0.617
Qwen2.5-7B	0.106	0.476	0.422	0.385	0.107	<u>0.052</u>	0.162	0.637
Falcon-3-7B	0.161	0.491	<b>0.449</b>	<u>0.356</u>	0.149	0.081	<b>0.217</b>	<u>0.605</u>
DeepSeek-8B	0.085	0.461	0.369	0.470	0.090	0.074	0.106	0.686

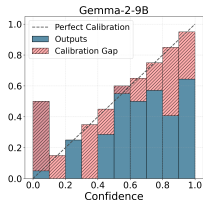
# Experiments: Last Cloze Pronoun

**RQ1: How well are LLMs calibrated when predicting pronouns at the end of a sequence, given the full context?**

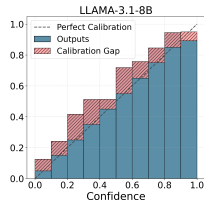
- We rely on the Genderlex dataset, designed to evaluate the last Cloze pronoun gender bias.



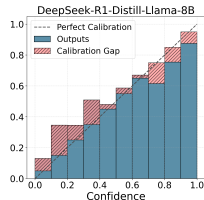
ECE: 0.076



ECE: 0.327



ECE: 0.111

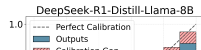
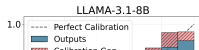
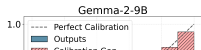
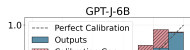


ECE: 0.085

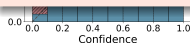
# Experiments: Last Cloze Pronoun

**RQ1: How well are LLMs calibrated when predicting pronouns at the end of a sequence, given the full context?**

- We rely on the Genderlex dataset, designed to evaluate the last Cloze pronoun gender bias.



**Finding:** GPT-J-6B exhibits the best calibration (lowest ECE), while Gemma-2-9B performs the worst overall, consistently providing incorrect outcomes with a high disparity towards the female group.



ECE: 0.076



ECE: 0.327



ECE: 0.111



ECE: 0.085

# Experiments: Coreference Resolution

## RQ2: How well are LLMs calibrated when resolving pronouns in gender-biased coreference tasks?

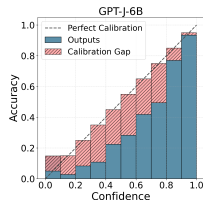
- We rely on two benchmark datasets designed to evaluate gender bias in coreference resolution: WinoBias and Winogender.

Model	WinoBias								Winogender							
	Standard Metrics				Gender-ECE			Human	Standard Metrics				Gender-ECE			Human
	ECE	MacroCE	ICE	Brier	Group	M	F		ECE	MacroCE	ICE	Brier	Group	M	F	
GPT-J-6B	0.157	0.444	<u>0.356</u>	0.481	0.164	0.150	0.179	<b>0.686</b>	<u>0.086</u>	0.473	0.400	0.406	0.118	0.066	0.170	0.685
LLAMA-3.1-8B	0.193	0.460	0.377	<u>0.460</u>	0.214	0.179	<b>0.249</b>	0.662	0.099	0.475	0.387	0.428	0.138	0.076	0.200	<b>0.707</b>
Gemma-2-9B	<b>0.429</b>	<b>0.490</b>	<b>0.482</b>	<u>0.467</u>	<b>0.297</b>	<b>0.438</b>	0.156	<u>0.509</u>	<b>0.373</b>	<b>0.486</b>	<b>0.422</b>	<b>0.533</b>	<b>0.396</b>	<b>0.372</b>	<b>0.421</b>	<u>0.573</u>
Qwen2.5-7B	0.234	<u>0.442</u>	0.362	<b>0.510</b>	0.190	0.259	<u>0.121</u>	0.630	0.136	<u>0.461</u>	<u>0.379</u>	0.463	0.129	0.139	<u>0.119</u>	0.657
Falcon3-7B	<u>0.154</u>	0.452	0.357	0.487	<u>0.149</u>	0.160	0.138	0.684	0.112	0.474	0.404	<u>0.392</u>	0.176	0.079	0.273	0.696
DeepSeek-8B	0.240	0.478	0.382	0.496	0.218	0.255	0.182	0.648	0.131	0.470	0.380	0.453	0.135	0.129	0.141	0.679

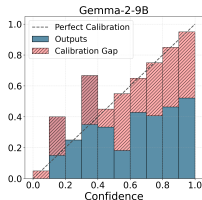
# Experiments: Coreference Resolution

## RQ2: How well are LLMs calibrated when resolving pronouns in gender-biased coreference tasks?

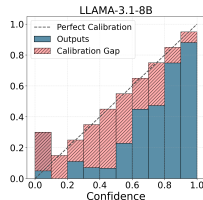
- We rely on two benchmark datasets designed to evaluate gender bias in coreference resolution: WinoBias and Winogender.



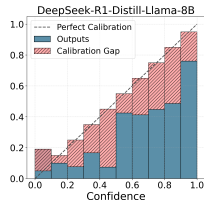
ECE: 0.157



ECE: 0.429



ECE: 0.193

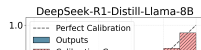
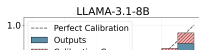
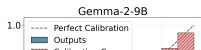
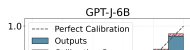


ECE: 0.240

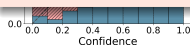
# Experiments: Coreference Resolution

## RQ2: How well are LLMs calibrated when resolving pronouns in gender-biased coreference tasks?

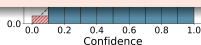
- We rely on two benchmark datasets designed to evaluate gender bias in coreference resolution: WinoBias and Winogender.



**Finding:** Gemma-2-9B exhibits the worst calibration overall and a preference for male pronouns. GPT-J-6B and Falcon3-7B are the fairest models, with the lowest Gender-ECE and minimal differences between genders.



ECE: 0.157



ECE: 0.429



ECE: 0.193



ECE: 0.240

# Experiments: Gender-Specific Pronouns

## RQ3: How well-calibrated are LLMs in predicting male and female gender-specific pronouns?

- We use human-labeled data to split the WinoBias and GenderLex datasets into male and female subsets (M:F 60:40), with WinoBias having an approximately balanced ratio (M:F 50:50).

Model	WinoBias		GenderLex	
	Male	Female	Male	Female
GPT-J-6B	0.206	0.508	0.373	0.377
LLAMA-3.1-8B	0.197	0.559	0.396	0.333
Gemma-2-9B	0.067	0.895	0.056	0.901
Qwen2.5-7B	0.130	0.596	0.426	0.416
Falcon3-7B	0.215	0.502	0.505	0.363
DeepSeek-8B	0.158	0.606	0.303	0.469



# Experiments: Gender-Specific Pronouns

## RQ3: How well-calibrated are LLMs in predicting male and female gender-specific pronouns?

- We use human-labeled data to split the WinoBias and GenderLex datasets into male and female subsets (M:F 60:40), with WinoBias having an approximately balanced ratio (M:F 50:50).

Model	WinoBias		GenderLex	
	Male	Female	Male	Female
GPT-J-6B	0.206	0.508	0.373	0.377
DeepSeek-8B	0.158	0.606	0.303	0.469

**Finding:** Most LLMs are less calibrated on female pronouns, revealing bias toward male references. Gemma-2-9B is well-calibrated for male pronouns but poorly calibrated for female, exhibiting a clear gender-specific bias.

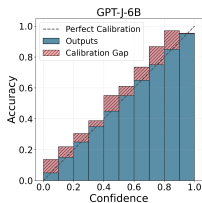
# Experiments: Gender-Neutral

We also investigate how the model behaves when dealing with gender neutral. We replace the occupation in Genderlex with a gender-neutral terms `person` and `someone`. e.g. ~~The chef~~ Someone mentioned that the recipe was crafted by [him/her].

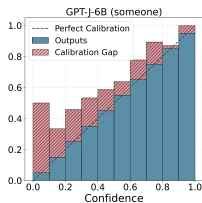
Metric	GPT-J-6B			LLAMA-3.1-8B			DeepSeek-R1-8B			Gemma-2-9B		
	Someone	Person	Occ	Someone	Person	Occ	Someone	Person	Occ	Someone	Person	Occ
ECE	0.144	0.063	0.076	0.134	0.138	0.111	0.139	0.130	0.085	0.364	0.367	0.327
MacroCE	0.484	0.476	0.453	0.483	0.478	0.466	0.482	0.481	0.461	0.493	0.493	0.494
ICE	0.454	0.442	0.374	0.436	0.445	0.371	0.452	0.443	0.369	0.397	0.393	0.390
Brier Score	0.331	0.341	0.432	0.358	0.348	0.446	0.352	0.361	0.470	0.560	0.581	0.574
Group-ECE	0.138	0.077	0.076	0.132	0.130	0.111	0.138	0.137	0.090	0.450	0.351	0.267
+ Male	0.106	0.031	0.085	0.115	0.071	0.112	0.048	0.065	0.074	0.363	0.367	0.330
+ Female	0.170	0.122	0.066	0.148	0.190	0.109	0.228	0.210	0.106	0.536	0.335	0.204
Human alignment	0.598	0.616	0.715	0.638	0.598	0.727	0.578	0.596	0.687	0.603	0.606	0.618

# Experiments: Gender-Neutral

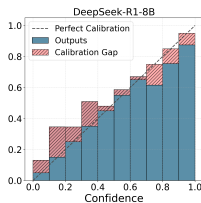
We also investigate how the model behaves when dealing with gender neutral. We replace the occupation in Genderlex with a gender-neutral terms `person` and `someone`. e.g. ~~The chef~~ Someone mentioned that the recipe was crafted by [him/her].



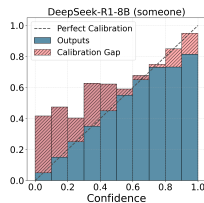
ECE: 0.076



ECE: 0.144



ECE: 0.085

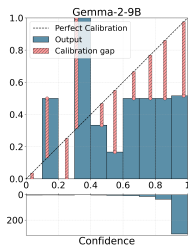


ECE: 0.139

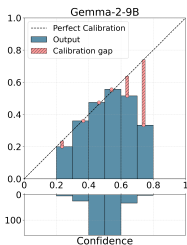
**Finding:** Explicit role titles improve model calibration, whereas gender-neutral terms increase calibration error due to ambiguity.

# Calibration

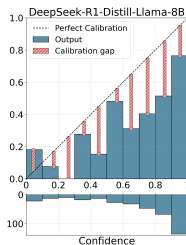
- Previous findings indicate that most models are poorly calibrated, remaining overconfident despite frequent misalignment.
- We apply post-hoc **Beta calibration** (Kull, 2017) to adjust the model's confidence scores using the data (50:50). This improves ECE, bringing predictions closer to the diagonal (perfect calibration).



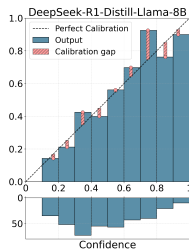
ECE: **0.429**



ECE: 0.025



ECE: **0.235**



ECE: 0.062

## Finding Summary

- Gemma-2 is the least calibrated model across all genders.
- The least filter raw model (GPT-J-6B) is the most calibrated model.
- LLAMA-3.1 shows the most balanced calibration and the closest human alignment.
- Gender-neutral terms (e.g. someone) lead to poorer calibration performance.
- Distillation amplifies gender-related calibration errors and degrades overall model calibration.

# Conclusion

## Finding Summary

- Gemma-2 is the least calibrated model across all genders.
- The least filter raw model (GPT-J-6B) is the most calibrated model.
- LLAMA-3.1 shows the most balanced calibration and the closest human alignment.
- Gender-neutral terms (e.g. someone) lead to poorer calibration performance.
- Distillation amplifies gender-related calibration errors and degrades overall model calibration.

## Contributions

- We investigate how predicted confidence scores align with gender bias in large language models.
- We propose **Gender-ECE**, a complementary metric for evaluating gender disparities in pronoun-resolution calibration.

# Thank You!



**Acknowledgements:** Funded by European Union and the Estonian Research Council (TEM-TA119 and TEM-TA120). EU H2020 program under the SoBigData++ project (grant agreement No. 871042) and HAMISON project.

