# Lost In Translation
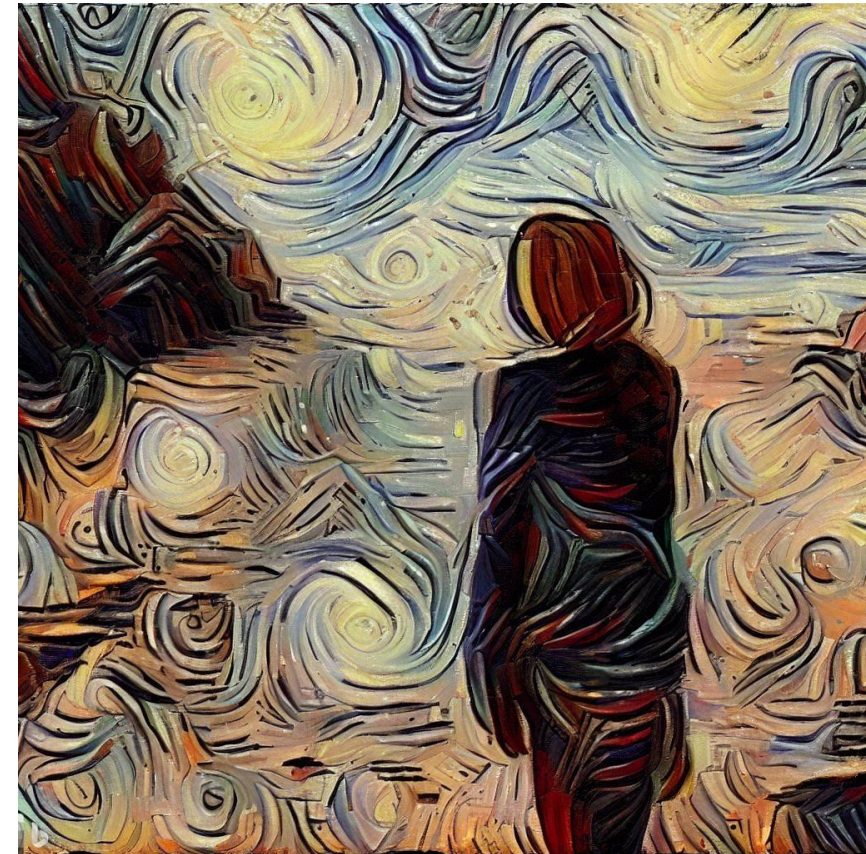## *The Difficulty of Evaluating Image Captioning*



## Hamid Palangi

**New Frontiers for Zero-Shot Image Captioning @CVPR2023**
**June 18, 2023, Vancouver, Canada**

# What I want to achieve in this session?

- **To convince you that Image Captioning Evaluation is Difficult**

- To discuss Possible Ways Forward

- [if time permits] To show relevant problems in text to image generation
    - How good Text-to-Image models are in terms of simple spatial relationships? https://github.com/microsoft/VISOR

# What do you see in this image?

# What do you see in this image?



**A group of men standing in a room**

# How about this one?

# How about this one?



**A group of people standing in a room**

# CodaLab

## Microsoft COCO Image Captioning Challenge

Organized by tl483 - Current server time: June 17, 2023, 2:05 a.m. UTC

▶ **Current**

Challenge

March 15, 2015, midnight UTC

**End**

Competition Ends

Never

Learn the Details    Phases    Participate    **Results**    Forums →

**Challenge**

**Phase description**
None

**Max submissions per day:** 1

**Max submissions total:** 5

⊞ Download CSV

### Results

| # | User | Entries | Date of Last Entry | BLEU-1 | | BLEU-2 | | BLEU-3 | | BLEU-4 | | METEOR | | ROUGE-L | | CIDEr-D | |
|---|------|---------|---------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| | | | | c5 ▲ | c40 ▲ | c5 ▲ | c40 ▲ | c5 ▲ | c40 ▲ | c5 ▲ | c40 ▲ | c5 ▲ | c40 ▲ | c5 ▲ | c40 ▲ | c5 ▲ | c40 ▲ |
| 1 | MS_Cog_Svcs-GIT2-Single_Model | 1 | 08/22/22 | 0.843 (2) | 0.981 (1) | 0.700 (2) | 0.944 (1) | 0.557 (2) | 0.876 (2) | 0.432 (3) | 0.783 (3) | 0.319 (2) | 0.421 (2) | 0.620 (2) | 0.784 (3) | 1.464 (2) | 1.498 (1) |
| 2 | OFA-Sys_OFA | 4 | 05/31/22 | 0.845 (1) | 0.981 (2) | 0.701 (1) | 0.944 (2) | 0.559 (1) | 0.878 (1) | 0.436 (1) | 0.787 (1) | 0.321 (1) | 0.427 (1) | 0.625 (1) | 0.790 (1) | 1.472 (1) | 1.496 (2) |
| 3 | MS_Cog_Svcs-GIT-Single_Model | 3 | 05/30/22 | 0.840 (5) | 0.979 (3) | 0.698 (3) | 0.944 (3) | 0.556 (3) | 0.876 (3) | 0.432 (2) | 0.783 (2) | 0.319 (3) | 0.420 (3) | 0.620 (3) | 0.784 (2) | 1.455 (3) | 1.488 (3) |

# OFA-Image_Caption

Gradio Demo for OFA-Image_Caption. Upload your own image or click any one of the examples, and click "Submit" and then wait for the generated caption.



**Caption**

a group of men in a room posing for a picture

Clear          Submit

# OFA-Image_Caption

Gradio Demo for OFA-Image_Caption. Upload your own image or click any one of the examples, and click "Submit" and then wait for the generated caption.



**Caption**

a group of people standing in a room

Clear          Submit

# Image Captioning Evaluation is Difficult

- Three important sources of difficulty:

    1. **Metrics**: How helpful are evaluation metrics for the usage goal?

    2. **Data**: Is the data a good representative of real world and/or target usage?

    3. **Models**: How much of the issues from different model components (e.g., object detector, tagger, etc) propagates in the output? Can model architectures game the metrics?

# Image Captioning Evaluation is Difficult

- BLEU, METEOR, ROUGE and CIDEr are doing evaluation based on n-grams focusing on lexical similarities

- SPICE uses scene graph of the image and runs dependency parser on the caption

- BERTScore and similar model based metrics use pretrained models to measure semantic similarity between reference captions (5 in the case of COCO) and generated caption

- …

# Image Captioning Evaluation is Difficult

- **Commonsense and Context Understanding** in the generated captions is complex and challenging to measure

- Evaluations are usually **Subjective** given the creativity of the task

- **Coverage** is an important factor, e.g., are the 5 reference captions in COCO covering all the aspects of the image?

- Measuring **Visual Grounding** is challenging when aspects beyond object name grounding in the generated captions are considered

- In object/attribute/relationship detection there are ground truth labels, in image captioning (and text generation in general) there is **No Gold Label**.

# Image Captioning Evaluation is Difficult

- Captions are **Challenging to Interpret**, e.g., out of all the objects/attributes/relations detected by a given detector, why the language model did or did not select specific set of them as salient to include in the caption?

- **User Studies** are essential but time consuming, how to address different aspects of this problem?

- Evaluation metrics should cover **Different Languages** and various aspects of each language.

- Metrics that account for various **Deployment Constraints**, what are different inference time shipping thresholds for image captioning systems in practice?

# Image Captioning Evaluation is Difficult

- How to measure **Creativity vs Factuality**?

- How to evaluate **Rare and Out Of Distribution (OOD)** scenarios? Their importance is not reflected in the aggregate metrics.

- How to evaluate **Robustness to Real World Perturbations**? E.g., does a JPEG compression and decompression have a huge impact on the performance?

- How to **Incorporate User Feedback** in the metrics and evaluation process?

- Do current metrics account for **Interactive and Dialogue Scenarios**? Is it trivial to extend them to measure quality of an image grounded conversation?

# Image Captioning Evaluation is Difficult

- Difficulty of creating metrics to measure **Fairness-Related Harms**:

    1. Generated captions reflect the most salient aspect of the image, which is by nature a subjective choice: the system might show differential treatment toward various social groups.

    2. Generated attributes and verbs for specific objects might be systematically assigned to specific social groups but not others.

    3. The issues mentioned in (1) and (2) above might be caused by the object detector, tagger, language model, or the complex interactions among them. What are good metrics to disentangle them?

    4. Would the language model erase objects and tags related to specific groups of people?

# What I want to achieve in this session?

- To convince you that Image Captioning Evaluation is Difficult

- **To discuss A Few Possible Ways Forward**

- [if time permits] To show relevant problems in text to image generation
  - How good Text-to-Image models are in terms of simple spatial relationships?
    https://github.com/microsoft/VISOR

# A Few Possible Ways Forward

- Enforcing the desired characteristics explicitly, e.g., **Relations**



a baseball player swinging a bat

a baseball player holding a bat

# A Few Possible Ways Forward
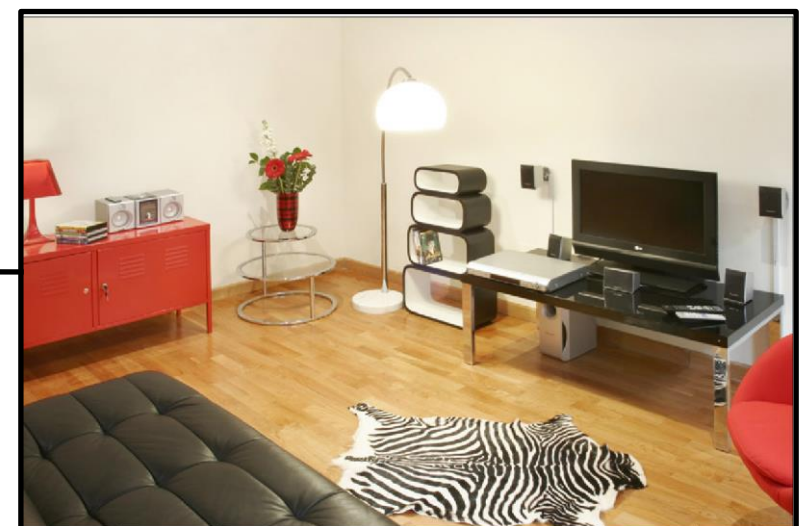
• Enforcing the desired characteristics explicitly, e.g., **Relations**



| Type | Examples | Classes | Instances |
|------|----------|---------|-----------|
| Entities | | | |
| Part | arm, tail, wheel | 32 | 200k (25.2%) |
| Artifact | basket, fork, towel | 34 | 126k (16.0%) |
| Person | boy, kid, woman | 13 | 113k (14.3%) |
| Clothes | cap, jean, sneaker | 16 | 91k (11.5%) |
| Vehicle | airplane, bike, truck, | 12 | 44k (5.6%) |
| Flora | flower, plant, tree | 3 | 44k (5.5%) |
| Location | beach, room, sidewalk | 11 | 39k (4.9%) |
| Furniture | bed, desk, table | 9 | 37k (4.7%) |
| Animal | bear, giraffe, zebra | 11 | 30k (3.8%) |
| Structure | fence, post, sign | 3 | 30k (3.8%) |
| Building | building, house | 2 | 24k (3.1%) |
| Food | banana, orange, pizza | 6 | 13k (1.6%) |
| Relations | | | |
| Geometric | above, behind, under | 15 | 228k (50.0%) |
| Possessive | has, part of, wearing | 8 | 186k (40.9%) |
| Semantic | carrying, eating, using | 24 | 39k (8.7%) |
| Misc | for, from, made of | 3 | 2k (0.3%) |

Table 1. Object and relation types in Visual Genome, organized by super-type. Most, 25.2% of entities are parts and 90.9% of relations are geometric or possessive.

# A Few Possible Ways Forward

- Enforcing the desired characteristics explicitly, e.g., **Relations**
  - Discarding relationships classified with high confidence using a simple prior network.
  - From top 1600 objects and 500 relations in VG
    - Showing each predicate by Glove, run clustering to remove duplicates, e.g., "wears" and "is wearing a" ➔ 180 rels
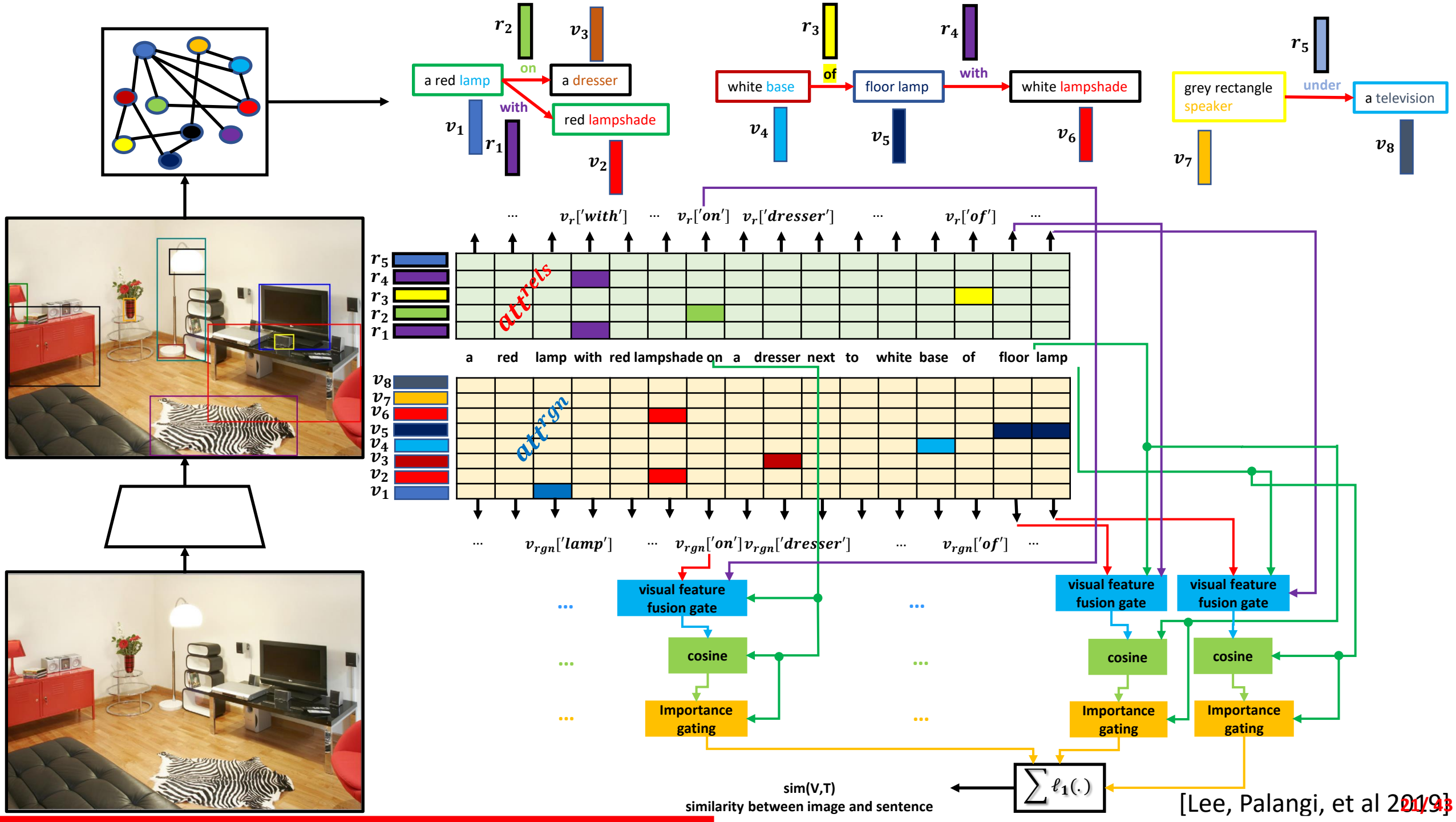    - Running prior network, removing relations that can be predicated with > 50% accuracy ➔ 117 rels
    - 58,983 images

# A Few Possible Ways Forward

- Enforcing the desired characteristics explicitly, e.g., **Relations**



| Metrics | SGDet | SGCls | PredCls | Metrics | SGDet | SGCls | PredCls |
|---------|-------|-------|---------|---------|-------|-------|---------|
| R50 | 27.2 | 35.8 | 65.2 | R50 | 14.8 | 16.5 | 46.7 |
| R100 | 30.3 | 36.5 | 67.1 | R100 | 17.4 | 19.2 | 52.5 |

Figures from https://arxiv.org/abs/1902.00313

# A Few Possible Ways Forward

- Enforcing the desired characteristics explicitly, e.g., **Relations** [Lee, Palangi, et al 2019]

[Lee, Palangi, et al 2019]

sim(V,T)
similarity between image and sentence

# A Few Possible Ways Forward

- Enforcing the desired characteristics explicitly, e.g., **Relations** [Lee, Palangi, et al 2019]



**Baseline**: a man standing on the side of a road
**New**: a man repairing a traffic light at an intersection



**Baseline**: a woman standing on a sidewalk talking on a cell phone
**New**: a woman standing on a sidewalk looking at her cell phone



**Baseline**: a man holding a nintendo wii game controller
**New**: a man sitting on a couch holding a wii remote

Figures from https://arxiv.org/pdf/1909.09953.pdf

# A Few Possible Ways Forward

- Having a unified architecture where different tasks help each other, e.g., **VQA helping Image Captioning** [Zhou, Palangi, et al 2020]



**Unified VLP:** A group of people standing under umbrellas in the rain.

**Question:** Are they dressed the same?

**Unified VLP:** Yes.
**Initialization from UniLM:** No.
**Initialization from BERT:** No.

**Image Captioning**

A girl with an upside-down umbrella.

**Unified Encoder-Decoder**

**Visual Question Answering**

A: Yes

**Unified Encoder-Decoder**

Is    the    ···    down

Q: Is the umbrella upside down?

Figures from https://ojs.aaai.org/index.php/AAAI/article/view/7005

# A Few Possible Ways Forward

- Collecting data that reflects a specific failure case, e.g., **NICE dataset** including about 2 million images and corresponding text that reflects qualities like empathy and emotion [Chen, Huang, et al 2021]



(a) NICE

(b) MS-COCO

**Image Captioning:**
Some houses are on the foot of mountains.

**Image Commenting:**
*Comment 1:* Gorgeous!
*Comment 2:* Where it is?
*Comment 3:* I don't know.
*Comment 4:* The scenery here is so beautiful.
*Comment 5:* I really want to go there.

# A Few Possible Ways Forward

- Collecting data that reflects a specific failure case, e.g., **AdaVision** that is a human in the loop tool for data collection [Gao, Ilharco, et al 2022]

# A Few Possible Ways Forward

- Collecting data that reflects a specific failure case, e.g., measuring specific **Representational Harms** [Wang, Barocas, et al 2022]
  - *Stereotyping*: which occurs when oversimplified beliefs about social groups reproduce harmful social hierarchies
  - *Demeaning*: which occurs when social groups are cast as being lower status and less deserving of respect



Figures from https://arxiv.org/abs/2206.07173

# A Few Possible Ways Forward

- Collecting data that reflects a specific failure case, e.g., using **explainability methods like SHAP** to interpret the captions w.r.t the input image [Lundberg, Lee, 2017]

# A Few Possible Ways Forward

- Leveraging the recent large scale language models and perform instruction tuning, e.g., using **LLaVA** to interact about the input image [Liu, Li, et al 2023]

# A Few Possible Ways Forward

- Let's take a closer look at **LLaVA** [Liu, Li, et al 2023]

- The red boxes show a combination of errors caused by the vision components (e.g., most of the people are not smiling or looking at the camera) and **Hallucination** caused by the LLM.

- The model still totally loses the point about only two people looking at each other in the group.



Figures from https://llava.hliu.cc/

# A Few Possible Ways Forward

- Let's take a closer look at **LLaVA** [Liu, Li, et al 2023]

- One more round to show case the **Hallucination**.



LLaVA Chatbot

The image features a group of people standing in a circle, all smiling and looking at the camera. They appear to be posing for a group photo, with some of them wearing ties. The group consists of both adults and children, suggesting it might be a family gathering or a social event. The people are positioned closely together, creating a sense of unity and camaraderie among the group members.

Where did the man and the woman looking at each other did their graduate school?

The man and the woman looking at each other did their graduate school at the University of Wisconsin.

# A Few Possible Ways Forward

- Let's take a closer look at **LLaVA** [Liu, Li, et al 2023]

- One more round to show case the **Hallucination**.



As a youth, Murray read children's biographies of American heroes like Kit Carson, Wild Bill Hickok, and Davy Crockett.[12] He attended St. Joseph's grade school and Loyola Academy. During his teen years, he worked as a golf caddy to fund his education at the Jesuit high school,[12][14] and was the lead singer of a rock band called the Dutch Masters and took part in high school and community theater.[15] One of his sisters had polio and his mother suffered several miscarriages.[12]

After graduating from Loyola Academy, Murray attended Regis University in Denver, Colorado, taking pre-medical courses, but quickly dropped out and returned to Illinois.[12] Decades later, in 2007, Regis awarded him an honorary Doctor of Humanities degree.[16] On September 21, 1970, his 20th birthday, the police

LLaVA Chatbot

The image features a group of people standing in a circle, all smiling and looking at the camera. They appear to be posing for a group photo, with some of them wearing ties. The group consists of both adults and children, suggesting it might be a family gathering or a social event. The people are positioned closely together, creating a sense of unity and camaraderie among the group members.

Where did the man and the woman looking at each other did their graduate school?

The man and the woman looking at each other did their graduate school at the University of Wisconsin.

Scarlett Johansson

Johansson at the 2019 San Diego Comic-Con

**Born** — Scarlett Ingrid Johansson November 22, 1984 (age 38) New York City, U.S.

**Citizenship** — United States · Denmark

**Education** — Professional Children's School

Figures from https://llava.hliu.cc/

# A Few Possible Ways Forward

- Let's take a closer look at **LLaVA** [Liu, Li, et al 2023]

- One more round to show case the **Hallucination**.



LLaVA Chatbot

The image features a group of people standing in a circle, all smiling and looking at the camera. They appear to be posing for a group photo, with some of them wearing ties. The group consists of both adults and children, suggesting it might be a family gathering or a social event. The people are positioned closely together, creating a sense of unity and camaraderie among the group members.

Where did the man and the woman looking at each other did their graduate school?

The man and the woman looking at each other did their graduate school at the University of Wisconsin.

who are they?

The man and the woman looking at each other are a couple, and they are both wearing ties. They could be the hosts of the gathering or possibly the parents of the child in the group.

# A Few Possible Ways Forward

- Let's take a closer look at **LLaVA** [Liu, Li, et al 2023]

- One more round to show case the **Hallucination**.



> **LLaVA Chatbot**
>
> together, creating a sense of unity and camaraderie among the group members.
>
> Where did the man and the woman looking at each other did their graduate school?
>
> The man and the woman looking at each other did their graduate school at the University of Wisconsin.
>
> who are they?
>
> The man and the woman looking at each other are a couple, and they are both wearing ties. They could be the hosts of the gathering or possibly the parents of the child in the group.
>
> how old is the child?
>
> The child is described as a "little girl" and a "small child," indicating that she is likely between the ages of 4 and 10 years old.

Figures from https://llava.hliu.cc/

# A Few Possible Ways Forward

- Let's take a closer look at **LLaVA** [Liu, Li, et al 2023]

- One more round to show case the **Hallucination**.



Figures from https://llava.hliu.cc/

# What I want to achieve in this session?

- To convince you that Image Captioning Evaluation is Difficult

- To discuss Possible Ways Forward

- [if time permits] **To show relevant problems in text to image generation**
  - **How good Text-to-Image models are in terms of simple spatial relationships? https://github.com/microsoft/VISOR**

# Text-to-Image (T2I) Models: Transforming Words into Images

**2016**



"GAN_INT_CLS" Reed et al. ICML 2016 https://arxiv.org/abs/1605.05396

"StackGAN" Zhang et al. ICCV 2017 https://arxiv.org/abs/1612.03242

# Text-to-Image (T2I) Models: Transforming Words into Images

## 2022



panda mad scientist mixing sparkling chemicals, artstation

a corgi's head depicted as an explosion of a nebula

a propaganda poster depicting a cat dressed as french emperor napoleon holding a piece of cheese

a teddy bear on a skateboard in times square

A giant cobra snake on a farm. The snake is made out of corn.

A chrome-plated duck with a golden beak arguing with an angry turtle in a forest.

A dog looking curiously in the mirror, seeing a cat.

An extremely angry bird.

An oil painting of two rabbits in the style of American Gothic, wearing the same clothes as in the original

a portrait of a statue of the Egyptian god Anubis wearing aviator goggles, white t-shirt and leather jacket. The city of Los Angeles is in the background. hi-res dslr photograph.

A dignified beaver wearing glasses, a vest, and colorful neck tie. He stands next to a tall stack of books in a library. dslr photo.

The buildings of downtown Manhattan situated below Mount Everest. The Great Pyramid is in the foreground. DSLR photograph.

DALLE-v2Ramesh et al. preprint 2022 https://arxiv.org/abs/2204.06125

"Imagen": Saharia et al. preprint 2022 https://arxiv.org/abs/2205.11487

"PARTI": Yu et al. preprint 2022 https://arxiv.org/abs/2206.10789

# Survey of Existing Metrics for T2I Evaluation

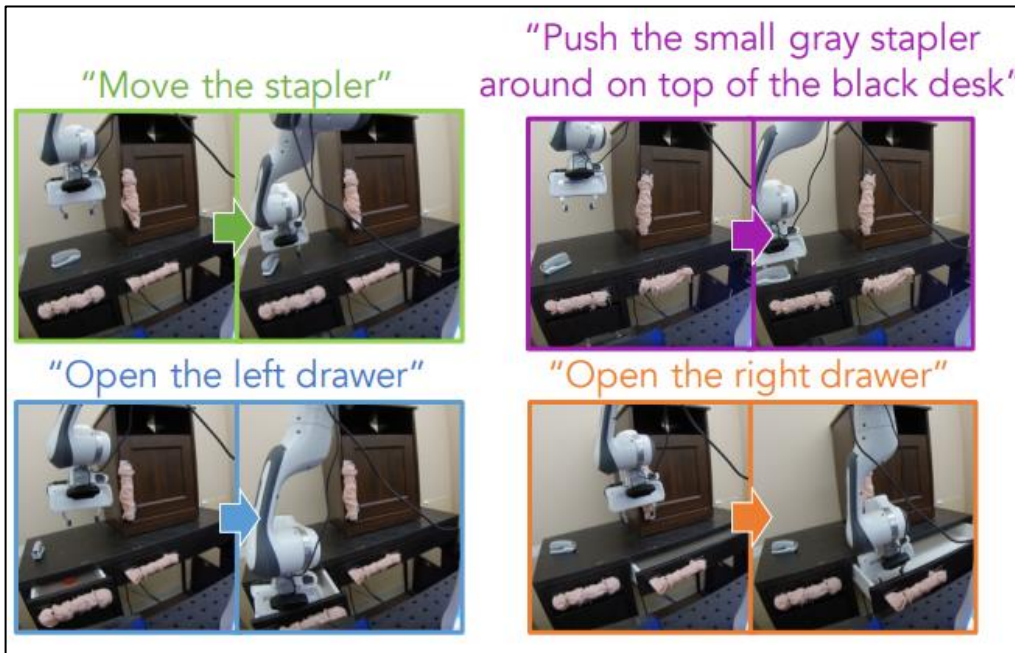| | StackGAN (Zhang et al. ICCV 2017) | DM-GAN (Zhu et al. CVPR 2019) | OP-GAN (Hinz et al. TPAMI 2020) | GLIDE (Nichol et al. NeurIPS 2021) | CogView-1/2 (Ding et al. NeurIPS 2021) | DALLE v1/v2 (Ramesh et al. 2021/2022) | Stable Diffusion (Rombach et al. CVPR 2022) |
|---|---|---|---|---|---|---|---|
| **IS:** Inception Score (Salimans et al. NeurIPS 2016) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| **FID:** Frechet Inception Distance (Heusel et al. NeurIPS 2017) | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| **R-Precision** (Xu et al. CVPR 2018) | | ✓ | ✓ | ✓ | | | |
| **Image Captioning Metrics** (Hong et al. CVPR 2018) | | | ✓ | | | | |
| **CLIPscore** (Hessel et al. EMNLP 2021) | | | | ✓ | | | |
| **SOA:** Semantic Object Accuracy (Hinz et al. T-PAMI 2020) | | | ✓ | | | | |
| Human Study | ✓ | | | ✓ | ✓ | ✓ | ✓ |

Four categories of existing evaluation metrics
1. Purely Visual Metrics for Photorealism : IS, FID
2. Image-text matching : Image Captioning / CLIPscore
3. Object-Level : SOA
4. Human study

# Spatial Relationships Directly Impact Practical Applications!

- Humans often communicate using spatial relationships between objects
  - "The Windows start button is on the bottom left of the screen"
  - "The toaster in my kitchen is to the left of the stove"
  - "Can you move the chart to the right of the text in your PPT?"

- Spatial relationships are also important for robotics and applications



Place the mug on the right of the box



"Move the stapler"

"Push the small gray stapler around on top of the black desk"

"Open the left drawer"

"Open the right drawer"



Standing in front of the family picture, turn left and walk straight through the bathroom past the tub and mirrors. Go through the doorway and stop when the door to the bathroom is on your right and the door to the closet is to your left.

Walk with the family photo on your right. Continue straight into the bathroom. Walk past the bathtub. Stop in the hall between the bathroom and toilet doorways.

Walk straight passed bathtub and stop with closet on the left and toilet on the right.

Mees et al. ICRA 2020 https://arxiv.org/abs/2001.08481v2

Nair et al. CORL 2020 https://proceedings.mlr.press/v164/nair22a/nair22a.pdf

Anderson et al. CVPR 2018
ttps://arxiv.org/pdf/1711.07280.pdf

# CLIPScore *is not effective*

CLIPscore: Hessel et al. EMNLP 2021 https://arxiv.org/abs/2104.08718

- CLIPscore is a recent T2I evaluation metric that uses CLIP to obtain scores for text—image pairs

- Consider text $t$

- Let $x$ be the image generated for $t$

- Let $t_{flipped}$ be the text with the same objects, but a flipped relation –

- Compute $CLIPscore(x, t)$ and $CLIPscore(x, t_{flipped})$

$t_{flipped}$ does not match with $x$ (it has the opposite relationship)
Yet, $CLIPscore(x, t) \sim CLIPscore(x, t_{flipped})$

➔ *CLIPscore is not a good metric for evaluating spatial relationships generated by T2I models*

t = "A backpack below a train"

CLIPscore(x, t) = 0.8691

t_flipped = "A backpack above a train"

CLIPscore(x, t_flipped) = 0.8779
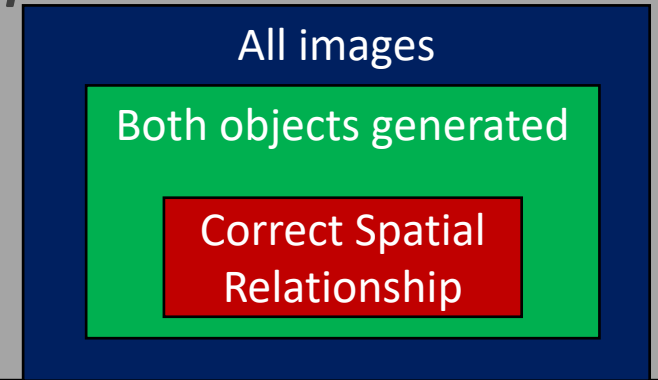
t="A toaster to the left of a fire hydrant"

CLIPscore(x, t) = 0.8052

t_flipped = "A toaster to the right of a fire hydrant"

CLIPscore(x, t_flipped) = 0.8086

40/ 43

# Verifying Spatial Object Relationships (VISOR)

$$VISOR = P(R_{gen} == R_{GT} \mid \exists A \cap \exists B) = \frac{\textcolor{red}{\blacksquare}}{\textcolor{green}{\blacksquare}}$$

All images

Both objects generated

Correct Spatial Relationship

| Model | OA | SRA | $VISOR_1$ | $VISOR_2$ | $VISOR_3$ | $VISOR_4$ |
|---|---|---|---|---|---|---|
| GLIDE | 3.36 | 1.98 | 6.72 | 1.02 | 0.17 | 0.03 |
| DALLE-mini | 27.1 | 16.17 | 38.31 | 17.5 | 6.89 | 1.96 |
| CogView2 | 18.47 | 12.17 | 33.47 | 11.43 | 3.22 | 0.57 |
| DALLE-v2 | 63.93 | 37.89 | 73.59 | 47.23 | 23.26 | 7.49 |
| Stable Diffusion | 29.86 | 18.81 | 40.6 | 20.11 | 6.89 | 1.63 |



For more comprehensive results please refer to https://arxiv.org/abs/2212.10015

a skateboard below a mouse

a mouse to the left of a couch

DALLE-v2



Stable Diffusion

**hpalangi@microsoft.com**  **www.hamidpalangi.com**