



10-423/10-623 Generative AI

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Vision-Language Models (VLMs)

Matt Gormley & Pat Virtue

Lecture 14

March 10, 2025

Reminders

- **Homework 3: Applying and Adapting LLMs**
 - Out: Sun, Feb 23
 - Due: Thu, Mar 13 at 11:59pm
- **Quiz 4**
 - in-class, Mon, Mar 17
 - lectures 12 – 15
- **Homework 4: Multimodal Foundation Models**
 - Out: Thu, Mar 13
 - Due: Mon, Mar 24 at 11:59pm

(Slides with blue titles from Henry Chai)

VISION LANGUAGE MODELS (VLMS)

Multimodal Models

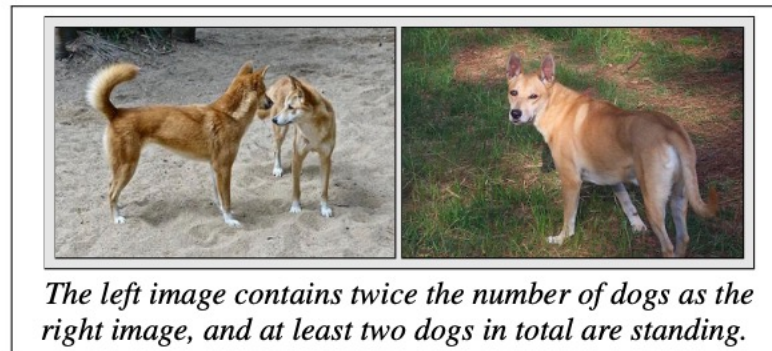
- Previously: Text-to-image models – adapt generative models for vision in order to guide their output toward some desired target using natural language
 - Output is still an image
- Today: visual language models (VLMs) – adapt generative models for text in order to allow them to interact with images (as well as text) as input
 - Output is (typically) still text

VLM: Tasks

- Common benchmarks for VLMs include
 - **Visual reasoning:** given an image (or a pair of images) determine if some natural language statement about the image(s) is true or false
 - **Visual grounding:** locate an object in some image given a natural language description
 - **Visual question answering:** given an image (or images), respond to arbitrary, potentially open-ended questions about the content.
 - **Caption generation:** create natural language descriptions of content of some image

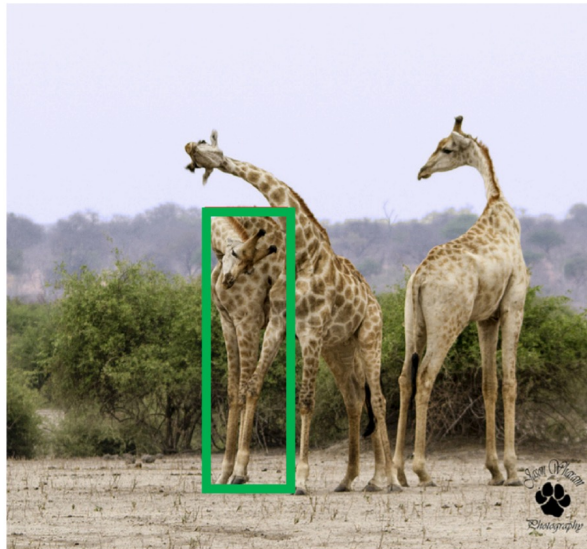
VLM: Tasks

- Common benchmarks for VLMs include
 - **Visual reasoning:** given an image (or a pair of images) determine if some natural language statement about the image(s) is true or false



VLM: Tasks

- Common benchmarks for VLMs include
 - **Visual reasoning:** given an image (or a pair of images) determine if some natural language statement about the image(s) is true or false
 - **Visual grounding:** locate an object in some image given a natural language description



RefCOCO:

1. giraffe on left
2. first giraffe on left

RefCOCO+:

1. giraffe with lowered head
2. giraffe head down

RefCOCOG:

1. an adult giraffe scratching its back with its horn
2. giraffe hugging another giraffe

VLM: Tasks

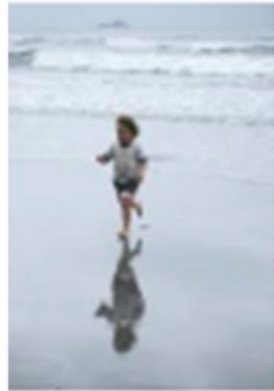
- Common benchmarks for VLMs include



- **Visual question answering:** given an image (or images), respond to arbitrary, potentially open-ended questions about the content.

VLM: Tasks

- Common benchmarks for VLMs include



Ground Truth Caption: A little boy runs away from the approaching waves of the ocean.

Generated Caption: A young boy is running on the beach.



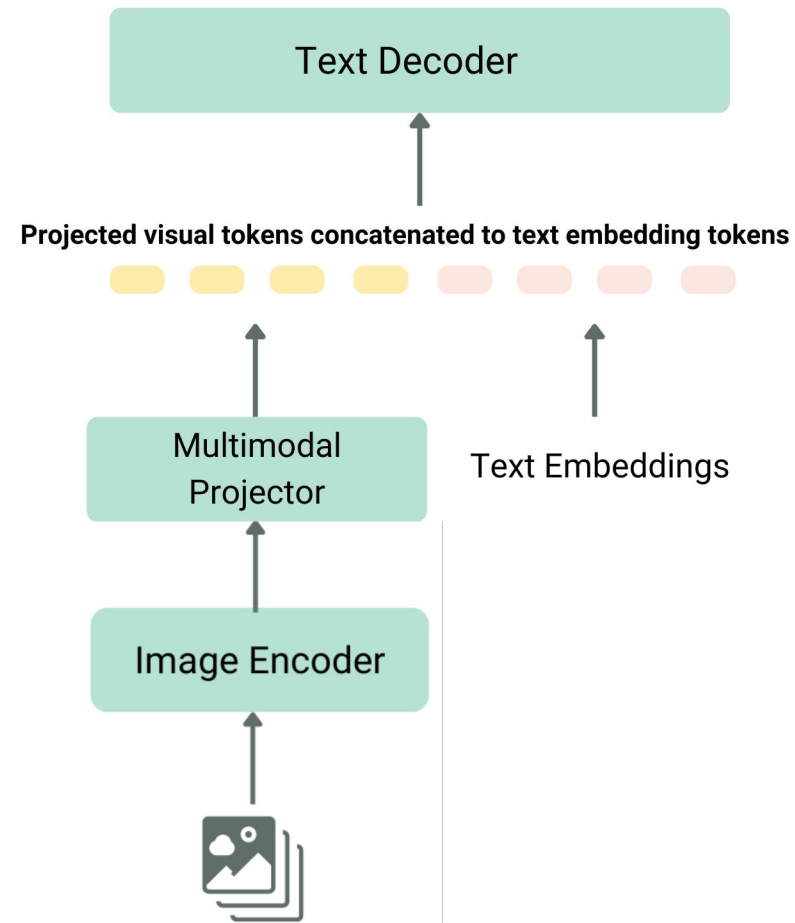
Ground Truth Caption: A brunette girl wearing sunglasses and a yellow shirt.

Generated Caption: A woman in a black shirt and sunglasses smiles.

- **Caption generation:** create natural language descriptions of content of some image

VLM: Architecture

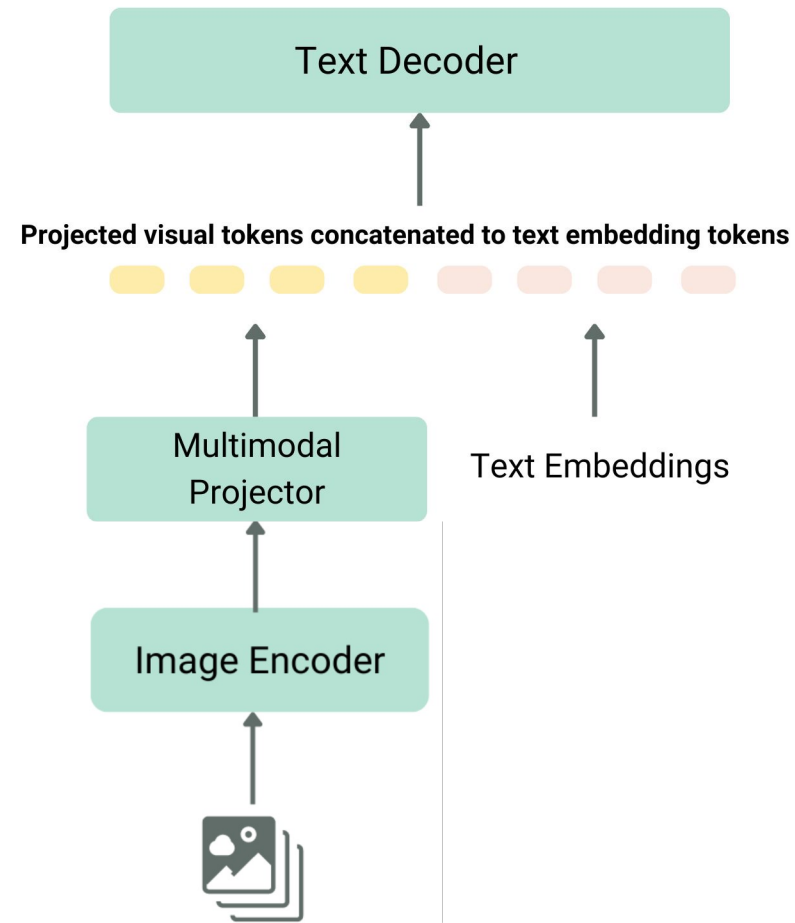
- High-level idea: convert both the image and the text inputs into embedding vectors, then pass those vectors into a decoder-only transformer and do next (text) token prediction



- Two common encoders:
 - VQ-VAE encoder followed by an embedding layer that converts the discrete tokens into dense numerical vectors
 - CLIP encoder, that directly learns an embedding vector using a contrastive pre-training objective

VLM: Architecture

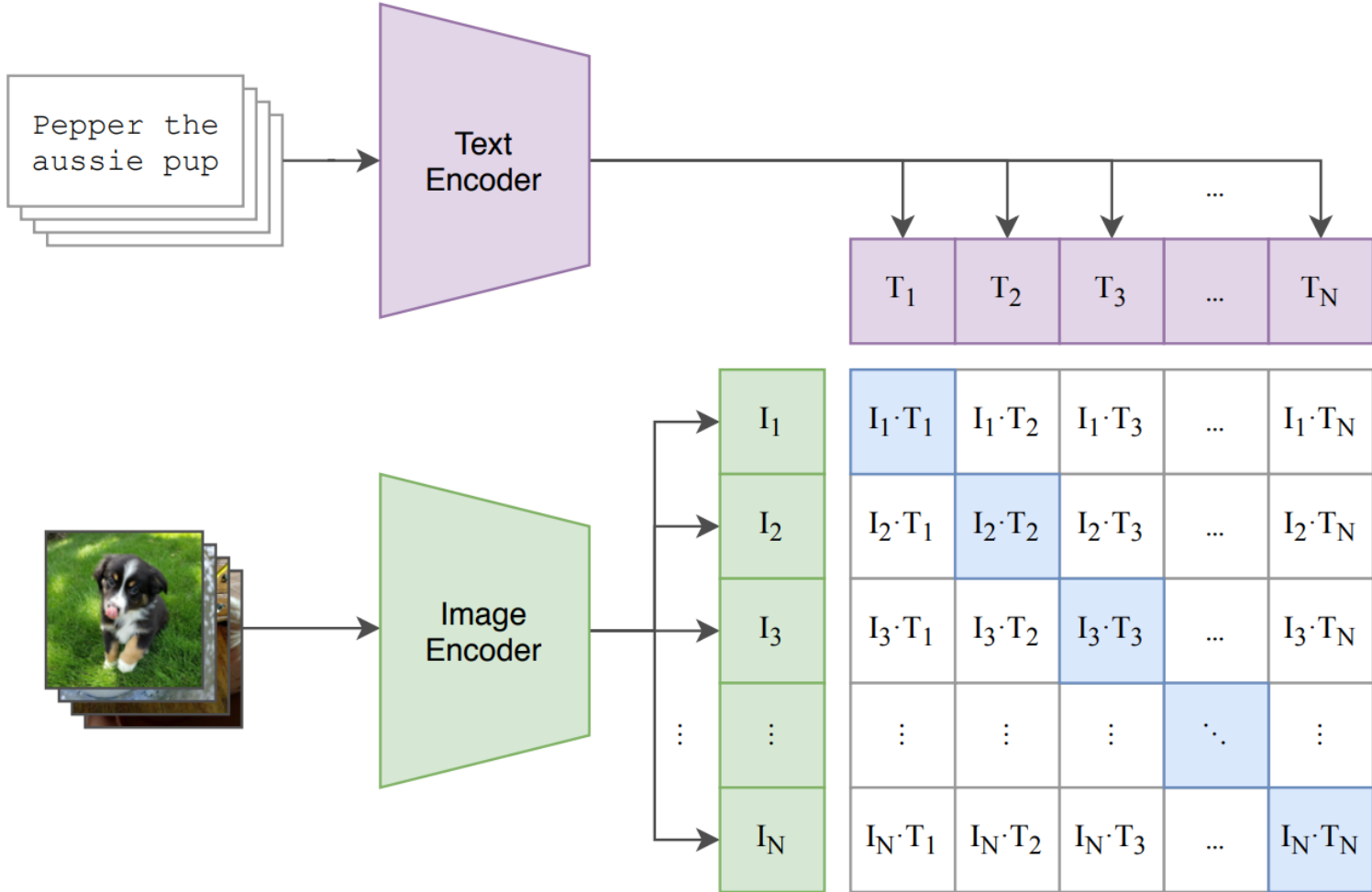
- High-level idea: convert both the image and the text inputs into embedding vectors, then pass those vectors into a decoder-only transformer and do next (text) token prediction



- Two common encoders:
 - VQ-VAE encoder followed by an embedding layer that converts the discrete tokens into dense numerical vectors
 - CLIP encoder, that directly learns an embedding vector using a contrastive pre-training objective

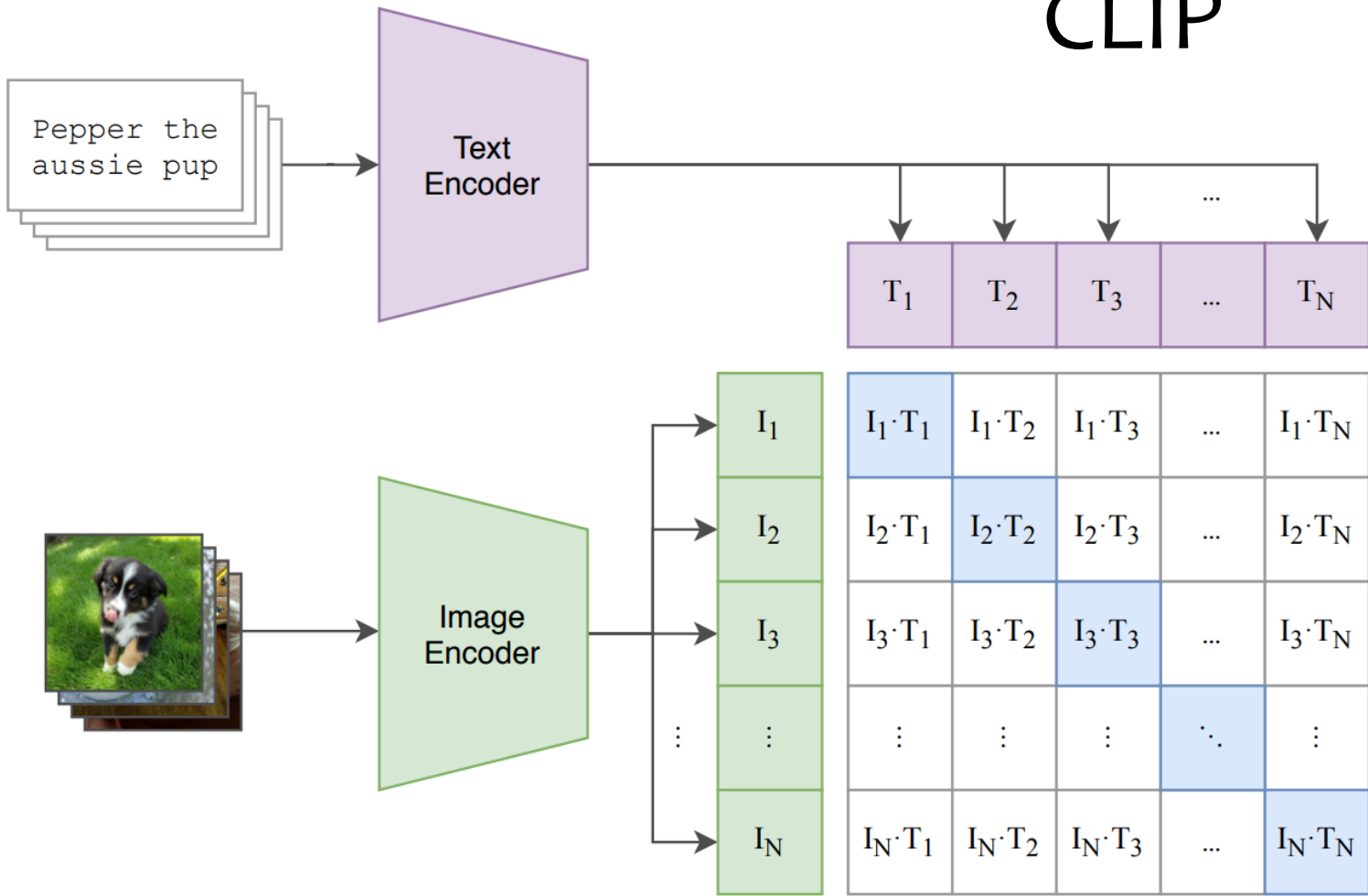
CLIP

CLIP



- The text encoder is, e.g., an encoder-only transformer
- The image encoder is, e.g., a ResNet-like CNN or ViT
- Both are linearly projected into same-dimensional vectors i.e., the multi-modal embedding space

CLIP

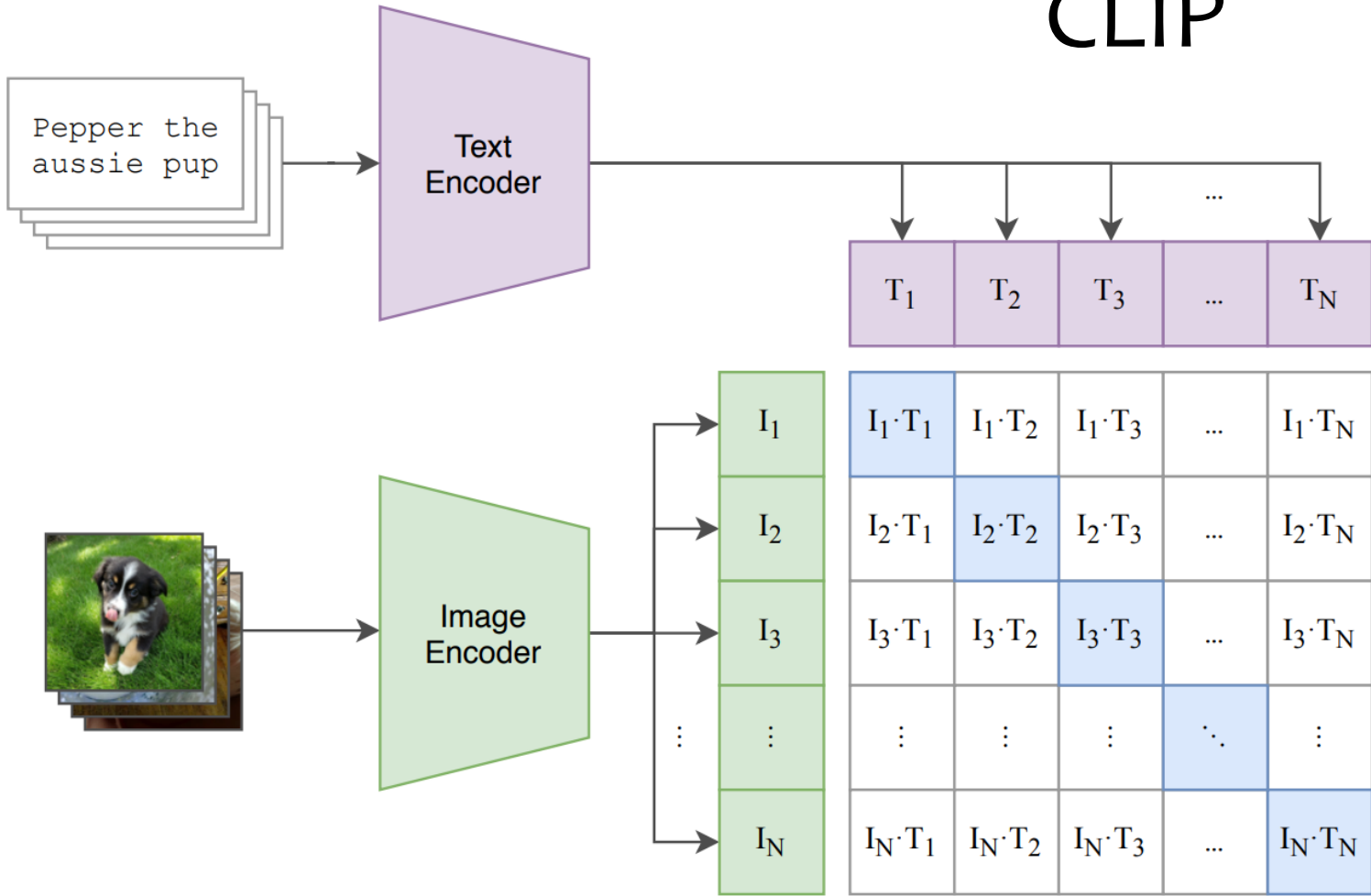


**Incorrect (but intuitive)
objective function:**

$$\max \left[\sum_{i=1}^N I_i^\top T_i - \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N I_i^\top T_j \right]$$

Given a mini-batch of N (image, caption) pairs, both encoders are simultaneously pre-trained to maximize the cosine similarity of corresponding image-caption embedding vectors and minimize all other pairwise cosine similarities

CLIP

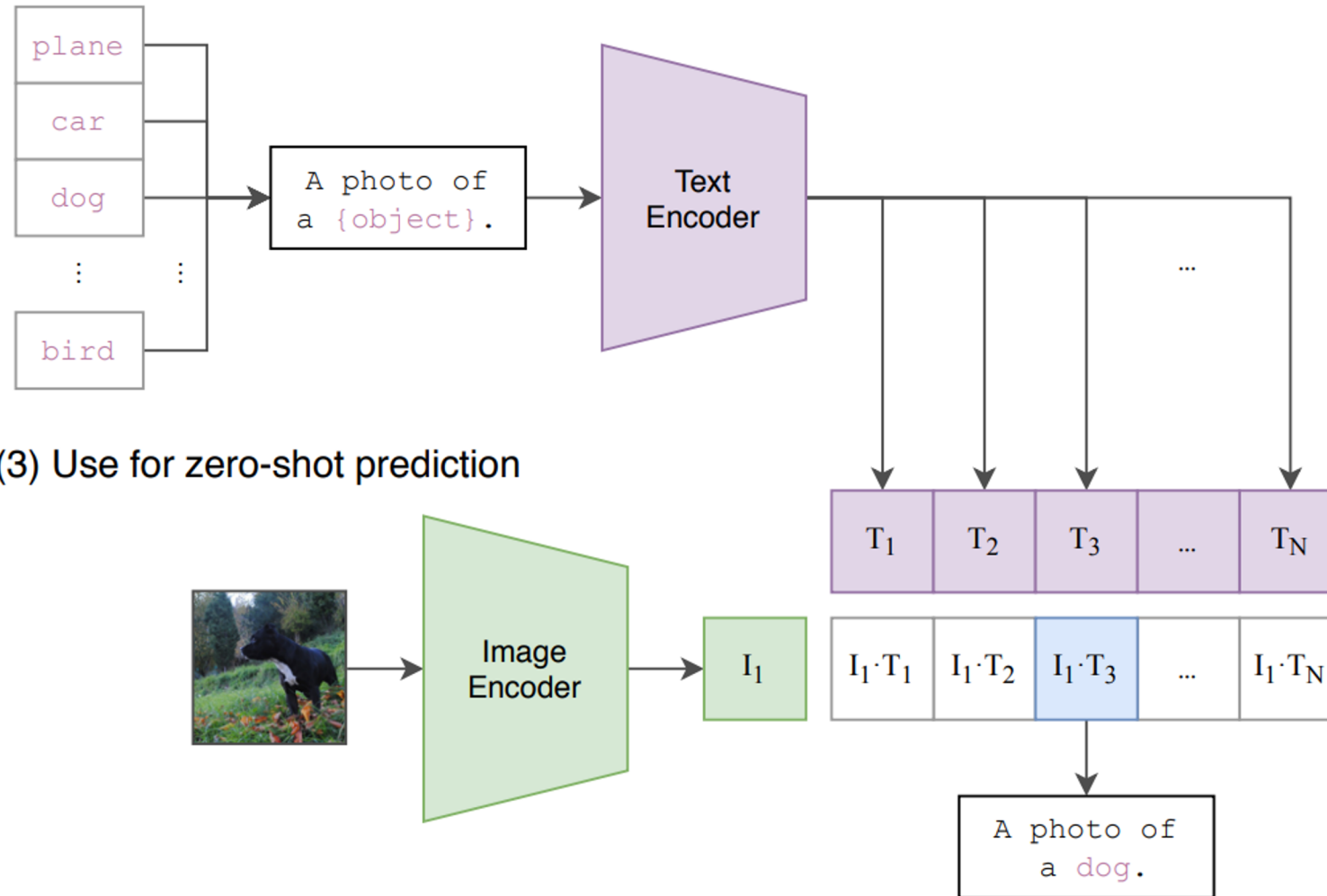


Correct objective function:

$$\max \sum_{i=1}^N \left[\log \frac{\exp \left(\frac{I_i^\top T_i}{\tau} \right)}{\sum_{j=1}^N \exp \left(\frac{I_i^\top T_j}{\tau} \right)} + \log \frac{\exp \left(\frac{I_i^\top T_i}{\tau} \right)}{\sum_{j=1}^N \exp \left(\frac{I_j^\top T_i}{\tau} \right)} \right]$$

Given a mini-batch of N (image, caption) pairs, both encoders are simultaneously pre-trained to maximize the cosine similarity of corresponding image-caption embedding vectors and minimize all other pairwise cosine similarities

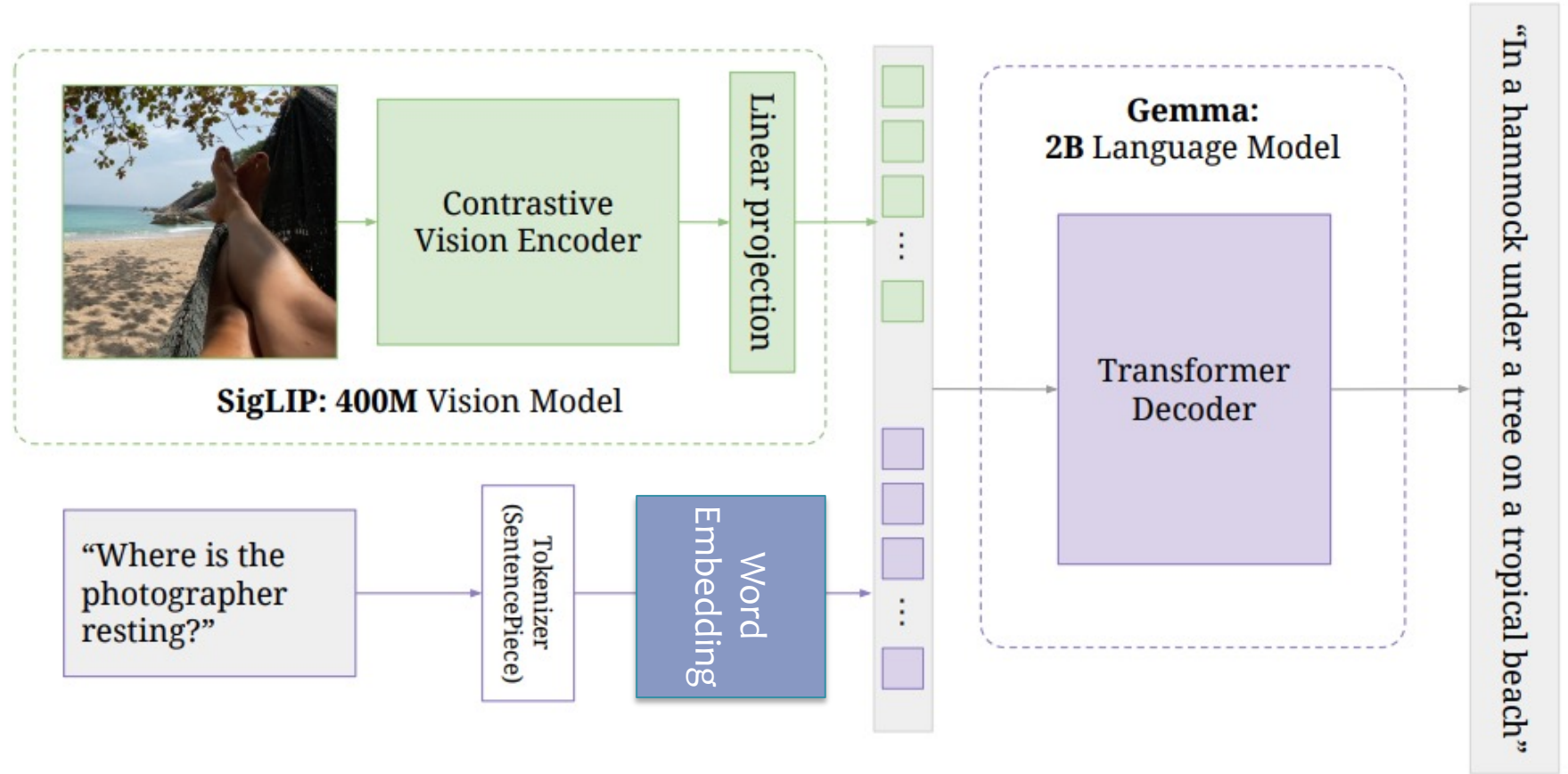
CLIP for Zero Shot Classification



VLMS WITH TEXT-ONLY DECODERS

PaliGemma

- SigLIP is a variant of CLIP
- Gemma is a 2B LLM (open source counterpart to Gemini)



Qwen-VL

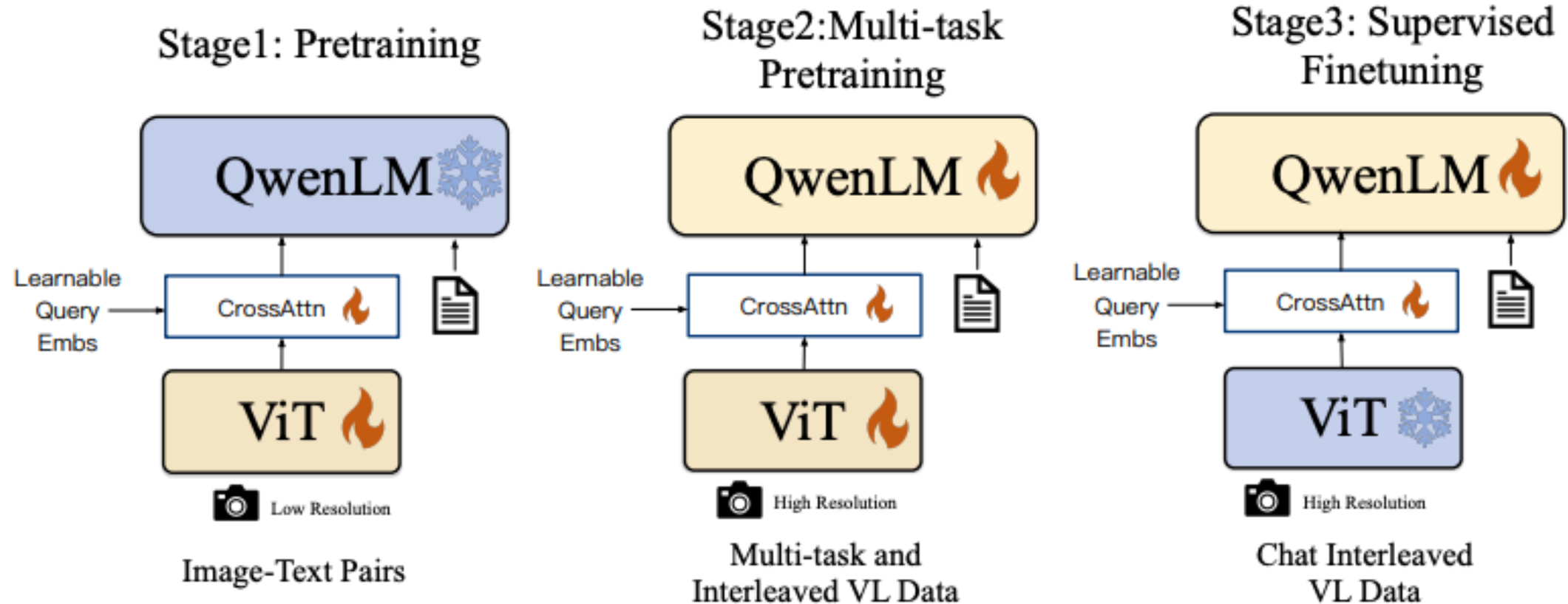
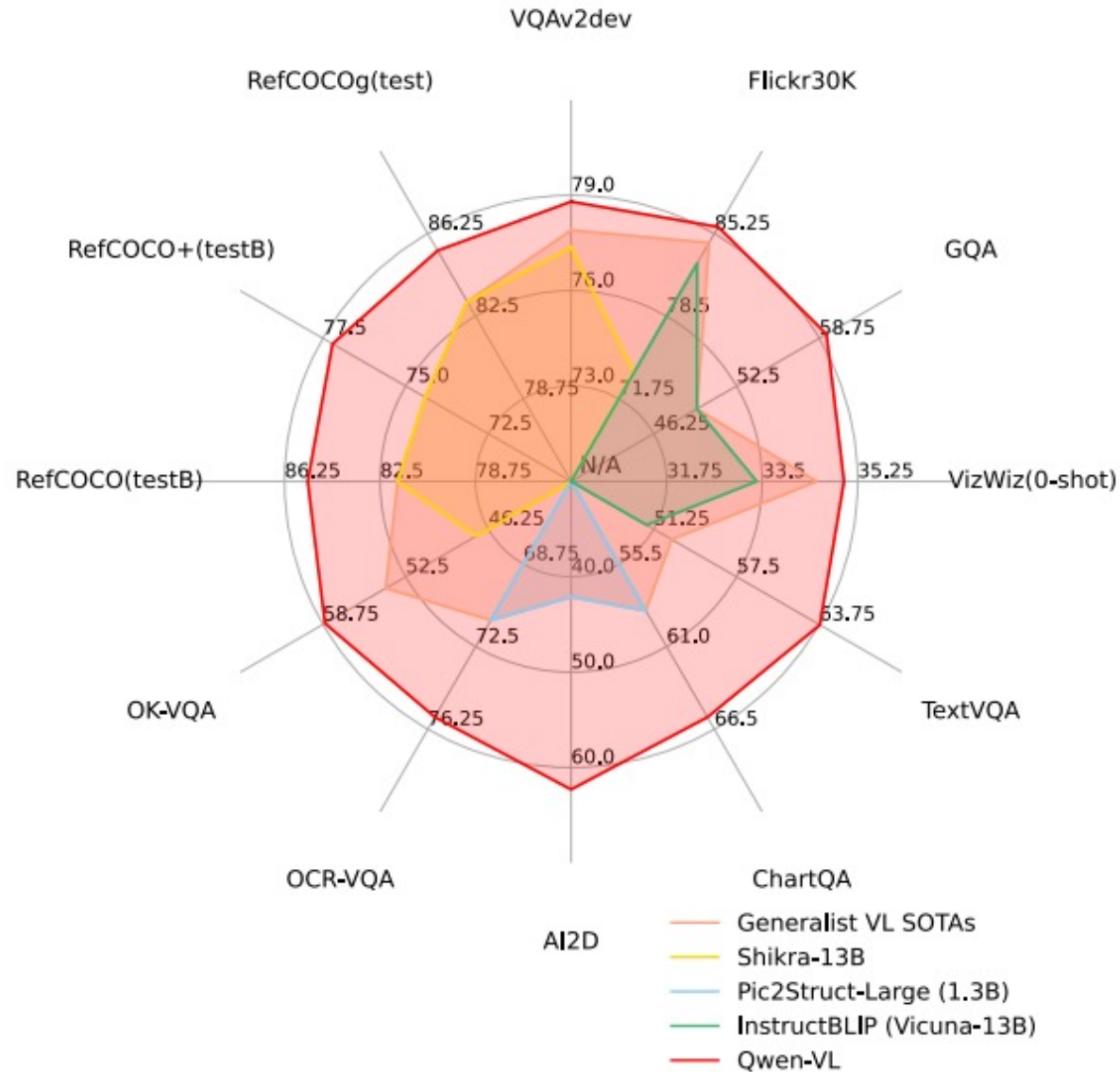


Figure 3: The training pipeline of the Qwen-VL series.

Qwen-VL



Llama 3.2 Vision

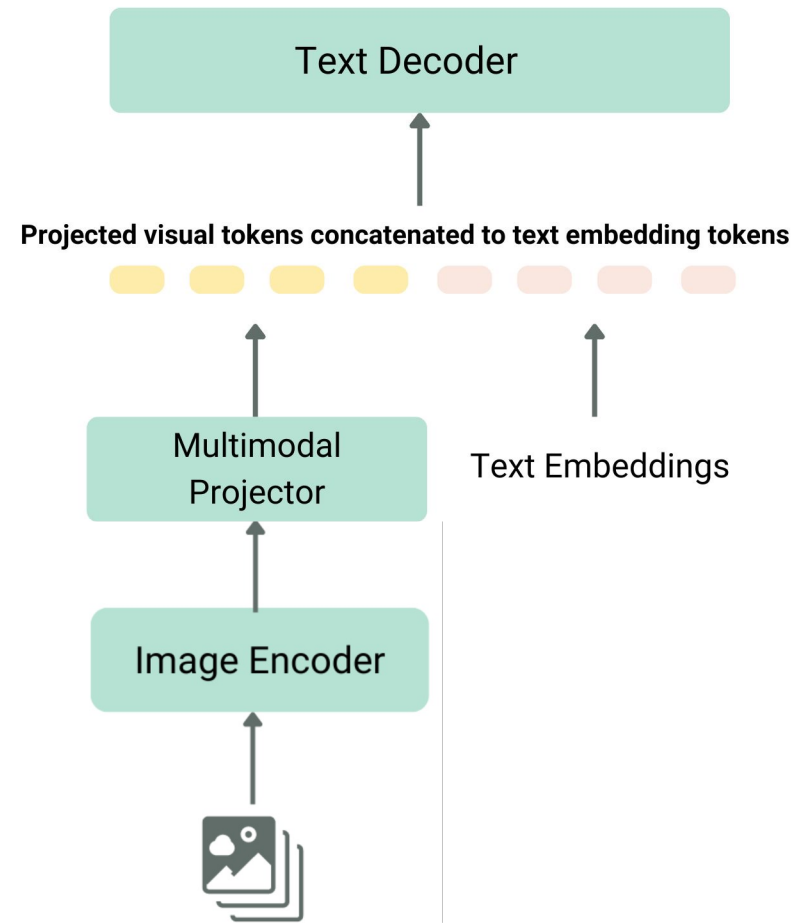
Vision instruction-tuned benchmarks

Modality	Category Benchmark	Llama 3.2 11B	Llama 3.2 90B	Claude 3 – Haiku	GPT-4o-mini
Image	College-level Problems and Mathematical Reasoning MMMU (val, 0-shot CoT, micro avg accuracy)	50.7	60.3	50.2	59.4
	MMMU-Pro, Standard (10 opts, test)	33.0	45.2	27.3	42.3
	MMMU-Pro, Vision (test)	23.7	33.8	20.1	36.5
	MathVista (testmini)	51.5	57.3	46.4	56.7
	Charts and Diagram Understanding ChartQA (test, 0-shot CoT relaxed accuracy)*	83.4	85.5	81.7	—
	A12 Diagram (test)*	91.1	92.3	86.7	—
	DocVQA (test, ANLS)*	88.4	90.1	88.8	—
	General Visual Question Answering VQAv2 (test)	75.2	78.1	—	—
Text	General MMLU (0-shot, CoT)	73.0	86.0	75.2 (5-shot)	82.0
	Math MATH (0-shot, CoT)	51.9	68.0	38.9	70.2
	Reasoning GPQA (0-shot, CoT)	32.8	46.7	33.3	40.2
	Multilingual MGSM (0-shot, CoT)	68.9	86.9	75.1	87.0

VISION LANGUAGE MODELS (VLMS)

VLM: Architecture

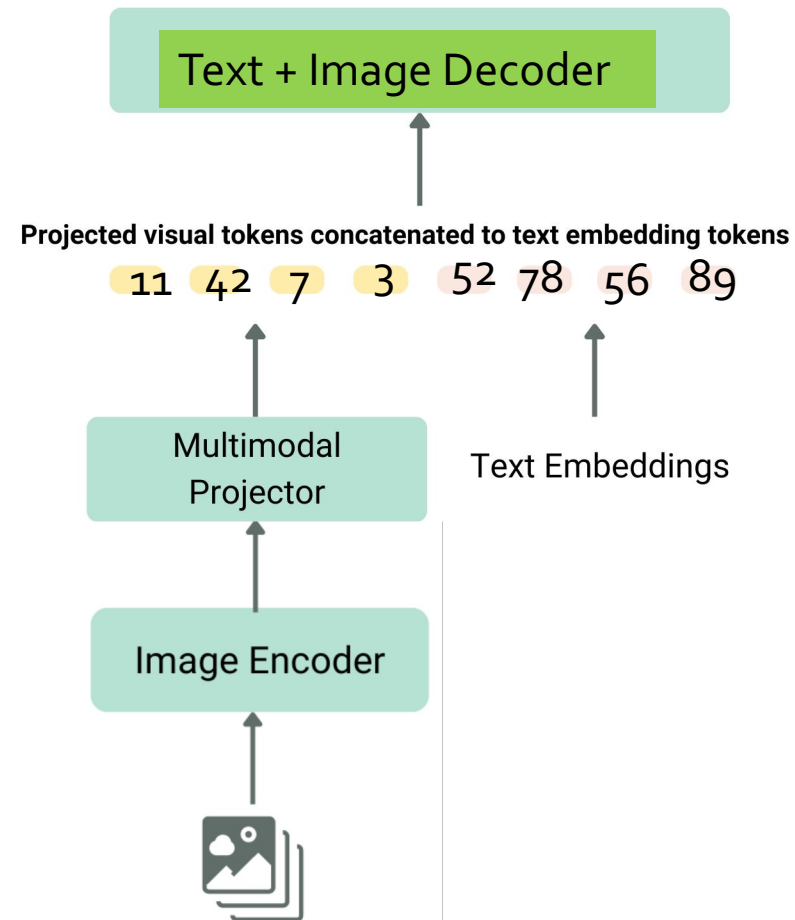
- High-level idea: convert both the image and the text inputs into embedding vectors, then pass those vectors into a decoder-only transformer and do next (text) token prediction



- Two common encoders:
 - VQ-VAE encoder followed by an embedding layer that converts the discrete tokens into dense numerical vectors
 - CLIP encoder, that directly learns an embedding vector using a contrastive pre-training objective

VLM: Architecture

- High-level idea: convert both the image and the text inputs into **integers**, then pass those **integers** into a decoder-only transformer and do next (**text or image**) token prediction

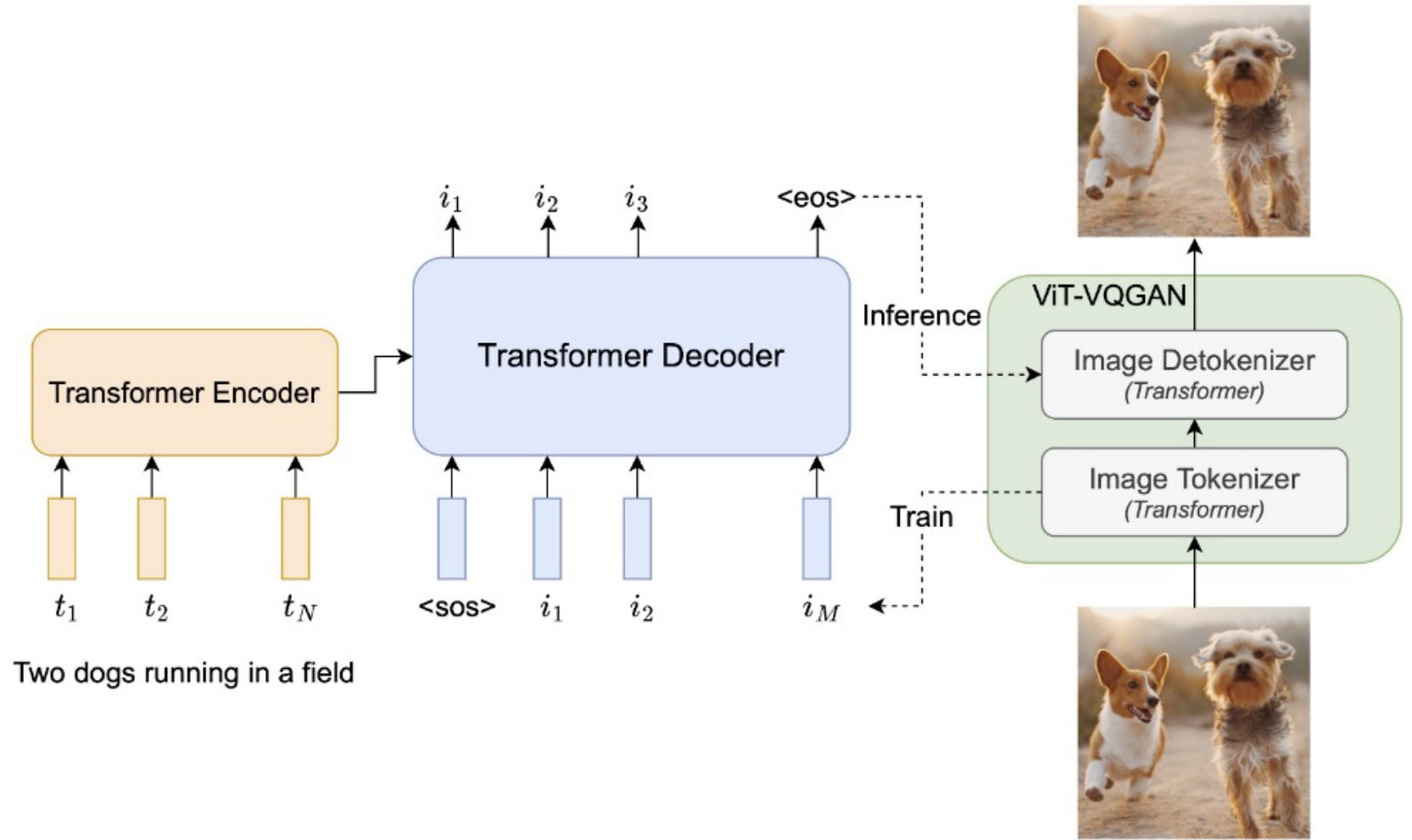


- Two common encoders:
 - VQ-VAE encoder followed by an embedding layer that converts the discrete tokens into dense numerical vectors
 - CLIP encoder, that directly learns an embedding vector using a contrastive pre-training objective

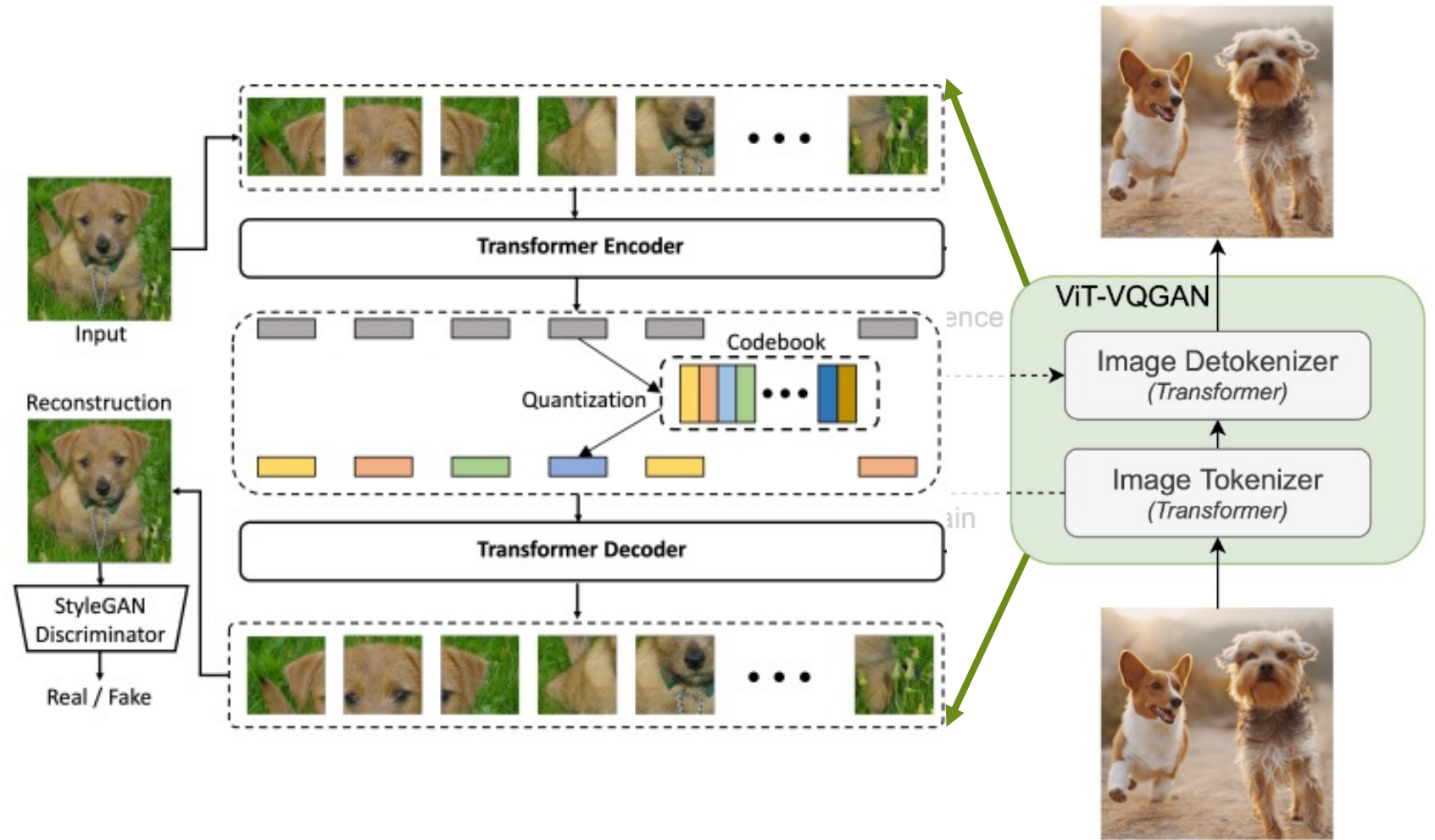
Why VLMs with Integer Tokens?

VQ-VAES

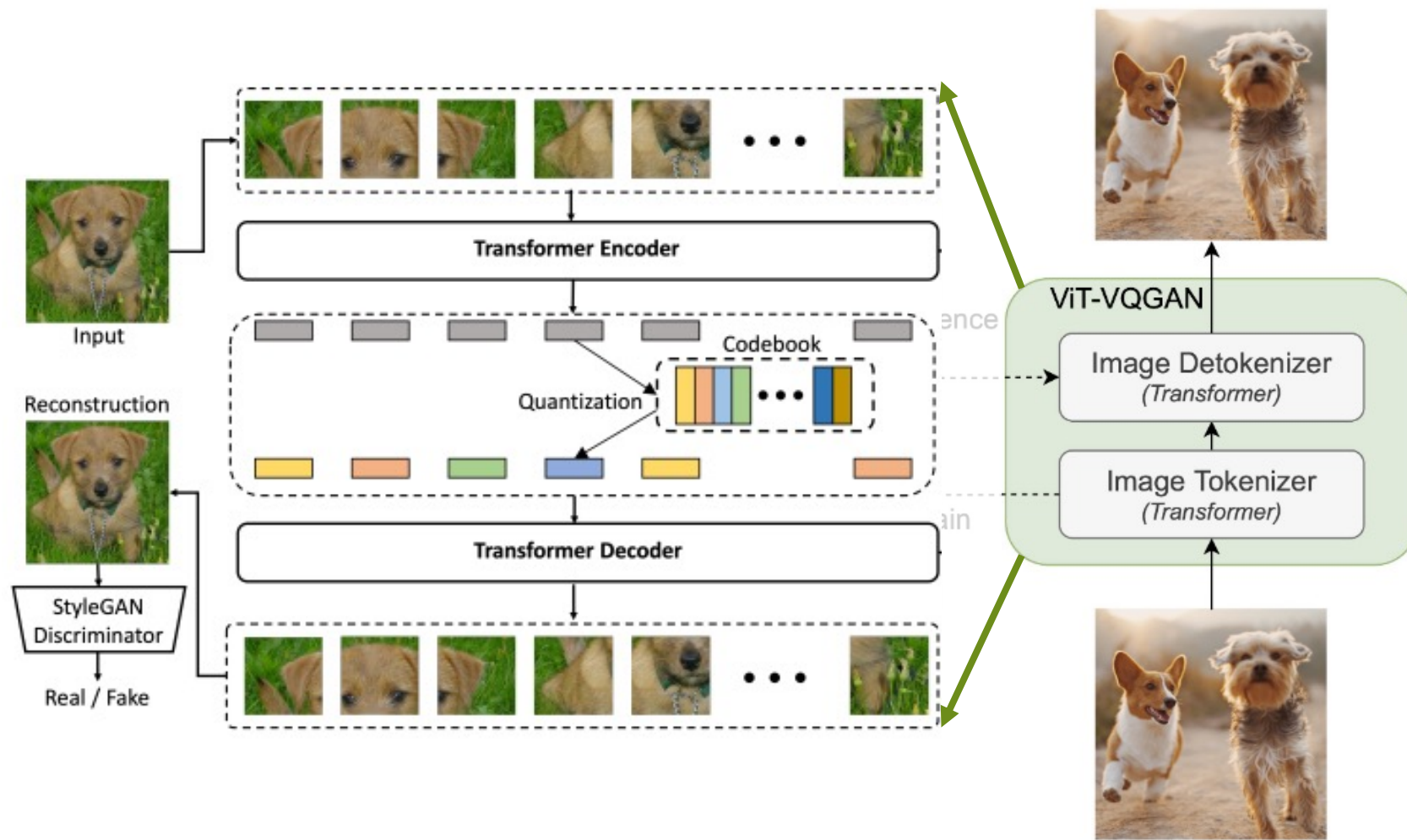
Recall: Parti

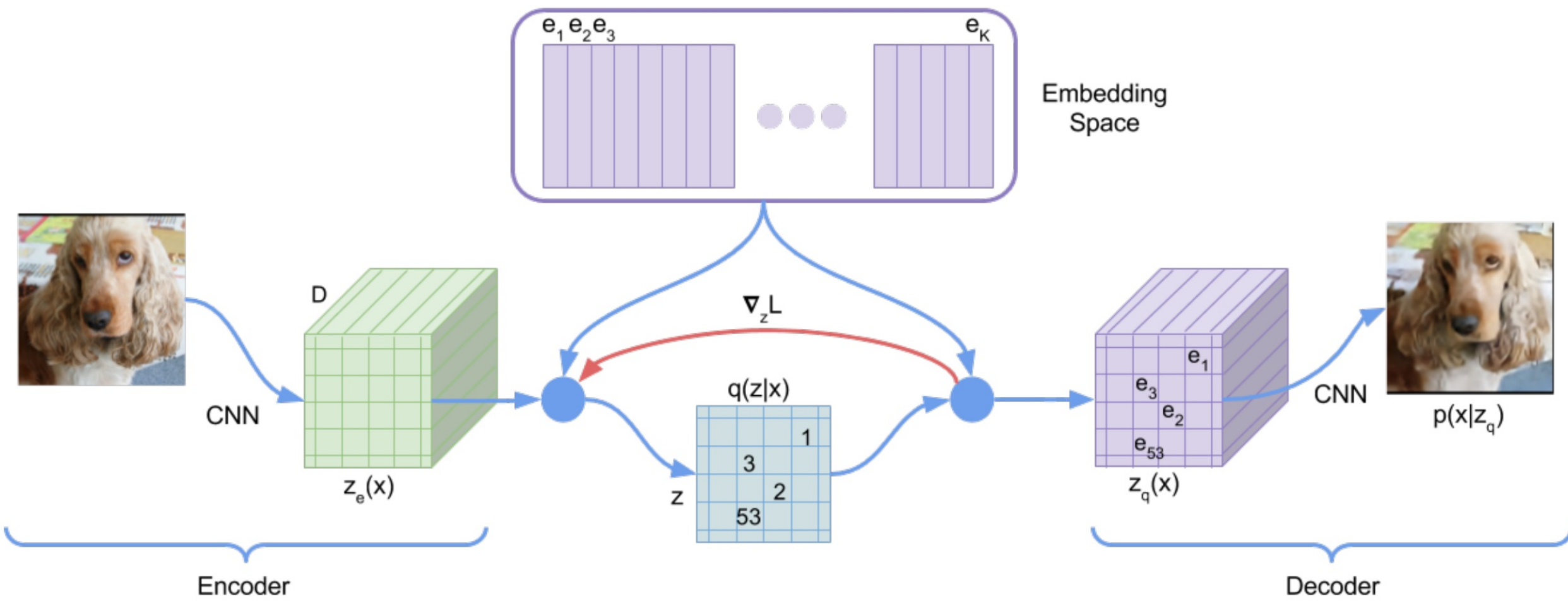


Recall: Image Tokenization

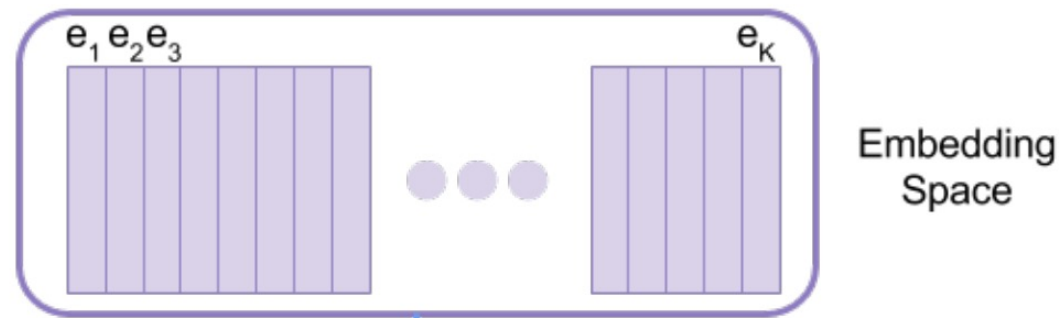


How can we (pre-)train these models given the non-differentiable quantization operation?





Vector-Quantized VAEs

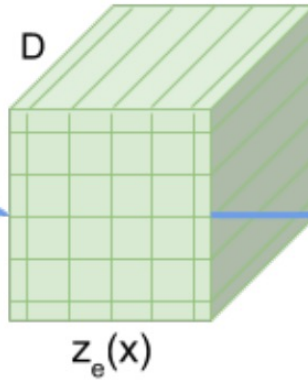


- Embedding space consists of K D -dimensional latent vectors $\{e_1, \dots, e_K\}$ which are learned during training
- The indices $[1, \dots, K]$ of each latent vector correspond to the “image tokens” in some fixed-length codebook

Vector-Quantized VAEs

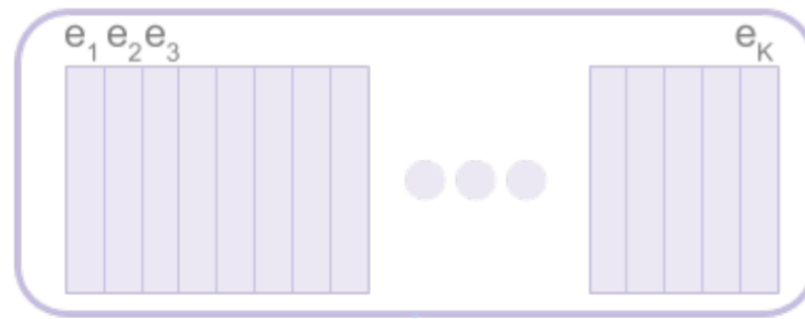


CNN



$z_e(x)$

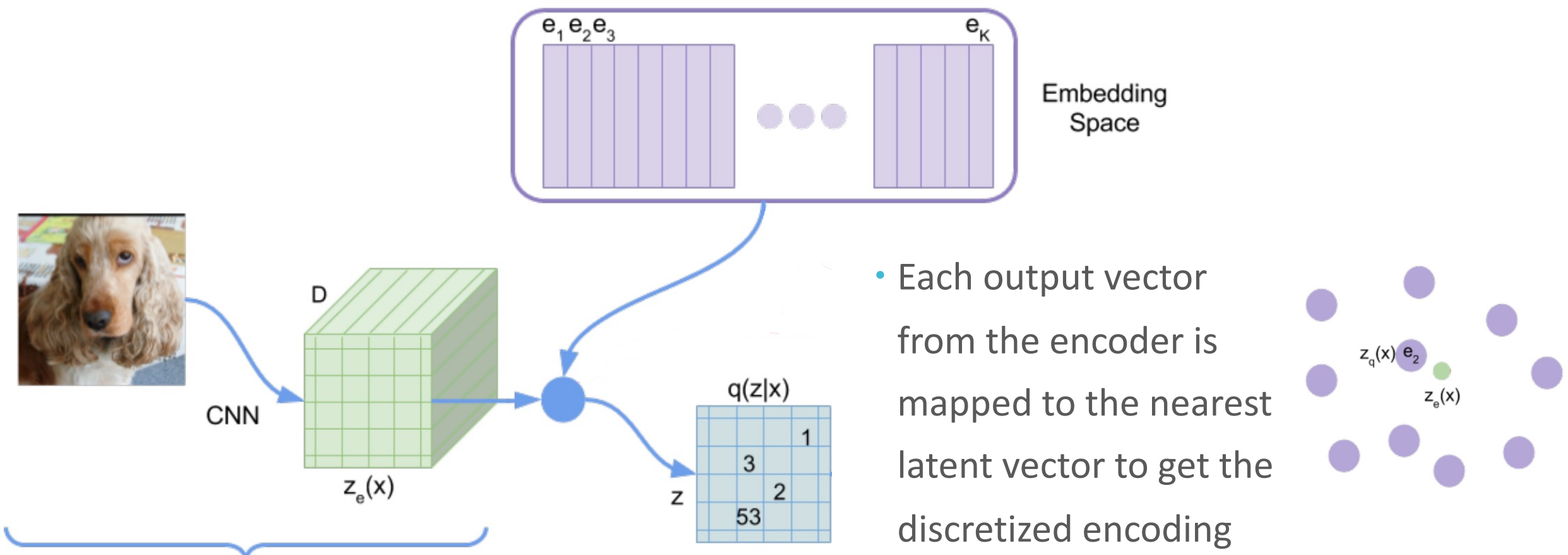
Encoder



Embedding Space

- The encoder (e.g., a ResNet-like CNN) maps images to N D -dimensional vectors

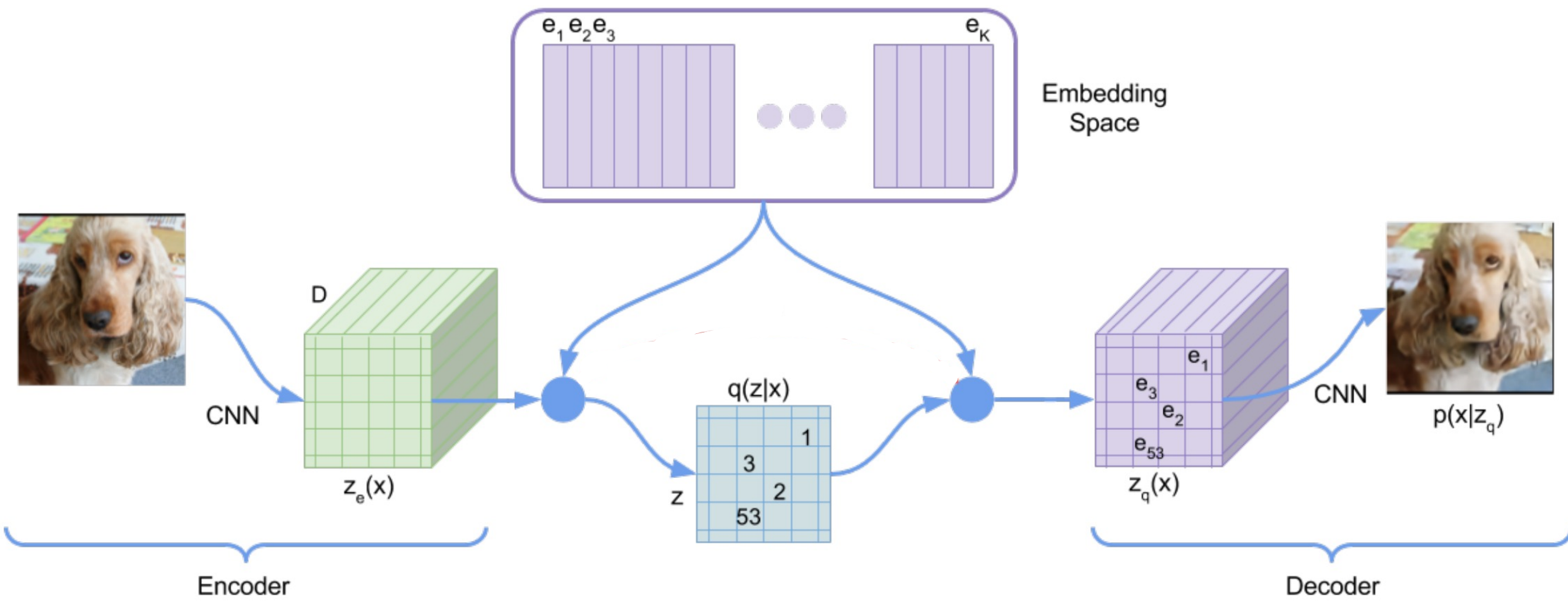
Vector-Quantized VAEs



- Each output vector from the encoder is mapped to the nearest latent vector to get the discretized encoding

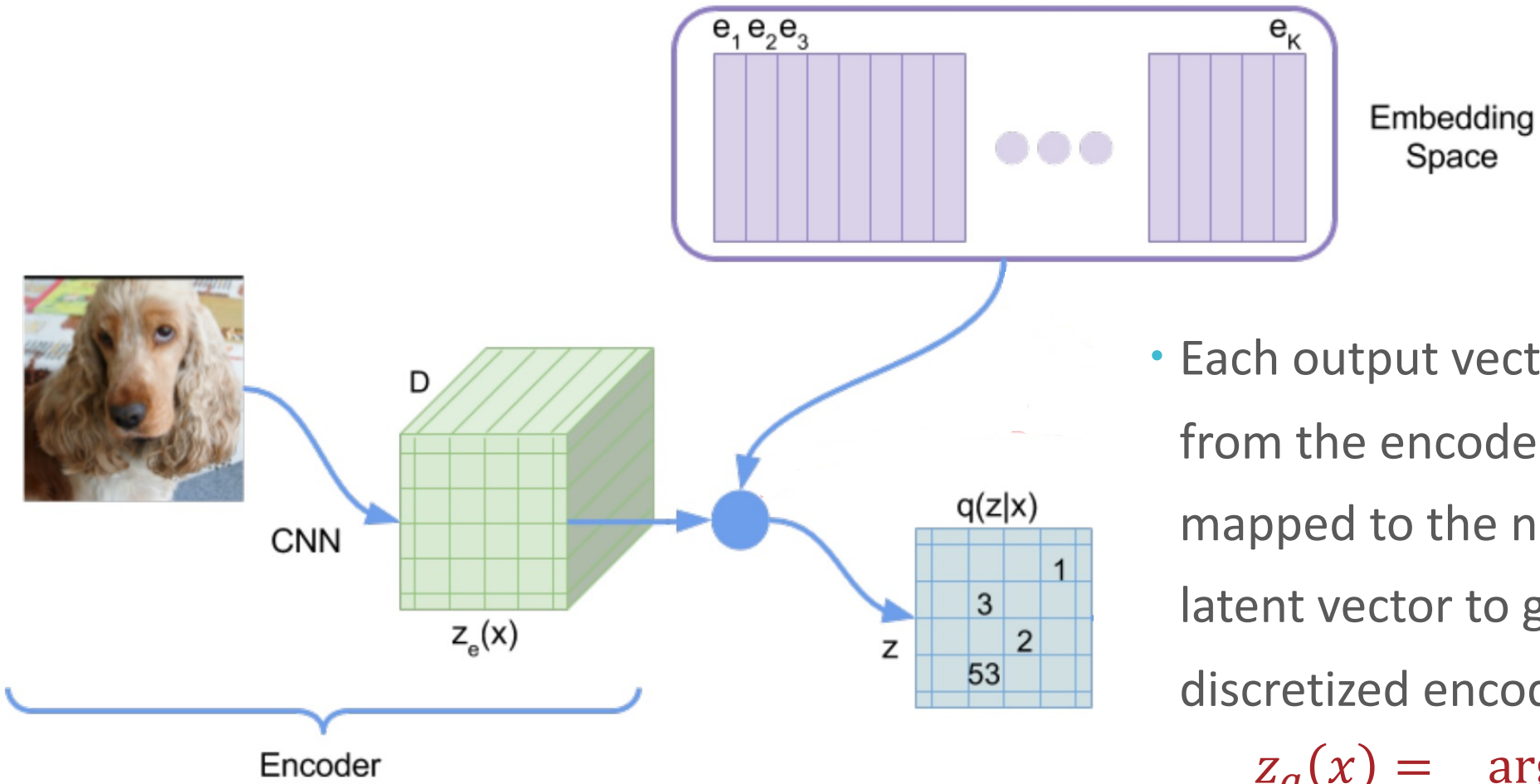
$$z_q(x) = \operatorname{argmin}_{e \in \{e_1, \dots, e_K\}} \|z_e(x) - e\|_2^2$$

Vector-Quantized VAEs

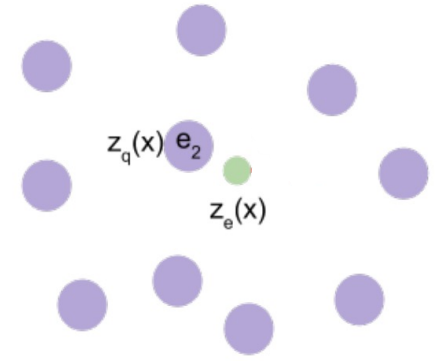


- The decoder takes the discretized representation and recreates the original image

Vector-Quantized VAEs

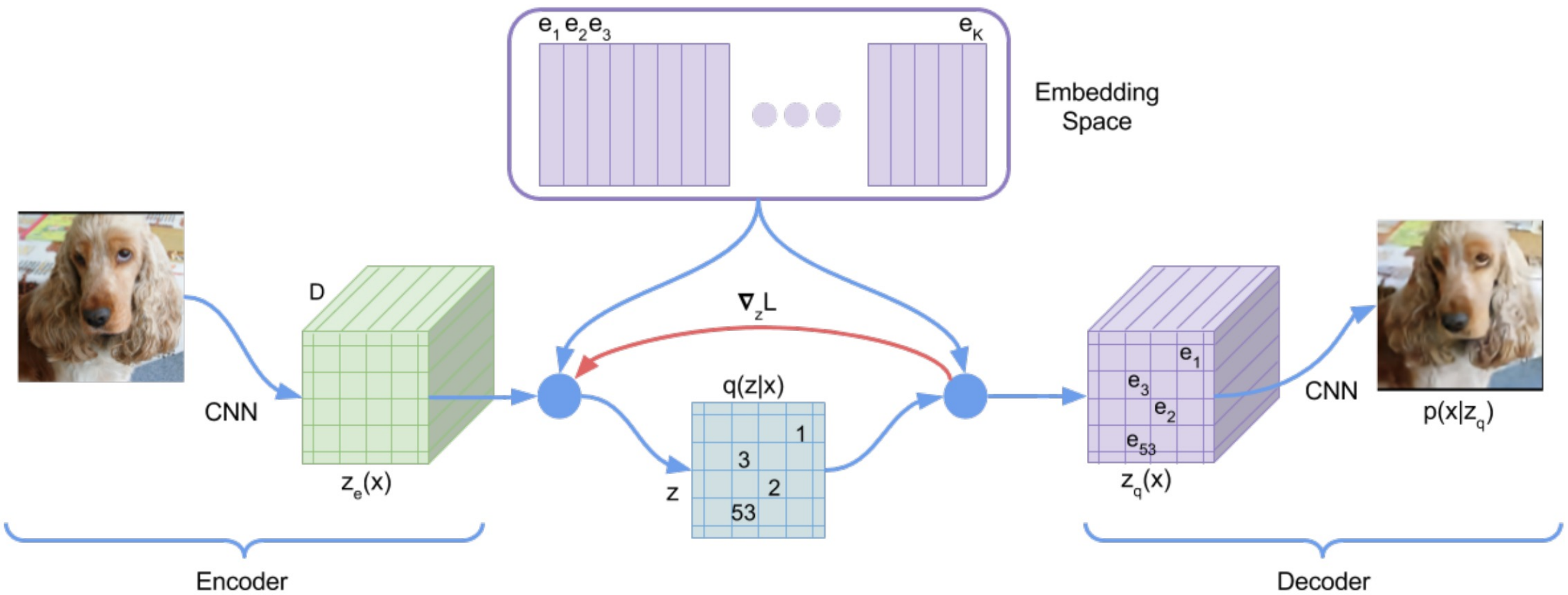


- Each output vector from the encoder is mapped to the nearest latent vector to get the discretized encoding



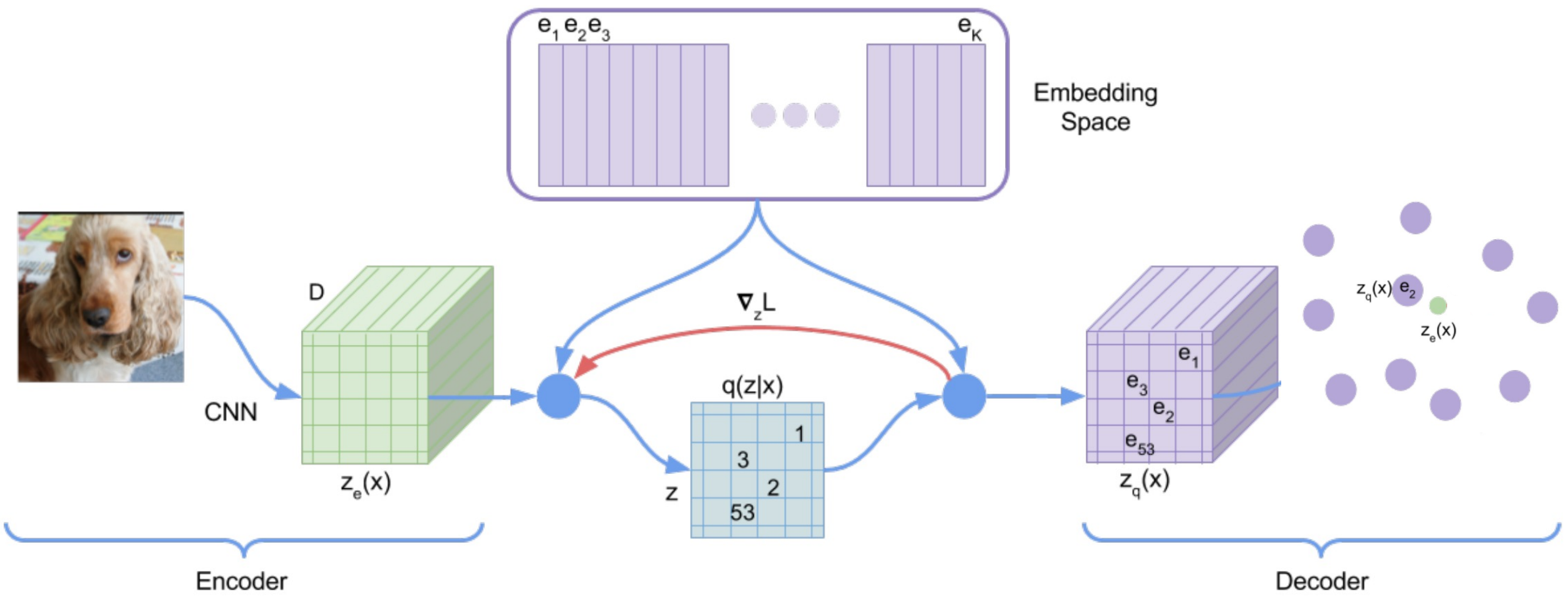
$$z_q(x) = \underset{e \in \{e_1, \dots, e_K\}}{\operatorname{argmin}} \|z_e(x) - e\|_2^2$$

Wait, how would we take the gradient through the argmin?



- Treat the gradient w.r.t. $z_q(x)$ as an estimate of the gradient w.r.t. $z_e(x)$

Straight-through Estimator



- Intuition: the closer $z_q(x)$ and $z_e(x)$, the better the estimate (under certain assumptions)

Straight-through Estimator

VQ-VAE Objective Function

- Intuition: we want the latent vectors to correspond to relevant points in the embedding space i.e., ones that are near the outputs of the encoder
- However, we also want the encoder to respect the latent vectors and not overfit to the training dataset
- Idea: augment the standard VAE objective with some regularizing terms that drive the two closer to each other

$$-\log p_{\theta}(x|z_q(x)) + \left\| \text{sg}[z_e(x)] - z_q(x) \right\|_2^2 + \beta \left\| z_e(x) - \text{sg}[z_q(x)] \right\|_2^2$$

where **sg** is the stop-gradient operator which fixes the argument to be non-updated constant

VQ-VAE Objective Function

- Intuition: we want the latent vectors to correspond to relevant points in the embedding space i.e., ones that are near the outputs of the encoder
- However, we also want the encoder to respect the latent vectors and not overfit to the training dataset
- Idea: augment the standard VAE objective with some regularizing terms that drive the two closer to each other

$$-\log p_{\theta}(x|z_q(x)) + \|sg[z_e(x)] - z_q(x)\|_2^2 + \beta \|z_e(x) - sg[z_q(x)]\|_2^2$$

- The first term is the typical reconstruction error objective

VQ-VAE Objective Function

- Intuition: we want the latent vectors to correspond to relevant points in the embedding space i.e., ones that are near the outputs of the encoder
- However, we also want the encoder to respect the latent vectors and not overfit to the training dataset
- Idea: augment the standard VAE objective with some regularizing terms that drive the two closer to each other

$$-\log p_{\theta}(x|z_q(x)) + \left\| \text{sg}[z_e(x)] - z_q(x) \right\|_2^2 + \beta \left\| z_e(x) - \text{sg}[z_q(x)] \right\|_2^2$$

- The second term drives the latent vector to be closer to the encoder output vector that was mapped to it

VQ-VAE Objective Function

- Intuition: we want the latent vectors to correspond to relevant points in the embedding space i.e., ones that are near the outputs of the encoder
- However, we also want the encoder to respect the latent vectors and not overfit to the training dataset
- Idea: augment the standard VAE objective with some regularizing terms that drive the two closer to each other

$$-\log p_{\theta}(x|z_q(x)) + \left\| \text{sg}[z_e(x)] - z_q(x) \right\|_2^2 + \beta \left\| z_e(x) - \text{sg}[z_q(x)] \right\|_2^2$$

- The third term drives the encoder to output vectors closer to the latent vectors

CLIP vs. VQ-VAEs

- VLMs with VQ-VAE encoders (or any vector quantized image model) can also generate images in addition to text by defining a loss over the image codebook tokens
- CLIP does not discretize its image embedding so VLMs with CLIP-based encoders cannot (naturally) define a loss over images and thus, can only output text
- However, CLIP embeddings are more expressive than the discrete VQ-VAE encodings so can lead to improved performance in some settings

VLMS WITH TEXT AND IMAGE DECODERS

Large World Model



Large World Model

H MORE IMAGE GENERATION EXAMPLES



A black dog



A blue colored pizza



A cube made of denim



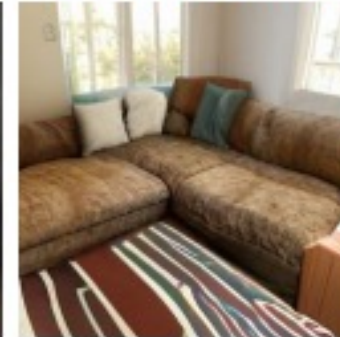
A glass of wine



*A yellow and black bus
cruising through a rainforest*



*Oil painting of a couple in
formal attire caught in the
rain without umbrellas*



*A couch in a cozy living
room*



*A carrot to the left of
broccoli*



*Fisheye lens of a turtle
in a forest*



A blue colored dog

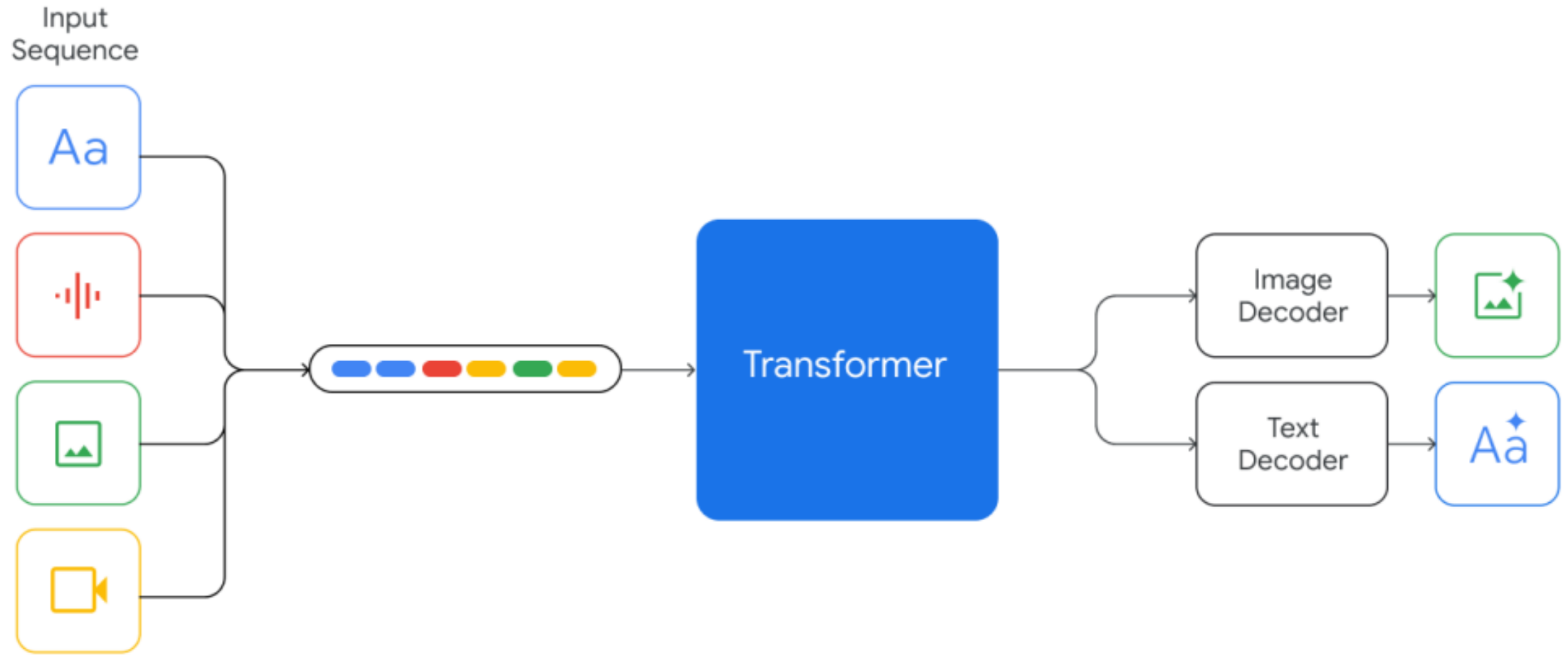


*Stained glass windows
depicting hamburgers and
french fries*





A pink car


Gemini





Gemini


 **Prompt**
Give me two ideas that I could do with these 2 colors


Colors: I see blue and yellow yarn 

How about a cute blue cat? 

Or a blue dog that would also have a yellow ear? 

Give me two ideas that I could do with these 2 colors 

 **Response**
Colors: I see green and pink yarn

Idea 1: How about a green avocado with pink seed? 


Idea 2: Or a green bunny with pink ears? 

Figure 6 | **Image Generation.** Gemini models can output multiple images interleaved with text given a prompt composed of image and text. In the left figure, Gemini Ultra is prompted in a 1-shot setting with a user example of generating suggestions of creating cat and dog from yarn when given two colors, blue and yellow. Then, the model is prompted to generate creative suggestions with two new colors, pink and green, and it generates images of creative suggestions to make a cute green avocado with pink seed or a green bunny with pink ears from yarn as shown in the right figure.