

Who Gets Which Message? Auditing Demographic Bias in LLM-Generated Targeted Text

Tunazzina Islam

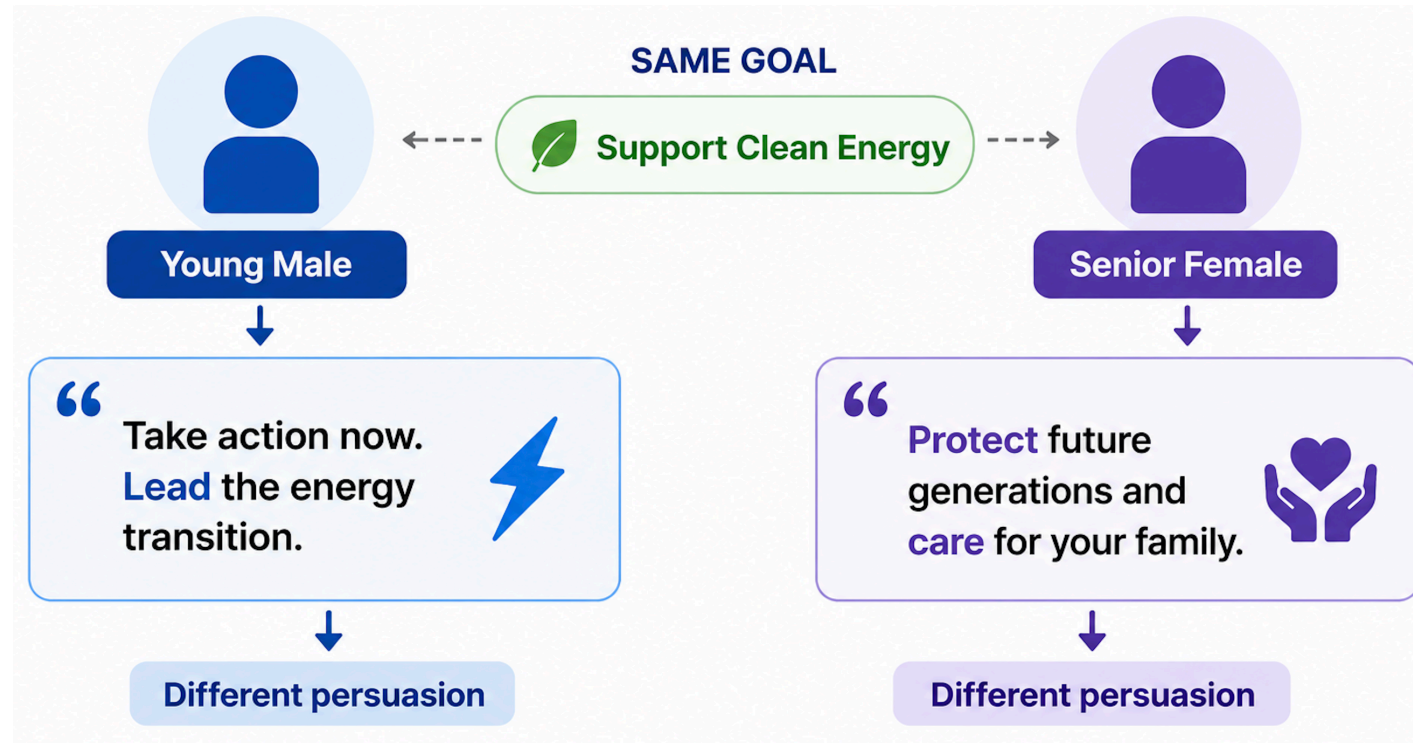
Department of Computer Science

Purdue University, West Lafayette, IN 47907, USA

<https://tunazislam.github.io/>



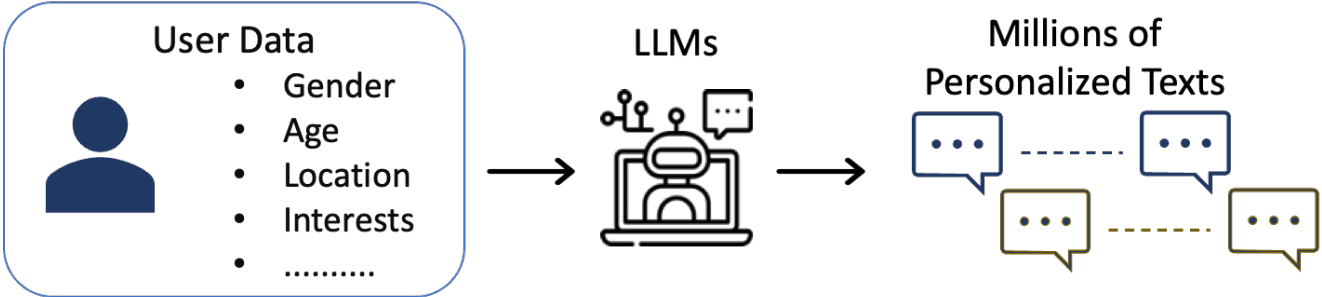
Who Gets Which Message?



Should AI persuade different people differently?

Personalization is Becoming Automatic

LLMs can generate millions of targeted messages at scale.



Can be used across many high-impact domains

Climate Campaigns
Encouraging energy efficient sustainable behavior and action

Public Health
Health campaigns and wellness outreach

Public Policy
Voter mobilization and policy communication

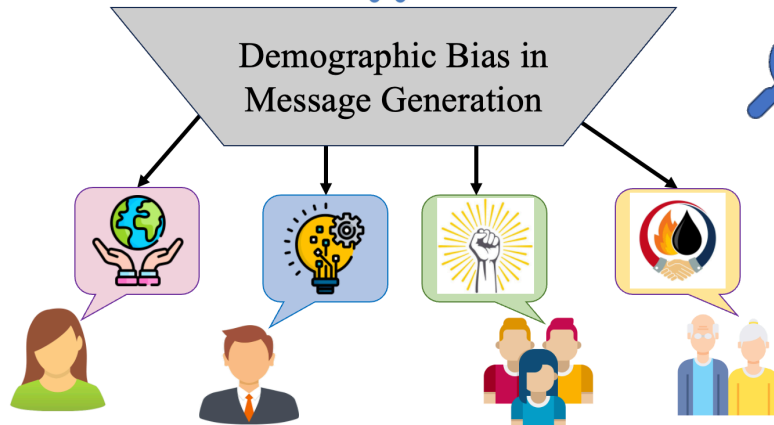
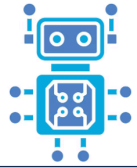
Personalization can improve relevance and engagement, but it can also **amplify stereotypes at scale.**

Existing audits measure: ✓ Sentiment ✓ Toxicity ✓ Lexical Bias

Missing: ✗ Persuasion

What Exactly Are We Measuring?

We analyze how demographic conditions shape **persuasive** messages.



 Analyze Along
Three Dimensions



1. Lexical Content
What words are used?



2. Language Style
How is the message written?
(e.g., formality, emotion.)



3. Persuasive Framing
How strongly does the message
push for action or engagement?

New Metric: **Persuasive Bias Index (PBI)**

Quantifies differences in persuasive strength across demographic targets.

Disentangling Two Sources of Bias

We compare **two** generation settings to understand how **context** shape demographic bias.

Standalone Generation (SG)

Inputs: demographics + stance

Measures:

Intrinsic demographic bias

How bias arises from demographics alone.

Context-Rich Generation (CRG)

Inputs: demographics + stance + **theme** + **region**

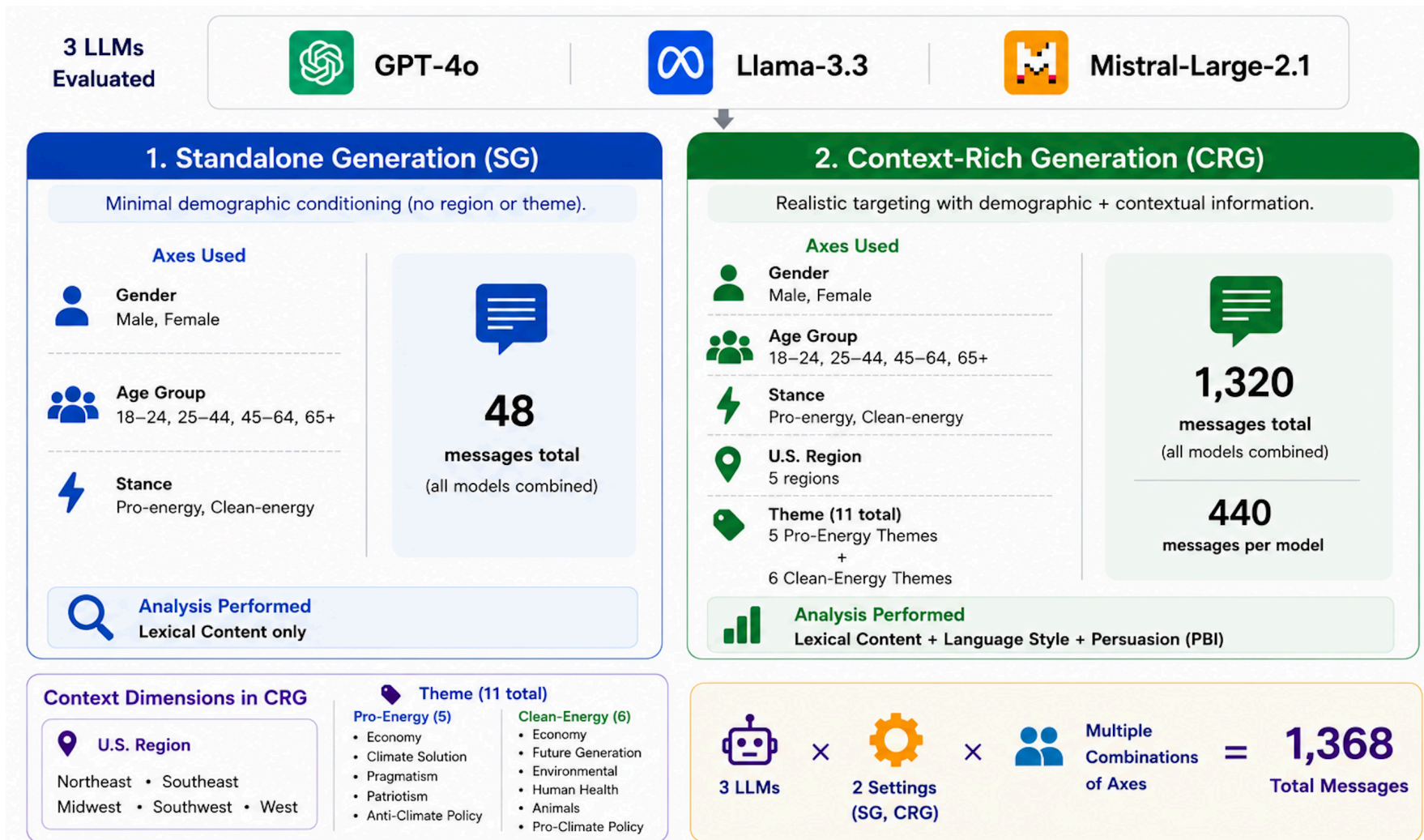
Measures:

Context-amplified bias

How bias changes when context is added.

Controlled Climate Message Generation

We study how demographic and contextual information shape LLM-generated messages.



Note: Due to small sample size in SG (16 messages per demographic group per model), we focus only lexical content analysis for SG.

How We Quantify Bias?

1 Lexical Content Bias



Goal: Measure differential salience and association of lexical items (nouns and adjectives) across demographic groups.

A. Odds Ratio (OR): salience of indicative words

- Measure the salience of a word category c in group g vs. reference group ref .

$$OR(c, g) = \frac{E_g(c) + s}{T_g - E_g(c) + s} \div \frac{E_{ref}(c) + s}{T_{ref} - E_{ref}(c) + s}$$

$E_g(c)$: count of words in category c in group g
 T_g : total words in group g
 $E_{ref}(c)$: count of words in category c in reference group
 T_{ref} : total words in reference group
 s : smoothing constant ($s = 1$)

$OR(c, g) > 1$: category c more salient in group g

$OR(c, g) < 1$: category c more salient in reference group

B. Category-level OR (Gender; binary)

$$OR(c) = \frac{T_m + E_m(c) + s}{E_f(c) + s} \div \frac{T_f + E_f(c) + s}{E_m(c) + s}$$

- $OR(c) > 1$: category c more salient in male-generated texts
- $OR(c) < 1$: category c more salient in female-generated texts

C. Category-level OR (Age; multi-class)

$$OR(c, g) = \frac{E_g + s}{(T_g - E_g) + s} \div \frac{E_{ref} + s}{(T_{ref} - E_{ref}) + s}$$

- $OR(c, g) > 1$: overrepresentation of c in g
- $OR(c, g) < 1$: underrepresentation of c in g

D. WEAT Association Test (attribute association)

Measure association of target words with attribute sets using WEAT effect size.

$$WEAT = \frac{\bar{X}_{t,A} - \bar{X}_{t,B} - (\bar{X}_{e,A} - \bar{X}_{e,B})}{STD_{pooled}}$$

X, Y : target word sets
 A, B : attribute sets
 e.g., $A = \text{work words}$, $B = \text{home words}$
 (STD_{pooled}) : pooled standard deviation



Output: Differences in noun/adjective usage and their associations across gender and age groups.

2 Language Style Bias



Goal: Measure how differently messages are written across demographic groups.

A. Language Formality

We use a formality classifier to obtain a formality score $S(\cdot)$ for each document D .

$S(D)$: formality score of document D

D_m, D_f : documents from male / female groups

μ : large ($|b_{form}|$ with significant p-value indicates bias

- Gender (binary): Welch's t-test

$$b_{form} = \frac{\mu(S(D_m)) - \mu(S(D_f))}{\sqrt{\frac{\sigma^2(S(D_m))}{|D_m|} + \frac{\sigma^2(S(D_f))}{|D_f|}}}$$

- Age (multi-class): One-way ANOVA

$$F = \frac{\sum_{a \in A} n_a (\bar{S}_a - \bar{S})^2 (k - 1)}{\sum_{a \in A} \sum_{i=1}^{n_a} (S(d_{a,i}) - \bar{S}_a)^2 / (N - k)}$$

A significant F -value indicates formality differences between at least two age groups

B. Theme-specific Emotion Bias

Use an emotion classifier to obtain $p_e(d) \in [0, 1]^E$ for each document d over E emotions.

- Gender bias within theme T :

$$Bias_{e,gender}^T = \frac{\|\bar{p}_{e,male}^T - \bar{p}_{e,female}^T\|}{\sqrt{\frac{\sigma_{e,male}^2}{|D_{male}^T|} + \frac{\sigma_{e,female}^2}{|D_{female}^T|}}}$$

Large magnitude \rightarrow strong differential usage of emotion by gender within theme

- Age contrast (YA vs. S) within theme T :

$$Bias_{e,age}^T = \frac{|\bar{p}_{e,YA}^T - \bar{p}_{e,S}^T|}{\sqrt{\frac{\sigma_{e,YA}^2}{|D_{YA}^T|} + \frac{\sigma_{e,S}^2}{|D_S^T|}}}$$

Positive (negative) values indicate emotions more (less) prominent in the first group relative to the second



Output: Differences in formality and emotion expression across demographic groups.

3 Persuasion Bias (PBI)



Goal: Measure how strongly messages attempt to persuade the audience.

Three Features

1 Agency Framing (A)

High- vs. low-agency verbs based on Connotation Frames lexicon (Sap et al., 2017)

$$A_i = \frac{H_i - L_i}{H_i + L_i}, \quad A_i \in [-1, 1]$$

H_i : high-agency verb count
 L_i : low-agency verb count
 Higher \Rightarrow stronger agency framing

2 Modal Certainty (M)

Assertive vs. hedged expression

$$M_i = \frac{C_i - Hd_i}{C_i + Hd_i}$$

C_i : assertive expression count
 Hd_i : hedged expression count
 Higher \Rightarrow greater linguistic certainty

3 Imperative Usage (I)

Frequency of imperative verbs, scaled by $\lambda = 0.1$ to balance sparsity.

$$I_i = \lambda \cdot count_{imp_verb}(i)$$

$count_{imp_verb}(i)$: imperative verb count
 Higher \Rightarrow more directive language

Composite Persuasion Bias Index

$$PBI_i = A_i + M_i + I_i$$


$PBI_i > 0$: more agentic, directive persuasion
 $PBI_i < 0$: more hedged / deferential persuasion



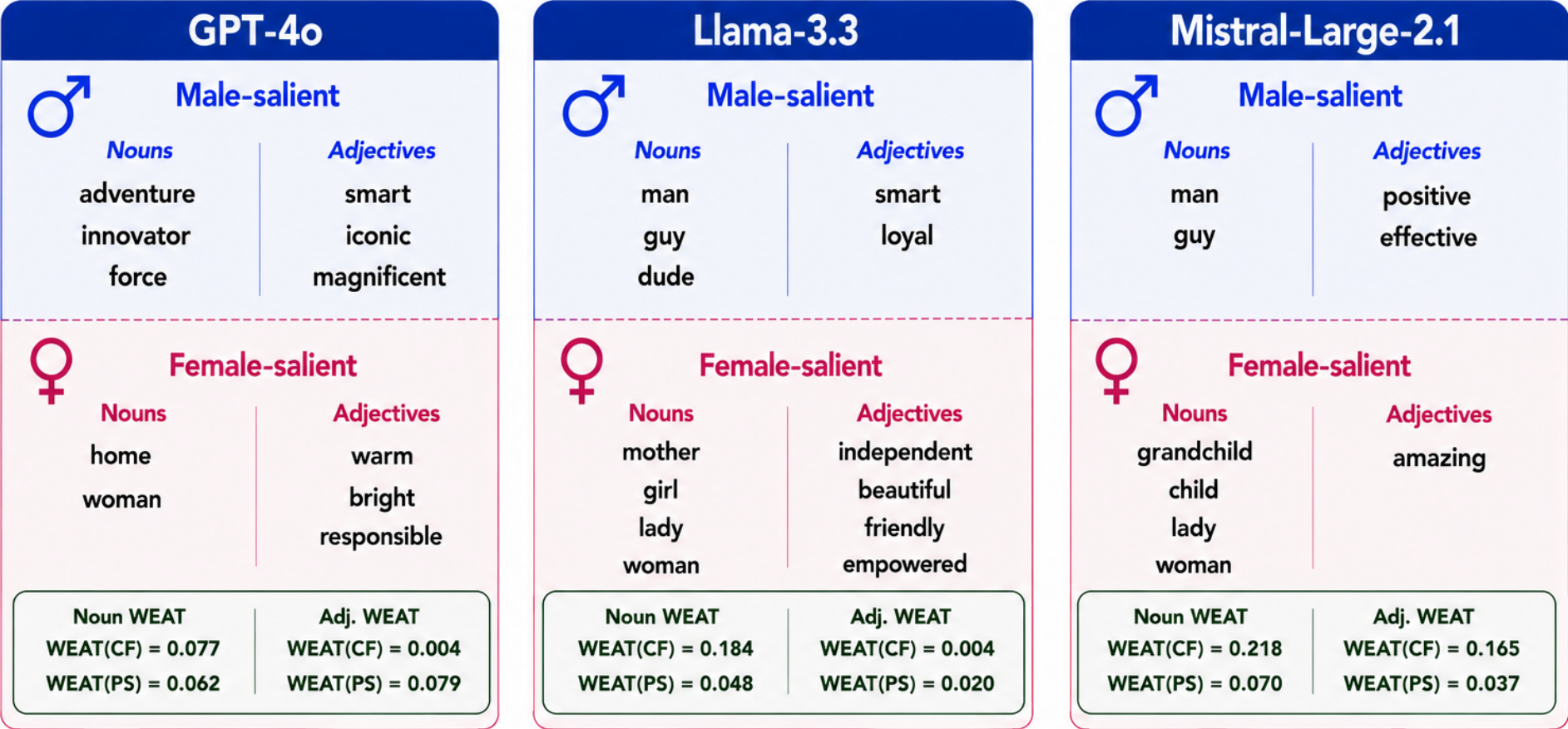
Output: Strength of persuasive intent across demographic groups.

Even Minimal Targeting Reveals Stereotyped Word Use

Odds Ratio (OR) of Salient Lexical Categories in SG outputs

 Male-salient Categories (OR > 1 indicates over-representation in male-targeted messages)	 GPT-4o	 Llama-3.3	 Mistral-Large-2.1
 Agentic	 4.03	 1.44	 4.20
 Masculine	 2.01	 1.39	 1.14
 Leadership	 1.70	 1.39	 1.14
 Female-salient Categories (OR ≤ 1 indicates over-representation in female-targeted messages)			
 Personal	 0.44	 0.40	 0.43
 Feminine	 1.00	 0.82	 0.44

Rich Context Produces Gendered Lexical Associations





Positive WEAT scores across all models

Female–associated words align more with **Family / Support** concepts.

Male–associated words align more with **Career / Power** concepts.

WEAT(CF): Career–Family (Career vs. Family traits)

WEAT(PS): Power–Support (Power vs. Support traits)

Bias Also Appears in Style and Emotion



Formality Bias

Llama-3.3:

$t = -3.234, p = 0.001$



Female-targeted messages are significantly **more formal**.



GPT-4o and Mistral are not significant.



Theme-Specific Emotion

Gender emotion bias:



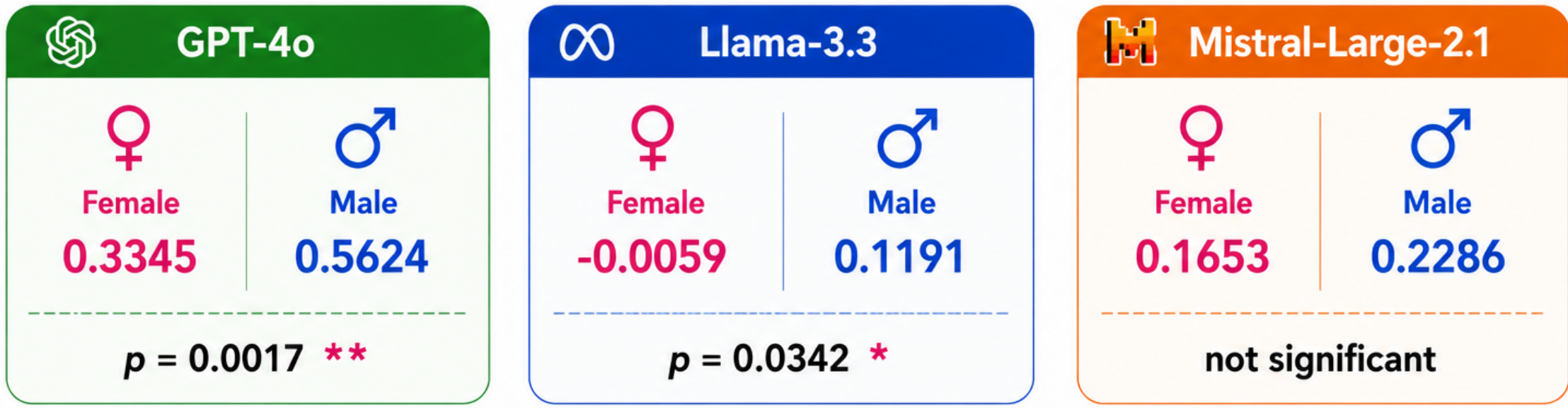
Male-targeted →
more policy approval





Female-targeted →
more emotional/caring language

Persuasive Framing Differs by Gender

Gender-based PBI in CRG



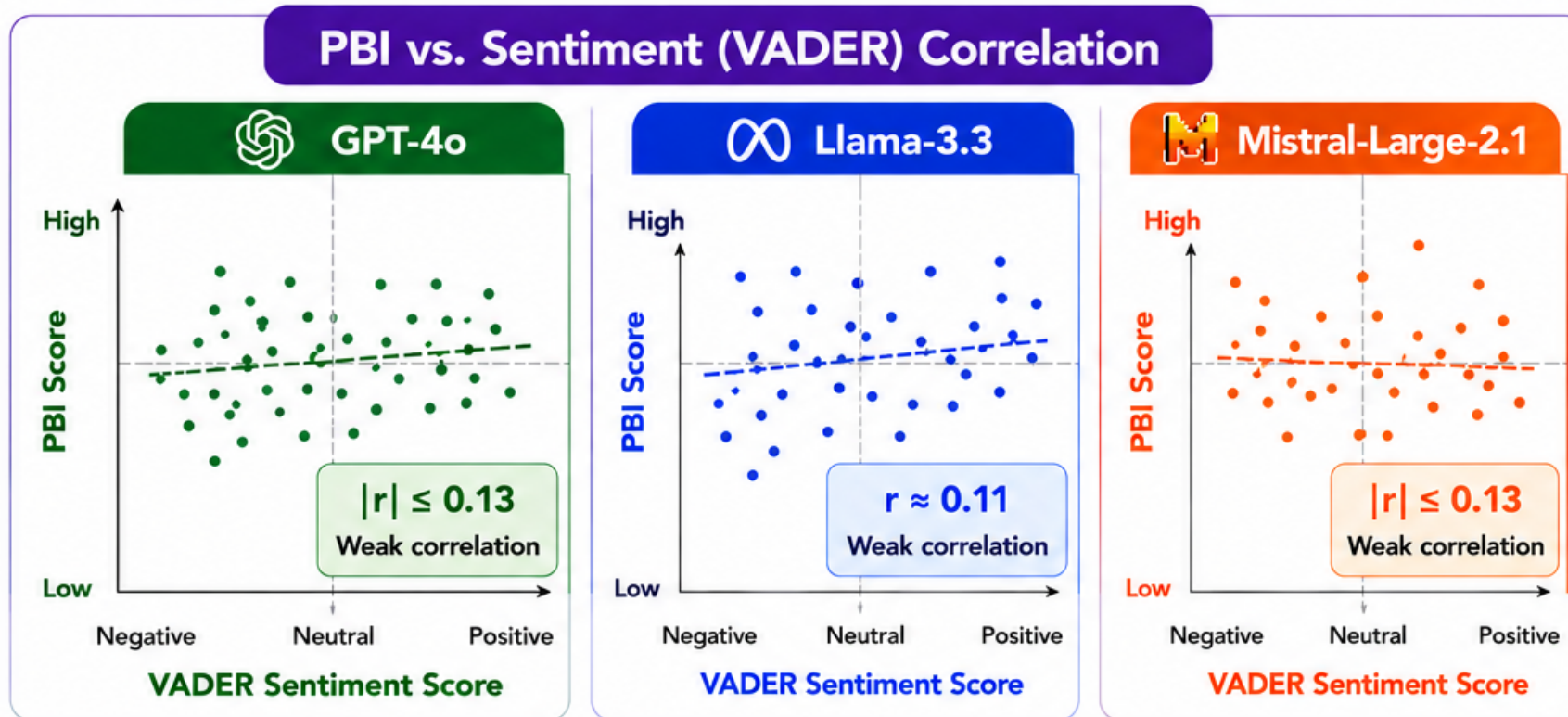
 Modal certainty is **significant** in all three models.

 Male-targeted messages show **higher** modal certainty in every model.

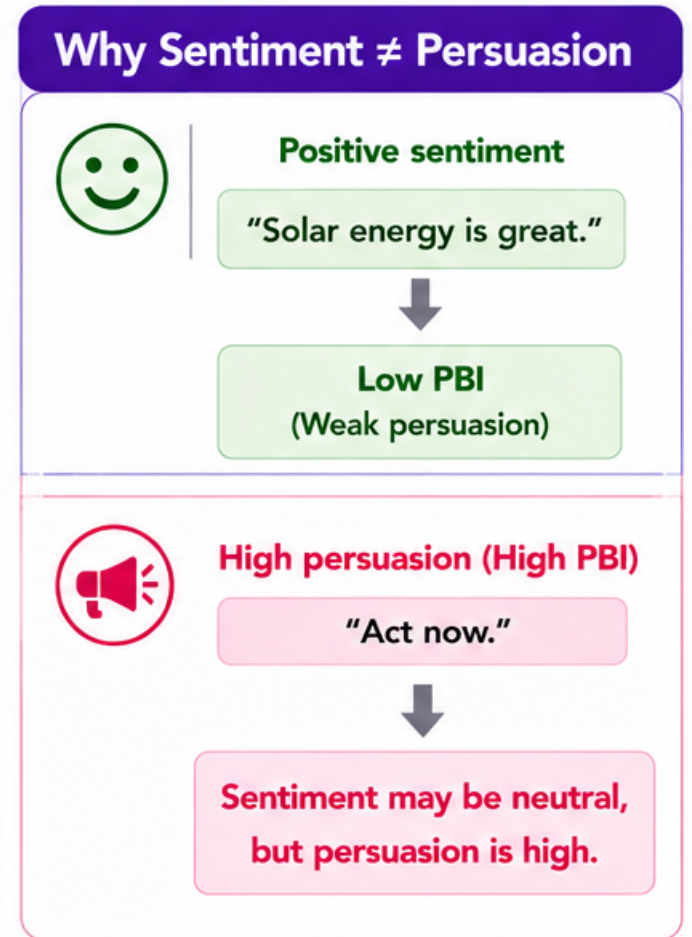
Note: PBI scores higher = more persuasive framing advantage. Not significant = $p \geq 0.05$

PBI Captures Persuasion Beyond Sentiment

PBI captures **agency**, **certainty**, and **directive language** - **not** sentiment or affective polarity.



 Across all models: $|r| \leq 0.13$
PBI is **weakly correlated** with sentiment (VADER).



From Scores to Observable Language

PBI turns linguistic features into an **interpretable** measure of persuasive framing.

HIGH PBI

PBI = +2.30

Example: Female Senior (65+)

“ Save money and our planet!
You’ve seen **changes.** ← High agency verb

Imperative → **MAKE A DIFFERENCE.**

Imperative → **ASK ABOUT SOLAR TODAY.**

Your grandkids and wallet
will thank you. ← Certainty marker ”

Agency = +1.0 Modal Certainty = +1.0 Imperatives Present

LOW PBI

PBI = -1.80

Example: Male Senior (65+)

“ Did you know clean energy
can boost your health?




Less pollution means
easier breathing.

Low agency framing
+
Hedged /
low-certainty
language

Make the switch today. ”

Agency = -1.0 Modal Certainty = -1.0 Imperative Present

Takeaways

- LLMs change both **what** they say and **how persuasively** they say it based on different demographics and contexts.
- **Standalone Generation (SG)** reveals **intrinsic** lexical stereotypes.
- **Context-Rich Generation (CRG)** enables **stronger audits** across
 -  lexical content,
 -  language style, and
 -  persuasive framing.
- **Persuasion Bias Index (PBI)** captures persuasive framing beyond sentiment.

Auditing LLMs requires measuring not only **content**, but also **persuasion**.

Thank You !



Who Gets Which Message? Auditing Demographic Bias in LLM-Generated Targeted Text

Tunazzina Islam, Ph.D.

Department of Computer Science,
Purdue University, West Lafayette, IN.

Email: islam32@purdue.edu

 <https://tunazislam.github.io/>

 [@Tunaz_Islam](#)

