

Generating Images with Multimodal Language Models

jykoh.com/gill

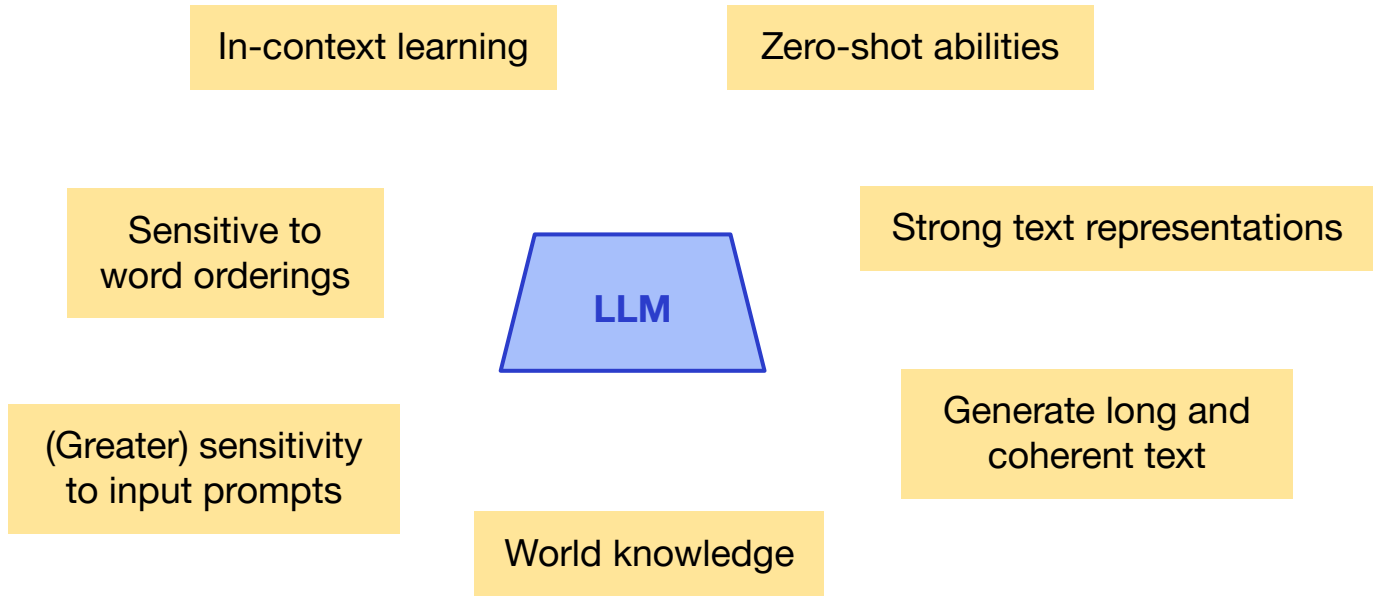
Jing Yu Koh, Daniel Fried, Ruslan Salakhutdinov

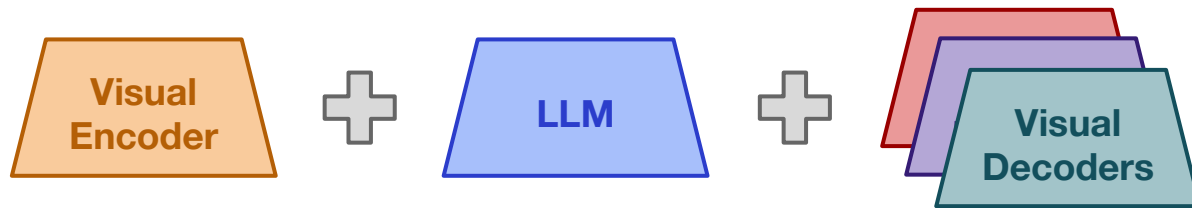
7 July 2023

Deep Learning: Classics and Trends (DLCT)

The logo for Carnegie Mellon University, featuring the text "Carnegie Mellon University" in white serif font on a red square background.


Carnegie
Mellon
University






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

Generating Images with Large Language Models



How can I make this more nutritious?


You can add vegetables to your ramen noodles, but you should be careful not to overdo it.




What are some vegetables I can add to it?

Broccoli, carrots, and green beans are all good choices.


I'm looking for some ideas for a new tattoo. What do you think looks good on a petite person?




I want a tattoo of a sunflower, but I'm not sure what style would look best on me.




That's beautiful. Show me one that would work well as a complementary piece





How should I publicize these at the market?

I would suggest a simple sign with the name of the business and a picture of the cupcakes.




GILL: A More General Multimodal LM

User prompts Retrieved Generated

I am thinking of an animal. It is brown and furry, and has a tail. What is it? 

 A squirrel. 

No, the animal I am thinking of lives in the river. 

 A beaver. 

Show me an illustration of a beaver 

A beaver. 



Yes, what about a pencil drawing of one? 




 What about a comic illustration of one? 










Ours


Retrieval and generation multimodal LM.
Decides when to retrieve or generate.


I am thinking of an animal. It is brown and furry, and has a tail. What is it? 


 A squirrel. 


No, the animal I am thinking of lives in the river. 


 A beaver. 



Show me an illustration of a beaver 

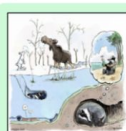
A beaver. 




Yes, what about a pencil drawing of one? 




 What about a comic illustration of one? 










FROMAGe


Retrieval only multimodal language model. Image outputs limited to the candidate retrieval set.


I am thinking of an animal. It is brown and furry, and has a tail. What is it? 


 A squirrel. 


No, the animal I am thinking of lives in the river. 


 A beaver. 



Show me an illustration of a beaver 


A beaver. 




Yes, what about a pencil drawing of one? 



 What about a comic illustration of one? 





Stable Diffusion

Generation only text-to-image model.
Less sensitive to longer text inputs (such as dialogue).

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



Generating Images with Large Language Models

- **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

Learning to *Process* Images



Image #1

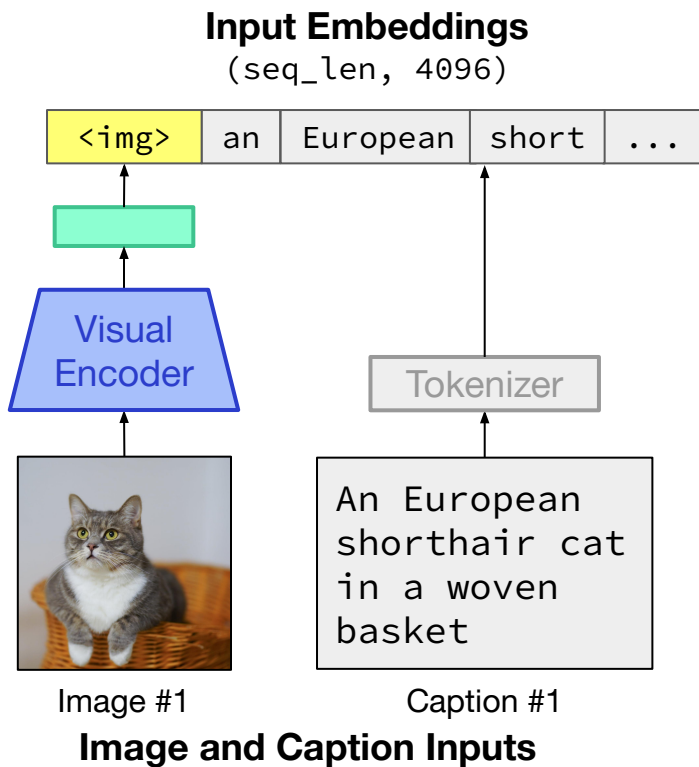
An European
shorthair cat
in a woven
basket

Caption #1

Image and Caption Inputs

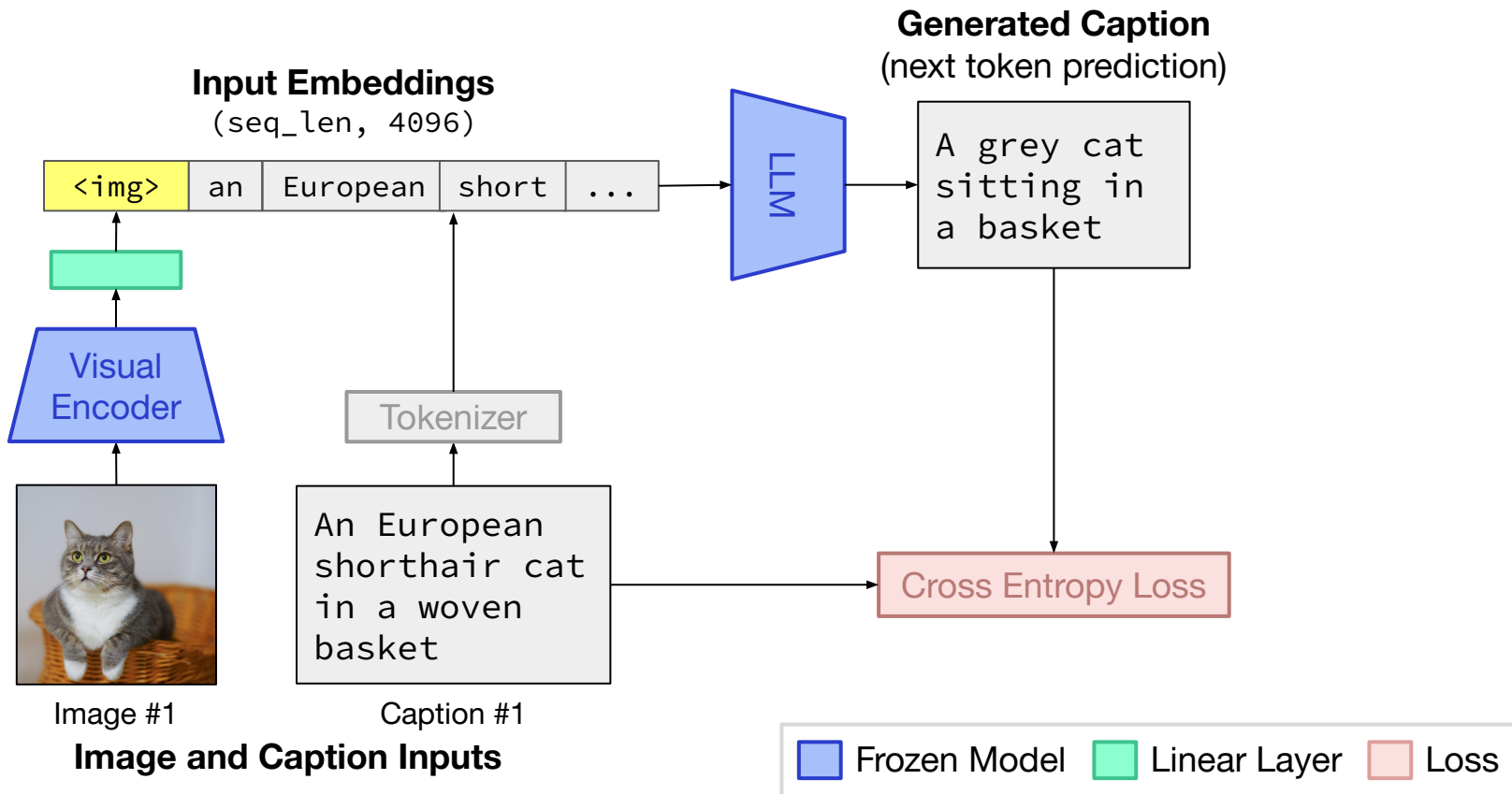
 Frozen Model  Linear Layer  Loss

Learning to *Process* Images



 Frozen Model  Linear Layer  Loss

Learning to *Process* Images



Learning to *Produce* Images

An European
shorthair cat in
a woven basket

[IMG1]...[IMG{r}]

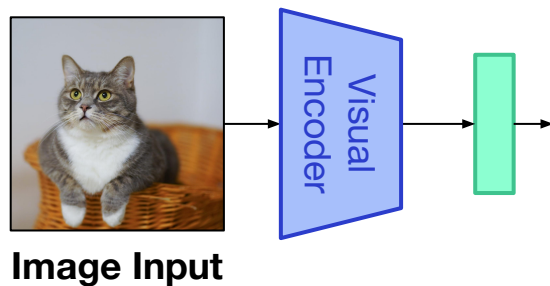
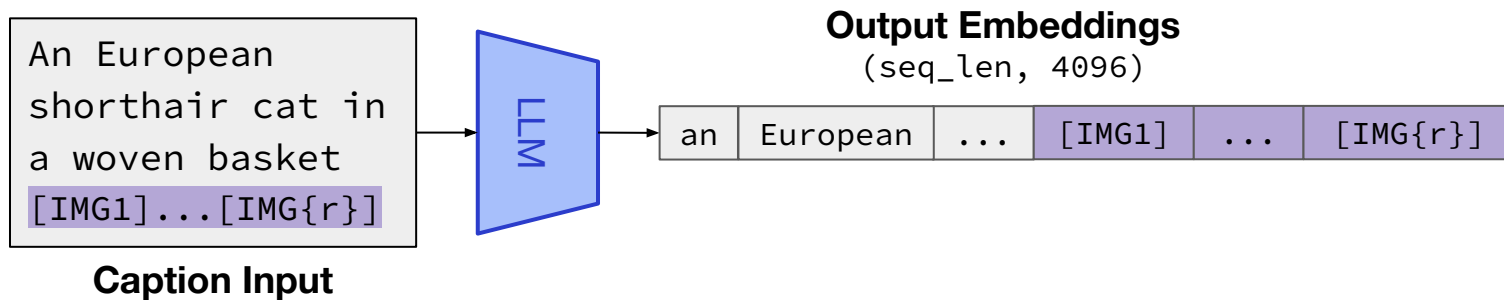
Caption Input



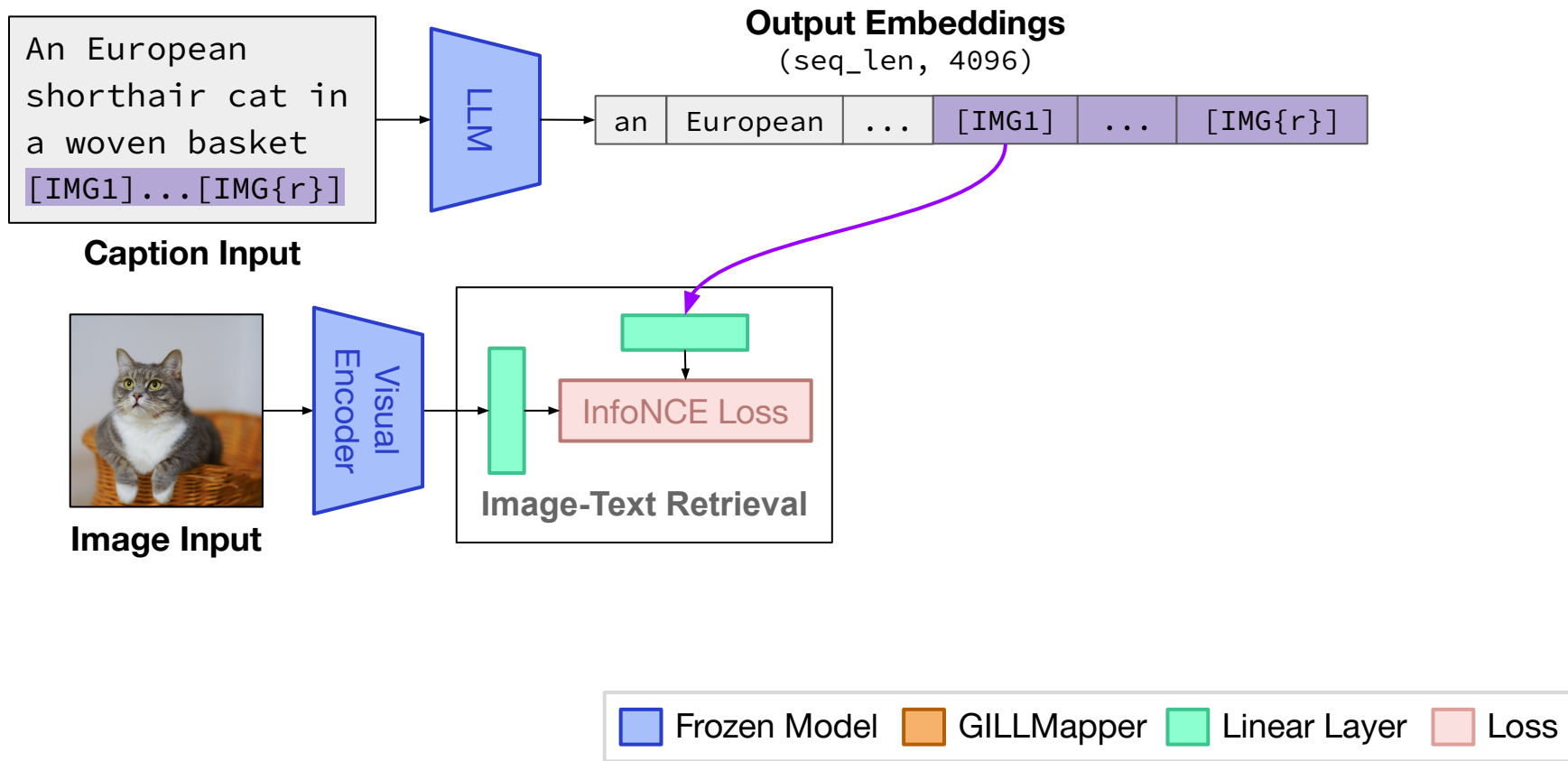
Image Input

 Frozen Model  GILLMapper  Linear Layer  Loss

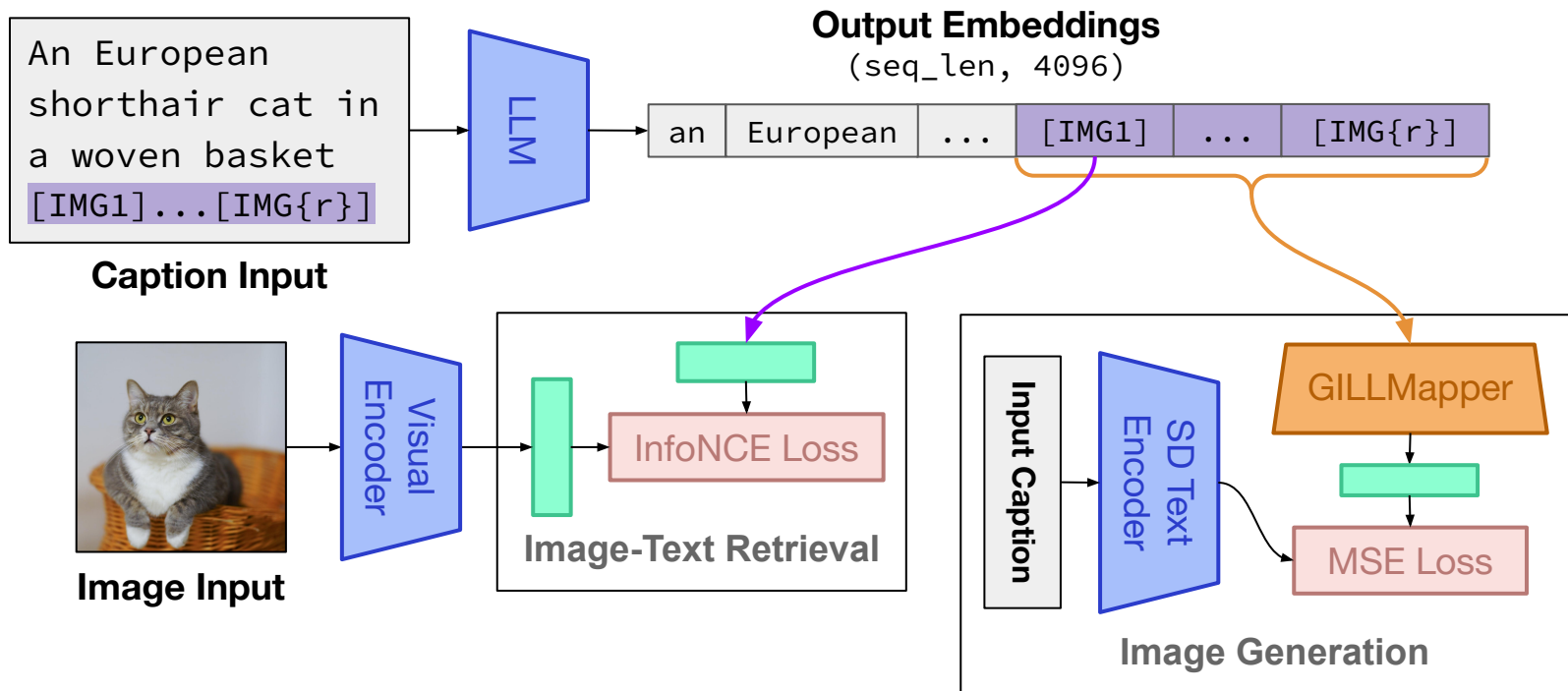
Learning to *Produce* Images



Learning to *Produce* Images



Learning to *Produce* Images



GILLMapper: An Improved LLM-to-Generator Map

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

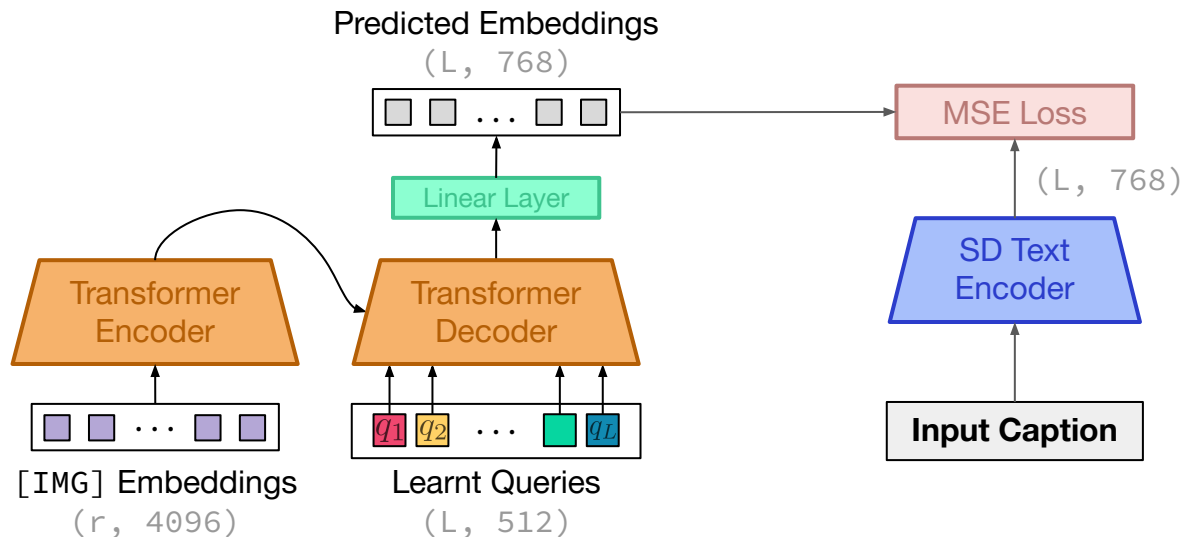


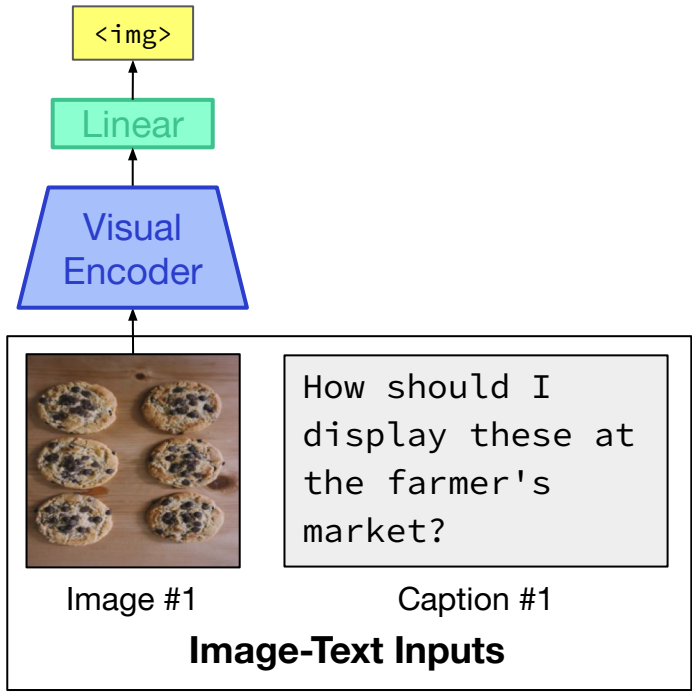


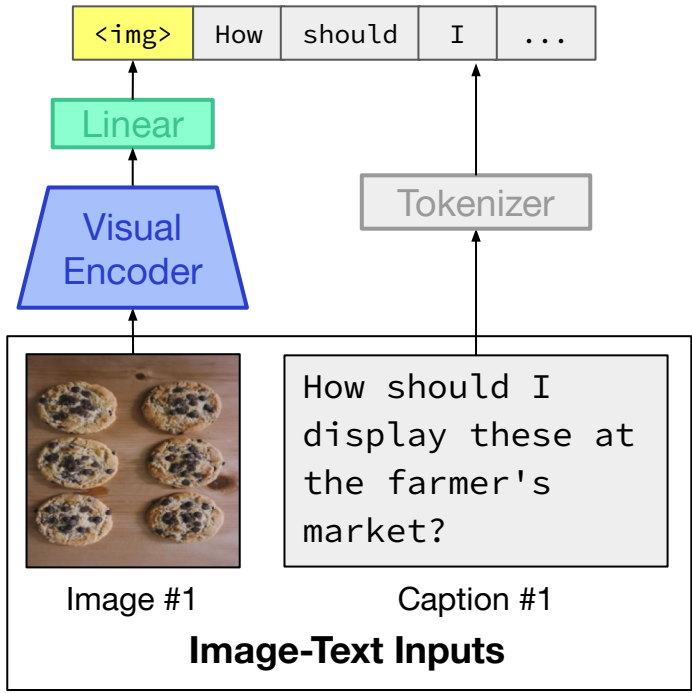
Image #1

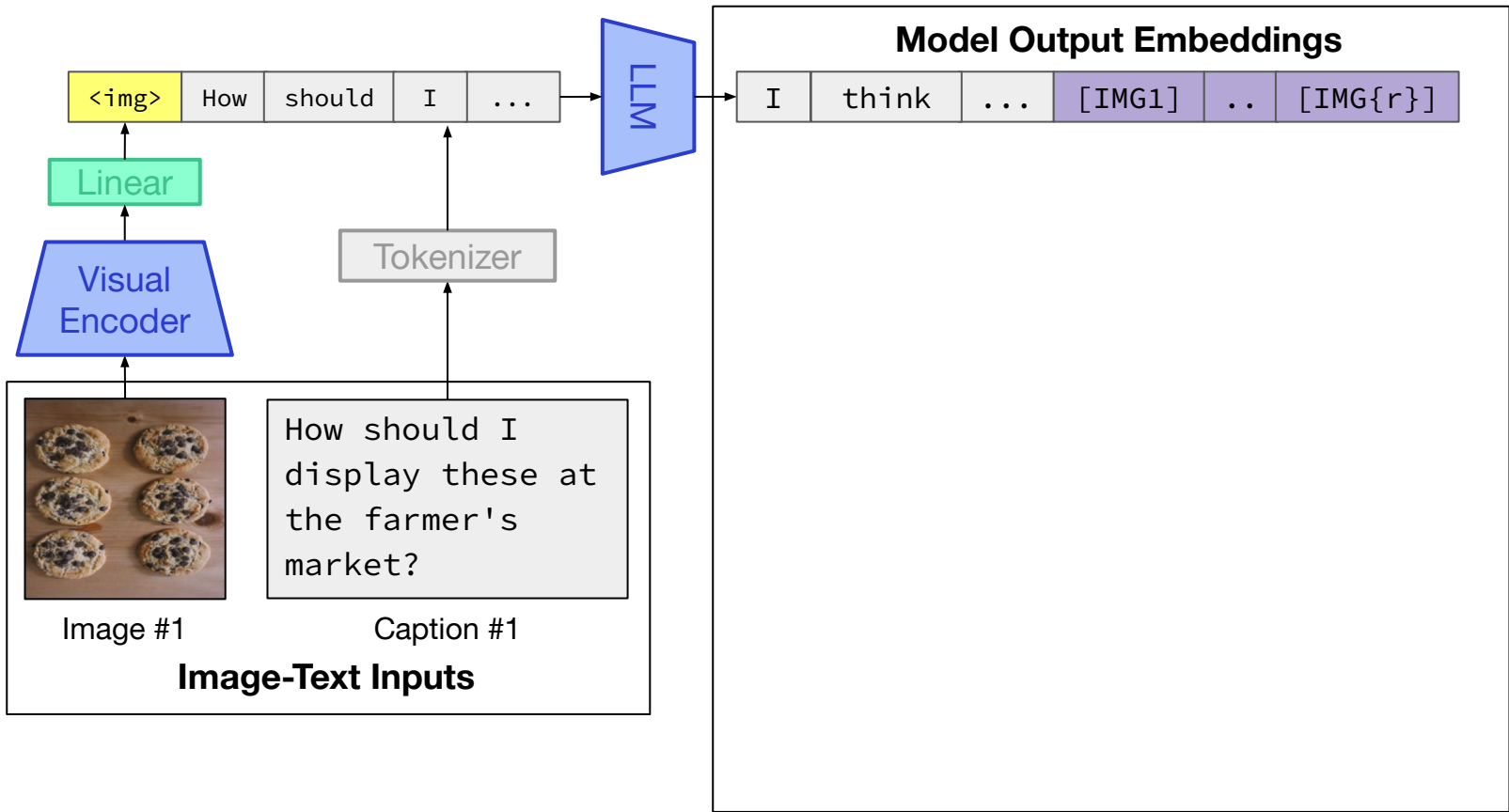
How should I
display these at
the farmer's
market?

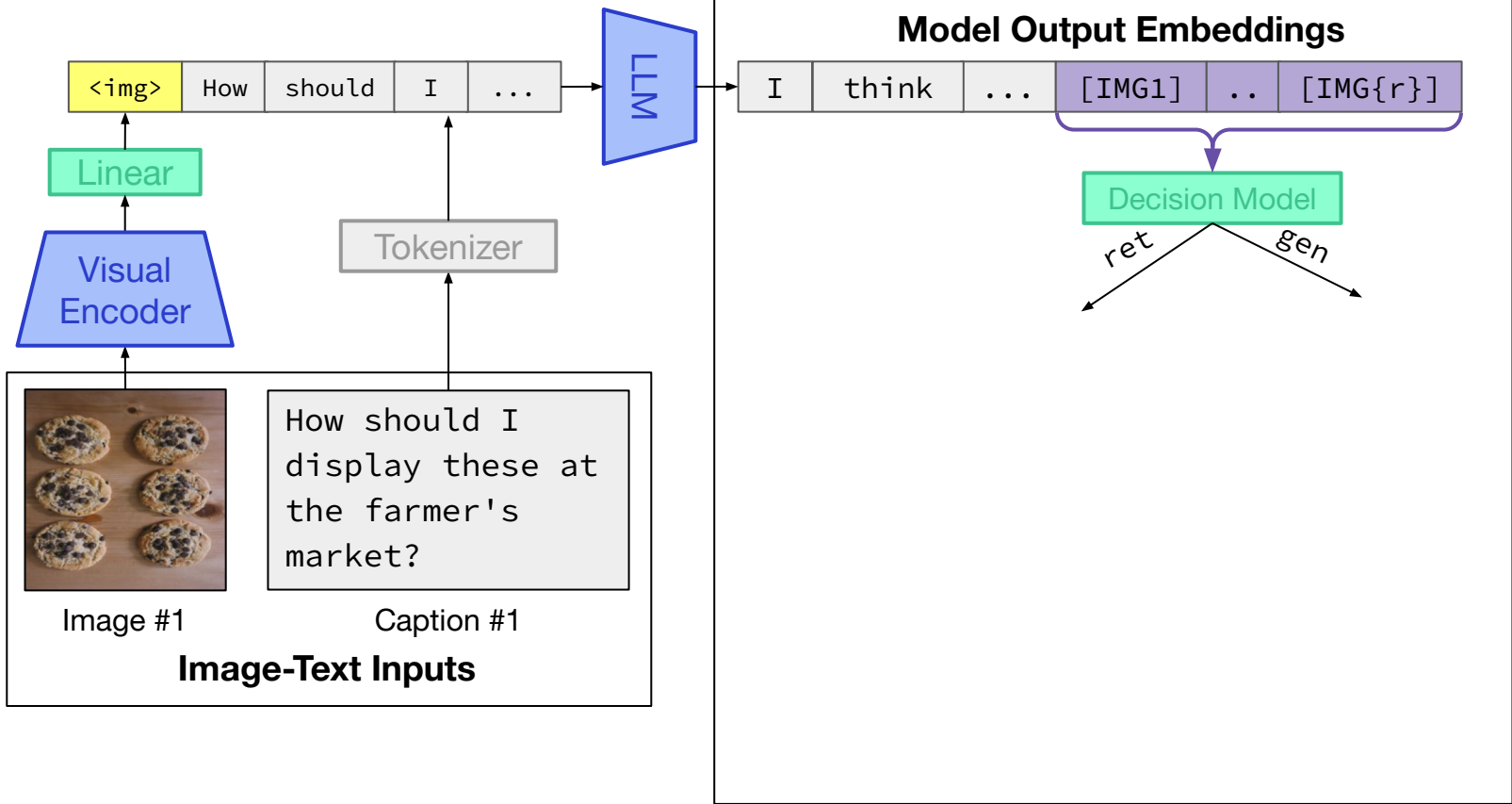
Caption #1

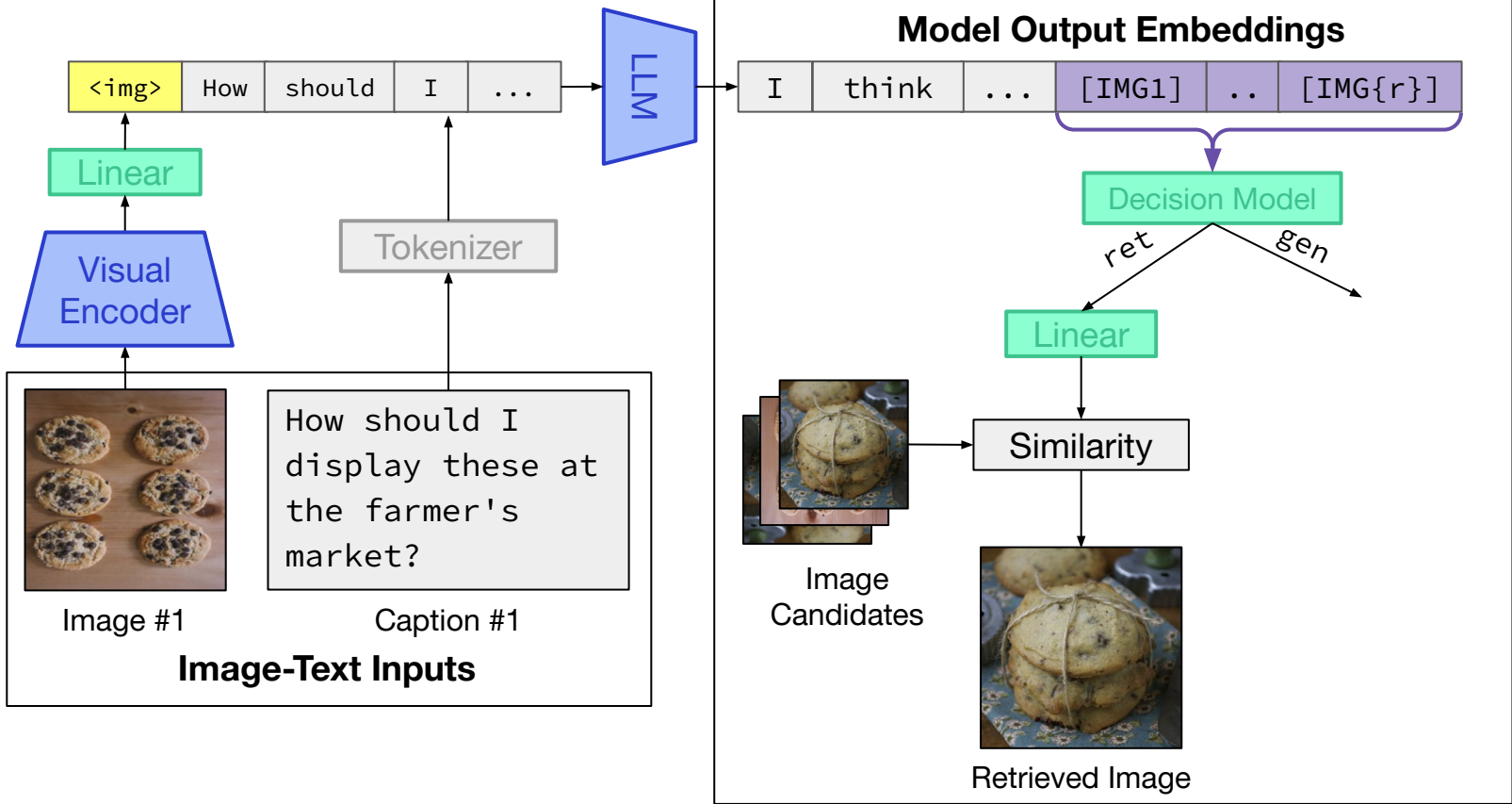
Image-Text Inputs

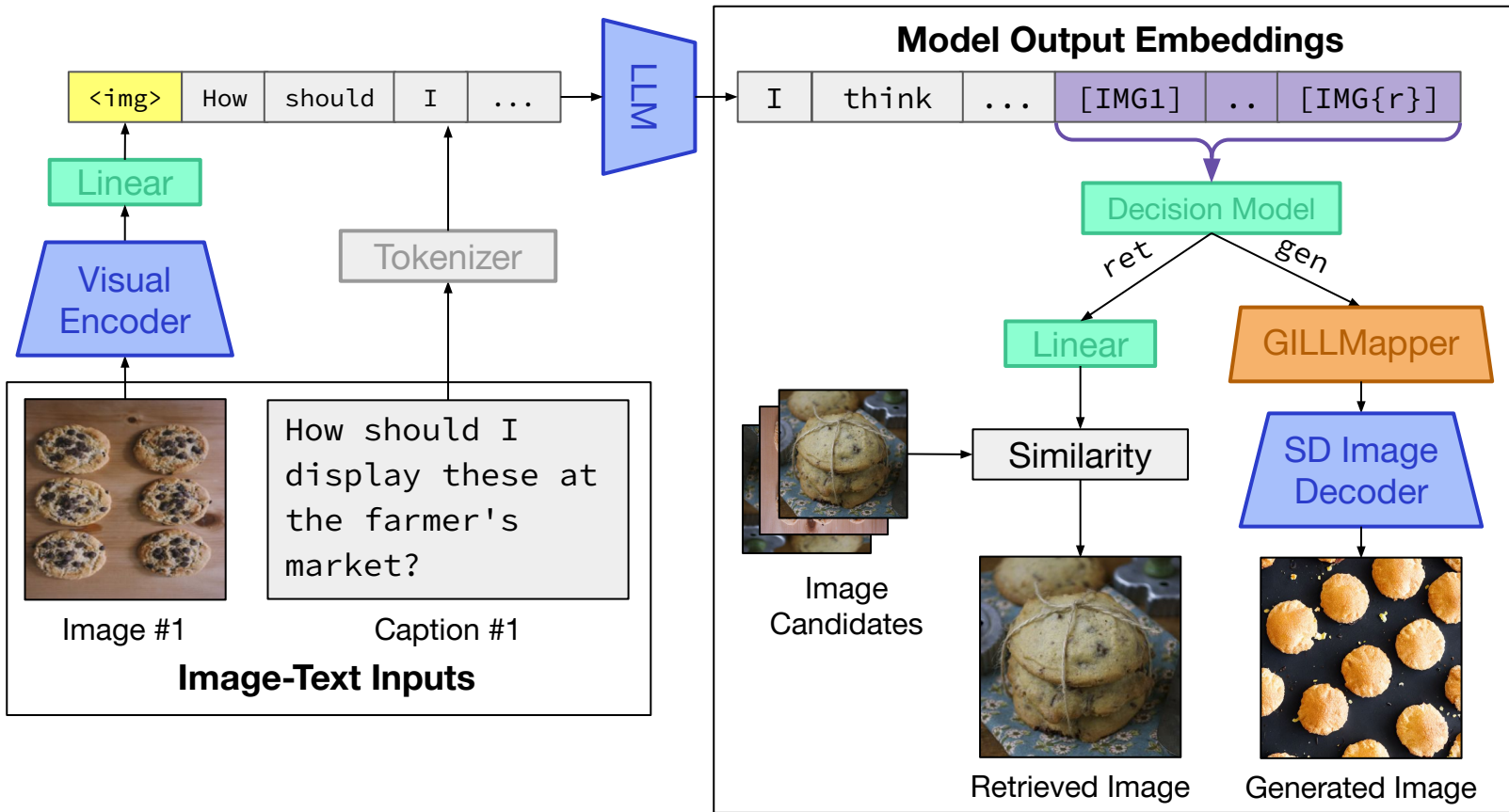













Final Model Outputs:

I think they look best when they are on a tray with a little bit of space between them. 

Evaluation: Contextual Image Generation

- Given a Visual Story, generate a relevant image



Image and Text Inputs

Evaluation: Contextual Image Generation

- 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

Once while I was on vacation in this nice brick hotel		I woke up and took my dog Trixie for a walk.		Trixie ran around and enjoyed the fresh air.		We had lots of fun playing fetch together		After a while she got tired and had to take a rest.			
	1		2		3		4	5			

Image and Text Inputs

Stable Diffusion Ours Groundtruth

Evaluation: Contextual Image Generation

Model	CLIP Similarity (\uparrow)			LPIPS (\downarrow)		
	1 caption	5 captions	5 caps, 4 images	1 caption	5 captions	5 caps, 4 images
GLIDE [34]	0.582	0.591	-	0.753	0.745	-
Stable Diffusion [43]	0.592 ± 0.0007	0.598 ± 0.0006	-	0.703 ± 0.0003	0.704 ± 0.0004	-
GILL	0.581 ± 0.0005	0.612 ± 0.0011	0.641 ± 0.0011	0.702 ± 0.0004	0.696 ± 0.0008	0.693 ± 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)

Evaluation: Contextual Image Generation

- Given a Visual Dialogue, generate a relevant image

Q: is the man alone? A: yes, the man is alone 1	Q: is it sunny outside? A: no, it is not sunny outside 2	Q: what color is the snowboard? A: the snowboard is grey in color 3	Q: is the man wearing a cap? A: the man is wearing a black cap 4	...	Q: what color are the glasses? A: the glasses are white in color 8	Q: can you see the sky? A: no it's totally dark 9	Q: does it look like he's having fun? A: he seems to be enjoying 10
---	--	---	--	-----	--	---	---

VisDial Inputs

Q: what color are the dogs? A: 1 of the dog is white and the other dog is light brown 1	Q: can you tell what breed they are? A: i can't really tell what breed they are, perhaps german shepherd 2	Q: are they both wearing a hat? A: only 1 is wearing a hat 3	...	Q: are they standing in grass? A: no, they are standing on dirt 8	Q: are they looking at each other? A: no, they are facing away from each other 9	Q: do they seem like they like each other? A: can't tell 10
---	--	--	-----	---	--	---

VisDial Inputs

Evaluation: Contextual Image Generation

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

Q: is the man alone? A: yes, the man is alone 1	Q: is it sunny outside? A: no, it is not sunny outside 2	Q: what color is the snowboard? A: the snowboard is grey in color 3	Q: is the man wearing a cap? A: the man is wearing a black cap 4	...	Q: what color are the glasses? A: the glasses are white in color 8	Q: can you see the sky? A: no it's totally dark 9	Q: does it look like he's having fun? A: he seems to be enjoying 10
---	--	---	--	-----	--	---	---

VisDial Inputs



Stable Diffusion



Ours



Groundtruth

Q: what color are the dogs? A: 1 of the dog is white and the other dog is light brown 1	Q: can you tell what breed they are? A: i can't really tell what breed they are, perhaps german shepherd 2	Q: are they both wearing a hat? A: only 1 is wearing a hat 3	...	Q: are they standing in grass? A: no, they are standing on dirt 8	Q: are they looking at each other? A: no, they are facing away from each other 9	Q: do they seem like they like each other? A: can't tell 10
---	--	--	-----	---	--	---

VisDial Inputs



Stable Diffusion



Ours



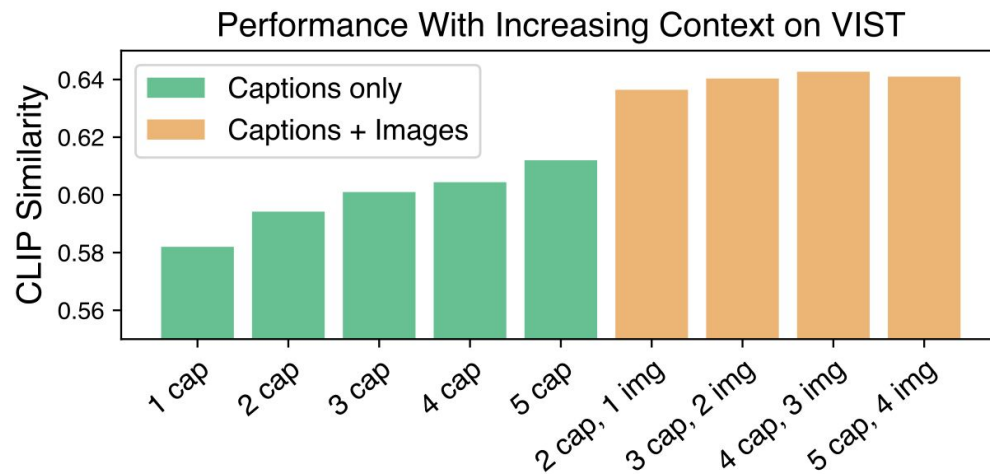
Groundtruth

Evaluation: Contextual Image Generation

Model	CLIP Similarity (\uparrow)			LPIPS (\downarrow)		
	1 round	5 rounds	10 rounds	1 round	5 rounds	10 rounds
GLIDE [34]	0.562	0.595	0.587	0.800	0.794	0.799
Stable Diffusion [43]	0.552 \pm 0.0015	0.629 \pm 0.0015	0.622 \pm 0.0012	0.742 \pm 0.0010	0.722 \pm 0.0012	0.723 \pm 0.0008
GILL	0.528 \pm 0.0014	0.621 \pm 0.0009	0.645 \pm 0.0010	0.742 \pm 0.0022	0.718 \pm 0.0028	0.714 \pm 0.0006

The Effect of Context

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



GILLMapper: A Stronger LLM-to-Generator Mapping

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

Model	CC3M	VIST
	FID (↓)	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

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



Ours

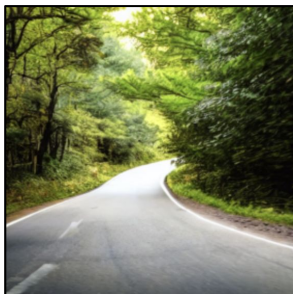


Stable Diffusion

“Snow mountain and tree reflection in the lake”



Ours



Stable Diffusion

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



Ours



Stable Diffusion

“a group of penguins in a snowstorm”




Ours


Other Abilities: Image Refinement

I want to paint a landscape with mountains and a river. Can you help me find some inspiration?

I'm not sure



Can you show me some with horses?



What about something more abstract?



Future Work

- **Train on more diverse data**
 - CC3M is small by modern standards – we would get a lot more from training on LAION
 - 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 8 🟢 Running on A10G 📄 Logs 📱 App 📁 Files 👤 Community 2 ⚙️ 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: [Generating Images with Multimodal Language Models](#)

Project Website: [GILL Website](#)

Code and Models: [GitHub](#)

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

How can I publicize these?

I would suggest you start with a local newspaper.



(Generated)

 Upload Image

Message

Type a message

Submit

Undo

Reset All

Frequency multiplier for returning images (higher means more frequent) 1.3

Max # of words 32

Temperature (0 for deterministic, higher for more randomness) 0

Example Conversations



How can I make this recipe healthier?



What are some ingredients I can add to it?

You can add any vegetables you like.



How can I make this recipe healthier?



What are some ingredients I can add to it?

You can add any vegetables you like.



How can I make this recipe healthier?



What are some ingredients I can add to it?

Thanks!

jykoh@cmu.edu

jykoh.com/gill