

# Multimodal AI

## Lecture 5.1 – Large Multimodal Models

**Paul Liang**

Assistant Professor

MIT Media Lab & MIT EECS



<https://pliang279.github.io>

[ppliang@mit.edu](mailto:ppliang@mit.edu)

 [@pliang279](https://twitter.com/pliang279)



# Assignments for This Coming Week

HW2 due tomorrow Wed (3/4).

HW2 presentation this Thurs + in-depth tutorial on implementing multimodal LLMs.

HW3 (multimodal LLMs) out later this week.

# Today's lecture

- 1 Multimodal transformers & foundation models
- 2 Adapting LLMs into multimodal LLMs
- 3 From text to multimodal generation
- 4 Latest directions

# Recap: Large Language Models

**Large Language Models**



Question answering

Open-ended dialog

Translate to different languages

Retrieve news

Solving math problems

Writing code

and more...

# From Large Language Models to Multimodal Models



**Classification:** What is the tone of the man in the grey shirt?

**Open-ended:** Describe the relationships between these 2 people.

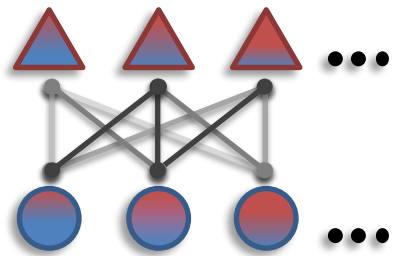
**Explanation:** Explain why, citing visual and verbal evidence.

**Generation:** Animate a story inspired by this short video clip.

**Counterfactual:** What if these people were from a different society or culture?

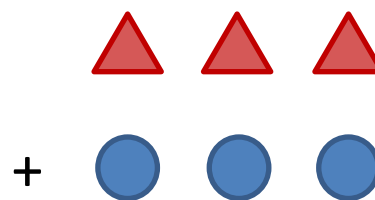
# From Large Language Models to Multimodal Models

*It's just a privilege to watch your mind at work.*



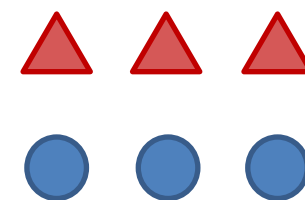
Multimodal representation

*This person is being sarcastic.  
They seem to be close friends.*



*(quote previous episodes)  
(highlight multimodal information)*

*Here's a story of them in  
a different culture...*

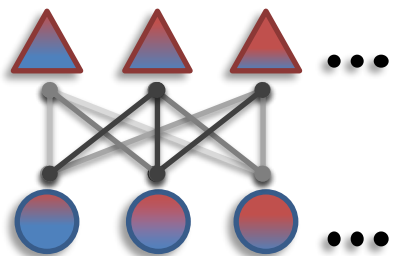


*(generate future episodes)*

# Lecture outline

## Part 1: Multimodal foundation model representations of text, video, audio

*It's just a privilege to  
watch your mind at work.*



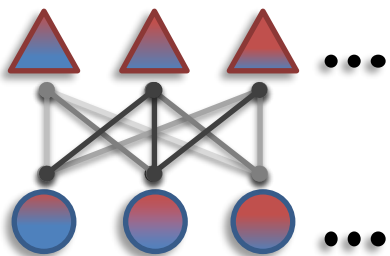
Multimodal  
representation



# Lecture outline

## Part 2: Adapting large language models for multimodal text generation

*It's just a privilege to  
watch your mind at work.*



Multimodal  
representation



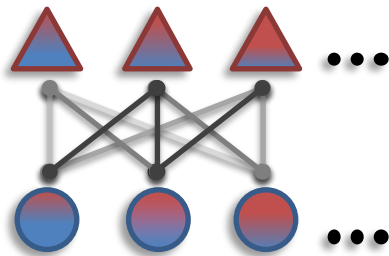
*This person is being sarcastic.  
They seem to be close friends.*



# Lecture outline

## Part 3: Enabling text and image generation

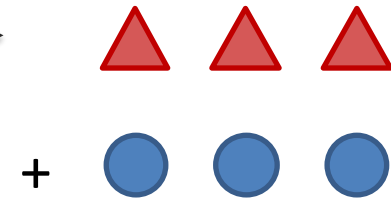
*It's just a privilege to  
watch your mind at work.*



Multimodal  
representation

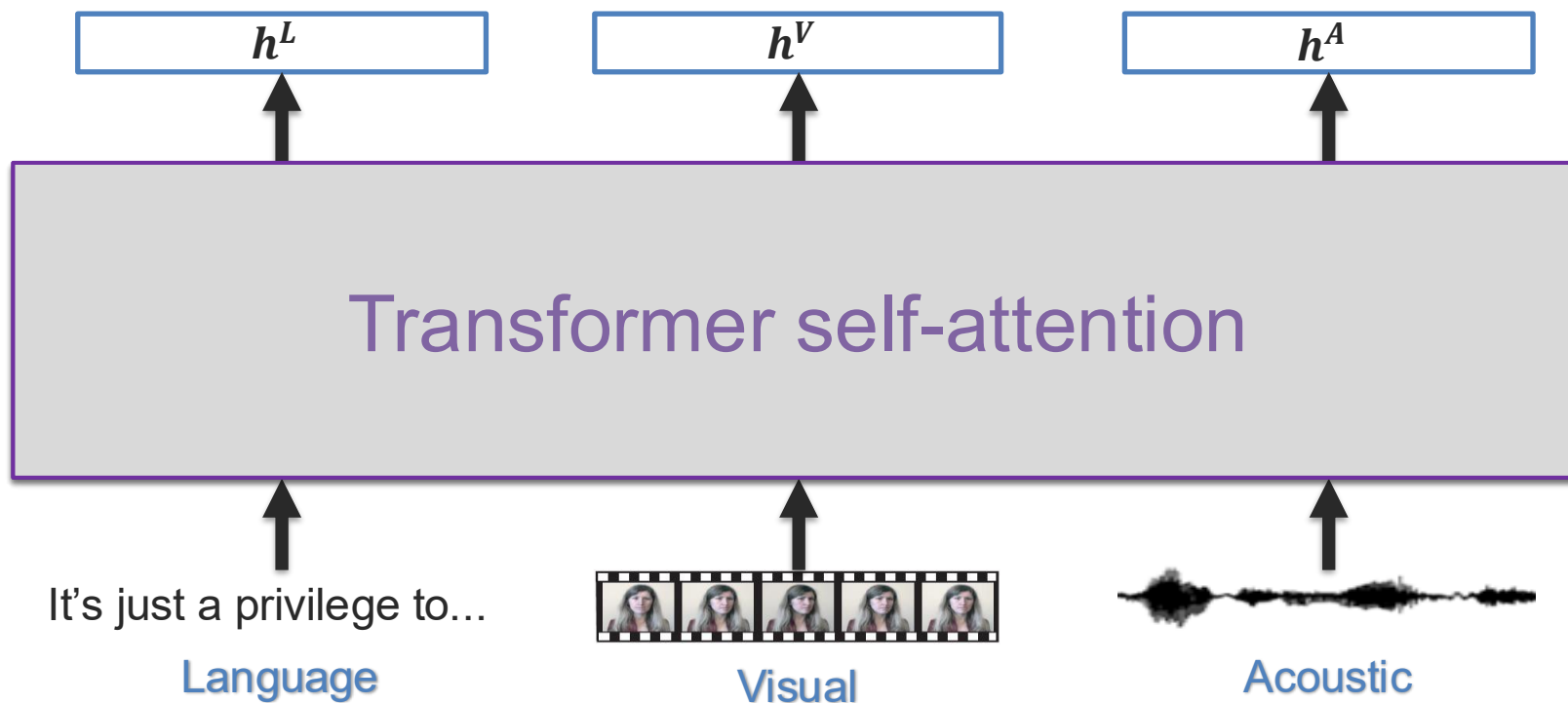


*This person is being sarcastic.  
They seem to be close friends.*

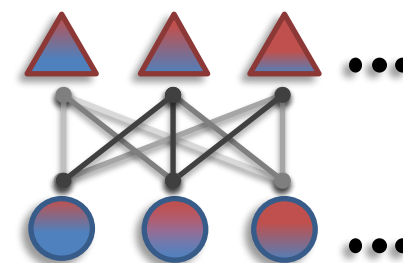


*(quote previous episodes)  
(highlight multimodal information)*

# Large Multimodal Models



**Contextualized representation**



Implicit alignment  
+ representation

# Multimodal Transformers

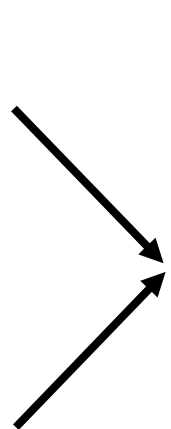
*It's just a privilege to watch your mind at work.*



privilege  
mind

vocal  
eye-roll  
emphasis

vocal  
emphasis  
eye-roll



Vision-to-language  
attention

privilege	0.7			0.3
mind	1.0			
	vocal		eye-roll	
	emphasis			

emphasis

privilege  
mind



New **language** representation  
contextualized with vision

(row) normalize to 0-1

$$h = \text{softmax} \left( \frac{X_1 W_q W_k^T X_2^T}{\sqrt{d}} \right) X_2 W_v$$

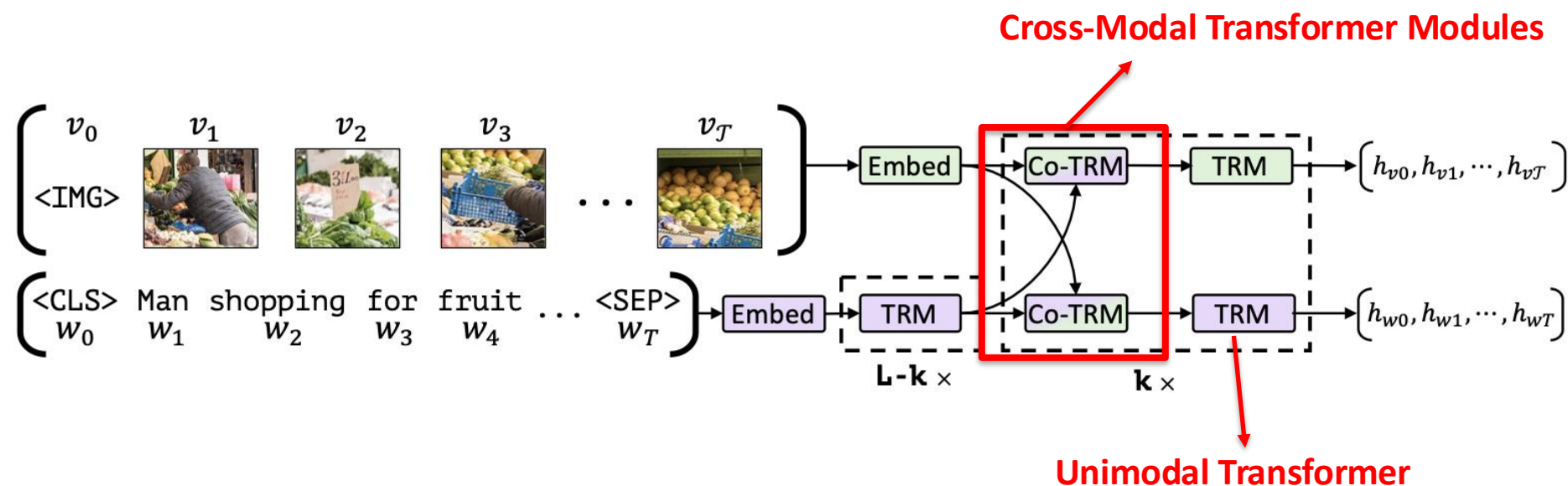
$3 \times d$      $d \times 4$      $4 \times d$   
 $X_1$      $W_q$      $W_k^T$      $X_2^T$      $X_2$      $W_v$

normalize wrt dimension  $d$

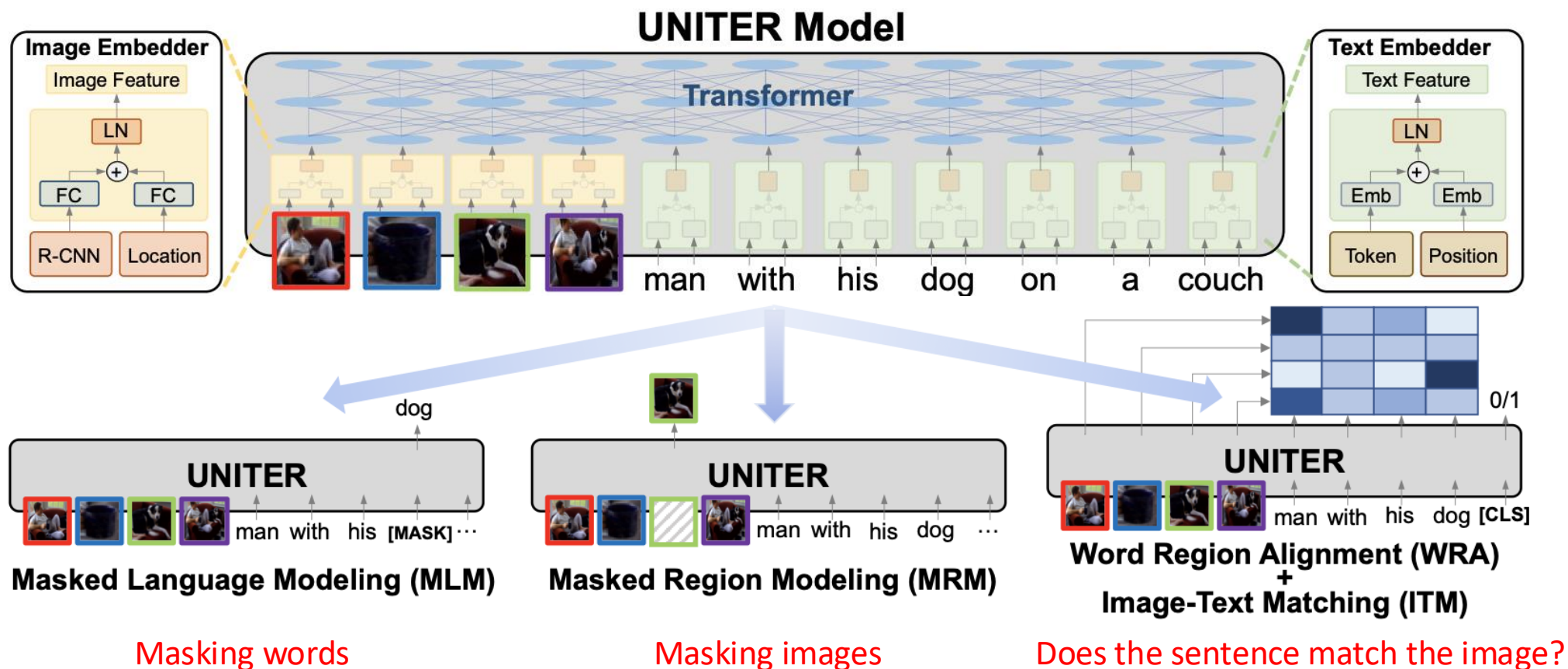
$3 \times 4$   
(weighted) outer product

Sarcasm

# Multimodal Cross-attention Transformers

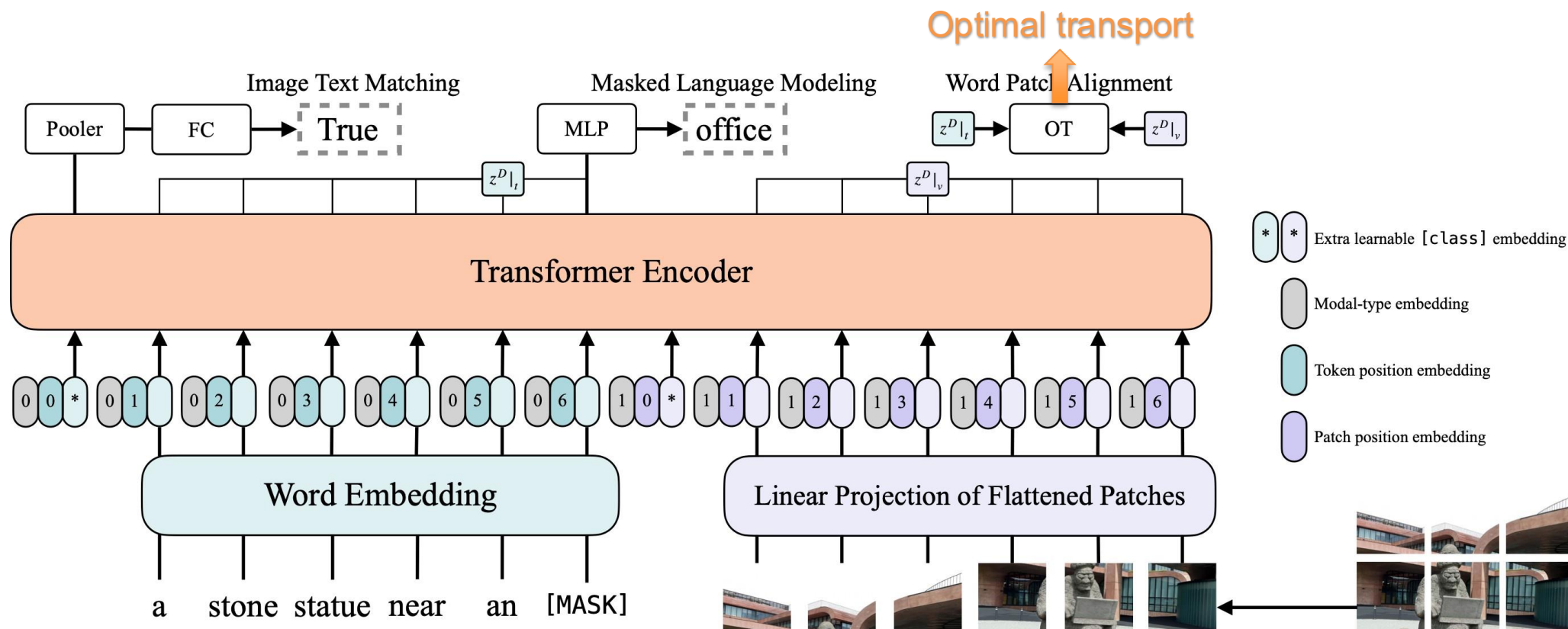


# Multimodal Cross-attention Transformers



# Visual-and-Language Transformer (ViLT)

( $\approx$  BERT + ViT)



# Visual-and-Language Transformer (ViLT)

Example of alignment between modalities:



a display of **flowers** growing out and over the retaining **wall** in front of **cottages** on a **cloudy** day.



flowers



wall



cottages



cloudy



a room with a **rug**, a **chair**, a **painting**, and a **plant**.



rug



chair

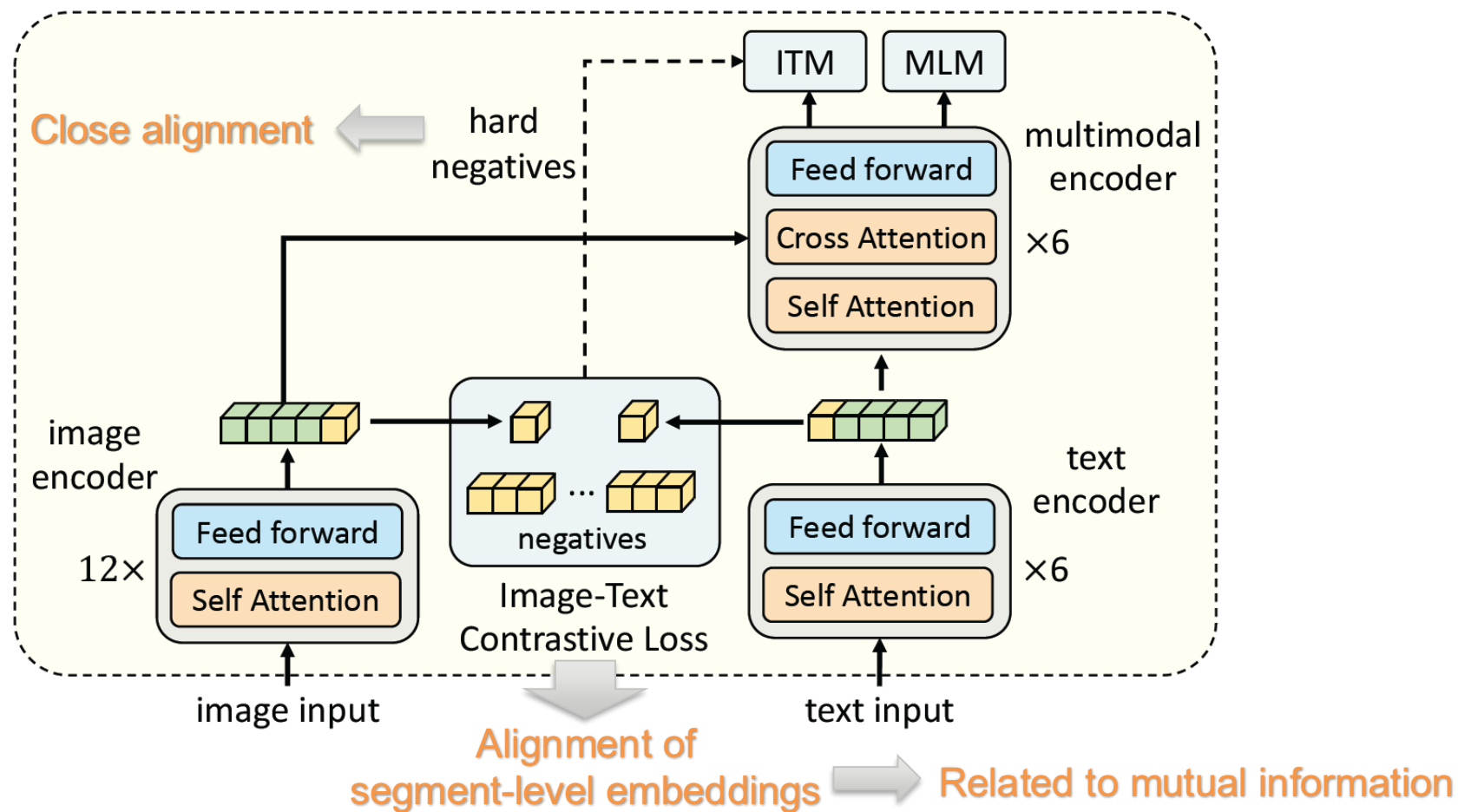


painting



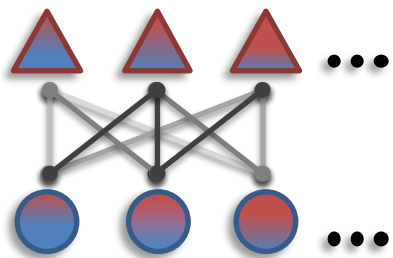
plant

# ALBEF: Align Before Fusion (≈ BERT + ViT + CLIP-ish)



# Adapting Large Language Models to Multimodal

*It's just a privilege to  
watch your mind at work.*



Multimodal  
representation



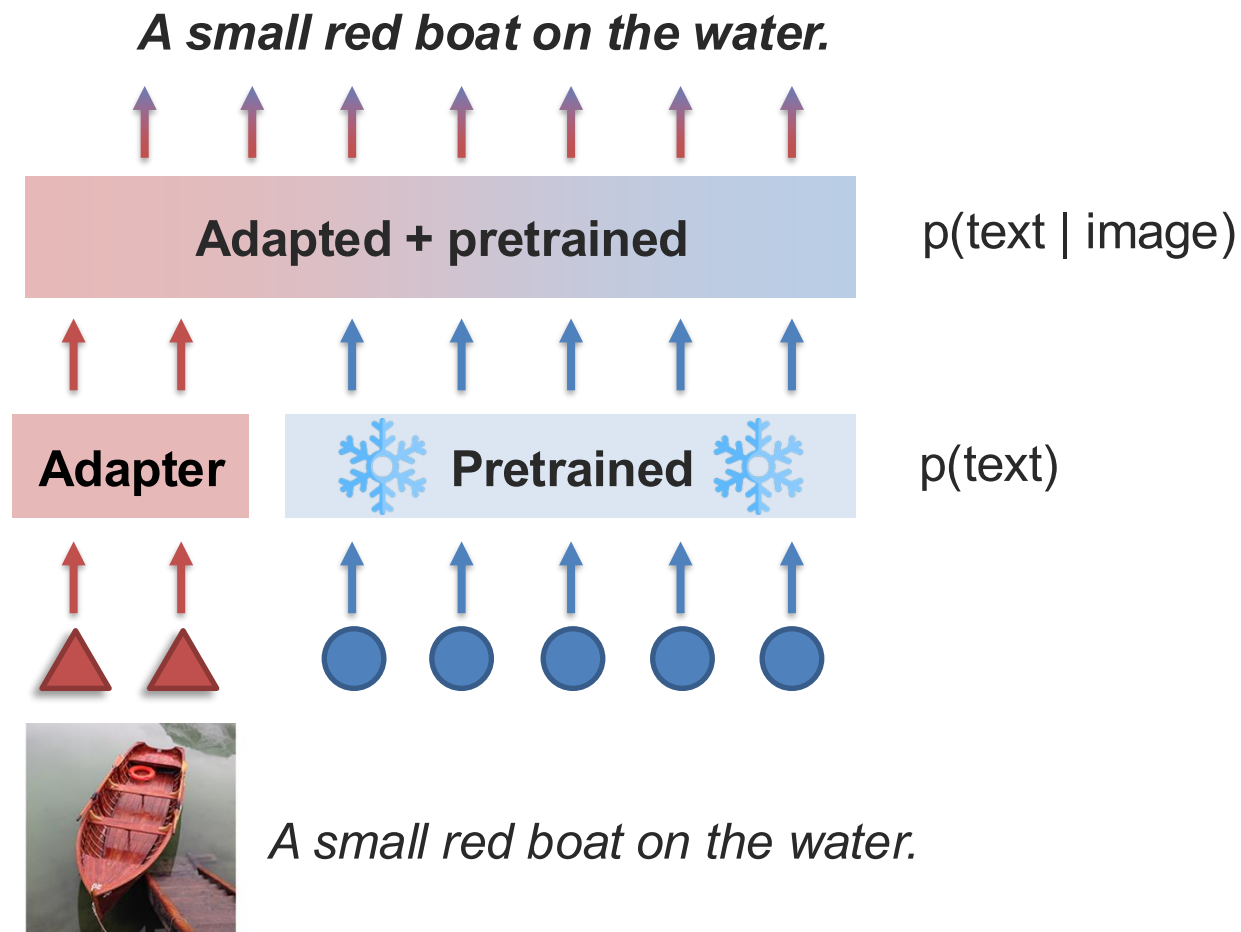
*This person is being sarcastic.  
They seem to be close friends.*



# Adapting Large Language Models

## Conditioning via prefix tuning

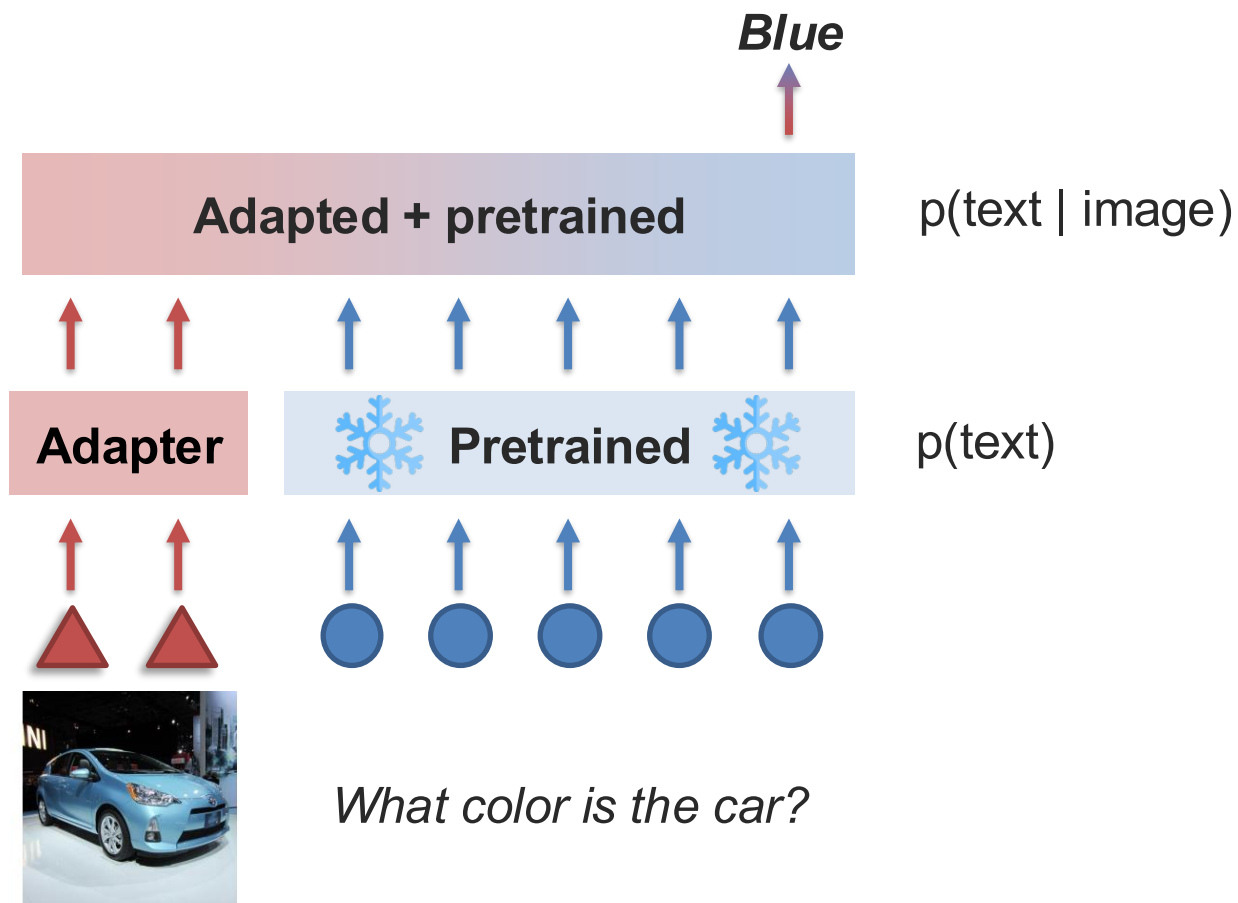
Modeling  $p(\text{text} \mid \text{image})$ :



# Adapting Large Language Models

## Conditioning via prefix tuning

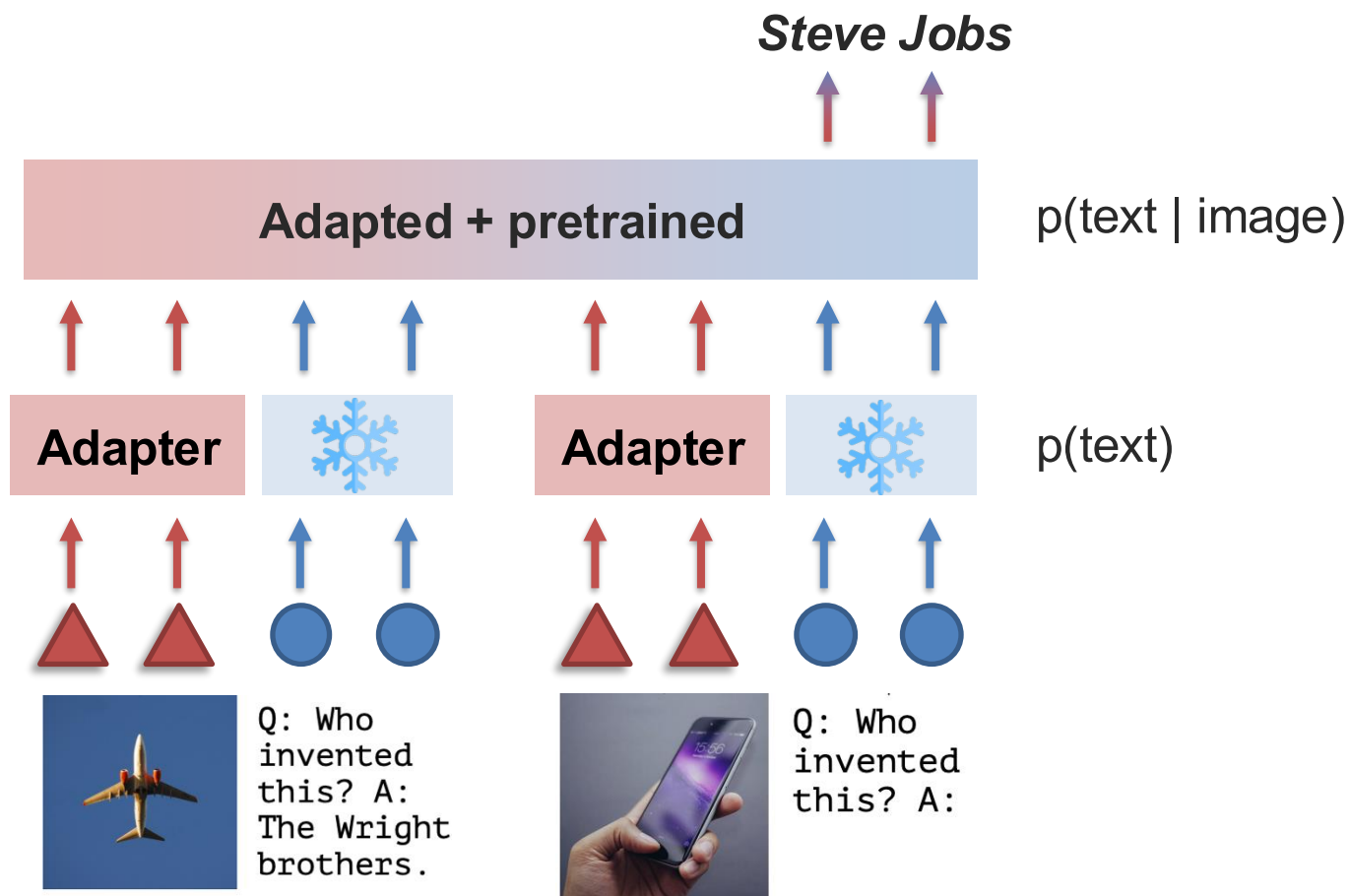
Modeling  $p(\text{text} \mid \text{image})$ :



# Adapting Large Language Models

## Conditioning via prefix tuning

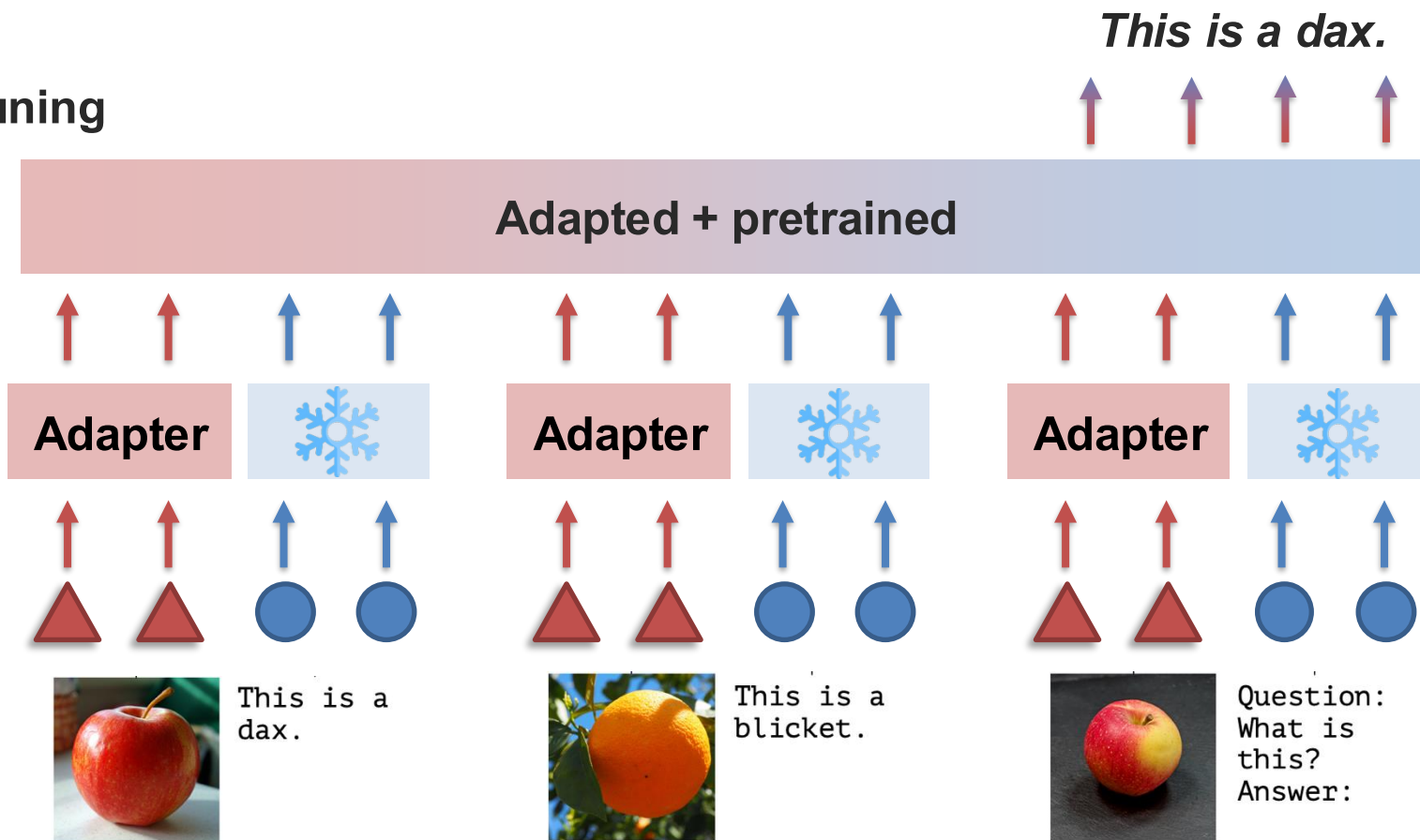
Modeling  $p(\text{text} \mid \text{image})$ :



# Adapting Large Language Models

## Conditioning via prefix tuning

Modeling  $p(\text{text} \mid \text{image})$ :



# Adapting Large Language Models

## Flamingo



PP This is a picture of two teddy bears on the moon.

What are they doing?

PP They are having a conversation.

What object are they using?

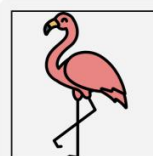
PP It looks like a computer.

Is this surprising?

PP Yes, it is surprising.

Why is this picture surprising to you?

PP I think it is surprising because teddy bears are not usually found on the moon.



What is the common thing about these three images?

PP They are all flamingos.

What is the difference between these three images?

PP The first one is a cartoon, the second one is a real flamingo, and the third one is a 3D model of a flamingo.



PP This is an apple with a sticker on it.

What does the sticker say?

PP The sticker says "iPod".

Where is the photo taken?

PP It looks like it's taken in a backyard.

Do you think it is printed or handwritten?

PP It looks like it's handwritten.

What color is the sticker?

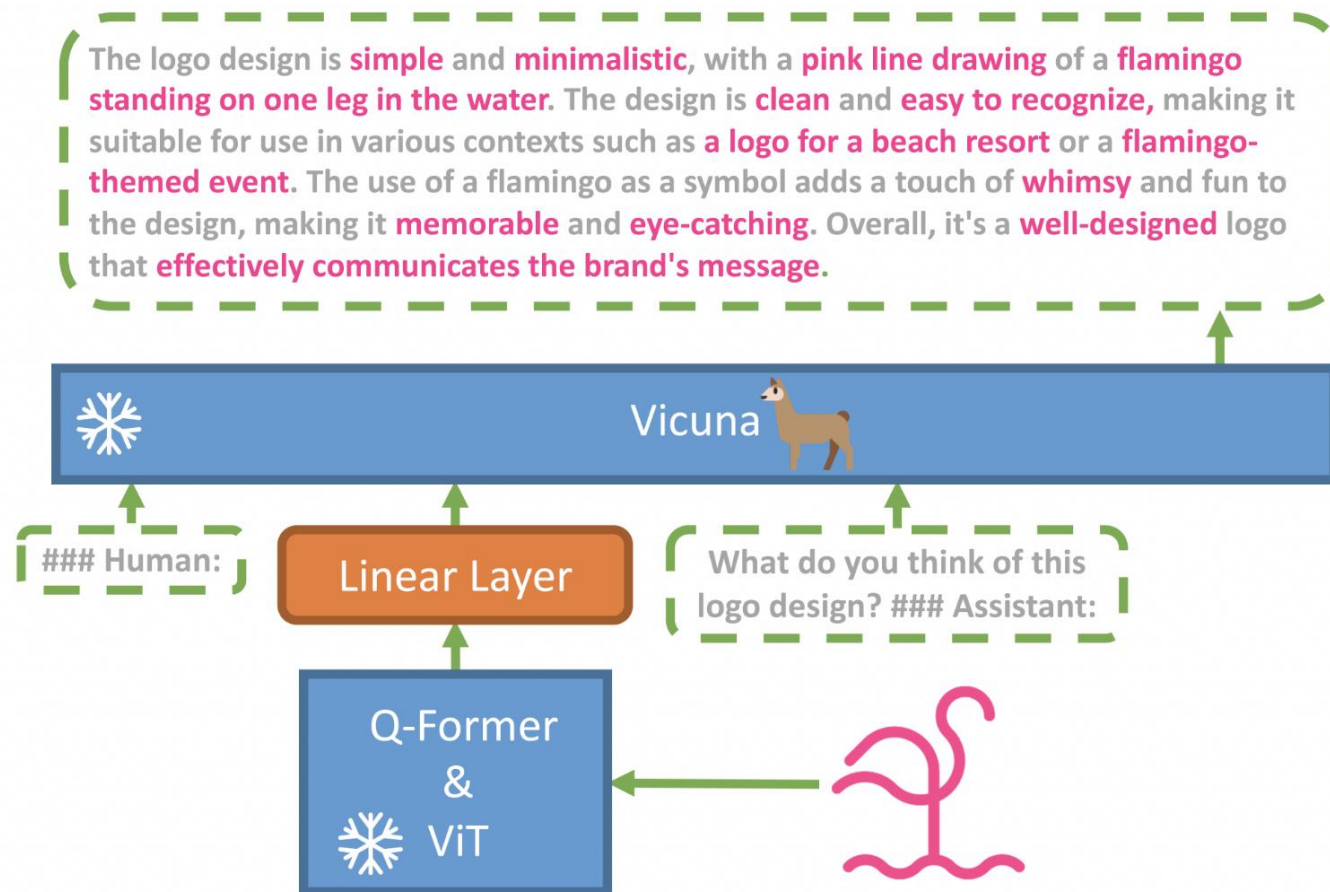
PP It's white.

# Scaling Large Multimodal Models

## Mini-GPT4

Stage 1: **Alignment** using paired image-text data.

Stage 2: **Instruction tuning** using image + text instructions and example completions.



The architecture of MiniGPT-4.

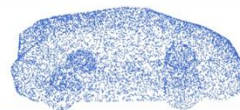
# Scaling Large Multimodal Models

## LLaMA-Adapter

LLaMA-Adapter:  
Bilingual Multi-modality  
Instruction Model



Example: 3D Point Cloud to Image (Bilingual)



Generate an image from the 3D point cloud.



Hello



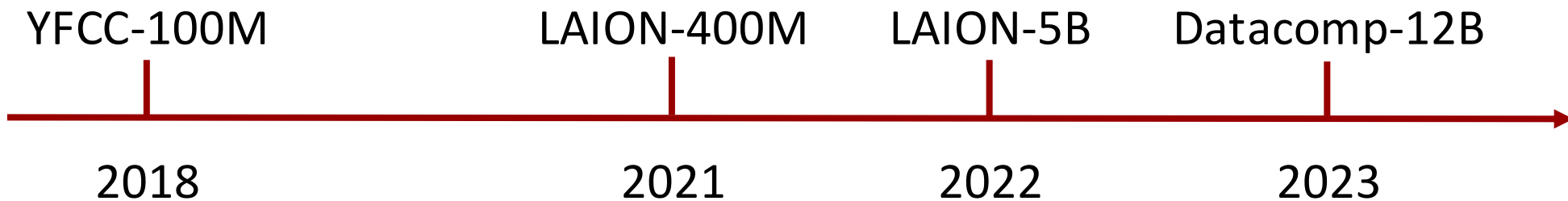
根据这个3D点云生成一张图片。

你好



# Multimodal Pre-training Datasets

- Largest dataset is DataComp. It has 12.8 billion image-text pairs.
- Recent efforts shifted more towards filtering for high quality multimodal data. Examples include DFN (2B), COYO (600M), and Obelics (141M)



Data type	dataset	#samples	sampling prob.
Image-Caption	DFN [Fang et al., 2023]	2B	27%
	COYO [Byeon et al., 2022]	600M	11.25%
	HQITP	400M	6.75%
Interleaved Text	Obelics [Laurençon et al., 2024a]	141M Docs	45%
	DCLM [Li et al., 2024b]	6.6T Toks	10%

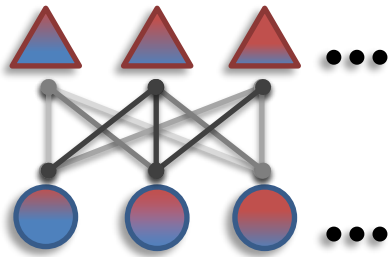
# Multimodal Instruction Tuning Datasets

- More scattered, smaller in nature
- General domain: Vision-Flan (187K), LLaVA-Instruct (150K), InstructBLIP (~1.6M), M3IT (2.4M)
- Clinical: CLIMB-QA (4.51M), BioMed-VITAL (210K), LLaVA-Med (60K)

Dataset	# Tasks	Multi-Lingual	# of Instances	Avg. # of Manual Instructions / Task	Open-Sourced
MiniGPT4	N / A	✗	5K	N / A	✓
LLaVA	3	✗	1.15M	N / A	✓
MultiModalGPT	3	✗	6K	5	✗
MultiInstruct	26	✗	~ 235K	5	✗
InstructBLIP	28	✗	~ 1.6M	9.7	✗
M <sup>3</sup> IT (Ours)	40	✓	2.4M	10	✓

# From Text to Multimodal Generation

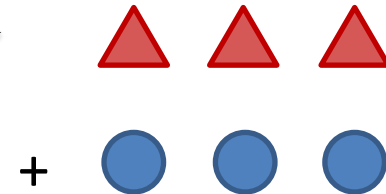
*It's just a privilege to  
watch your mind at work.*



Multimodal  
representation



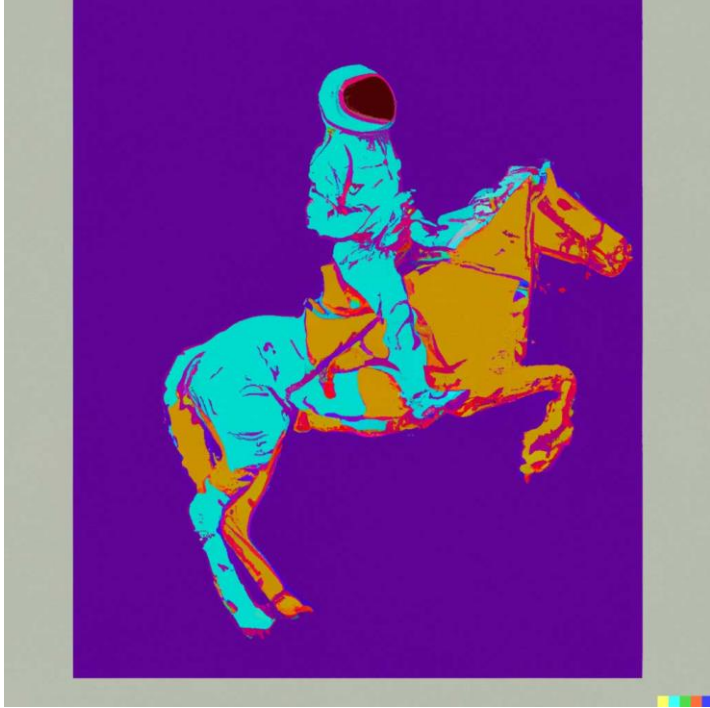
*This person is being sarcastic.  
They seem to be close friends.*



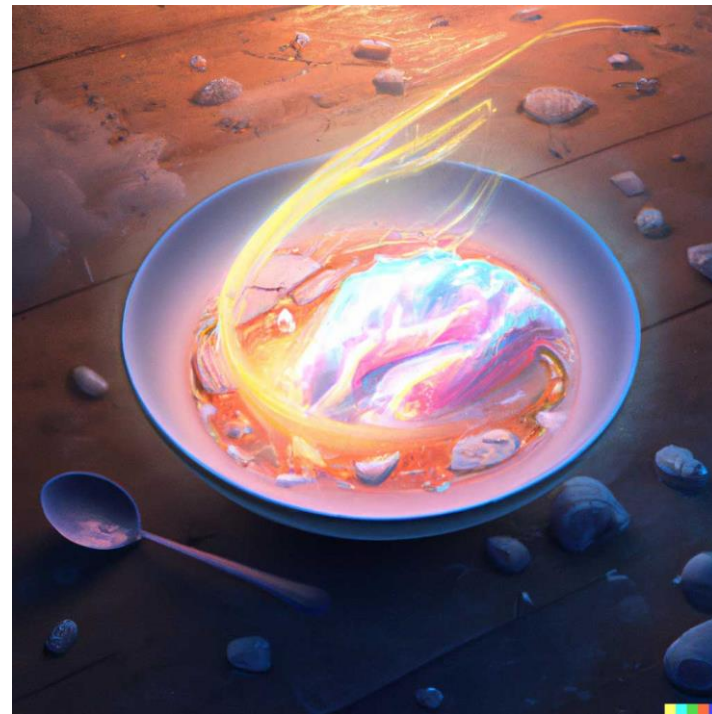
*(retrieve next episode)  
(highlight multimodal evidence)*

# From Text to Multimodal Generation

*An astronaut riding a horse in the style of Andy Warhol.*



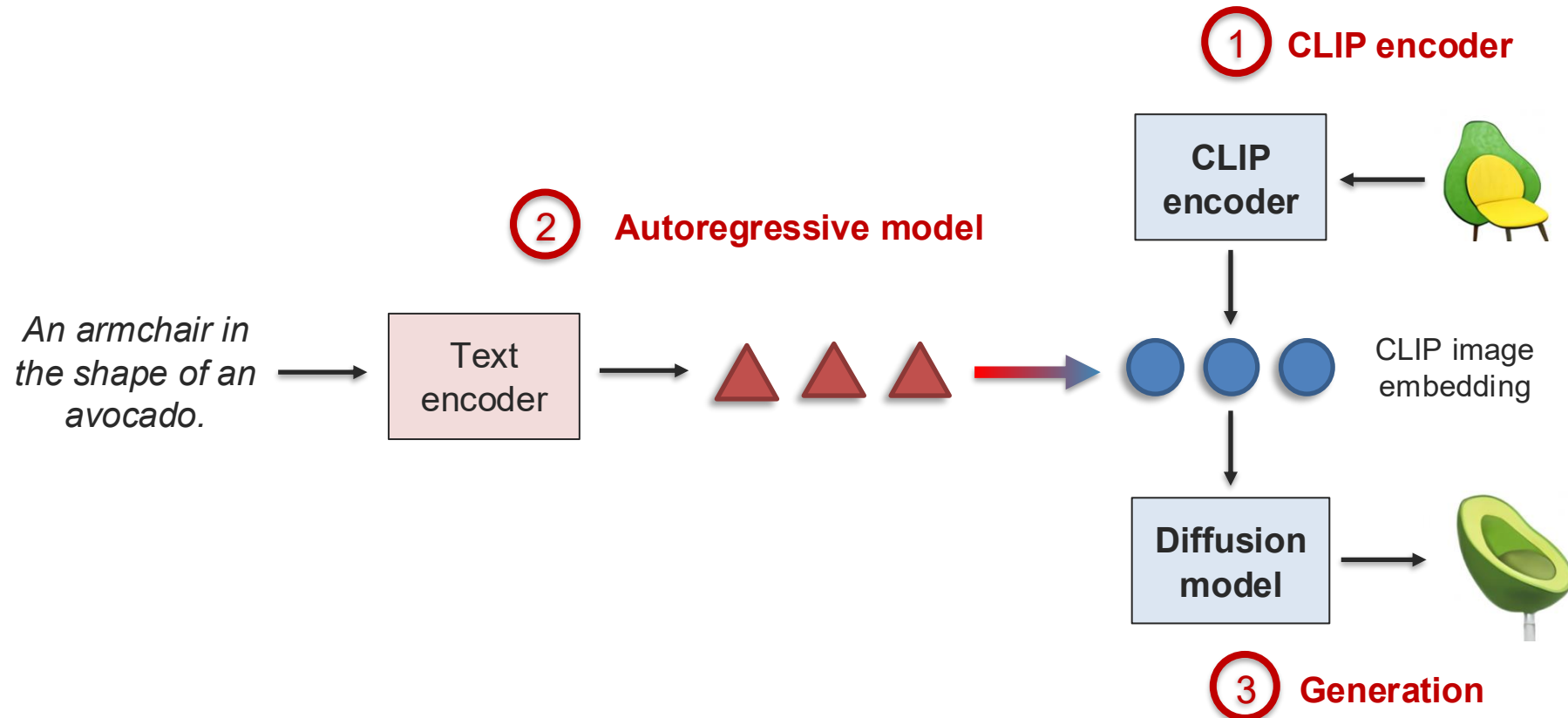
*A bowl of soup that is a portal to another dimension as digital art*



# From Text to Multimodal Generation

Directly training diffusion models with conditional information

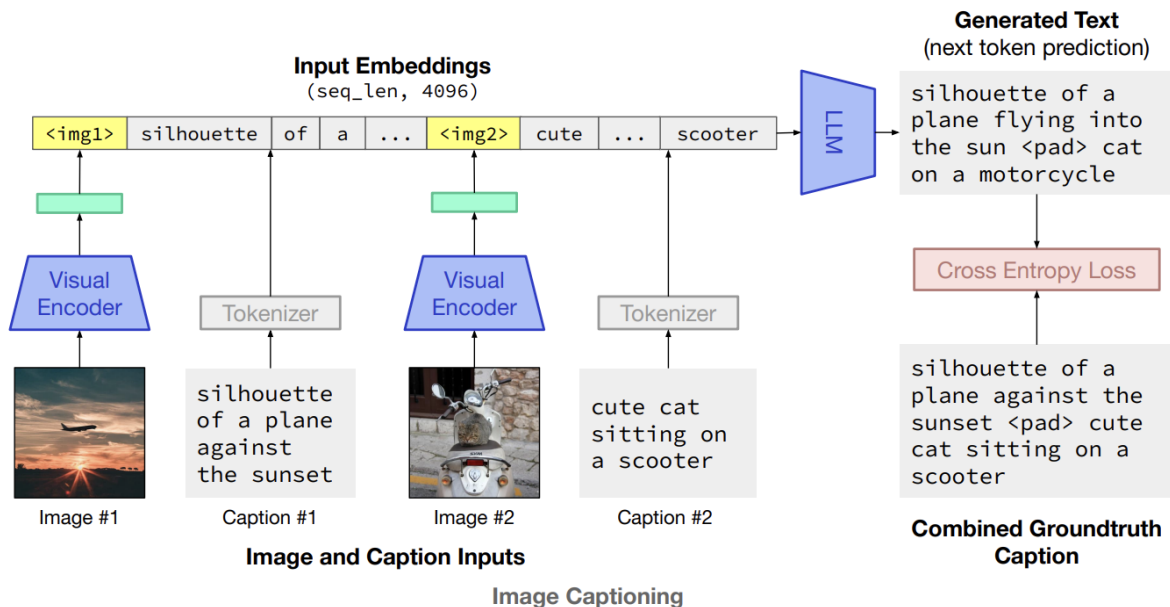
Conditional latent variables are pretrained CLIP embeddings, then diffusion model to generate image.



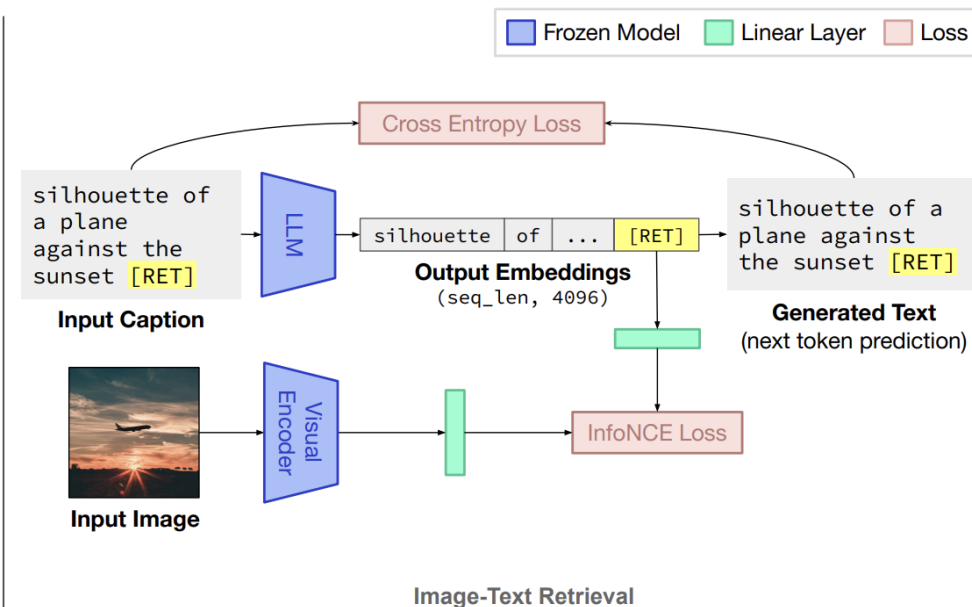
# Grounding LMs for Image Retrieval

LIMBeR + CLIP. Trainable in 1 day on 1 GPU


## Interleaved images and text





## CLIP, with a frozen LLM





# Grounding LMs for Image Retrieval


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. 


That's right, what do beavers look like? 


They have a tail. 


Show me an illustration of a beaver 


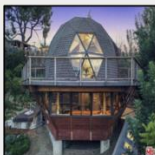
A beaver. 





Yes, what about a pencil drawing of one? 





What about a photo of one? 






What is this? 


This is a house with a unique design. 


What is its shape? 


The house is shaped like a dome. 


Can you show me more houses like this? 


The house is designed by architect. 





That's nice, any others? 


Yes, there are many houses like this in the world. 

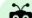



Provide instructions on how to cook. 


How do I make macarons? 


Macarons are a French dessert that are made by combining egg whites and sugar. The egg whites are beaten until stiff and then folded into the sugar. 


What does the stiff egg white look like? 

The egg white should be stiff and glossy. 

Show me a picture of what that looks like 



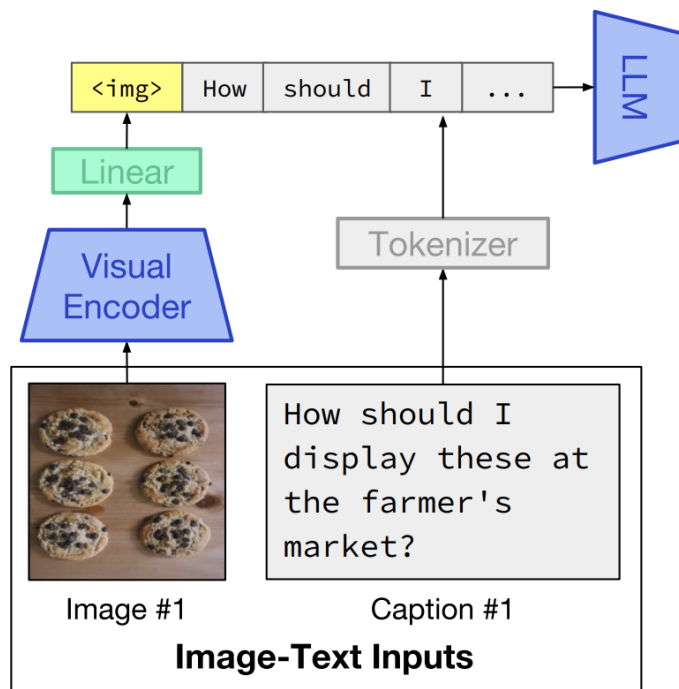
After this, what do I do? 

The macaron is then baked in a hot oven. 

# Grounding LMs for Multimodal Generation

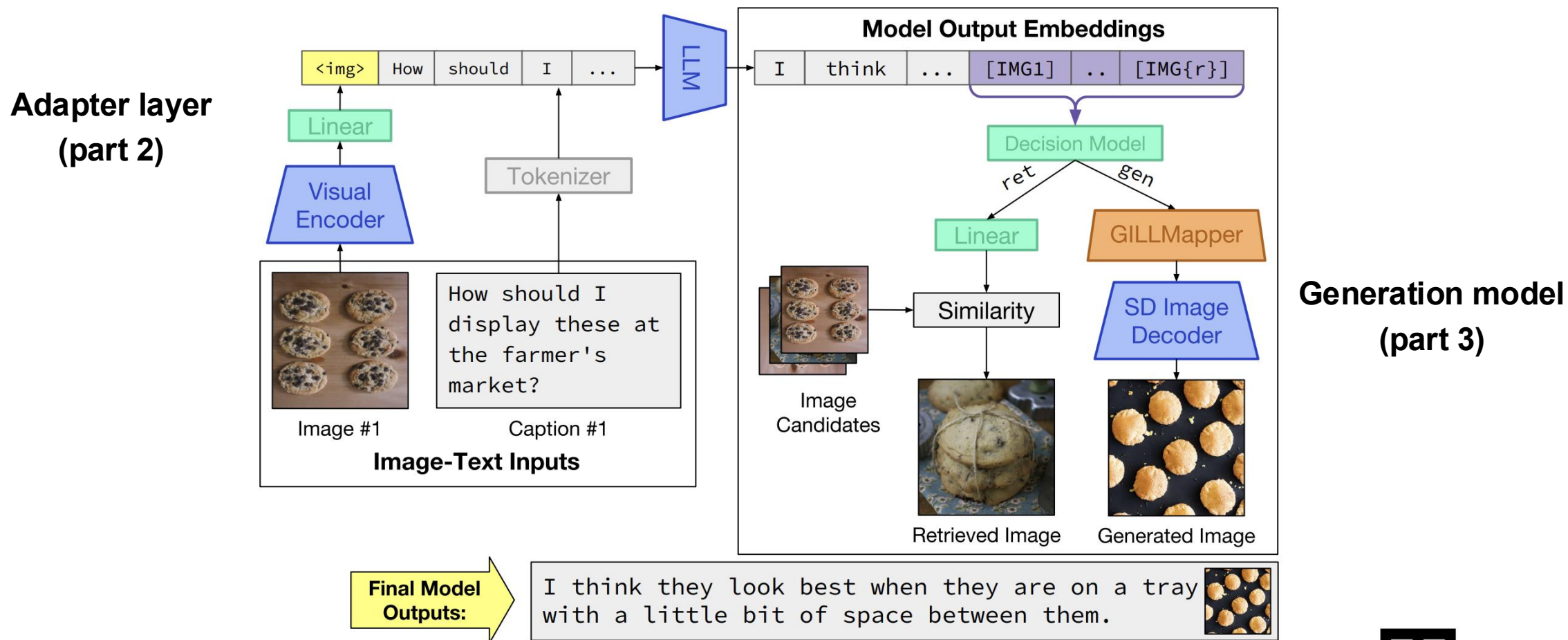
## Large multimodal models with image generation

Adapter layer  
(part 2)



# Grounding LMs for Multimodal Generation


## Large multimodal models with image generation



# Grounding LMs for Multimodal Generation


## Visual Storytelling

Our model can condition on interleaved image-and-text inputs to generate more relevant images compared to non-LLM based text-to-image generation 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.

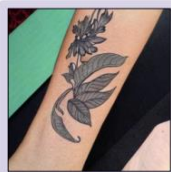
I'm looking for some ideas for a new tattoo. What do you think looks good on a petite person?




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.




Hi, I'm looking for ideas for a customised cake design for my daughter's birthday.

I'm not sure what you mean by customised.




I want a cake with a unicorn theme, with pastel colors and a whimsical feel.

I think



Can you also show me some decorations that would work well for the party?

I can show you some ideas



## Multimodal Dialogue

Our model can generate multimodal dialogue, weaving together text, retrieved images, and generated images.

User prompts
  Retrieved
  Generated

# Summary

- 1 Multimodal transformers & foundation models
- 2 Adapting LLMs into multimodal LLMs
- 3 From text to multimodal generation
- 4 Latest directions: natively multimodal, multimodal MoE, real-world modalities

# Assignments for This Coming Week

HW2 due tomorrow Wed (3/4).

HW2 presentation this Thurs + in-depth tutorial on implementing multimodal LLMs.

HW3 (multimodal LLMs) out later this week.