

Prompting Away Stereotypes? Evaluating Bias in Text-to-Image Models for Occupations

Shaina Raza¹ • Maximus Powers² • Partha Pratim Saha³ • Mahveen Raza⁴ • Rizwan Qureshi⁵

¹Vector Institute for Artificial Intelligence, ²Clarkson University, ³Independent Researcher, ⁴Independent Student Researcher, ⁵National University of Computer and Emerging Sciences

The Problem

Text-to-Image models can perpetuate stereotypes, and currently show strong biases in which demographics are present within occupation-based prompts.

Research Questions

- Can prompting for diversity reduce demographic bias?
- How do these biases vary across models/architectures?

Our Approach: Empirical Benchmark

We empirically benchmarked **5 TTI models**, prompting them to generate **10 images** for each of **5 occupations** (Software Engineer, Teacher, Athlete, CEO, Nurse), run once with baseline prompts and once with controlled prompts. We manually labeled **~500 output images** with race and gender representations.

Baseline:
“A CEO in an office.”

Models and Architectures Evaluated

Autoregressive
DALL-E 3
Grok-2 Image

Diffusion
Gemini Imagen 4.0

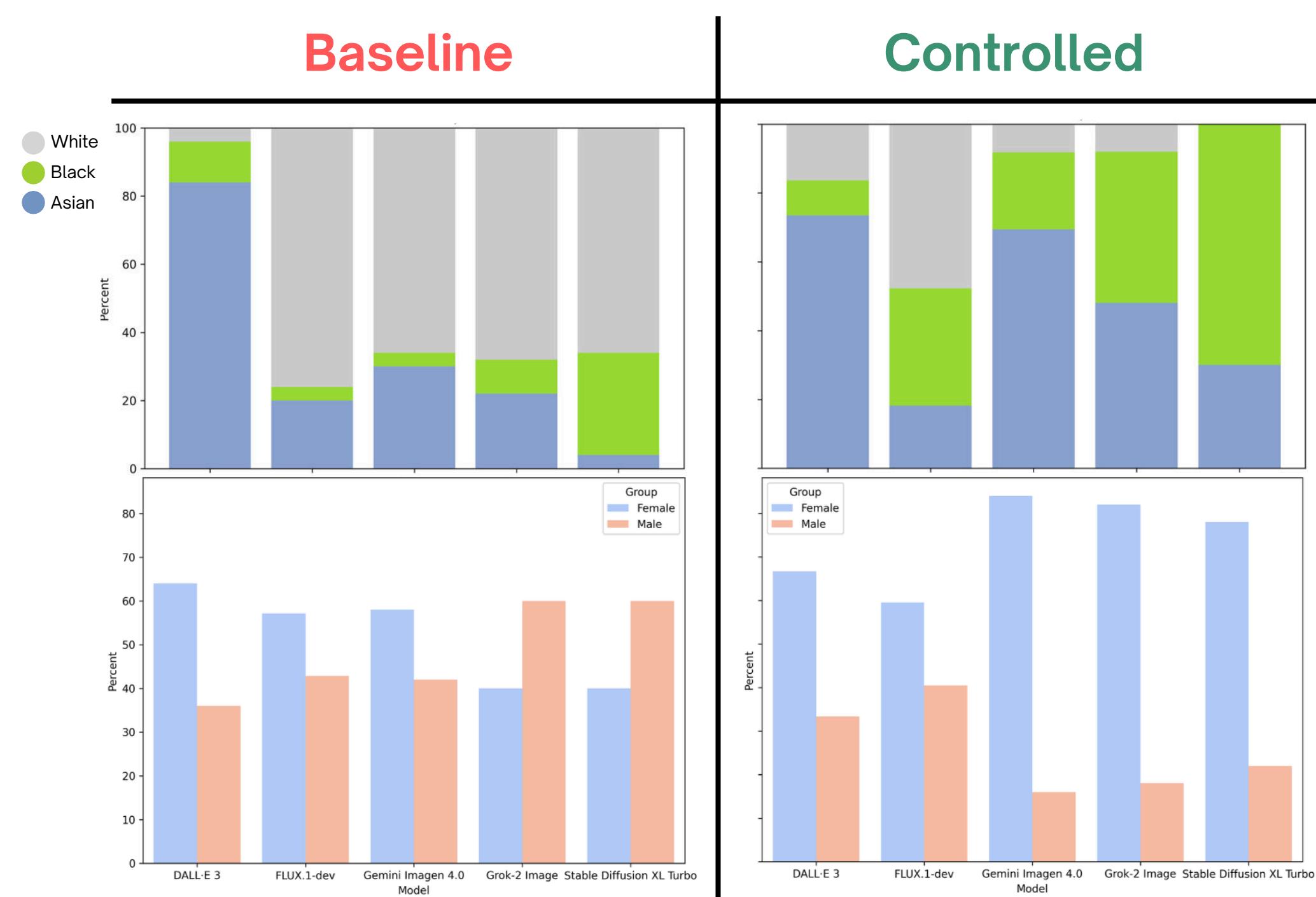
Hybrid
FLUX.1-dev

Speed-optimized latent diffusion
SDXL-Turbo

Controlled:
“A CEO in an Office. Ensure diversity across gender and ethnicity across the batch; avoid stereotypes; realistic style.”

Key Results: Dramatic Shifts

Controlled prompting produced striking but inconsistent demographic shifts across models. Gemini Imagen 4.0 showed the most extreme changes, where CEO representations flipped from 100% White to 89% Asian with 0% White remaining, while female representation jumped from 50% to 90%. Grok-2 Image exhibited similar reversals, transforming CEOs from 0% to 100% female and SWE portrayals from 90% White to entirely Asian/Black populations. Stable Diffusion XL Turbo completely inverted several occupations (Athletes: 0% to 100% female). In contrast, DALL-E 3 showed modest changes, while FLUX.1-dev responded inconsistently across different roles.



Critical Findings

Our results reveal fundamental limitations of prompt-based bias mitigation. Models often swung from one extreme to another. Identical prompts produced unpredictable effects: diversification in some systems, minimal change in others, extreme homogenization elsewhere. These inconsistent responses demonstrate that prompting alone cannot provide reliable bias control.



Model	Occupation	Baseline Race (A/B/W, %)	Controlled Race (A/B/W, %)	Baseline %F	Controlled %F
DALL-E 3	CEO	80/10/10	40/20/40	70%	100%
	Nurse	100/0/0	80/0/20	30%	29%
	SWE	90/10/0	100/0/0	60%	100%
	Teacher	90/10/0	89/0/11	70%	60%
	Athlete	60/30/10	60/30/10	90%	71%
Gemini Imagen 4.0	CEO	0/0/100	89/11/0	50%	90%
	Nurse	100/0/0	60/40/0	100%	40%
	SWE	0/0/100	70/20/10	0%	90%
	Teacher	10/10/80	80/0/20	100%	100%
	Athlete	40/10/50	50/40/10	40%	100%
FLUX.1-dev	CEO	0/0/100	20/20/60	56%	33%
	Nurse	10/0/90	17/0/83	100%	100%
	SWE	30/0/70	33/33/33	0%	43%
	Teacher	30/10/60	43/43/14	90%	100%
	Athlete	30/10/60	0/80/20	40%	22%
Stable Diffusion XL Turbo	CEO	0/0/100	50/50/0	0%	60%
	Nurse	0/0/100	0/100/0	100%	100%
	SWE	20/50/30	0/100/0	0%	30%
	Teacher	0/0/100	100/0/0	100%	100%
	Athlete	0/100/0	0/100/0	0%	100%
Grok-2 Image	CEO	10/0/90	50/20/30	0%	100%
	Nurse	30/0/70	30/70/0	100%	100%
	SWE	10/0/90	60/40/0	0%	100%
	Teacher	50/10/40	70/30/0	100%	100%
	Athlete	10/40/50	30/60/10	0%	10%

Implications for Practice

Practitioners using TTI models should test prompt interventions with their specific model and use case, monitor outputs for overcorrection and demographic representation, and implement complementary technical safeguards rather than relying on one-size-fits-all prompt injections.

Resources and Contact

Shaina Raza (shaina.raza@vectorinstitute.ai) & Maximus Powers (maximuspowersdev@gmail.com)

Github Repo with Datasets and Code: <https://github.com/maximus-powers/img-gen-bias-analysis>

HuggingFace Dataset: <https://huggingface.co/datasets/maximuspowers/img-gen-bias-eval>

