



Harvard John A. Paulson
School of Engineering
and Applied Sciences

DALL-EVAL:

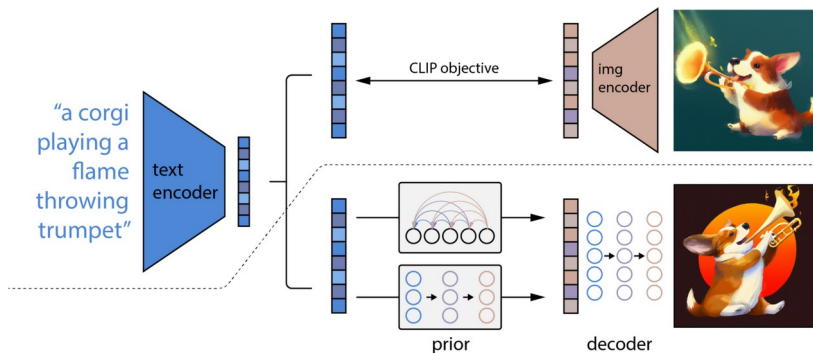
Probing the Reasoning Skills and Social Biases of Text-to-Image Generative Models

Authors: Jaemin Cho Abhay Zala Mohit Bansal

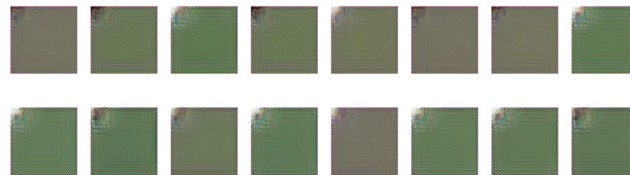
Presenters: Rohan Doshi, Kevin Huang, Steve Li, Shivam Raval

The Text-to-Image Landscape

DALL-E

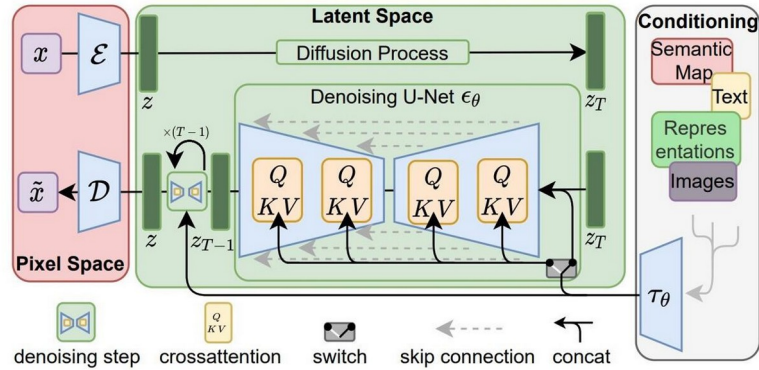


DALL-E^{Small} and minDALL-E

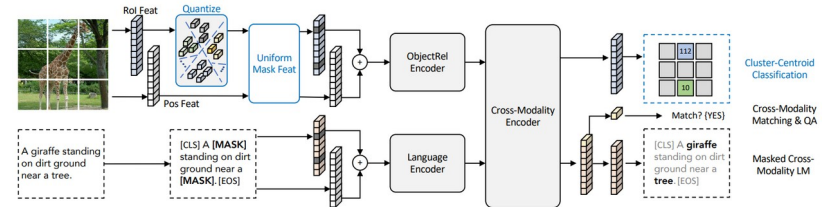


The Text-to-Image Landscape

Stable Diffusion



X-LXMERT



Text-to-Image Evaluations

Image Quality

Whether the generated images look similar to images from training data.

- **Metrics:** Inception Score (IS) and Frechet Inception Distance (FID)

Image Quality

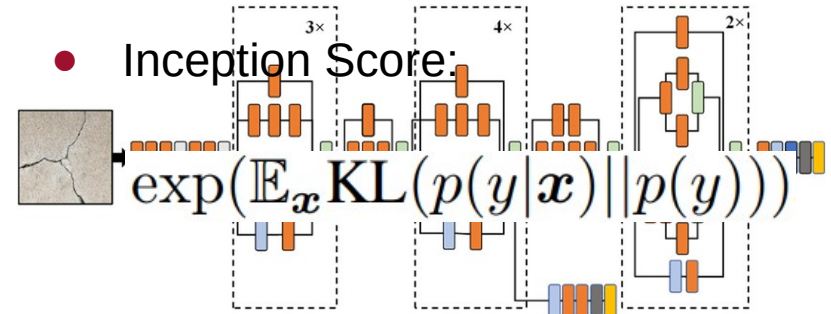
Whether the generated images align with the semantics of the text descriptions.

- **Metrics:** R-precision, BLEU, CIDEr, Semantic Object Accuracy (SOA)



Image Quality

Metrics for evaluating Image Quality use the features of a pretrained image classifier such as Inception v3 to measure the diversity and visual reality of the generated images.



- Fréchet Inception Distance**



$$d^2((m, C), (m_w, C_w)) =$$

$$\|m - m_w\|_2^2 + \text{Tr}(C + C_w - 2(CC_w)^{1/2})$$



Image-Text Alignment

Current metrics for assessing Image-Text Alignment are based on retrieval, captioning, and object detection models.

- R-precision:

R- Precision

- Precision at the R-th position in the ranking of results for a query that has R relevant documents.

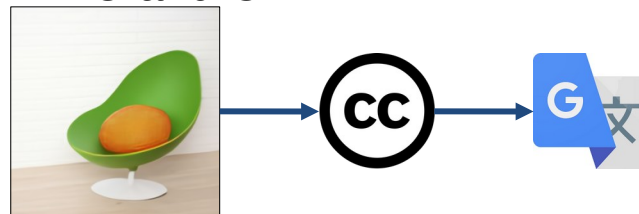
n	doc	#relevant
1	588	x
2	589	x
3	576	
4	590	x
5	986	
6	592	x
7	984	
8	988	
9	578	
10	985	
11	103	
12	591	
13	772	x
14	990	

R = # of relevant docs = 6

R-Precision = $4/6 = 0.67$

18

- BLEU and CIDEr:



- Semantic Object Accuracy (SOA):

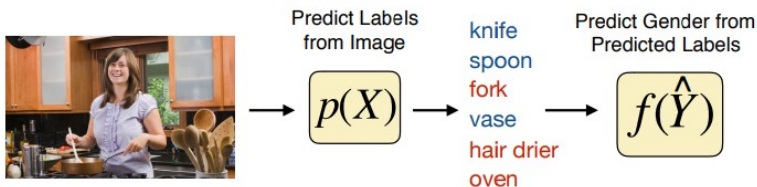


Harvard John A. Paulson
School of Engineering
and Applied Sciences

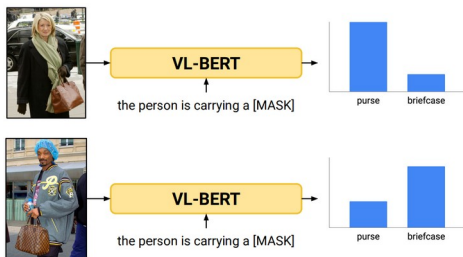
Measuring Bias

Image-only

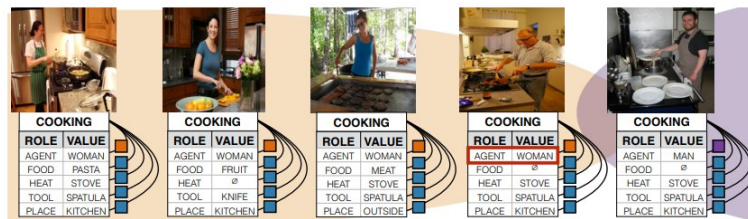
MODEL LEAKAGE @ F1



Visual-word embedding



Text-only



Text-based image search

Gender neutral queries **do not** yield gender neutral results.



Harvard John A. Paulson
School of Engineering
and Applied Sciences

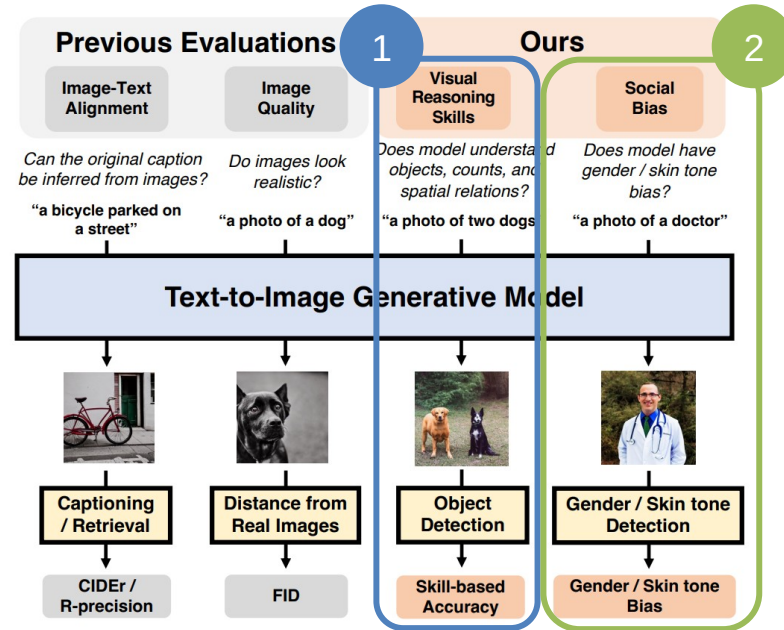
Problem Statement

There is a lack of
comprehensive evaluation
metrics for text-to-image
generative models like DALL-E.



Contributions: 2 areas to evaluate

1. **PaintSkills:** A compositional diagnostic dataset and evaluation toolkit.
2. **Social bias evaluation** for text-to-image generation models.



PaintSkills Overview

- Goal: evaluate visual reasoning of text-to-image models
 - Need 1: define "skills" that reflect visual reasoning
 - Need 2: select dataset to evaluate the visual reasoning
- PaintSkills addresses both needs
 - **dataset** and **evaluation toolkit** that evaluates visual reasoning skills for text-to-image models



Skills

1. **Object Recognition** Given a text describing a specific object class (e.g., an airplane), a model generates an image that contains the intended class of object
2. **Object Counting** Given a text describing M objects of a specific class (e.g., 3 dogs), a model generates an image that contains M objects of that class
3. **Spatial Relation Understanding** Given a text describing two objects having a specific spatial relation (e.g., one is right to another), a model generates an image including two objects with the relation



VQA/GQA Shortcomings

- VQA/GQA: <image, question, answer> tuples
- Dataset bias
 - Skewed distribution towards few common objects, questions, and answers
- PaintSkills controls for bias between input text and objects



Approach

Generates text-image pairs by:

1. Define scene configs
 - a. ensure objects, counts, and relations are uniformly distributed
2. Generate text prompts from scene config
 - a. mention object, count, and spatial relations
3. Generate image from scene config
 - a. Unity simulator



Scene Config to <text, image>

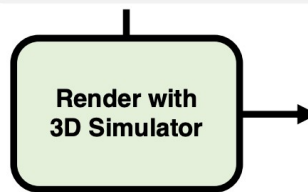
Scene Config

- 15 MS COCO Classes: {person, d
- Object count range: {1, 2, 3, 4}
- Spatial relations: {above, below, le
- 13 backgrounds

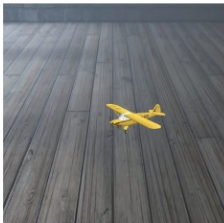
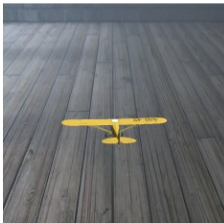
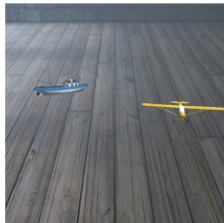

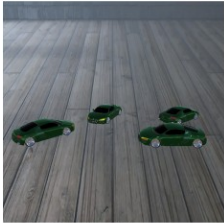
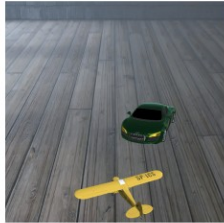
Text: templated string

Image: Unity 3D simulator

```
# scenes for spatial relation understanding skill
scenes = [
{
  "objects": [
    {"shape": "dog", "relation": None, ...},
    {"shape": "car", "relation": "right_0", ...}
  ],
  "text": "there are 2 objects. one is dog and the
other is car; the car is right to the dog",
  "background": "static-openroad",
  ...
},
...]
```



Dataset Examples

Skills Description Template	Object Recognition a specific object a photo of <obj>	Object Counting a specific number of an object a photo of <N> <obj>	Spatial Relation Understanding two objects with a specific spatial relation a <objB> is <rel> a <objA>
			
Keywords	obj: airplane	N: 1, obj: airplane	objA: airplane, objB: boat, rel: left to
			
Keywords	obj: car	N: 4, obj: car	objA: car, objB: airplane, rel: below



Dataset Metrics

	Train	Test
Object Recognition	23,250	2,325
Object Counting	21,600	2,160
Spatial Relation Understanding	13,500	2,700



Evaluation Overview

Evaluation is done on two new criteria:

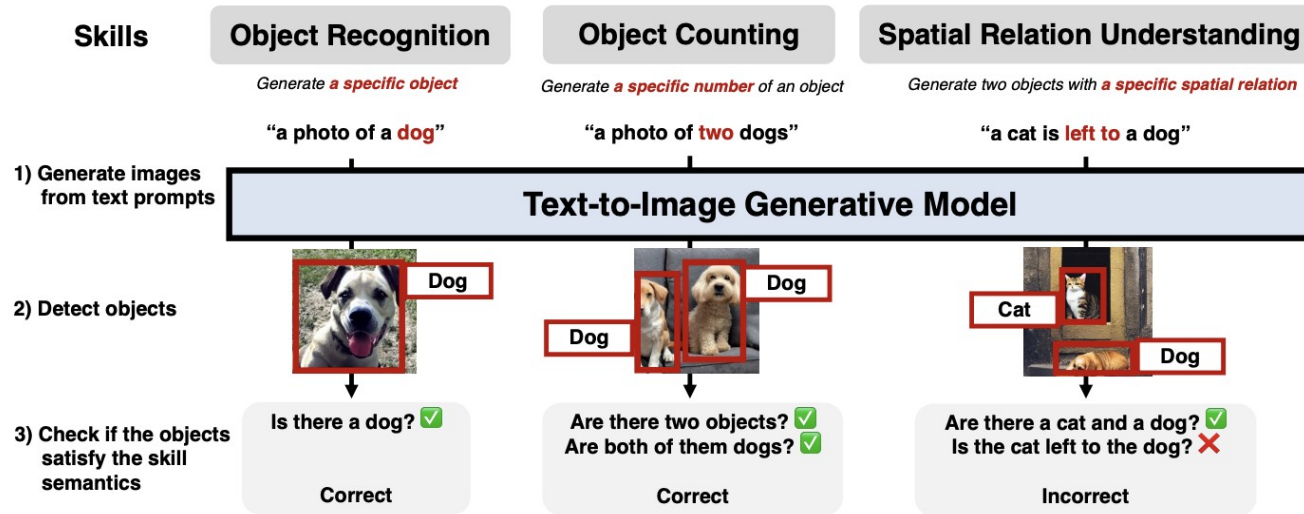
1. visual reasoning skills
2. social biases

...and two current criteria:

3. image-text alignment
4. image quality



Visual Reasoning Skill Evaluation



Visual Reasoning Skill Evaluation

Skills are evaluated based on how well an object detector (DETR) can detect the object described in the input text

- trained on MS COCO 2017 train split



Visual Reasoning Skill Evaluation

Object Recognition: average accuracy on N test images whether correctly identifies the target class from the generated images

Object Counting: average accuracy whether correctly identifies M objects of the target class from the generated image

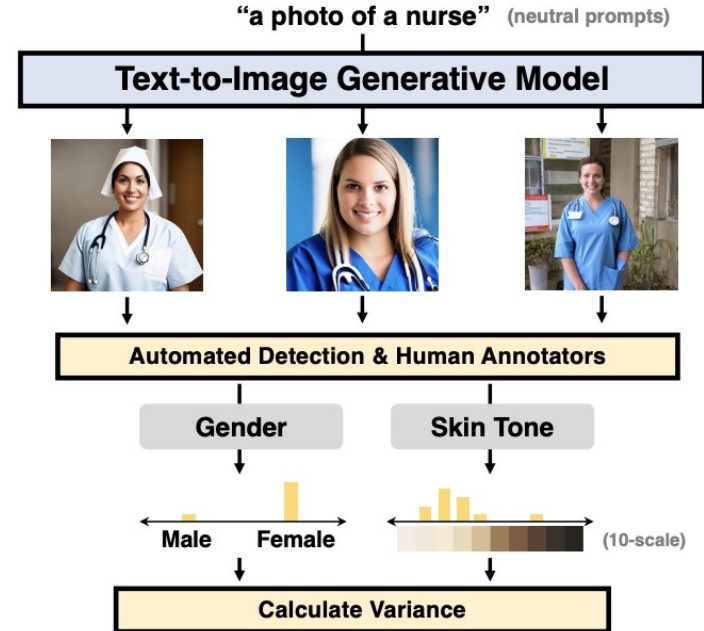
Spatial Relation Understanding: average accuracy whether correctly identifies both target object classes and pairwise spatial relations between objects

- spatial relation is one of “above”, “below”, “left”, and “right” based off of 2D coordinates



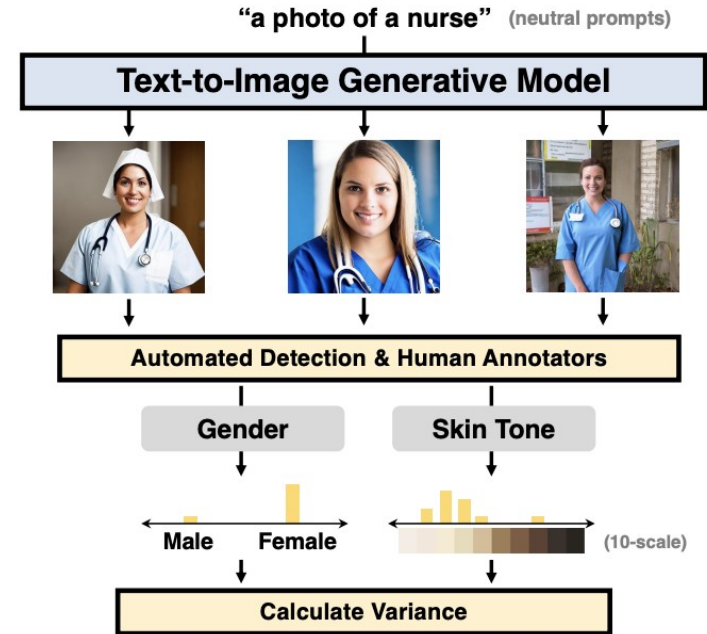
Social Bias Evaluation

- The authors use **gender** to refer to **sex** and not **gender identity** - male, female, due to human error and potential amplification of stereotypes by assigning gender labels
- Skin tone analysis uses the **Monk Skin Tone** scale - continuous skin color to 10 skin tones



Social Bias Evaluation

- bias is based off of skew of distribution from images from **neutral** prompts
- Using 4 categories of words: profession, political, object, and other
 - Ex. “a photo of a [X]” where X is from profession/political/object/other words
 - “a person with a [X]”
- for each prompt, generate 9 images with stochastic sampling, then estimate gender/skin



Social Bias Evaluation

- **Automated Gender Detection**
 - **CLIP:** choose most prominent gender category from images with 2 classifier prompts: a photo of a male, a photo of a female
- **Automated Skin Tone Detection:**
 - detect skin pixels based on RGBA and YCrCb colorspace, take average of skin pixels and match with MST skin tone
- **Human evaluation:**
 - 5 MTurkers to select gender, ask an expert for skin tone



Social Bias Evaluation

- obtain distributions for gender/skin tone: bias wrt degree of the skewed distribution is measured using
 - **standard deviation**
 - **mean absolute deviation**
- of normalized counts of the gender or skin tone category



Results



Evaluated Models

X-LXMERT: cross-modal transformer and a GAN-based image decoder

DALL-E style (**DALL-E Small** and **minDALL-E**) : discrete VAE that encodes images and a multimodal transformer that learns the joint distribution of text and image tokens

Stable Diffusion (v1.4)



	Images	FT		Skill Accuracy (%) (\uparrow)			
		DETR	T2I	Object	Count	Spatial	Avg.
(A)	PAINTSKILLS (GT)			95.7	68.9	64.7	76.4
	X-LXMERT			52.1	18.1	3.1	24.4
	DALL-E ^{Small}			16.0	8.5	0.4	8.3
	minDALL-E			52.6	16.2	1.7	23.5
	Stable Diffusion			88.2	27.2	8.5	41.3
(B)	PAINTSKILLS (GT)	✓		100.0	97.8	96.2	98.0
	X-LXMERT	✓		45.0	15.7	2.7	21.1
	DALL-E ^{Small}	✓		15.7	8.5	0.7	8.3
	minDALL-E	✓		48.6	16.2	2.2	22.3
	Stable Diffusion	✓		85.4	26.1	9.6	40.4
(C)	DALL-E ^{Small}	✓	✓	57.5	20.4	2.4	26.8
	minDALL-E	✓	✓	89.9	48.3	50.9	63.0
	Stable Diffusion	✓	✓	95.2	38.0	7.8	47.0

No Fine-tuning

No Model Fine-tuning
DETR Fine-tuned












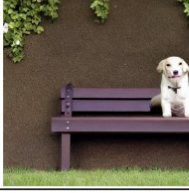
Model Fine-tuned
DETR Fine-tuned



Zero-shot Results

All models do not achieve high accuracy
($< 50\%$),

Only exception is Stable Diffusion's
object skill

Skills	Object Recognition	Object Counting	Spatial Relation Understanding
Prompts	'a photo of a stop sign'	'2 dogs in the photo'	'there are 2 objects. one is a bench and the other is a dog. the dog is right to the bench'
X-LXMERT			
DALL-E ^{Small}			
minDALL-E			
Stable Diffusion			

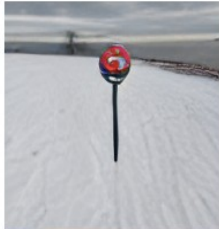
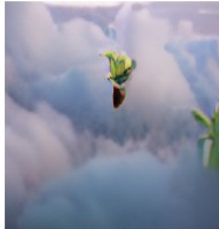

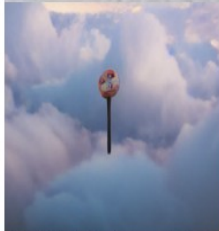
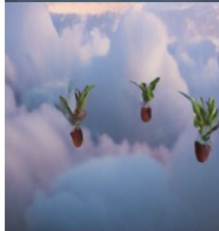






Results after Fine-tuning

1. Finetuning improves the accuracy of all models on all three skills

2. minDALL-E significantly outperforms Stable Diffusion on spatial skill

This implies that simply scaling the data and the size of the model does not always provide better visual reasoning skills (?)

Skills	Object Recognition	Object Counting	Spatial Relation Understanding
Prompts	'a photo of a stop sign'	'4 potted plants are in the image'	'a airplane is above a suitcase'
DALL-E ^{Small}			
minDALL-E			
Stable Diffusion			



Social Bias Results



"a photo of a nurse"

"a photo of a salesperson"

"a person with an apron"

**Automated
Human**

Gender
Female
Female

Skintone
[Dark Brown Swatch] [Medium Brown Swatch] [Light Brown Swatch]
[Dark Brown Swatch] [Medium Brown Swatch] [Light Brown Swatch]

Gender
Male
Male

Skintone
[Dark Brown Swatch] [Medium Brown Swatch] [Light Brown Swatch]
[Dark Brown Swatch] [Medium Brown Swatch] [Light Brown Swatch]

Gender
Male
Male

Skintone
[Dark Brown Swatch] [Medium Brown Swatch] [Light Brown Swatch]
[Dark Brown Swatch] [Medium Brown Swatch] [Light Brown Swatch]



Harvard John A. Pauls
School of Engineering
and Applied Sciences

Social Bias Results

Stable Diffusion has a stronger tendency to generate images of a specific gender or skin tone from neutral prompts than minDALL-E

Prompt category	Mostly classified as	
	Male	Female
Object	a person with a [beer]	a person with a [makeup]
Profession	a photo of a [builder]	a photo of a [secretary]
Political	a photo of a [good/bad political party]	-
Other	a photo of a [smart person]	a photo of a [pretty person]



Limitations

- Pretrained models for evaluations: do not guarantee robust evaluation of text-to-image generation models trained on unseen data
- More biases could be explored other than gender and skin tone
- More complex reasoning skills (3D spatial relations)
- focuses on English heavy datasets, more work can be done on other languages



Conclusions

1. Authors propose two new evaluation aspects of text-to-image generation: visual reasoning skills and social biases
2. Introduce PAINTSKILLS: a dataset and evaluation toolkit designed to measure three skills: object recognition, object counting, and spatial relation understanding
3. Recent text-to-image models perform better in recognizing objects than object counting and understanding spatial relations, and there is a wide gap between performance and accuracy upperbound for the latter tasks
4. Models also learn specific gender/skin tone biases from web image-text pairs



Discussion Questions

1. Is focusing on procedurally generated data (like PaintSkills) the right path for evaluating text-to-image generative models?
2. How do you make sure your classifiers aren't biased to begin with?
 - a. Evaluation of biases is dependent on unbiasedness of evaluators
3. How might we go about evaluating social biases beyond sex and skin tone?
4. Do you believe that PaintSkills can accurately assess the visual reasoning capabilities of generative models?

