# Generating Images with Multimodal Language Models

**Jing Yu Koh**, Daniel Fried, Ruslan Salakhutdinov 6 September 2023 Cohere for Al







Can we ground text-only LLMs to pretrained visual encoders and decoders?



<u>Generating Images with Large Language Models</u>



### **GILL: A More General Multimodal LM**



### **GILL: A More General Multimodal LM**

- Frozen (Tsimpoukelli et al., 2021)
   Flamingo (Alayrac et al., 2022)
   BLIP-2 (Li et al., 2023)
  - Process image + text, generate text only
- FROMAGe (Koh et al., 2023)
  - Process image + text, generate text + retrieve images
- **GILL** (this work)
  - Process image + text, generate text + retrieve images + generate images
  - Decides whether to retrieve images or generate from scratch
  - Resource efficient: trained on 2 GPUs for 2 days



<u>Generating Images with Large Language Models</u>

#### • Capable of retrieving images, generating images, and generating text

- Can condition on arbitrarily interleaved image + text inputs
- Generate text, generate images, and retrieve images as part of the output

#### • Leverage the learnt abilities of pre-trained text-only LLMs

- In-context learning
- Sensitivity to input prompts
- Generate long and coherent dialogue

#### Model agnostic

- We use a 7B LLM, the CLIP encoder, and the Stable Diffusion image generator
- Likely benefits from using larger and stronger LLMs in the future
- Can be applied with other visual models (e.g., OCR) to introduce new abilities



An European shorthair cat in a woven basket

Image #1 Caption #1
Image and Caption Inputs
(from CC3M)





Frozen Model

Linear Layer

Loss



An European shorthair cat in a woven basket [IMG1]...[IMG{r}]

#### **Caption Input**



Image Input







#### **Caption Input**



Image Input









Frozen Model 📃 GILLMapper 📃 Linear Layer



### **GILLMapper: An Improved LLM-to-Generator Map**

- Previous approaches use <u>linear mappings</u> between LLMs and visual models
- This is insufficient for image generation: decoders require <u>dense</u> information



Multimodal Few-Shot Learning with Frozen Language Models (<u>Tsimpoukelli et al., 2021</u>) Linearly Mapping from Image to Text Space (<u>Merullo et al., 2023</u>) Grounding Language Models to Images for Multimodal Inputs and Outputs (<u>Koh et al., 2023</u>)

















• Given a Visual Story, generate a relevant image



Image and Text Inputs

- Given a Visual Story, generate a relevant image
- Need to condition on long, temporally dependent text
- (Optionally) Condition on image inputs interleaved within the text



Image and Text Inputs

Stable Diffusion

Ours

Groundtruth

|                                     | CLIP Similarity (↑)           |  |                             | LPIPS $(\downarrow)$                                     |  |                             |
|-------------------------------------|-------------------------------|--|-----------------------------|--|--|-----------------------------|
| Model                               | 1 caption                     | 5 captions   | 5 caps, 4 images            | 1 caption  | 5 captions   | 5 caps, 4 images            |
| GLIDE [34]<br>Stable Diffusion [43] | 0.582<br><b>0.592</b> ±0.0007 | $\begin{array}{c} 0.591 \\ 0.598 \pm 0.0006 \end{array}$ | -                           | $\begin{array}{c} 0.753 \\ 0.703 \pm 0.0003 \end{array}$ | $\begin{array}{c} 0.745 \\ 0.704 \pm 0.0004 \end{array}$ | -                           |
| GILL                                | $0.581 \pm 0.0005$            | $\textbf{0.612} \pm 0.0011$                              | $\textbf{0.641} \pm 0.0011$ | $0.702 \pm 0.0004$                                       | $\textbf{0.696} \pm 0.0008$                              | $\textbf{0.693} \pm 0.0008$ |

- Our model outperforms Stable Diffusion on longer input contexts
- This is despite GILL (essentially) distilling from SD!
- GILL benefits from the abilities of the LLM (sensitivity to longer inputs, word orderings, in-context learning)

• Given a Visual Dialogue, generate a relevant image



- Given a Visual Dialogue, generate a relevant image
- Need to condition on long dialogue-like text (OOD with finetuning data)



|                                     | CLIP Similarity (↑)           |                               | LPIPS $(\downarrow)$                                     |                               |  |  |
|-------------------------------------|-------------------------------|-------------------------------|--|-------------------------------|--|--|
| Model                               | 1 round                       | 5 rounds                      | 10 rounds  | 1 round                       | 5 rounds   | 10 rounds  |
| GLIDE [34]<br>Stable Diffusion [43] | <b>0.562</b><br>0.552 ±0.0015 | 0.595<br><b>0.629</b> ±0.0015 | $\begin{array}{c} 0.587 \\ 0.622 \pm 0.0012 \end{array}$ | 0.800<br><b>0.742</b> ±0.0010 | $\begin{array}{c} 0.794 \\ 0.722 \pm 0.0012 \end{array}$ | $\begin{array}{c} 0.799 \\ 0.723 \pm 0.0008 \end{array}$ |
| GILL                                | $0.528 \pm 0.0014$            | $0.621 \pm 0.0009$            | <b>0.645</b> ±0.0010                                     | $\textbf{0.742} \pm 0.0022$   | $\textbf{0.718} \pm 0.0028$                              | $\textbf{0.714} \pm 0.0006$                              |

#### **The Effect of Context**

Multi-modal context is **worth more** than uni-modal context, producing more relevant generation results.

Performance With Increasing Context on VIST



#### **GILLMapper: A Stronger LLM-to-Generator Mapping**

Image generators require **denser** input sequences. Linear mappings are insufficient.

|                            | CC3M                      | VIST         |
|----------------------------|---------------------------|--------------|
| Model                      | <b>FID</b> $(\downarrow)$ | CLIP Sim (†) |
| Stable Diffusion [43]      | 13.94                     | 0.598        |
| Ours + Linear              | 15.50                     | 0.500        |
| Ours + 3-layer MLP         | 15.33                     | 0.502        |
| Ours + Transformer Encoder | 16.30                     | 0.605        |
| Ours + GILLMapper          | 15.31                     | 0.641        |

### **Other Abilities: Text-to-Image Generation**





**Stable Diffusion** 

Ours

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





**Stable Diffusion** 

"A drop-top sports car coming around a bend in the road"





**Stable Diffusion** 

Ours

"Snow mountain and tree reflection in the lake"





**Stable Diffusion** 

Ours

"a group of penguins in a snowstorm"

### **Other Abilities: Image Refinement**



### **Future Work**

#### • Train on more diverse data

- CC3M is small by modern standards we would get a lot more from training on LAION
- Training on interleaved image-text data would also likely help a lot
- GILLMapper will likely be more aligned to SD
- Apply to even larger LLMs and stronger visual models
  - We use a 7B LLM, but you can likely train a 13B LLM with a few A6000 GPUs

#### • Use a finetuned LLM

• For example, instruction finetuned, or dialogue finetuned

#### • Perform more explicit image conditioning

• May allow the model to be better at tasks such as image editing

### Try the model!

#### huggingface.co/spaces/jykoh/gill

😕 Spaces 🔎 jykoh/gill 🖆 🎔 like 💈 💠 Running on A10G 🚍 Logs

🥪 App 📲

🖗 App 📲 Files 🥔 Community 🛛 🌼 Settings 📑 🖉 🛛

#### 🔍 GILL

This is the official Gradio demo for the GILL model, a model that can process arbitrarily interleaved image and text inputs, and produce image and text outputs.

Paper: <u>Generating Images with Multimodal Language Models</u> Project Website: <u>GILL Website</u> Code and Models: <u>GitHub</u>

#### Tips:

- Start by inputting either image or text prompts (or both) and chat with GILL to get image-and-text replies.
- Tweak the level of sensitivity to images and text using the parameters on the right.
- Check out cool conversations in the examples or community tab for inspiration and share your own!
- If the model outputs a blank image, it is because Stable Diffusion's safety filter detected inappropriate content. Please try again with a different prompt.
- Outputs may differ slightly from the paper due to slight implementation differences. For reproducing paper results, please use our official code.
- For faster inference without waiting in queue, you may duplicate the space and use your own GPU: 🔘 Duplicate Space

| 🕫 🕸 GILL Chatbot                          | Frequency multiplier for returning images (higher means 1.3 (*)<br>more frequent) |           |  |   |
|---|---|-----------|--|---|
| How can I publicize these?                |   |           |  |   |
| I would suggest you start with a local ne | ewspaper.   |           | Max # of words   | 32 🗘  |
| (Generated)                               |   |           | Temperature (0 for deterministic, higher for<br>randomness)                        |   |
|   | Message   | Submit    |  |   |
|   | Type a message  | Unde      | When any server upgetables ( can add as 17<br>Vous can add any segmables you like, |   |
|   |   | Undo      |  | Involve<br>Booking for some ideast for a tattoo. What do you think would<br>good on a partice person? |
| C Upload Image                            |   | Reset All |  |   |

# Thanks!

jykoh@cmu.edu jykoh.com/gill