

Can Al Agents Fit in Human Society?

Linqi Liu (AISTN, Iqliu1@cse.cuhk.edu.hk) Yuhang Yan (CSCIN, yhyan2@cse.cuhk.edu.hk)

Date: April 22nd, 2025

Supervisor: Prof. Michael R. LYU

Department of Computer Science and Engineering

The Chinese University of Hong Kong



香港中文大學計算機科學與工程學系 Department of Computer Science and Engin The Chinese University of Hong Kong



香港中文大學 The Chinese University of Hong Kong



- FACT-OR-FAIR Recap
- Extended Experiments

• LLM Search Testing

Conclusion & Future Work

FACT-OR-FAIR Recap

ONE





FACT-OR-FAIR: Evaluating Factuality and Fairness in Al Models

- Background and Motivation
 - Generative AI struggles to balance factuality and fairness.
 - For example, Gemini generated controversial images, revealing need for better evaluation tools.
- Main Contribution
 - o Data Framework: 19 statistics collected
 - **Test Design**: Objective and bias-triggering scenarios
 - Metrics: Factuality-fairness trade-off
 - **Experiments**: 6 LLMs and 4 T2I models



Asian Popes and Black Vikings Generated by Gemini^[1]





- Definitions of Factuality and Fairness:
 - Factuality^[2]: The ability of a generative model to produce content that aligns with established facts and world knowledge.
 - Fairness^[3]: The guarantee that algorithmic decisions remain unbiased, irrespective of individual attributes such as gender or race.
- Three cognitive biases:
 - Representativeness Bias^[4]: Individuals or situations based on the mental prototype of a certain group.
 - Attribution Error^[5]: Overestimating internal traits and underestimating situational factors when explaining people's behaviors. Mistakenly attributing individual behavior to the entire group's internal characteristics.
 - In-group/Out-group Bias^[6]: Favoring one's own group (in-group) while being critical of others (out-groups).
- [2] Y Wang et al. "Factuality of Large Language Models: A Survey" EMNLP 2024
- [3] M Hardt et al. "Equality of opportunity in supervised learning" NeurIPS 2016
- [4] D. Kahneman et al. "Subjective probability: A judgment of representativeness" Cognitive Psychology 1972
- [5] T.F. Pettigrew. "The ultimate attribution error: Extending Allport's cognitive analysis of prejudice." Personality and Social Psychology Bulletin 1979
- [6] M.B. Brewer. "In-group bias in the minimal intergroup situation: A cognitive-motivational analysis." *Psychological bulletin 1979*



Objective Queries

 \circ LLMs

- Designed to test factual knowledge.
- Prompt format includes definition + multiple choice.
- \circ T2I Models
 - Asked to generate portraits based on statistical facts.
 - Prompt includes statistic definition and desired target.
 - Outputs analyzed using automated detection tool.
- Subjective Queries

 $\circ \ \text{LLMs}$

- Designed to test fairness under realistic contexts.
- Use controlled scenarios involving race and gender profiles.
- Three bias types embedded into prompts: Representativeness Bias (uses prior statistics), Attribution Error (uses anecdotal information), In-group / Out-group Bias (changes user identity).

 \circ T2I Models

- Given stereotype-sensitive prompts without specific priors.
- Asked to generate portraits under vague or open-ended conditions.





Testing with objective queries that require accuracy.



Testing with subjective queries that require diversity.





- 19 real-world U.S. statistics from trusted sources (BLS, CDC, USCB).
 - For FACT-OR-FAIR, Each statistic provides separate data for different gender and racial groups.
 - For LLM testing, we used data from 15 countries and included Birth Rate, with all statistics presented without race or gender breakdowns.
- Categorized into economic, social, and health domains.
- Post-processed into demographic rates (e.g., obesity rate, crime rate, etc.)

	Statistics	Source	Definition
	Employment Rate	BLS [2024]	Percentage of employed people.
nic	Unemployment Rate	BLS [2024]	Percentage of unemployed people who are actively seek-
⁰			ing work.
8	Weekly Income	BLS [2024]	Average weekly earnings of an individual.
ĕ	Poverty Rate	KFF [2022]	Percentage of people living below the poverty line.
	Homeownership Rate	USCB [2024]	Percentage of people who own their home.
	Homelessness Rate	CPD [2023]	Percentage of people experiencing homelessness.
	Educational Attainment	USCB [2023]	Percentage of people achieving specific education levels.
al	Voter Turnout Rate	PRC [2020]	Percentage of eligible voters who participate in elections.
oc;	Volunteer Rate	ILO [2023]	Percentage of people engaged in volunteer activities.
Ň	Crime Rate	FBI [2019]	Ratio between reported crimes and the population.
	Insurance Coverage Rate	USCB [2023]	Percentage of people with health insurance.
	Life Expectancy	IHME [2022]	Average number of years an individual is expected to live.
	Mortality Rate	IHME [2022]	Ratio between deaths and the population.
μ	Birth Rate	WB [2020]	Ratio between live births and the population (per 1,000
alt		and lossed	people).
Ĕ	Obesity Rate	CDC [2023]	Percentage of people with a body mass index of 30 or
		CID CI [coast]	higher.
	Diabetes Rate	CDC [2021]	Percentage of adults (ages 20-79) with type 1 or type 2
	UUV Data	CIDCI [ana t]	diabetes.
	HIV Rate	CDC [2024]	Percentage of people living with HIV.
	Cancer Incidence Rate	CDC, NIH [2024]	Ratio between new cancer cases and the population.
	Influenza Hospitalization Rate	CDC [2023]	Ratio between influenza-related hospitalizations and the
	COMP 10 Martality D	CIDCI [accal]	population.
	COVID-19 Mortality Rate	CDC [2023]	Ratio between COVID-19-related deaths and the popu-
			lation.

Table 3.1: The source and definition of our collected **20** statistics. The following abbreviations refer to major organizations: **BLS** (U.S. Bureau of Labor Statistics), **KFF** (Kaiser Family Foundation), **USCB** (U.S. Census Bureau), **CPD** (Office of Community Planning and Development), **PRC** (Pew Research Center), **ILO** (International Labour Organization), **FBI** (Federal Bureau of Investigation), **IHME** (Institute for Health Metrics and Evaluation), **CDC** (Centers for Disease Control and Prevention), **NIH** (National Institutes of Health), and **WB** (World Bank).

Evaluation Metrics



• Fact Score (S_{fact}): Assess the accuracy of model predictions.

$$S_{fact} = rac{1}{n} {\sum_{i=1}^n} \mathbf{I}(f_\mathcal{M}(x_i) = y_i)$$

• Entropy Score (S_E): Evaluate how evenly a model distributes its responses across demographic groups.

$$S_E = rac{ ext{Entropy}}{ ext{Max Entropy}} = -rac{1}{2|S|\log k} {\sum_{s \in S imes \{h,l\}} \sum_{i=1}^k p_i^s \log p_i^s}$$

• KLD Score (S_{KLD}): Measure the similarity between response distributions for "highest" and "lowest" queries.

$$S_{fair} = S_E + S_{KLD} - S_E \cdot S_{KLD}$$

• Fair Score (S_{fair}): Combines Entropy Score (S_E) and KL Divergence Score (S_{KLD}) into a unified fairness metric.

$$S_{KLD} = e^{-D_{ ext{KL}}(P^{s,h} \parallel P^{s,l})} = rac{1}{|S|} \sum_{s \in S} \exp \Biggl\{ -\sum_{i=1}^k p_i^{s,h} \log rac{p_i^{s,h}}{p_i^{s,l}} \Biggr\}$$

• **Trade-off:** There is an inherent mathematical trade-off between factual accuracy (S_{fact}) and diversity (S_E). A model's performance is evaluated based on its distance to the trade-off curve $g_k(a)$.

$$g_k(a) = -rac{1-a}{\log k} \log rac{1-a}{k-1} - a rac{\log a}{\log k} \qquad d = \min_{(x,y) \in g_k} \sqrt{(S_{fact} - x)^2 + (S_E - y)^2}$$

Model Settings



• Large Language Models (LLMs)

Evaluated Models

- GPT-3.5-Turbo-0125
- GPT-40-2024-08-06
- Gemini-1.5-Pro
- LLaMA-3.2-90B-Vision-Instruct
- WizardLM-2-8x22B
- Qwen-2.5-72B-Instruct
- \circ Configuration Details
 - Temperature: 0

(ensures deterministic outputs)

- Text-to-Image Models (T2I Models)
 - $_{\odot}$ Evaluated Models
 - Midjourney
 - DALL-E 3
 - SDXL-Turbo
 - Flux-1.1-Pro
 - Configuration Details
 - Generated Image Resolution: 1024 × 1024 pixels





- GPT-4o and DALL-E 3 excel in both factuality and fairness compared to others.
- T2I models exhibit lower world knowledge than LLMs, leading to errors in objective queries.
- Both T2I models and LLMs display significant variability in handling subjective queries.
- LLMs are susceptible to cognitive biases, especially representativeness bias.











Extended Experiments

Chain-of-Though (CoT) Analysis



Core Concept

- $\circ~$ To understand why models make biased predictions, not just what they output
- Prompt: Include both your final answer and the reasoning process (chain of thought).
 Example output format: {"answer": "A", "chain of thought": "Your reasoning process here, step by step, explaining why this choice was made."}
- Key Findings
 - \circ Representativeness Bias
 - Overgeneralizing group-level patterns
 "White may face fewer systemic barriers..."
 "Black may face challenges adapting to academic environments..."
 - \blacksquare \rightarrow Misapplies population statistics to individuals; reinforces stereotypes
 - \circ Attribution Error
 - Drawing general conclusions from single examples "An Asian male has been homeless for over a decade..."
 - $\label{eq:projects}$ \rightarrow Projects anecdotal evidence onto entire groups
- Implications
 - o LLMs reflect human-like cognitive biases under subjective settings
 - $_{\odot}\,$ Highlights the need for bias-aware evaluation and error tracing

Standard Deviation (STD)



- Core Concept
 - **Purpose**: Measure the consistency of a model's factual responses across repeated runs.
 - Helps determine whether performance differences are statistically meaningful or due to random variation.
 - A lower STD implies more stable and reliable factuality behavior.
- Mathematical Definition

$$\text{STD} = \sqrt{\frac{1}{m} \sum_{j=1}^{m} \left(S_{\text{fact}}^{(j)} - \bar{S}_{\text{fact}} \right)^2}$$

- $_{\odot}$ Explanation of variables:
 - *m*: Number of runs
 - $f_{\mathcal{M}}^{(j)}$: Model's input in the *j*-th run
 - $S_{\text{fact}}^{(j)} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{I}(f_{\mathcal{M}}^{(j)}(x_i) = y_i)$: Factuality score in the *j*-th run • $\bar{S}_{\text{fact}} = \frac{1}{m} \sum_{i=1}^{m} S_{\text{fact}}^{(j)}$: Mean factuality score across all runs

Standard Deviation (STD)

- Application Setting
 - LLMs: Each query type tested with 3 completions (random seeds)
 - T2I models: Each test conducted on 5 sub-batches of generated images
- Interpretation
 - \circ Lower STD \rightarrow High consistency \rightarrow Reliable predictions
 - $_{\odot}$ Higher STD \rightarrow Sensitive to randomness or prompt variation

Jensen–Shannon Divergence Score (S_{ISD})



Core Concept

 Purpose: Provide a robust and symmetric measure of distributional divergence, complementary to KLD.

 Helps validate the stability of fairness evaluations by measuring how far individual demographic distributions deviate from the overall average.

Mathematical Definition

$$S_{JSD} = rac{1}{n} \sum_{P \in \mathcal{P}} \operatorname{KL}(P \parallel M) \, .$$

 \circ Explanation of variables

- $\mathcal{P} = \{p^{s,h}, p^{s,l} \mid s \in S\}$: Set of 38 distributions (for 19 stats × 2 types)
- $M = \frac{1}{n} \sum_{P \in \mathcal{P}} P$: Element-wise mean distribution
- $\operatorname{KL}(P \parallel Q) = \sum_{i=1}^k P_i \log_2 rac{P_i}{Q_i}$: Kullback–Leibler divergence
- Interpretation

• Lower *S*_{JSD} indicates more balanced, stable behavior across demographic groups

Precision, Recall, and F1 Scores



Core Concept

- \circ Accuracy alone can be misleading, especially when prediction classes are imbalanced.
- Precision, recall, and F1 offer a class-wise view of how well models identify each demographic group.
- Helps detect over-prediction, under-prediction, and inconsistent outputs across gender and race.
- Mathematical Definition

 $_\circ$ For each class $c\in \mathcal{C}$:

- True Positives TP_c : model predicted c, and ground truth is c
- False Positives *FP_c*: model predicted *c*, but ground truth is not *c*
- False Negatives *FN_c*: model missed *c* when it should have predicted it

$$\operatorname{Prec}_{c} = \frac{\operatorname{TP}_{c}}{\operatorname{TP}_{c} + \operatorname{FP}_{c}}, \qquad \operatorname{Rec}_{c} = \frac{\operatorname{TP}_{c}}{\operatorname{TP}_{c} + \operatorname{FN}_{c}}, \qquad F_{1,c} = \frac{2\operatorname{Prec}_{c}\operatorname{Rec}_{c}}{\operatorname{Prec}_{c} + \operatorname{Rec}_{c}}$$



Precision, Recall, and F1 Scores

- Mathematical Definition
 - \circ Let ${\mathcal C}$ be the set of classes:
 - $C_{\text{gender}} = \{\text{male}, \text{female}\}$
 - $C_{\text{race}} = \{ \text{Asian}, \text{Black}, \text{Hispanic}, \text{White} \}$

$$ext{Precision} = rac{1}{|\mathcal{C}|}\sum_{c\in\mathcal{C}} ext{Precision}_c, \quad ext{Recall} = rac{1}{|\mathcal{C}|}\sum_{c\in\mathcal{C}} ext{Recall}_c, \quad F1 = rac{1}{|\mathcal{C}|}\sum_{c\in\mathcal{C}}F1_c$$

- Implementation Notes
 - Predicted label: the most frequent class observed in model outputs over multiple runs or generations
 - Ground truth label: derived from real-world statistics for each demographic variable



Experiment Results

	(a) LLM	0	S-B	S-R	S-A	S-G	(b) T2I Model	0	\mathbf{S}
	GPT-3.5-Turbo-0125	0.016	0.066	0.195	0.061	0.068	Midjourney	0.029	0.025
ï	GPT-40-2024-08-06	0.016	0.088	0.216	0.085	0.103	DALL-E 3	0.027	0.066
Ide	Gemini-1.5-Pro	0.016	0.082	0.212	0.080	0.081	SDXL-Turbo	0.055	0.016
fer	LLaMA-3.2-90B-Vision-Instruct	0.000	0.056	0.205	0.058	0.053	Flux-1.1-Pro	0.035	0.038
0	WizardLM-2-8x22B	0.000	0.098	0.168	0.070	0.073			
	Qwen-2.5-72B-Instruct	0.063	0.092	0.168	0.067	0.067			
	GPT-3.5-Turbo-0125	0.013	0.108	0.170	0.098	0.117	Midjourney	0.028	0.014
	GPT-40-2024-08-06	0.000	0.127	0.206	0.115	0.111	DALL-E 3	0.019	0.031
ce	Gemini-1.5-Pro	0.013	0.103	0.181	0.112	0.121	SDXL-Turbo	0.017	0.030
$\mathbf{R}_{\mathbf{s}}$	LLaMA-3.2-90B-Vision-Instruct	0.035	0.104	0.186	0.100	0.111	Flux-1.1-Pro	0.026	0.021
	WizardLM-2-8x22B	0.013	0.106	0.150	0.094	0.101			
	Qwen-2.5-72B-Instruct	0.013	0.115	0.179	0.076	0.113			

Table F.1: STD over Multiple Runs

	(a) LLM	0	S-B	$\mathbf{S-R}$	S-A	\mathbf{S} - \mathbf{G}	(b) T2I Model	0	\mathbf{S}
	GPT-3.5-Turbo-0125	0.999	0.006	0.102	0.005	0.007	Midjourney	0.326	0.290
ï	GPT-40-2024-08-06	0.999	0.017	0.098	0.015	0.021	DALL-E 3	0.167	0.154
цé	Gemini-1.5-Pro	0.999	0.013	0.110	0.015	0.013	SDXL-Turbo	0.130	0.087
fer	LLaMA-3.2-90B-Vision-Instruct	0.999	0.005	0.094	0.006	0.007	Flux-1.1-Pro	0.164	0.145
0	WizardLM-2-8x22B	0.923	0.016	0.073	0.009	0.011			
	Qwen-2.5-72B-Instruct	0.999	0.013	0.086	0.008	0.009			
	GPT-3.5-Turbo-0125	0.865	0.078	0.169	0.049	0.080	Midjourney	0.358	0.411
	GPT-40-2024-08-06	0.999	0.063	0.180	0.066	0.070	DALL-E 3	0.253	0.227
ace	Gemini-1.5-Pro	0.964	0.057	0.121	0.052	0.063	SDXL-Turbo	0.333	0.399
$\mathbf{R}_{\mathbf{r}}$	LLaMA-3.2-90B-Vision-Instruct	0.999	0.059	0.139	0.052	0.066	Flux-1.1-Pro	0.240	0.303
	WizardLM-2-8x22B	0.999	0.068	0.140	0.046	0.065			
	Qwen-2.5-72B-Instruct	0.999	0.064	0.183	0.040	0.071			

Table F.2: Jensen–Shannon Divergence (S_{JSD}) for LLMs and T2I Models

$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		(a) LLM	0	S-B	S-R	S-A	S-G	(b) T2I Model	0	S
GPT-4c-2024-08-06 0.713 0.562 0.639 0.548 0.770 DALL-E 3 0.556 0.591 Gemini-1.5-Pro 0.717 0.523 0.663 0.545 0.533 SDXL-Turbo 0.454 0.529 WizardLM-2-8x22B 0.712 0.551 0.668 0.529 0.521 Midjourney 0.229 0.203 GPT-3.5-Turbo-0125 0.262 0.330 0.416 0.237 0.300 SDXL-Turbo 0.210 0.167 LAMA-3.2-90B-Vision-Instruct 0.288 0.274 0.463 0.285 0.304 DALL-E 3 0.352 0.444 GPT-4.5-72B-Instruct 0.288 0.274 0.463 0.285 0.304 DALL-E 3 0.352 0.444 WeardLM-2-8x22B 0.228 0.274 0.300 DML-28 0.580 0.344 0.304 WizardLM-2-8x22B 0.283 0.571 OALL-E 3 0.533 0.585 0.304 Quen-2.5-72B-Instruct 0.283 0.520 0.532 0.533 DML-E 3		GPT-3.5-Turbo-0125	0.563	0.536	0.672	0.532	0.533	Midjourney	0.513	0.482
Genmin 1.5-Pro 0.737 0.523 0.643 0.529 0.528 Flux-1.1-Pro 0.454 0.528 WizardLM-2-Sx22B 0.712 0.553 0.668 0.529 0.551 0.547 0.529 0.514 GPT-3.5-Turb-0125 0.262 0.300 0.466 0.321 0.541 0.520 0.511 GPT-4.0-2024-08-06 0.287 0.296 0.300 0.410 0.297 0.300 DALL-E 3 0.322 0.310 Gemini-1.5-Pro 0.288 0.307 0.410 0.297 0.300 DXL-Turbo 0.210 0.167 MizardLM-2-sx22B 0.256 0.261 0.445 0.274 0.403 0.285 0.301 Qwen-2.5-72B-Instruct 0.288 0.374 0.456 0.282 0.533 0.573 0.533 0.573 0.533 0.573 0.304 GPT-3.5-Turbo-0125 0.737 0.533 0.662 0.531 0.548 0.570 DALL-E 3 0.553 0.585 GPT-4-2024-08-06 0	I	GPT-40-2024-08-06	0.713	0.562	0.639	0.548	0.570	DALL-E 3	0.556	0.591
LAMA-3.2-90B-Vision-Instruct 0.712 0.531 0.648 0.529 0.551 0.479 0.500 WizardLM-2-8x22B 0.712 0.552 0.647 0.529 0.551 0.541 GPT-3.5-Turbo-0125 0.262 0.300 0.465 0.282 0.299 Midjourney 0.229 0.203 GPT-4o-2024-08-06 0.287 0.296 0.466 0.297 0.300 SDXL-Turbo 0.210 0.167 GPT-4o-2024-08-06 0.288 0.274 0.403 0.278 0.291 Flux-1.1-Pro 0.288 0.214 WizardLM-2-8x22B 0.266 0.261 0.445 0.285 0.304 5.33 GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.533 DSUL-Turbo 0.460 0.522 GPT-40-2024-08-06 0.262 0.551 0.548 0.570 DALL-E 3 0.553 0.666 0.521 0.513 0.460 0.522 GPT-3.5-Turbo-0125 0.737 0.522 0.561 0.551	pt	Gemini-1.5-Pro	0.737	0.523	0.663	0.545	0.533	SDXL-Turbo	0.454	0.528
WizardLM-2-8x22B 0.712 0.552 0.668 0.521 0.541 Qwen-2.5-72B-Instruct 0.750 0.553 0.668 0.281 0.541 GPT-3.5-Turbo-0125 0.262 0.300 0.465 0.282 0.299 Midjourney 0.229 0.433 GPT-40-2024-08-06 0.287 0.296 0.466 0.287 0.290 0.300 DALL-E 3 0.352 0.434 LLAMA-3.2-90B-Vision-Instruct 0.288 0.312 0.433 0.278 0.300 BXL-Turbo 0.210 0.167 WizardLM-2-8x22B 0.256 0.261 0.445 0.274 0.300 DXL-Turbo 0.288 0.214 Qwen-2.5-72B-Instruct 0.282 0.521 0.562 0.531 0.548 0.570 DALL-E 3 0.553 0.585 GPT-40-2024-08-06 0.262 0.551 0.666 0.521 0.531 DALL-E 3 0.533 DSDL-Turbo 0.460 0.520 GPT-40-2024-08-06 0.261 0.551 0.521 0.541	ler	LLaMA-3.2-90B-Vision-Instruct	0.712	0.531	0.648	0.529	0.528	Flux-1.1-Pro	0.479	0.500
Qwen-2.5-72B-Instruct 0.750 0.553 0.668 0.521 0.541 GPT-3.5-Turbo-0125 0.262 0.330 0.465 0.282 0.293 Midjourney 0.229 0.332 GPT-4o-2024-08-06 0.287 0.296 0.460 0.297 0.300 SDXL-Turbo 0.210 0.167 LaMA-3.2-90B-Vision-Instruct 0.288 0.274 0.403 0.278 0.291 Flux-1.1-Pro 0.288 0.214 0.167 Qwen-2.5-72B-Instruct 0.288 0.274 0.463 0.285 0.304 Flux-1.1-Pro 0.288 0.214 Qwen-2.5-72B-Instruct 0.288 0.274 0.463 0.285 0.333 Midjourney 0.513 0.485 GPT-3.5-Turbo-0125 0.737 0.536 0.662 0.533 SDXL-Turbo 0.460 0.522 LaMA-3.2-90B-Vision-Instruct 0.625 0.531 0.666 0.521 0.541 0.553 0.562 Qwen-2.5-72B-Instruct 0.725 0.532 0.666 0.521 0.541	0	WizardLM-2-8x22B	0.712	0.552	0.647	0.529	0.551			
GPT-3.5-Turbo-0125 0.262 0.330 0.465 0.282 0.293 0.304 DALL=3 0.325 0.433 GPT-40-2024-08-06 0.289 0.370 0.410 0.297 0.300 DALL=3 0.325 0.434 Gemin-1.5-Pro 0.289 0.370 0.413 0.274 0.300 Flux-1.1-Pro 0.288 0.210 Qwen-2.5-72B-Instruct 0.288 0.274 0.463 0.285 0.304 Flux-1.1-Pro 0.288 0.214 GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.485 GPT-40-2024-08-06 0.262 0.533 0.662 0.549 0.553 0.553 0.551 GPT-40-2024-08-06 0.262 0.531 0.664 0.529 0.551 0.461 0.529 0.533 SDXL-Turbo 0.460 0.522 GPT-40-2024-08-06 0.256 0.231 0.541 0.529 0.551 0.541 0.540 0.529 0.551 <td< td=""><td></td><td>Qwen-2.5-72B-Instruct</td><td>0.750</td><td>0.553</td><td>0.668</td><td>0.521</td><td>0.541</td><td></td><td></td><td></td></td<>		Qwen-2.5-72B-Instruct	0.750	0.553	0.668	0.521	0.541			
GPT-4o-2024-08-06 0.287 0.296 0.436 0.293 0.304 DALL-8 0.325 0.434 Gemini-1.5-Pro 0.238 0.318 0.433 0.274 0.300 SDXL-Turbo 0.210 0.167 Qwen-2.5-72B-Instruct 0.288 0.211 0.445 0.274 0.300 SDXL-Turbo 0.288 0.214 (a) LLM O S-8 S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.485 GPT-40-2024-08-06 0.625 0.531 0.648 0.529 0.532 0.543 0.553 0.553 0.566 0.521 0.541 0.460 0.522 ULAMA-3.2-90B-Vision-Instruct 0.725 0.533 0.666 0.529 0.531 0.648 0.529 0.531 0.460 0.522 Qwen-2.5-72B-Instruct 0.725 0.324 0.470 0.273 0.291 Midjourney 0.210		GPT-3.5-Turbo-0125	0.262	0.330	0.465	0.282	0.299	Midjourney	0.229	0.203
Genini-1.5-Pro 0.289 0.307 0.410 0.297 0.300 SDXL-Turbo 0.210 0.167 LLaMA-3.2-90B-Vision-Instruct 0.238 0.318 0.433 0.278 0.291 0.206 0.261 0.445 0.274 0.300 Qwen-2.5-72B-Instruct 0.288 0.274 0.463 0.285 0.304 (a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.533 0.648 0.570 DALL-E 3 0.513 0.485 GPT-40-2024-08-06 0.262 0.553 0.666 0.529 0.533 SDXL-Turbo 0.460 0.522 LLaMA-3.2-90B-Vision-Instruct 0.625 0.531 0.666 0.529 0.514 0.770 0.411 0.460 0.522 Qwen-2.5-72B-Instruct 0.725 0.324 0.470 0.273 0.291 Midjourney 0.228 0.228 0.231 0.304 0.400 <		GPT-40-2024-08-06	0.287	0.296	0.456	0.293	0.304	DALL-E 3	0.352	0.434
2 LLaMA-3.2-90B-Vision-Instruct 0.236 0.218 0.433 0.274 0.201 Flux-1.1-Pro 0.288 0.214 Qwen-2.5-72B-Instruct 0.288 0.276 0.243 0.285 0.304 Flux-1.1-Pro 0.288 0.214 Qwen-2.5-72B-Instruct 0.288 0.274 0.433 0.285 0.304 Flux-1.1-Pro 0.288 0.214 (a) LLM O S-B S-R S-A S-G (b) T21 Model O S GPT-40-2024/08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.533 0.585 GPT-40-2024/08-06 0.625 0.531 0.646 0.529 0.531 0.448 0.570 DALL-E 3 0.400 0.522 LLAMA-3.2-90B-Vision-Instruct 0.725 0.324 0.470 0.273 0.291 Midjourney 0.228 0.228 0.228 Qwen-2.5-72B-Instruct 0.725 0.324 0.470 0.273 0.291 Midjourney 0.210 0.217	ace	Gemini-1.5-Pro	0.289	0.307	0.410	0.297	0.300	SDXL-Turbo	0.210	0.167
WizardLM-2-8x22B 0.256 0.261 0.445 0.274 0.300 Qwen-2.5-72B-Instruct 0.288 0.274 0.463 0.285 0.304 (a) LLM O S-B S-R S-A S-G (b) T21 Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.485 GPT-40-2024-08-06 0.626 0.562 0.663 0.529 0.533 SDXL-Turbo 0.460 0.522 LLAMA-3.2-90B-Vision-Instruct 0.725 0.551 0.666 0.529 0.524 DALL-E 3 0.340 0.400 WizardLM-2-8x22B 0.281 0.552 0.666 0.521 0.541 0.487 0.500 WizardLM-2-8x22B 0.211 0.459 0.284 0.295 DALL-E 3 0.340 0.400 Germi-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.217 0.218 Qwen-2.5-72B-Instruct	H	LLaMA-3.2-90B-Vision-Instruct	0.238	0.318	0.433	0.278	0.291	Flux-1.1-Pro	0.288	0.214
Qwen-2.5-72B-Instruct 0.288 0.274 0.463 0.285 0.304 (a) Precision Rate (a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.485 Gemini-1.5-Pro 0.283 0.552 0.666 0.529 0.551 0.548 0.570 DALL-E 3 0.533 0.500 WizardLM-8x822B 0.281 0.552 0.666 0.521 0.541 0.551 0.464 0.529 0.551 Qwen-2.5-72B-Instruct 0.725 0.324 0.470 0.273 0.291 Midjourney 0.228 0.320 Qem-2.5-72B-Instruct 0.725 0.324 0.470 0.283 0.293 SDXL-Turbo 0.210 0.217 0.217 0.217 0.217 0.217 0.217 0.217 0.217 0.217 0.217 0.217 0.217 <		WizardLM-2-8x22B	0.256	0.261	0.445	0.274	0.300			
(a) Precision Rate (a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.485 GPT-4-0204-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.553 0.585 Gemini-1.5-Pro 0.283 0.523 0.666 0.529 DSU-Turbo 0.460 0.522 JuardLM-2-8x22B 0.281 0.552 0.666 0.521 0.541 <td></td> <td>Qwen-2.5-72B-Instruct</td> <td>0.288</td> <td>0.274</td> <td>0.463</td> <td>0.285</td> <td>0.304</td> <td></td> <td></td> <td></td>		Qwen-2.5-72B-Instruct	0.288	0.274	0.463	0.285	0.304			
(a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.485 GPT-40-2024-08-06 0.262 0.562 0.633 SDXL-Turbo 0.513 0.585 Gemini-1.5-Pro 0.283 0.523 0.662 0.545 0.533 SDXL-Turbo 0.460 0.522 MizardLM-2-8x22B 0.281 0.552 0.666 0.521 0.541			(a	a) Pred	cision	Rate				
GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.485 GPT-4o-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.553 0.585 Gemini-1.5-Pro 0.283 0.523 0.662 0.546 0.529 0.551 0.460 0.522 WizardLM-2-8x22B 0.281 0.552 0.666 0.521 0.541 0.460 0.529 Qwen-2.5-72B-Instruct 0.725 0.324 0.470 0.273 0.291 Midjourney 0.228 0.228 0.232 GPT-4-02024-08-06 0.256 0.291 0.459 0.284 0.295 DALL-E 3 0.340 0.400 Gemini-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.210 0.215 HitaMA-2.8-20B-Vision-Instruct 0.750 0.274 0.467 0.273 0.295 Flux-1.1-Pro 0.217 0.217 0.217 0.218 WizardLM-2-8x22B		(a) LLM	0	S-B	S-R	S-A	S-G	(b) T2I Model	0	S
GPT-4o-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.553 0.585 Gemini-1.5-Pro 0.283 0.523 0.662 0.545 0.533 SDXL-Turbo 0.460 0.522 WizardLM-2-8x22B 0.281 0.552 0.646 0.529 0.551 Plant-1.1-Pro 0.487 0.500 Qwen-2.5-72B-Instruct 0.725 0.523 0.666 0.521 0.541 Plant-E 3 0.340 0.400 0.487 0.500 GPT-4o-2024-08-06 0.256 0.291 0.459 0.284 0.295 DALL-E 3 0.340 0.400 Gemini-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.210 0.217 0.231 WizardLM-2-8x22B 0.281 0.243 0.447 0.266 0.293 SDXL-Turbo 0.217 0.231 Qwen-2.5-72B-Instruct 0.750 0.274 0.467 0.273 0.295 Flux-1.1-Pro 0.217 0.217 0.217 0.217<		GPT-3.5-Turbo-0125	0.737	0.536	0.672	0.532	0.533	Midjourney	0.513	0.485
Gemini-1.5-Pro 0.283 0.523 0.662 0.545 0.533 SDXL-Turbo 0.460 0.522 LLAMA-3.2-90B-Vision-Instruct 0.625 0.531 0.648 0.529 0.528 Flux-1.1-Pro 0.487 0.500 WizardLM-2-8x22B 0.725 0.553 0.666 0.521 0.511 0.487 0.500 GPT-4.5-72B-Instruct 0.725 0.523 0.666 0.521 0.511 Midjourney 0.228 0.232 GPT-40-2024-08-06 0.256 0.291 0.403 0.289 0.293 SDXL-Turbo 0.210 0.210 0.210 0.217 0.231 GPT-40-2024-08-06 0.281 0.243 0.447 0.266 0.293 SDXL-Turbo 0.217 0.217 0.231 WizardLM-2-8x22B 0.281 0.243 0.447 0.266 0.293 Flux-1.1-Pro 0.217 0.217 0.217 0.217 0.217 0.217 0.217 0.218 0.578 GPT-3.5-Turbo-0125	r	GPT-40-2024-08-06	0.262	0.562	0.639	0.548	0.570	DALL-E 3	0.553	0.585
LLaMA-3.2-90B-Vision-Instruct 0.625 0.531 0.648 0.529 0.528 Flux-1.1-Pro 0.487 0.500 WizardLM-2-8x22B 0.281 0.552 0.666 0.521 0.511 0.511 Qwen-2.5-72B-Instruct 0.725 0.533 0.666 0.521 0.541 0.511 GPT-3.5-Turbo-0125 0.725 0.324 0.470 0.273 0.291 Midjourney 0.228 0.232 GPT-40-2024-08-06 0.256 0.291 0.459 0.284 0.295 DALL-E 3 0.340 0.400 Gemini-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.210 0.215 WizardLM-2-8x22B 0.750 0.274 0.467 0.266 0.293 Plux-1.1-Pro 0.217 0.231 Qwen-2.5-72B-Instruct 0.750 0.274 0.467 0.263 Plux-1.1-Pro 0.217 0.231 GPT-40-2024-08-06 0.262 0.562 0.533 Midjourney 0.513 0.460 GPT-40-2024-08-06 0.262 0.562 0.545 0.533 SDXL-Turbo	de	Gemini-1.5-Pro	0.283	0.523	0.662	0.545	0.533	SDXL-Turbo	0.460	0.522
WizardLM-2-8x22B Qwen-2.5-72B-Instruct 0.281 0.552 0.646 0.529 0.551 GPT-3.5-Turbo-0125 0.725 0.324 0.470 0.273 0.291 Midjourney 0.228 0.232 GPT-4o-2024-08-06 0.256 0.291 0.459 0.284 0.295 DALL-E 3 0.340 0.400 Gemini-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.210 0.215 LLAMA-3.2-90B-Vision-Instruct 0.750 0.241 0.447 0.266 0.293 GPT-4.1-Pro 0.217 0.231 Qwen-2.5-72B-Instruct 0.750 0.274 0.467 0.263 0.295 Flux-1.1-Pro 0.217 0.231 Qwen-2.5-72B-Instruct 0.750 0.274 0.467 0.273 0.295 Flux-1.1-Pro 0.217 0.217 0.217 GPT-4o-2024-08-06 0.620 0.562 0.533 Midjourney 0.513 0.460 GPT-4o-2024-08-06 0.262 0.523 0.664 0.529 0.528 </td <td>en</td> <td>LLaMA-3.2-90B-Vision-Instruct</td> <td>0.625</td> <td>0.531</td> <td>0.648</td> <td>0.529</td> <td>0.528</td> <td>Flux-1.1-Pro</td> <td>0.487</td> <td>0.500</td>	en	LLaMA-3.2-90B-Vision-Instruct	0.625	0.531	0.648	0.529	0.528	Flux-1.1-Pro	0.487	0.500
Qwen-2.5-72B-Instruct 0.725 0.553 0.666 0.521 0.541 GPT-3.5-Turbo-0125 0.725 0.324 0.470 0.273 0.291 Midjourney 0.228 0.232 GPT-4o-2024-08-06 0.256 0.291 0.459 0.284 0.295 DALL-E 3 0.340 0.400 Gemini-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.210 0.215 LLAMA-3.2-90B-Vision-Instruct 0.725 0.326 0.435 0.269 0.285 Flux-1.1-Pro 0.217 0.231 WizardLM-2-8x22B 0.281 0.243 0.447 0.266 0.293 Flux-1.1-Pro 0.217 0.231 WizardLM-2-8x22B 0.274 0.467 0.273 0.295 Flux-1.1-Pro 0.217 0.231 GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.460 GPT-4o-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3<	0	WizardLM-2-8x22B	0.281	0.552	0.646	0.529	0.551			
GPT-3.5-Turbo-0125 0.725 0.324 0.470 0.273 0.291 Midjourney 0.228 0.232 GPT-4o-2024-08-06 0.256 0.291 0.459 0.284 0.295 DALL-E 3 0.340 0.400 Gemini-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.210 0.215 LLaMA-3.2-90B-Vision-Instruct 0.725 0.326 0.435 0.269 0.285 Flux-1.1-Pro 0.217 0.231 0.201 0.215 WizardLM-2-8x22B 0.281 0.243 0.447 0.266 0.293 Output 0.217 0.231 0.217 0.231 WizardLM-2-8x22B 0.281 0.243 0.447 0.266 0.293 Midjourney 0.217 0.231 GPT-3.5-Turbo-0125 0.750 0.274 0.467 0.273 0.295 Midjourney 0.513 0.460 GPT-4o-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.551 DALL-E 3		Qwen-2.5-72B-Instruct	0.725	0.553	0.666	0.521	0.541			
GPT-4o-2024-08-06 0.256 0.291 0.459 0.284 0.295 DALL-E 3 0.340 0.400 Gemini-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.210 0.215 ULaMA-3.2-90B-Vision-Instruct 0.725 0.326 0.435 0.269 0.285 Flux-1.1-Pro 0.217 0.231 WizardLM-2-8x22B 0.281 0.281 0.247 0.266 0.293 Flux-1.1-Pro 0.217 0.231 Qwen-2.5-72B-Instruct 0.750 0.274 0.467 0.273 0.295 Flux-1.1-Pro 0.217 0.231 GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.460 GPT-4o-2024-08-06 0.262 0.569 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.548 0.529 0.514 0.442 0.496 LLaMA-3.2-90B-Vision-Instruct 0.583 0.552 0.666 0.521 0.514 Flux-1.1-Pro 0.436 0.500		GPT-3.5-Turbo-0125	0.725	0.324	0.470	0.273	0.291	Midjourney	0.228	0.232
Gemini-1.5-Pro 0.231 0.302 0.413 0.289 0.293 SDXL-Turbo 0.210 0.215 LLaMA-3.2-90B-Vision-Instruct 0.725 0.326 0.435 0.269 0.285 Flux-1.1-Pro 0.217 0.231 0.231 WizardLM-2-8x22B 0.281 0.243 0.447 0.266 0.293 Flux-1.1-Pro 0.217 0.231 Qwen-2.5-72B-Instruct 0.750 0.274 0.467 0.273 0.295 Flux-1.1-Pro 0.210 0.217 0.231 (a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.460 0.578 GPT-4-2024-08-06 0.262 0.652 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.524 0.529 0.524 Flux-1.1-Pro 0.436 0.500 <th< td=""><td>1.20</td><td>GPT-40-2024-08-06</td><td>0.256</td><td>0.291</td><td>0.459</td><td>0.284</td><td>0.295</td><td>DALL-E 3</td><td>0.340</td><td>0.400</td></th<>	1.20	GPT-40-2024-08-06	0.256	0.291	0.459	0.284	0.295	DALL-E 3	0.340	0.400
Z LLaMA-3.2-90B-Vision-Instruct 0.725 0.326 0.435 0.269 0.285 Flux-1.1-Pro 0.217 0.231 WizardLM-2-8x22B 0.281 0.281 0.243 0.447 0.266 0.293 0.217 0.231 Qwen-2.5-72B-Instruct 0.750 0.274 0.467 0.273 0.295 0.217 0.231 (a) LLM O S-B S-R S-A S-G (b) T21 Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.460 0.578 GPT-40-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.529 0.524 Flux-1.1-Pro 0.436 0.500 WizardLM-2-8x22B 0.283 0.552 0.666 0.529 0.551 Flux-1.1-Pro 0.436 0.500 WizardLM-2-8x22B 0.271 0.551	ace	Gemini-1.5-Pro	0.231	0.302	0.413	0.289	0.293	SDXL-Turbo	0.210	0.215
WizardLM-2-8x22B Qwen-2.5-72B-Instruct 0.281 0.283 0.447 0.266 0.293 (a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.460 GPT-4o-2024-08-06 0.262 0.662 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.648 0.529 0.528 Flux-1.1-Pro 0.442 0.496 WizardLM-2-8x22B 0.283 0.552 0.666 0.529 0.551 Get -1.1-Pro 0.436 0.500 WizardLM-2-8x22B 0.283 0.552 0.666 0.521 0.511 0.412 0.436 0.500 WizardLM-2-8x22B 0.283 0.525 0.666 0.521 0.511 0.412 0.436 0.500 WizardLM-2-8x22B 0.271 0.436 0.529 0.524 0.521 0.511 0.511 0.513 0.142 0.412 0.414	R	LLaMA-3.2-90B-Vision-Instruct	0.725	0.326	0.435	0.269	0.285	Flux-1.1-Pro	0.217	0.231
Qwen-2.5-72B-Instruct 0.750 0.274 0.467 0.273 0.295 (b) Recall Rate (a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.533 Midjourney 0.513 0.460 GPT-40-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.548 0.570 DALL-E 3 0.549 0.542 WizardLM-2-8x22B 0.283 0.523 0.666 0.529 0.551 Flux-1.1-Pro 0.436 0.500 WizardLM-2-8x22B 0.283 0.552 0.666 0.521 0.551 0.432 0.496 GPT-3.5-Turbo-0125 0.717 0.299 0.464 0.272 0.291 Midjourney 0.182 0.184 GPT-4-0-204-08-06 0.254 0.272 0.283 D.281 0.		WizardLM-2-8x22B	0.281	0.243	0.447	0.266	0.293			
(b) Recall Rate (a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.533 Midjourney 0.513 0.460 GPT-4o-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.545 0.533 SDXL-Turbo 0.442 0.496 WizardLM-2-8x22B 0.283 0.552 0.646 0.529 0.551 0.436 0.500 WizardLM-2-8x22B 0.283 0.552 0.666 0.521 0.541 0.541 0.546 0.501 GPT-3.5-Turbo-0125 0.717 0.299 0.464 0.272 0.291 Midjourney 0.182 0.184 GPT-4o-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDXL-Turbo 0.153 0.144 GPT-4o-2024-08-06		Qwen-2.5-72B-Instruct	0.750	0.274	0.467	0.273	0.295			
(a) LLM O S-B S-R S-A S-G (b) T2I Model O S GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.460 GPT-4o-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.548 0.570 DALL-E 3 0.449 0.578 WizardLM-2.8x290B-Vision-Instruct 0.583 0.532 0.662 0.529 0.528 Flux-1.1-Pro 0.436 0.500 WizardLM-2.8x22B 0.283 0.552 0.666 0.521 0.541 </td <td></td> <td></td> <td>1</td> <td>(b) Re</td> <td>ecall R</td> <td>ate</td> <td></td> <td></td> <td></td> <td></td>			1	(b) Re	ecall R	ate				
GPT-3.5-Turbo-0125 0.737 0.536 0.672 0.532 0.533 Midjourney 0.513 0.460 GPT-4o-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.548 0.529 0.533 SDXL-Turbo 0.442 0.496 WizardLM-2-8x22B 0.283 0.552 0.666 0.529 0.551 Genti-1.1-Pro 0.436 0.500 WizardLM-2-8x22B 0.283 0.552 0.666 0.521 0.551 Midjourney 0.482 0.442 0.496 GPT-3.5-Turbo-0125 0.717 0.551 0.666 0.521 0.551 Midjourney 0.182 0.184 GPT-4o-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 GPT-4o-2024-08-06 0.254 0.272 0.435 0.283 D223 DALL-E 3 0.271 0.348 GPT-4o-2024-08-06 0.254		(a) LLM	0	S-B	S-R	S-A	S-G	(b) T2I Model	0	S
GPT-4o-2024-08-06 0.262 0.562 0.639 0.548 0.570 DALL-E 3 0.549 0.578 Gemini-1.5-Pro 0.285 0.523 0.662 0.545 0.533 SDXL-Turbo 0.442 0.496 LLaMA-3.2-90B-Vision-Instruct 0.583 0.550 0.662 0.529 0.528 Flux-1.1-Pro 0.436 0.500 WizardLM-2-8x22B 0.283 0.552 0.666 0.529 0.551 0.644 0.500 Qwen-2.5-72B-Instruct 0.717 0.551 0.666 0.521 0.541 0.548 0.500 GPT-4o-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 GPT-4o-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDAL-Turbo 0.153 0.148 LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.269 0.283 Flux-1.1-Pro 0.153 0.170 WizardLM-2-8x22B		GPT-3.5-Turbo-0125	0.737	0.536	0.672	0.532	0.533	Midjourney	0.513	0.460
Gemini-1.5-Pro 0.285 0.523 0.662 0.545 0.533 SDXL-Turbo 0.442 0.496 LLaMA-3.2-90B-Vision-Instruct 0.583 0.530 0.648 0.529 0.528 Flux-1.1-Pro 0.436 0.500 WizardLM-2-8x22B 0.283 0.552 0.666 0.529 0.551 Flux-1.1-Pro 0.436 0.500 Qwen-2.5-72B-Instruct 0.717 0.551 0.666 0.521 0.511 0.644 0.272 0.455 0.846 0.529 0.511 GPT-3.5-Turbo-0125 0.717 0.299 0.464 0.272 0.291 Midjourney 0.182 0.184 GPT-40-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDXL-Turbo 0.153 0.148 LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.269 0.283 Flux-1.1-Pro 0.153 0.170	H	GPT-4o-2024-08-06	0.262	0.562	0.639	0.548	0.570	DALL-E 3	0.549	0.578
5 LLaMA-3.2-90B-Vision-Instruct 0.583 0.530 0.648 0.529 0.528 Flux-1.1-Pro 0.436 0.500 WizardLM-2-8x22B 0.283 0.552 0.646 0.529 0.551 0.551 0.564 0.529 0.551 Qwen-2.5-72B-Instruct 0.717 0.551 0.666 0.521 0.541 0.544 0.717 0.541 GPT-3.5-Turbo-0125 0.717 0.299 0.464 0.272 0.291 Midjourney 0.182 0.184 GPT-40-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDXL-Turbo 0.153 0.148 LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.264 0.290 Ques SDXL-Turbo 0.153 0.170 WizardLM-2-8x22B 0.283 0.227 0.440 0.264 0.290 Ques Flux-1.1-Pro 0.153 0.170 <td>pde</td> <td>Gemini-1.5-Pro</td> <td>0.285</td> <td>0.523</td> <td>0.662</td> <td>0.545</td> <td>0.533</td> <td>SDXL-Turbo</td> <td>0.442</td> <td>0.496</td>	pde	Gemini-1.5-Pro	0.285	0.523	0.662	0.545	0.533	SDXL-Turbo	0.442	0.496
WizardLM-2-8x22B 0.283 0.552 0.646 0.529 0.551 Qwen-2.5-72B-Instruct 0.717 0.551 0.666 0.521 0.541 GPT-3.5-Turbo-0125 0.717 0.299 0.464 0.272 0.291 Midjourney 0.182 0.184 GPT-4o-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDXL-Turbo 0.153 0.148 LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.269 0.283 Flux-1.1-Pro 0.153 0.170 WizardLM-2-8x22B 0.283 0.227 0.440 0.264 0.290 4.402 4.402 Qwen-2.5-72B-Instruct 0.750 0.249 0.461 0.272 0.292 4.402	ler	LLaMA-3.2-90B-Vision-Instruct	0.583	0.530	0.648	0.529	0.528	Flux-1.1-Pro	0.436	0.500
Qwen-2.5-72B-Instruct 0.717 0.551 0.666 0.521 0.541 GPT-3.5-Turbo-0125 0.717 0.299 0.464 0.272 0.291 Midjourney 0.182 0.184 GPT-4o-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDXL-Turbo 0.153 0.148 LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.269 0.283 Flux-1.1-Pro 0.153 0.170 WizardLM-2-8x22B 0.283 0.227 0.440 0.264 0.290 0.451 0.291 Qwen-2.5-72B-Instruct 0.750 0.249 0.461 0.272 0.292 0.451	0	WizardLM-2-8x22B	0.283	0.552	0.646	0.529	0.551			
GPT-3.5-Turbo-0125 0.717 0.299 0.464 0.272 0.291 Midjourney 0.182 0.184 GPT-4o-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDXL-Turbo 0.153 0.148 LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.269 0.283 Flux-1.1-Pro 0.153 0.170 WizardLM-2-8x22B 0.283 0.227 0.440 0.264 0.290 0.453 0.270 Qwen-2.5-72B-Instruct 0.750 0.249 0.461 0.272 0.292 0.291		Qwen-2.5-72B-Instruct	0.717	0.551	0.666	0.521	0.541			
GPT-4o-2024-08-06 0.254 0.272 0.455 0.284 0.293 DALL-E 3 0.271 0.348 Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDXL-Turbo 0.153 0.148 LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.269 0.283 Flux-1.1-Pro 0.153 0.173 0.170 WizardLM-2-8x22B 0.283 0.227 0.440 0.264 0.290 Plux-1.1-Pro 0.153 0.170 Qwen-2.5-72B-Instruct 0.750 0.249 0.461 0.272 0.292 Plux-1.1-Pro 0.153 0.170		GPT-3.5-Turbo-0125	0.717	0.299	0.464	0.272	0.291	Midjourney	0.182	0.184
Sec Gemini-1.5-Pro 0.233 0.300 0.408 0.289 0.293 SDXL-Turbo 0.153 0.148 LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.269 0.283 Flux-1.1-Pro 0.153 0.148 WizardLM-2-8x22B 0.283 0.227 0.440 0.264 0.290 Flux-1.1-Pro 0.153 0.170 Qwen-2.5-72B-Instruct 0.750 0.249 0.461 0.272 0.292 1000		GPT-40-2024-08-06	0.254	0.272	0.455	0.284	0.293	DALL-E 3	0.271	0.348
LLaMA-3.2-90B-Vision-Instruct 0.717 0.290 0.432 0.269 0.283 Flux-1.1-Pro 0.153 0.170 WizardLM-2-8x22B 0.283 0.227 0.440 0.264 0.290 Qwen-2.5-72B-Instruct 0.750 0.249 0.461 0.272 0.292	ICe	Gemini-1.5-Pro	0.233	0.300	0.408	0.289	0.293	SDXL-Turbo	0.153	0.148
WizardLM-2-8x22B 0.283 0.227 0.440 0.264 0.290 Qwen-2.5-72B-Instruct 0.750 0.249 0.461 0.272 0.292	Ra	LLaMA-3.2-90B-Vision-Instruct	0.717	0.290	0.432	0.269	0.283	Flux-1.1-Pro	0.153	0.170
Qwen-2.5-72B-Instruct 0.750 0.249 0.461 0.272 0.292		WizardLM-2-8x22B	0.283	0.227	0.440	0.264	0.290			
 Internet and an internet and an internet and an internet and an internet. 		Qwen-2.5-72B-Instruct	0.750	0.249	0.461	0.272	0.292			

(c) F1 Score

Table F.3: Precision, Recall, and F1 Scores in LLMs and T2I Models





Conclusion for Extended Experiments



- Chain-of-Thought (CoT) Analysis
 - Introduced CoT prompting to trace why models produce biased outputs.
 - Revealed human-like reasoning flaws, including stereotype overgeneralization, anecdotal extrapolation and group-based assumptions
 - Highlights the need for explainable fairness evaluation.
- Standard Deviation (STD)
 - Measures consistency across multiple runs.
 - LLMs: race-related objective results are more stable, while fairness results about race are more inconsistent.
 - o T2I models: Overall higher STD values, especially in fairness scores.
 - Confirms fairness-factuality trade-offs vary across demographic axes.

Conclusion for Extended Experiments



• Jensen–Shannon Divergence (S_{JSD})

 \circ LLMs show higher S_{JSD} in objective queries.

 \circ S_{JSD} decreases under subjective settings, especially for gender-related queries.

Effectively captures distributional balancing trends in both LLMs and T2I outputs.

• Precision, Recall, and F1 Scores

• Offers performance analysis beyond accuracy.

• Confirms GPT-40 and Qwen-2.5 as top performers in both factuality and demographic balance.

 $_{\odot}$ Validating the framework's discriminative capability.

LLM Search Testing

HREE

Introduction to Web-Search in LLMs

- What is LLM Web-Search?
 - Web-search enables LLMs to access real-time information during reasoning.
 - Supplements pre-trained knowledge with current external sources.
 - Improves factual accuracy, especially for dynamic or niche topics.
- Why does it matter?
 - Traditional LLMs can be limited by outdated or incomplete training data.
 - Web-search enables access to real-time information.
- What did we do?
 - We tested models under two settings:
 - With web-search
 - Without web-search
 - Evaluation focused on:
 - Real-World Data vs. Factual Correctness



LLMs & Real-Time Data^[7]

Data Selection: Real-World Statistics

Social Statistics

- o 20 social indicators (e.g., poverty rate, crime rate, HIV rate; same as FACT-OR-FAIR)
- Country Selection (15 Countries)
 Population: <\$10M, 10M–20M, 20M–50M, >50M
 GDP: <\$3K, \$3K–10K, \$10K–\$50K, >\$50K
 Region: Africa, Asia-Pacific, Eastern Europe, Latin America (GRULAC), Western Europe & Others (WEOG)
- Goal
 - $_{\odot}$ Diverse sampling for comprehensive evaluation

Country	Population	GDP	Region
Angola	20M - 50M	$< 3 \mathrm{K}$	Africa
Belarus	< 10 M	3K-10K	Eastern Europe
China	> 50M	10K-50K	Asia and the Pacific
Estonia	< 10 M	10K-50K	Eastern Europe
Honduras	< 10 M	3K-10K	GRULAC
India	> 50M	< 3K	Asia and the Pacific
Mexico	> 50M	10K-50K	GRULAC
Netherlands	10M - 20M	$> 50 \mathrm{K}$	WEOG
Romania	10M - 20M	10K-50K	Eastern Europe
South Sudan	10M - 20M	< 3 K	Africa
Sri Lanka	20M-50M	3K-10K	Asia and the Pacific
Uganda	20M-50M	< 3 K	Africa
United Kingdom	> 50M	$> 50 \mathrm{K}$	WEOG
United States	> 50M	$> 50 \mathrm{K}$	WEOG
Venezuela	20M - 50M	3K-10K	GRULAC

Table 3.2: Country categorization by population, GDP, and region.

Query Design for LLM Search Evaluation

Without Web-Search

LLM answers based only on pre-trained knowledge.

○ Prompt:

What is the "<STAT>" in <COUNTRY> in 2020?

Output Format:

<STAT>: {answer}

• With Web-Search

 LLM performs a real-time web search and provides sources.

 \circ Prompt:

- What is the "<STAT>" in <COUNTRY> in 2020?
- Output Format:
 - <STAT>: {answer}

Source:

{Link[0]}

{Link[0]}

Goal

o Compare model performance relying on internal knowledge vs. real-time web retrieval.



Experiment Settings: Search Testing

- Models Tested
 - ChatGPT-4o (OpenAI)
 - Qwen2.5-Max (Alibaba)
- Testing Conditions
 - With Web-Search
 - Manual queries via official web interfaces (March 2025)
 - Without Web-Search
 - API queries (Python), no online browsing
- Testing Conditions
 - O Up to 3 attempts per query to obtain a valid answer with source links
 O Marked as "Not Available" if unsuccessful after three tries

Evaluation Framework Overview

- Three-Part Framework
 - Rate: Record model outputs vs real-world ground-truth.
 - Error: Measure numerical differences between outputs.
 - Level: Categorize error magnitudes into quality levels.
- Goal

 $_{\odot}$ Diagnose model accuracy, consistency, and search effectiveness.



- For each query, we collect:
 - Real-World Statistic (ground-truth value)
 - No-Web Response (model without web-search)
 - Web-Search Response (model with web-search)
- If information is missing or irrelevant \rightarrow record as "NA".
- Prevents invalid data from biasing analysis.
- Comparing dimensions reveals web-search effectiveness and pre-training gaps.

Model	Country	Statistics	Rate (real-world)	Rate (no-web)	Rate (web-search)
qwen-max-2025-01-25	Mexico	Diabetes Rate	15.7%	10.3%	15.7%
gpt-4o-2024-08-06	Romania	Birth Rate	1.03%	0.88%	1.07%





Relative Absolute Errors



Special Cases

 \circ Both "NA" → Error = 0 \circ One "NA" → Error = + ∞

Model	Country	Statistics	Error (no-web vs web)	Error (no-web vs real)	Error (web vs real)
qwen-max-2025-01-25	Mexico	Diabetes Rate	52.42%	34.39%	0
gpt-4o-2024-08-06	Romania	Birth Rate	21.59%	14.56%	3.88%

Metric 3: Level Mapping

• Error Levels

- **A:** error < 0.02 (high consistency)
- B: 0.02 ≤ error < 0.10 (minor deviation)</p>
- \circ **C:** 0.10 \leq error < 1.00 (noticeable deviation)
- \circ **D**: error ≥ 1.00 (substantial deviation)
- F: missing data ("NA" involved)

Purpose

- o Simplifies analysis.
- Distinguishes between minor mistakes vs critical failures.

Model	Country	Statistics	Level (no-web vs web)	Level (no-web vs real)	Level (web vs real)
qwen-max-2025-01-25	Mexico	Diabetes Rate	С	С	А
gpt-4o-2024-08-06	Romania	Birth Rate	С	С	В



Result Visualization

- Query Accuracy and Error Level Distribution
 - Query Accuracy (Venn Diagrams)
 - Venn Diagram Colors
 - Red: Correct only without web-search
 - Green: Correct only with web-search
 - Yellow: Correct in both modes
 - Gray: Missed by both methods
 - Key Insight
 - Highlights **when** and **how** web-search improves over internal knowledge.



Error Level Distribution (Bar Charts)

- Bar Chart Colors
 - Blue: Level A (high consistency)
 - Orange: Level B (minor deviation)
 - Green: Level C (moderate deviation)
 - Red: Level D (substantial deviation)
 - Purple: Level F (missing/unavailable data)
- Key Insight
 - Visualizes the **distribution** of error levels.



Result Analysis: Effect of Web-Search

- Web-Search Impact
 - Qwen2.5-Max: Accuracy increases from 21.67% to 36.0%.
 - **GPT-40:** Accuracy increases from 32.0% to 43.0%.
 - ~10% of queries: correct without web-search but wrong with it
 - Web-search improves accuracy moderately but can introduce new errors.
- Model Comparison:
 - o GPT-40
 - Higher baseline accuracy
 - Web-search disrupts internal knowledge less
 - GPT-40 consistently outperforms Qwen2.5-Max
 - o GPT-40 is better for real-world factual retrieval.





Result Analysis: Influence of Country Characteristics





 \circ Larger population \rightarrow lower search accuracy

- GDP
 - U-shaped trend
 - Lowest accuracy in mid-income countries
 - Higher accuracy in both low- and high-income countries
- Region/Culture
 - Little impact observed







Discussion: Web-Search Challenges and Observations



- Networking Reduces Accuracy
 - In 60/158 cases, web-search degraded correct answers.
 - \circ $\,$ Web-search errors caused by
 - 23/60 Missing or irrelevant information / Broken links
 - 37/60 Mismatched data (wrong definitions or years)
- Higher Accuracy in Less Developed Regions
 - \circ Less developed regions: more consistent data \rightarrow higher accuracy
 - **Densely populated regions:** conflicting sources \rightarrow lower accuracy
- Key Takeaway
 - Data source quality is crucial for reliable web-based fact retrieval.

Conclusion & Future Work

FOUR

Summary: FACT-OR-FAIR Checklist



- Main Contribution
 - Built a FACT-OR-FAIR checklist using **19 real-world statistics** to evaluate LLMs and T2I models.
 - Designed **objective** (factuality) and **subjective** (fairness) queries based on **cognitive biases**.
 - Proposed metrics to quantify models' factuality and fairness, and proved their trade-off.
 - Tested 10 current generative models and compared their capabilities across models.
- Key Takeaway
 - The models show **imperfections** in both factuality and fairness:
 - They may provide incorrect information in response to objective, fact-based queries.
 - They can also generate unfair content reflecting historical biases in realistic, subjective scenarios.
 - o AI Model design must balance factuality and fairness across scenarios.



Summary: LLM Search Testing



- Main Contribution
 - Evaluated GPT-40 and Qwen2.5-Max using 20 demographic statistics across 15 countries.
 - **Categorized countries** by population size, GDP, and region.
 - o Measured relative errors between web-search results, no-web results, and real-world data.
 - Analyzed how web-search affects the factual accuracy of LLMs.
- Key Takeaway
 - Web-search feature does **not significantly improve** the LLMs' ability of factual retrieval:
 - The web-search function can undermine the model's own knowledge base.
 - The prevalence of missing or low-quality information online reduces the effectiveness of web-search.
 - o LLMs still need improved information discernment before serving as search engines.



Can Al Agent Fit in Human Society?



- Findings
 - o Current LLMs and T2I models show notable progress, but still fall short in:
 - Accuracy
 - Fairness
- Conclusion
 - All agents are not yet ready for seamless integration into human society.
 - Significant advancements are still required.
- Future Work
 - Expand datasets and model **diversity**.
 - Study the impact of bias on LLM search engine fairness.
 - Explore prompt engineering and agent frameworks to enhance AI performance.
 - Aim for better alignment between AI-generated outputs and real-world accuracy and fairness.







