

INTRODUCTION

-
- A horizontal row of four portrait photographs of male doctors. From left to right: a white man with grey hair and glasses, a Black man with a beard and glasses, an Asian man with glasses, and a man with dark hair and glasses. All four are wearing white lab coats over collared shirts and dark ties, with stethoscopes around their necks. They are all smiling and standing in a brightly lit hospital corridor.

- Profession: A photo of a doctor
- Race 7: A photo of a white doctor
- Race 4: A photo of a white doctor
- Gender: A photo of a male doctor
- Age: A photo of a 0-2 year old doctor

	white	black	latino/hispanic	east asian	southeast asian	indian	middle eastern
"A photo of a <race> person"	9.40	61.94	15.38	66.24	5.88	60.61	26.48
<race> + Profession	68.92	70.90	15.38	43.46	20.59	57.58	13.24
<race> + Gender	20.72	67.91	15.38	44.73	11.76	51.52	29.86
<race> + Age	11.57	65.67	11.54	59.92	14.71	69.70	25.07

FRAMEWORK

The diagram illustrates the SDXL framework for generating images from both text and image prompts. It is divided into two main sections: a top section for text prompts and a bottom section for image prompts.

Top Section (Text Prompts):

- Inputs:** A list of text prompts: "A photo of a <demographic> doctor ...", "doctor", "engineer", and "chef".
- Encoding:** Each text prompt is processed by a "Text Encoder" (represented by a trapezoid) to produce a sequence of tokens (represented by rounded rectangles). The "doctor" prompt produces a red token, while "engineer" and "chef" produce grey tokens.
- SDXL Model:** These tokens are fed into the "SDXL" model (represented by a dashed box).
- Outputs:** The SDXL model generates three "Image Encoders" (represented by trapezoids) which output three "Image Encoders" (represented by rounded rectangles) containing the text "<avg>".

Bottom Section (Image Prompts):

- Inputs:** An image of three people (two men and one woman) is processed by an "Image Encoder" (represented by a trapezoid) to produce a sequence of tokens (represented by rounded rectangles).
- SDXL Model:** These tokens are fed into the "SDXL" model.
- Outputs:** The SDXL model generates three "Image Encoders" (represented by trapezoids) which output three "Image Encoders" (represented by rounded rectangles) containing the text "<avg>".

The diagram shows that the SDXL model can take either text or image prompts as input and generate a corresponding image output.

For each of these tests, the aim is to classify the profession of the person within the query image. Generating images helped most when the text utilized race as the demographic. It helped the least when gender was the demographic.

	“A photo of a”	Profession	Race 7	Race 4	Gender	Age
CLIP, Classify Profession	95.14		94.73	95.22	96.52	94.81
D3G, Classify Profession	95.54		95.22	95.30	96.52	95.06
Avg D3G, Classify Profession	95.87		95.62	95.38	96.76	95.54
CLIP, Classify Race7	28.20	44.65			28.61	25.69
D3G, Classify Race7	31.85	45.38			32.90	30.96
Avg D3G, Classify Race7	32.33	45.46			33.55	32.25

- Generate images based on demographics within the dataset (i.e. , generate more images of Hispanic people if they are underrepresented)
- Try modifying demographics of the person within the existing image

- This technique does not remove demographic biases but rather modifies the direction of the model bias.
- The images generated by the model can often reinforce certain demographic biases, and this method should only be used within appropriate contexts.