

INF721

2024/2

UFV

Deep Learning

L19: Multimodal Learning

Logistics

Announcements

- ▶ Midterm Exam II has been pushed back to Nov 27th

Last Lecture

- ▶ GPT
- ▶ BERT

Lecture Outline

- ▶ Multimodal Learning
- ▶ Vision Transformers (ViT)
- ▶ Audio Data
- ▶ Text-to-Speech
- ▶ Text-to-Image
- ▶ CLIP (Contrastive Language-Image Pre-training)

Multimodal Learning

Transformers make very few assumptions about input data, so they have become state-of-the-art in many different modalities:

- ▶ Text, Image , Video, Audio...

The core architecture has remained relatively constant across applications, what has changed is the **representation** and **encoding** of inputs and outputs

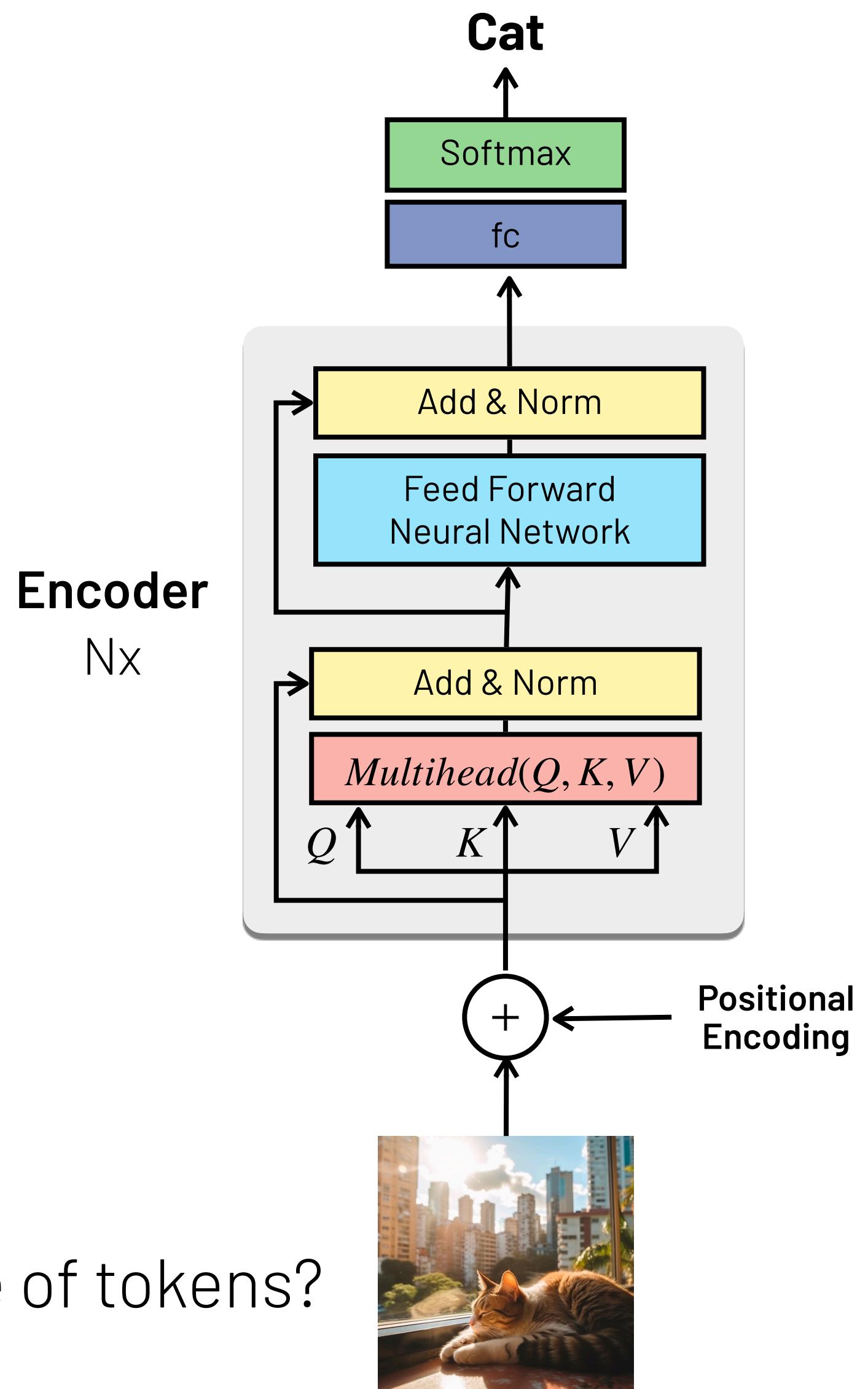
Vision Transformers

Vision transformers typically use an **transformer encoder** for image classification tasks

The overall architecture remains the same:

Main problems:

1. How to transform an image into a sequence of tokens?
2. How to encode positional information?



Images as Sequences (Classification): Naive

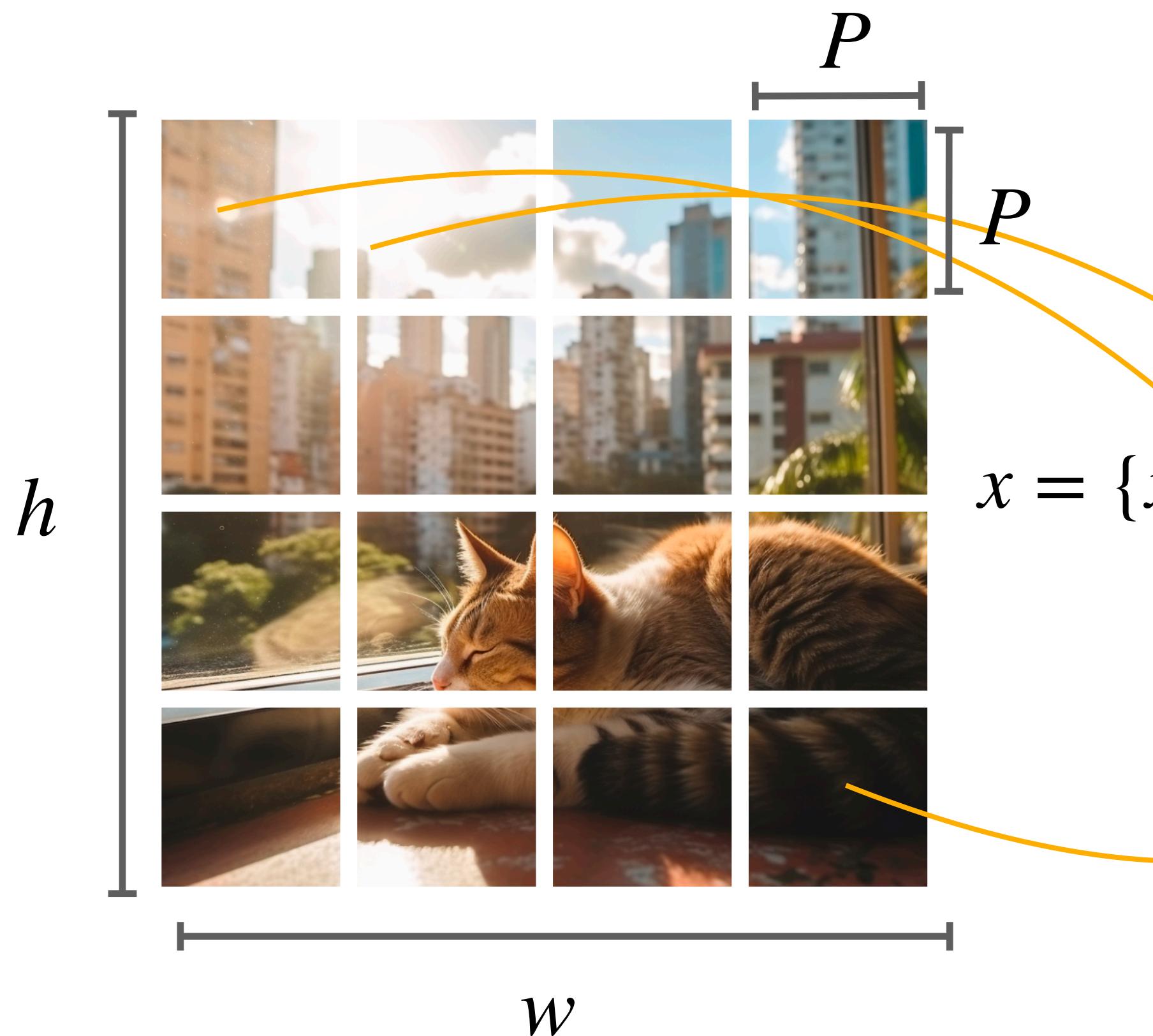
The simplest approach to represent images as sequences is to consider every pixel as an element:



The problem is that the memory required by the transformer grows quadratically $O(n^2)$ with the number of input tokens $n = wh$

Images as Sequences (Classification): Patches

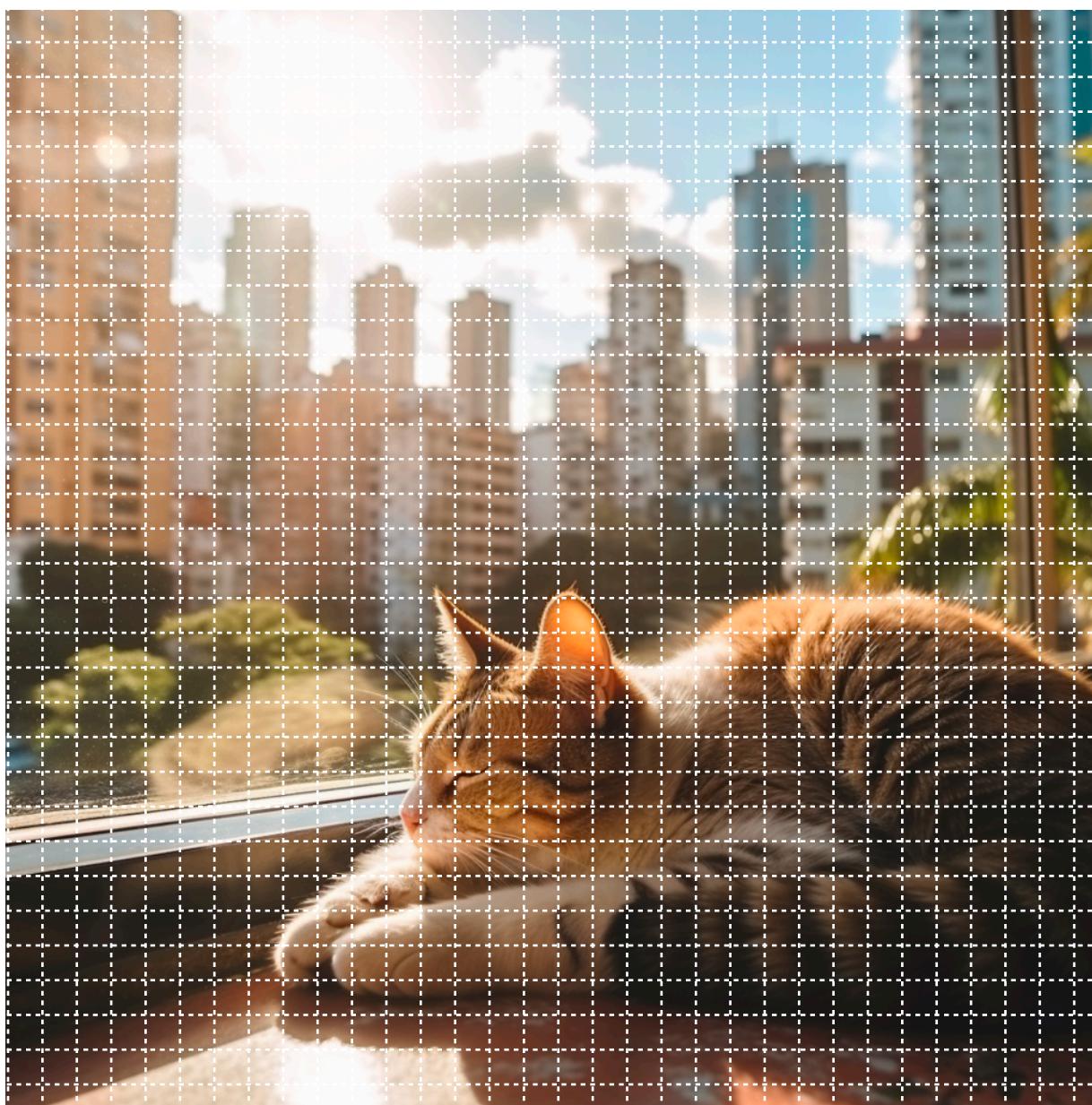
The most common approach is to split the input image into a sequence of non-overlapping patches of size $P \times P$



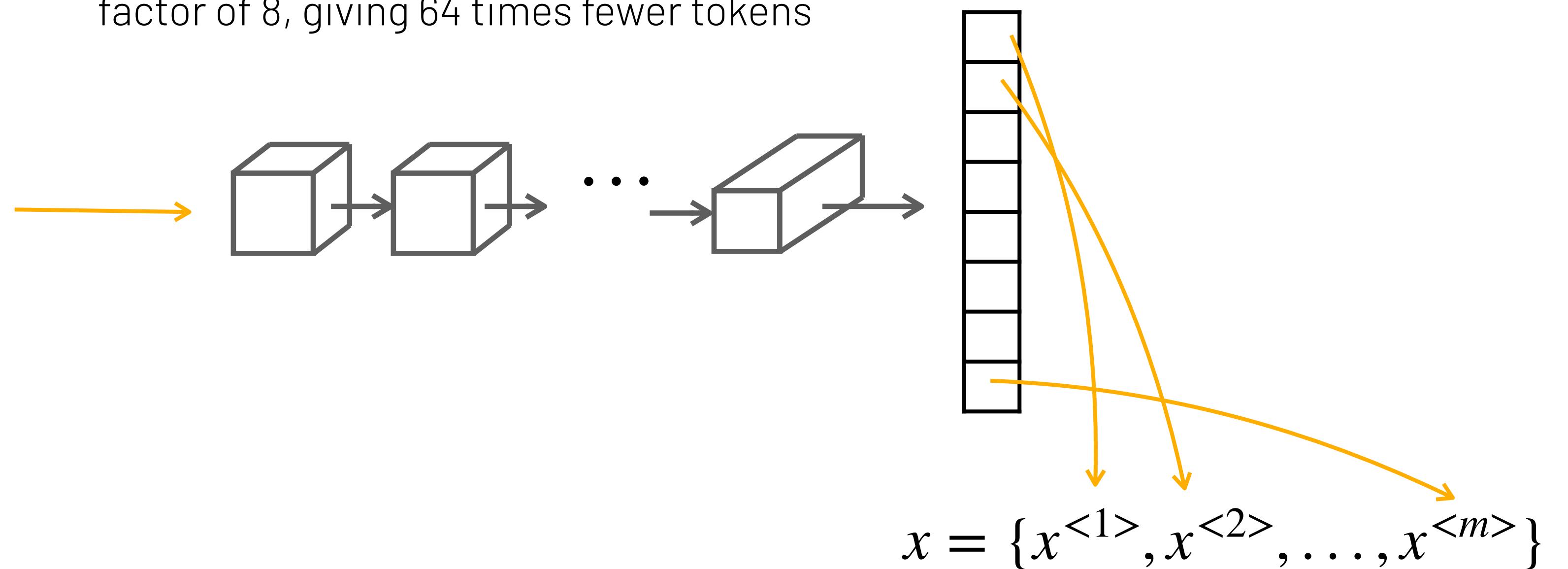
This approach reduces the size of the sequence to $m = \frac{hw}{P^2}$

Images as Sequences (Classification): CNN

Another approach to transform images in sequences is to use a Convolutional Neural Network (CNN) to extract a feature vector from the image:



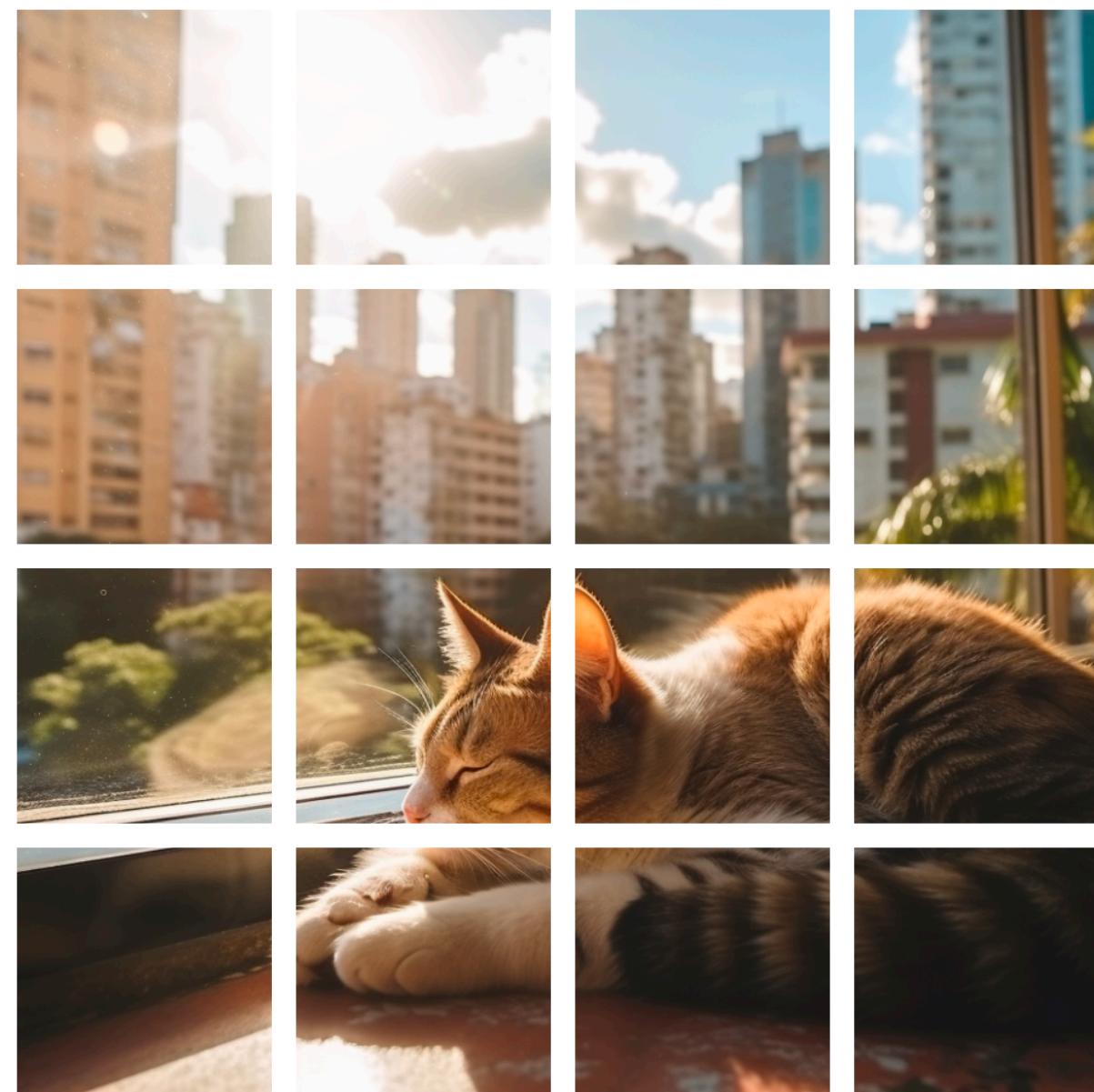
ResNet18 would downsample an image by a factor of 8, giving 64 times fewer tokens



Positional Information (Classification)

It is possible to define explicit positional encoding to a sequence of patches, but the most common approach is to **learn** these positional embeddings:

1. Create one hot encoding based on position indices:



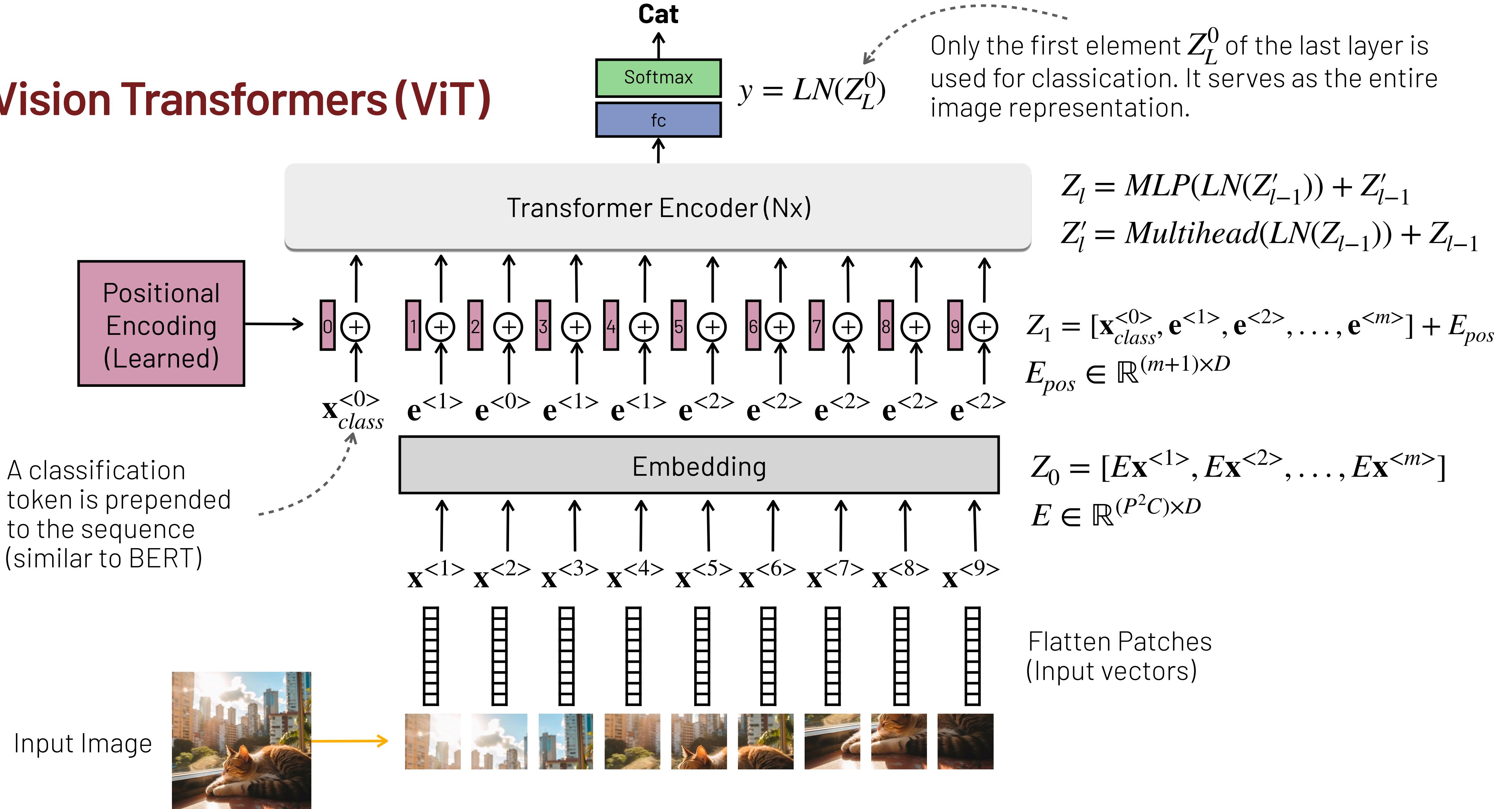
$$x = \{x^{<1>}, x^{<2>}, \dots, x^{<m>}\}$$
$$\begin{array}{llll} 0: & \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} & 1: & \begin{matrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} \\ 2: & \begin{matrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} & \dots & \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{matrix} \\ \vdots & & & \\ m: & \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} & & \end{array}$$
$$o_1 \quad o_2 \quad \dots \quad o_m$$

2. Multiply as an embedding matrix:

$$\begin{array}{cccccc} 0 & 1 & 2 & \cdots & m \\ \begin{matrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} & \times & \begin{matrix} 0 & 1 & 2 & \cdots & m \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{matrix} & & \begin{matrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} & = o_1 \end{array}$$

E_{pos} (weights initialized randomly)

Vision Transformers (ViT)



Generative Image Transformers

Vision transformers typically use an **transformer decoder** for image generation tasks

The overall architecture remains the same:

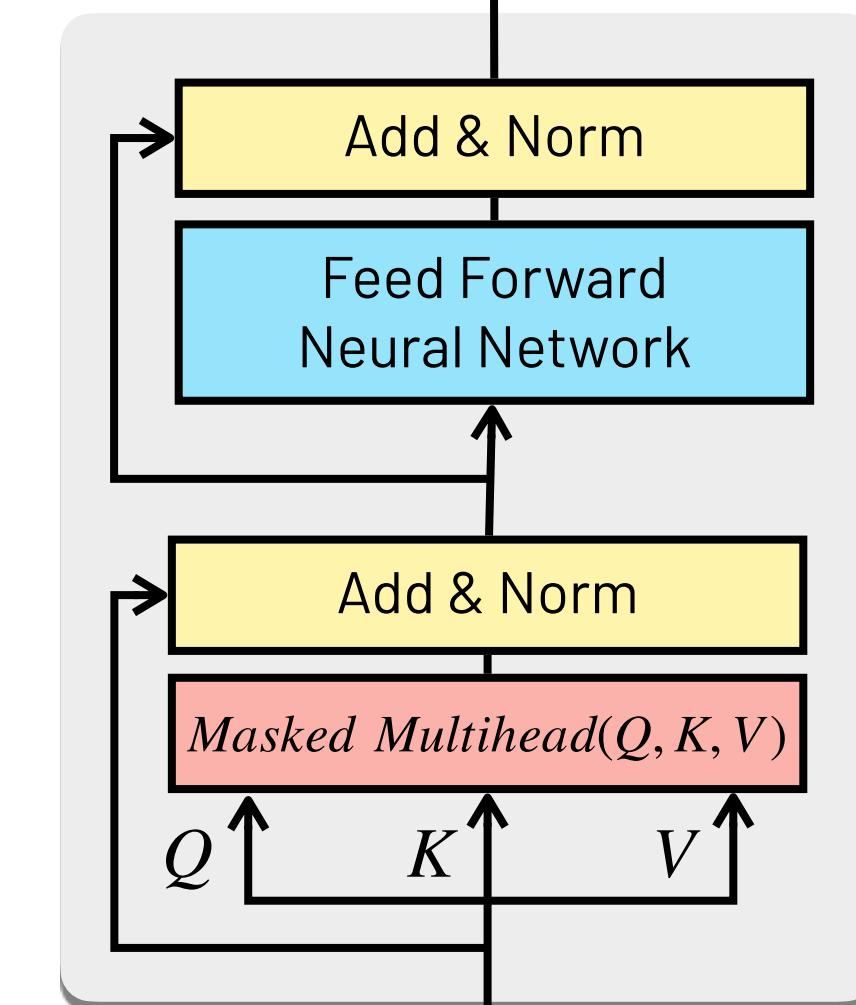
Explicit positional embedding (sinusoidal)

Represent image as a sequence of pixels or patches

Main problem:

rgb
 $(0, 0, 255)$ $\rightarrow 255^3 \approx 16M$ tokens

How to compress the color vocabulary?



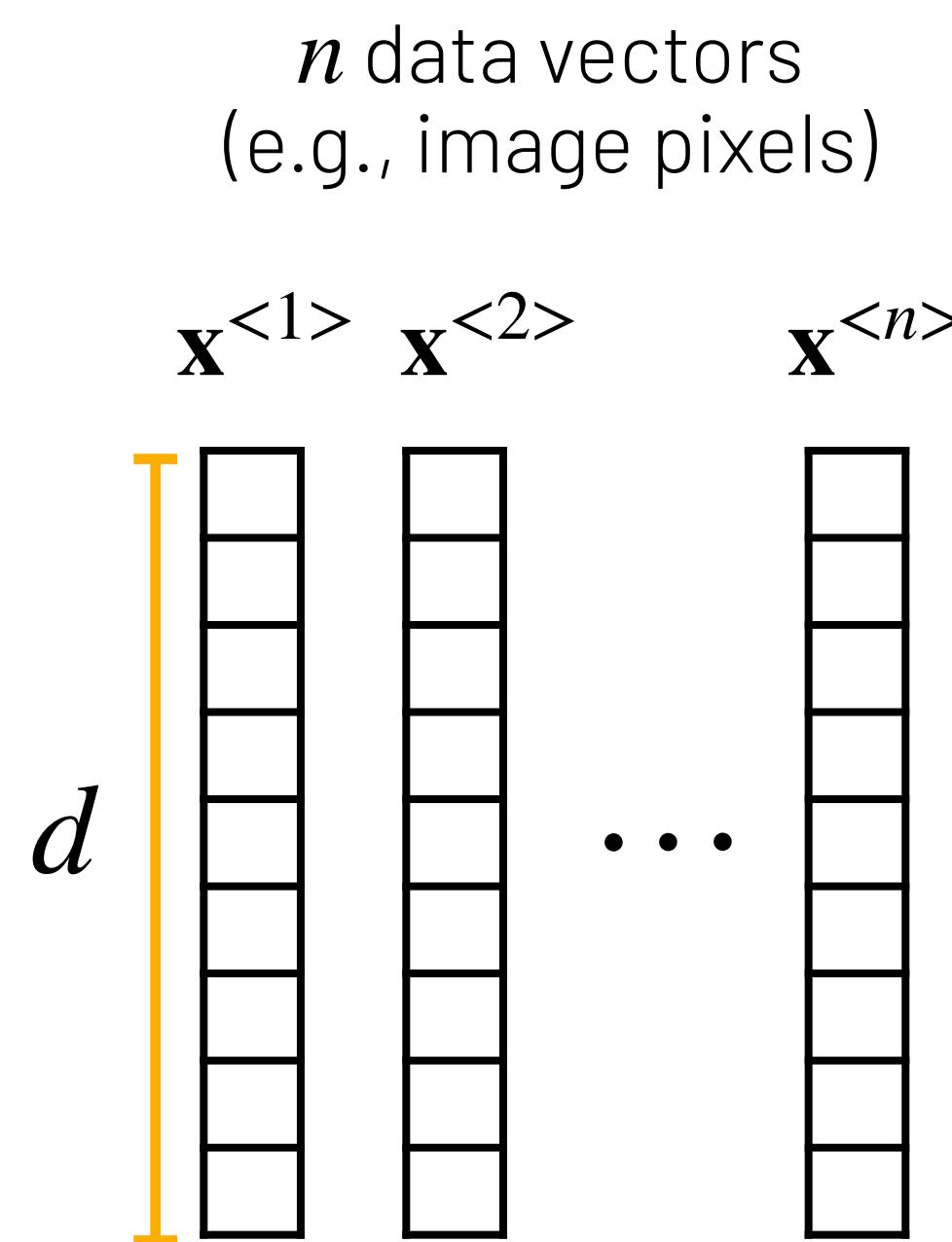
Decoder
Nx



Positional
Encoding

Vector Quantization

One way to address the problem of the high color vocabulary dimensionality is to use a technique called **Vector Quantization**, which can be viewed as a form of **data compression**.

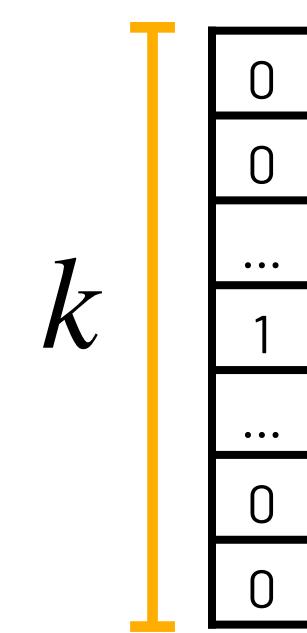


1. Approximate $\mathbf{x}^{<i>}$ by its nearest $\mathbf{c}^{<j>}$ according to some similarity metric

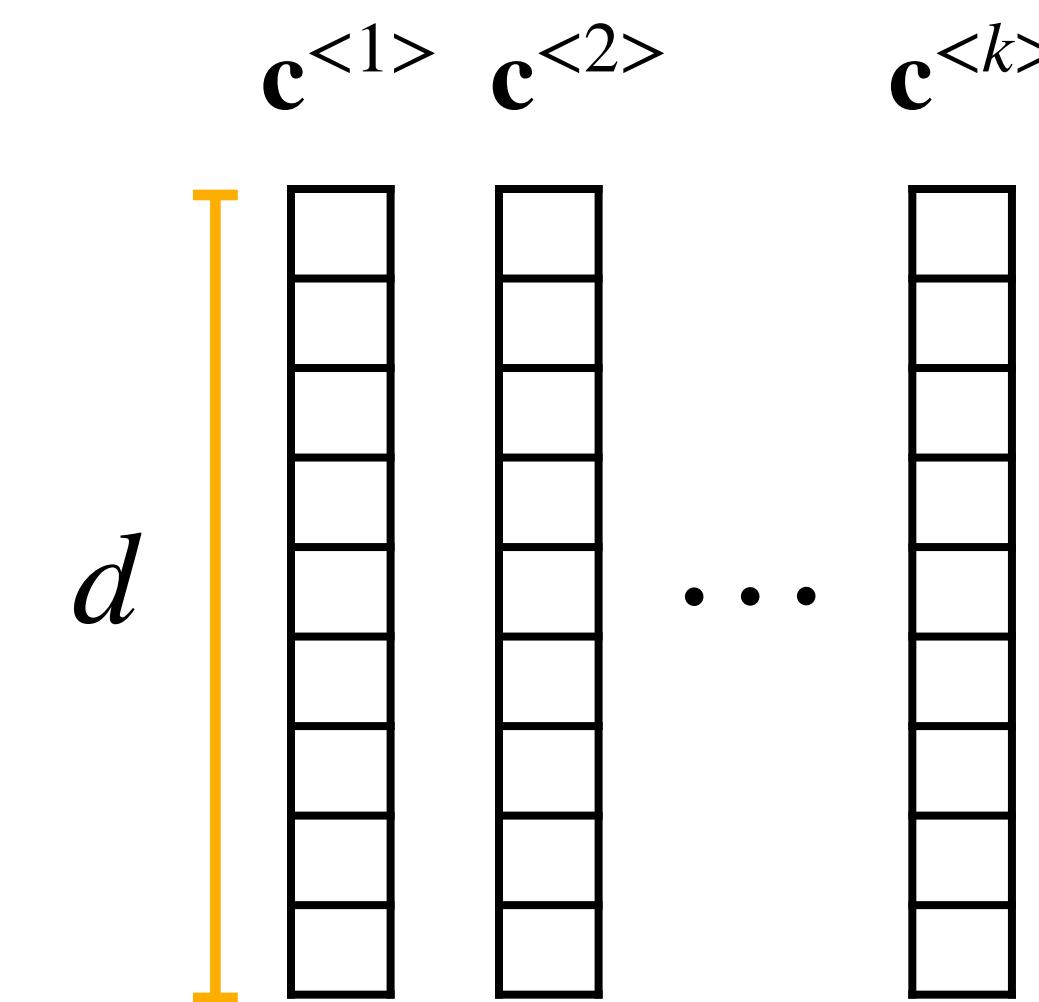
$$\mathbf{x}^{<i>} \rightarrow \operatorname{argmin}_{\mathbf{c}^{<j>}} \| \mathbf{x}^{<i>} - \mathbf{c}^{<j>} \| ^2$$

(usually Euclidian distance)

2. Represent each $\mathbf{x}^{<i>}$ as one-hot vector with k dimensions



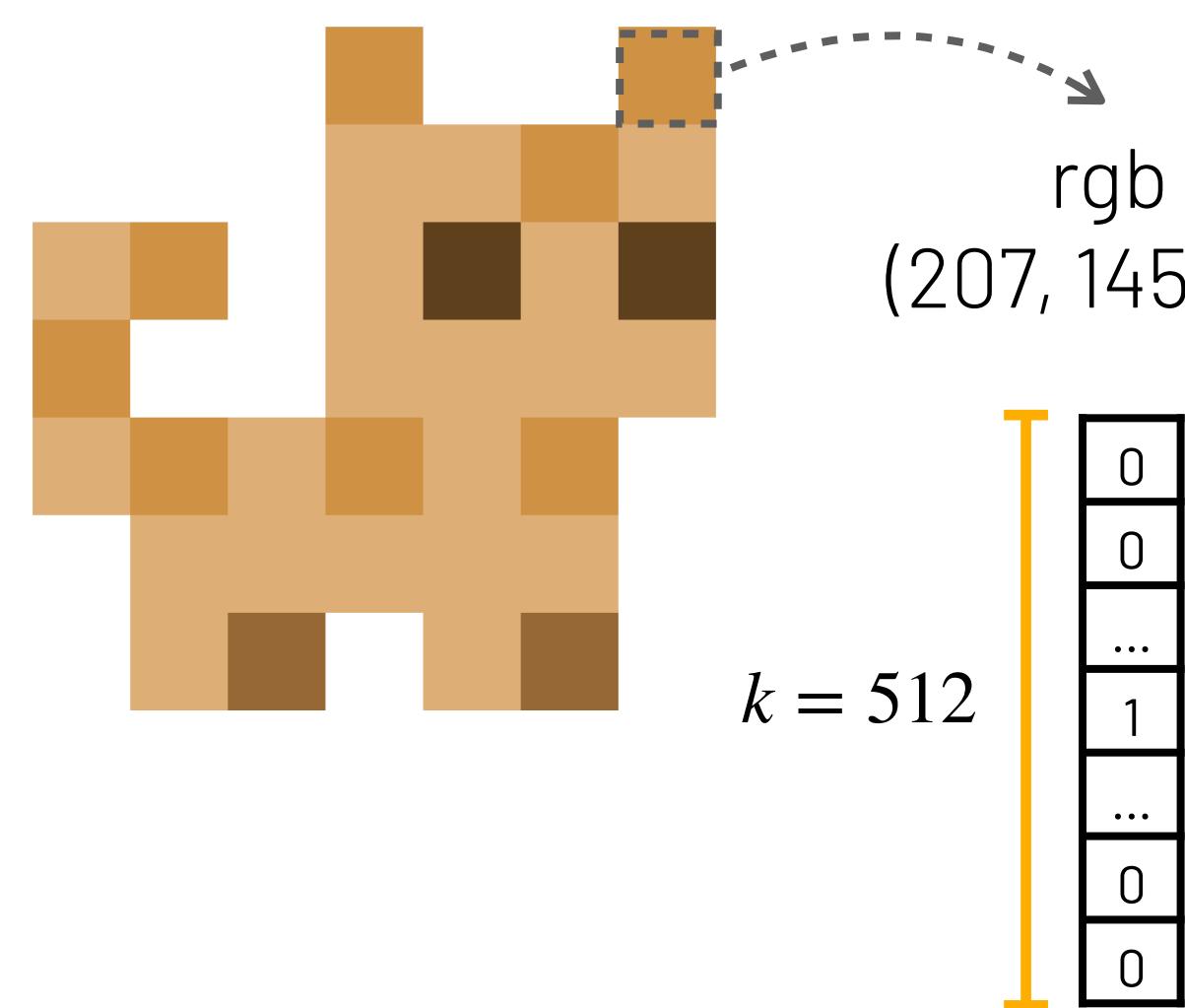
k codebook vectors $k \ll d$



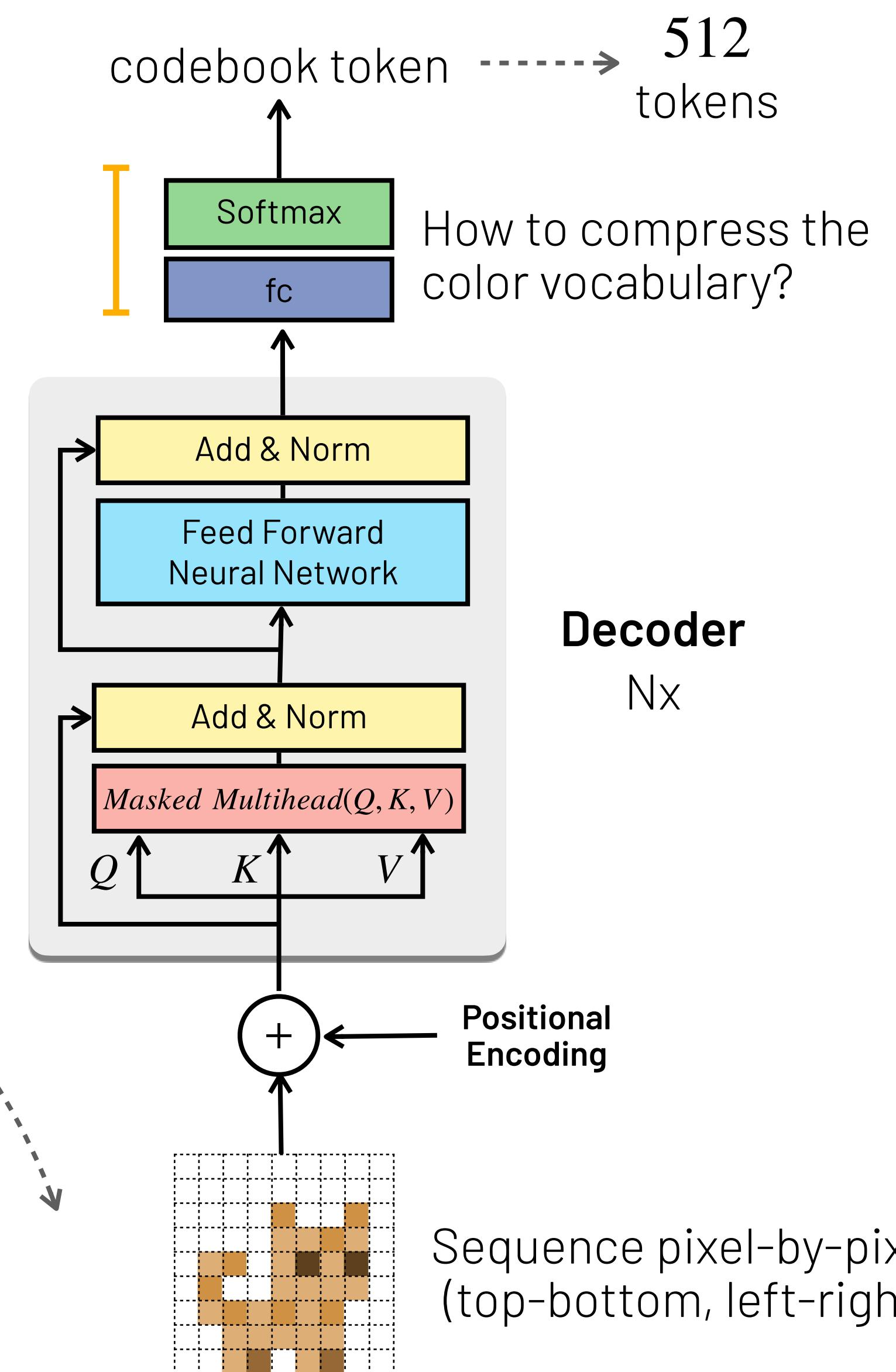
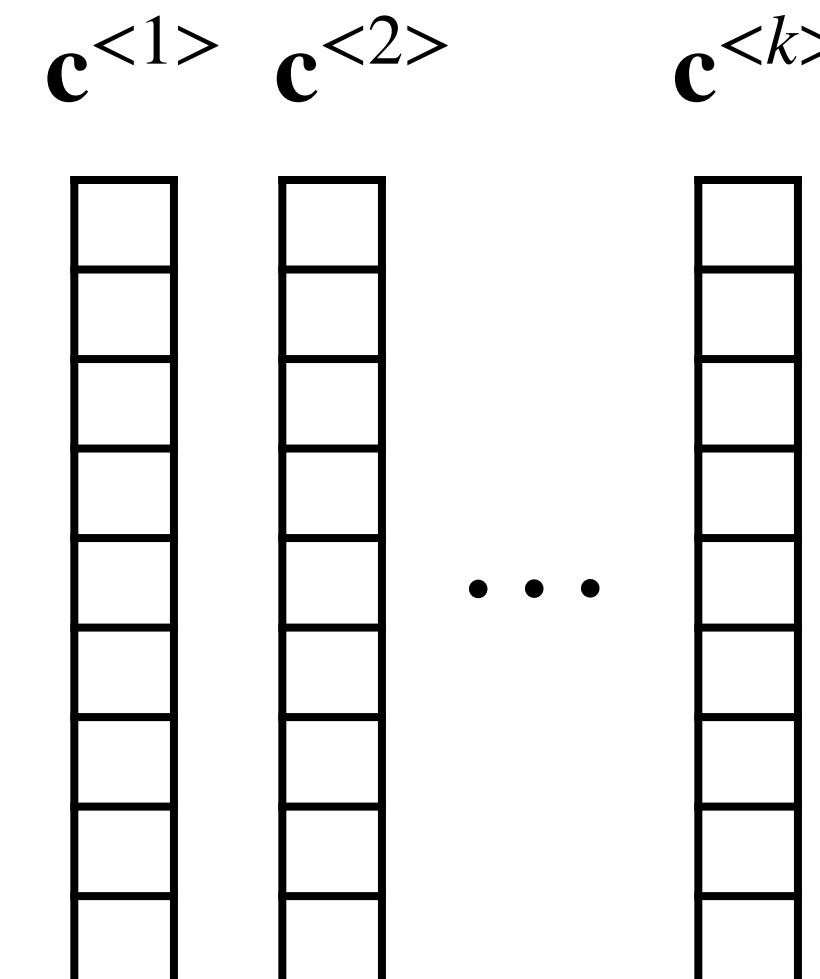
Use this one-hot vector to retrieve
the codebook vector

ImageGPT

ImageGPT was the first autoregressive transformers to generate images.

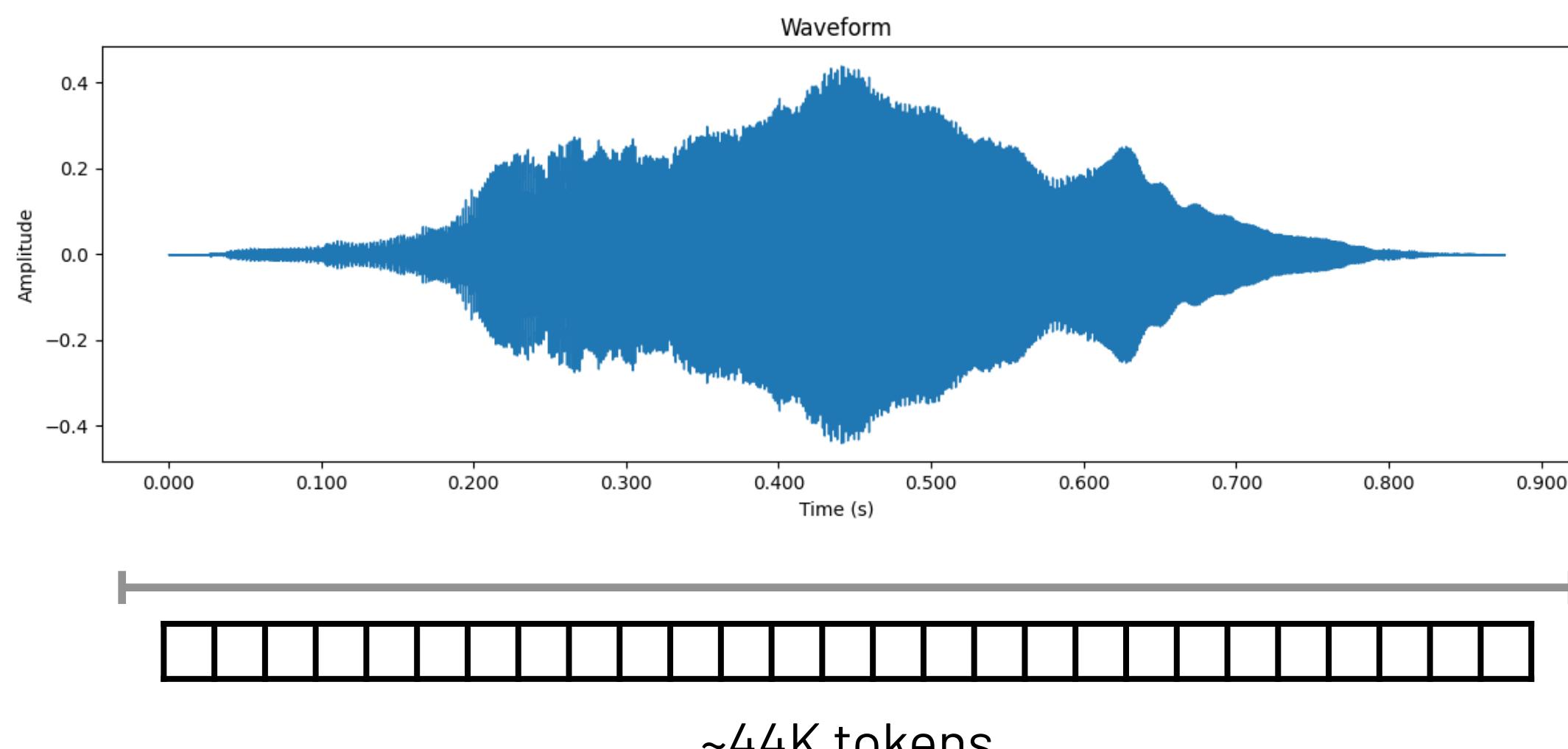


$k = 512$ codebook vectors

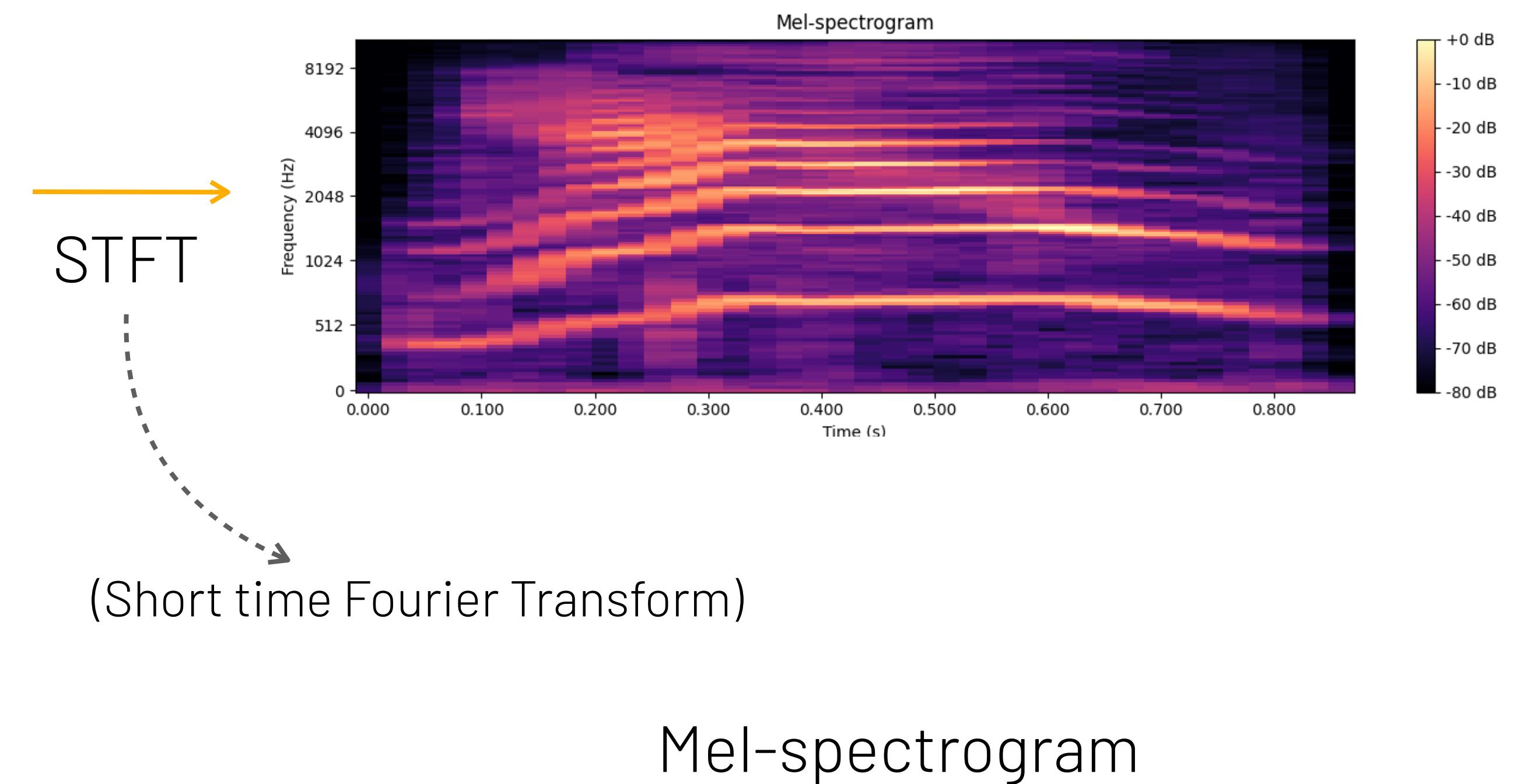


Audio Classification

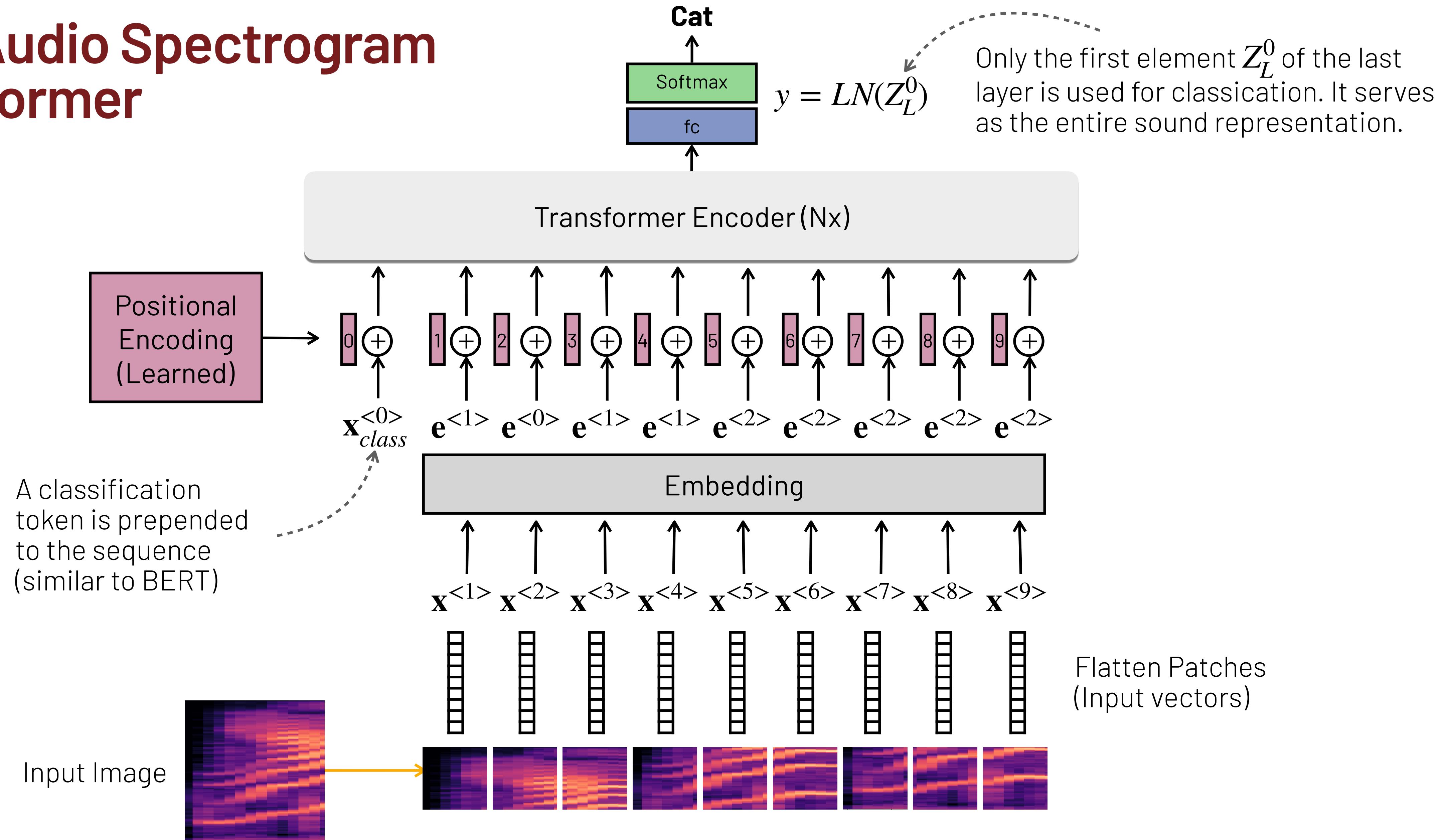
Audio signals are represented as waveforms, which are measures of amplitude of the air pressure at regular time intervals.



Audio signals are represented as waveforms, which are measures of amplitude of the air pressure at regular time intervals.



AST: Audio Spectrogram Transformer

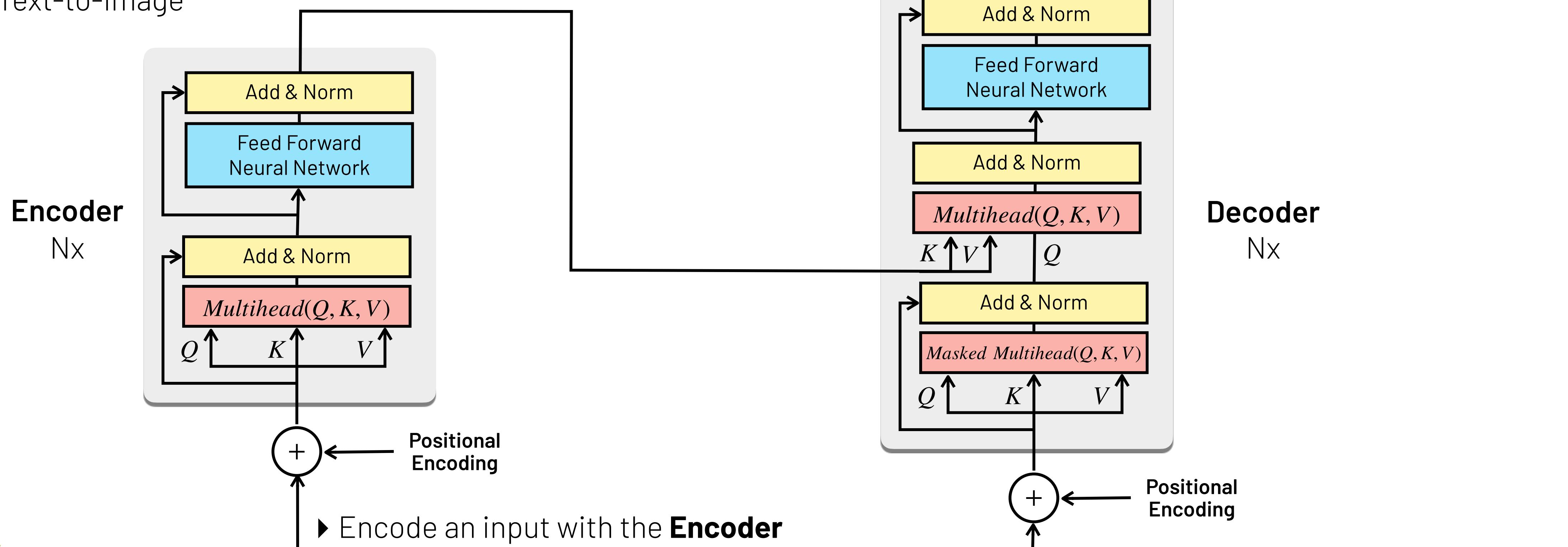


Multimodal Learning

Since transformers can process many different kinds of data, we can create multimodal models relatively straightforward with a "machine translation" framework

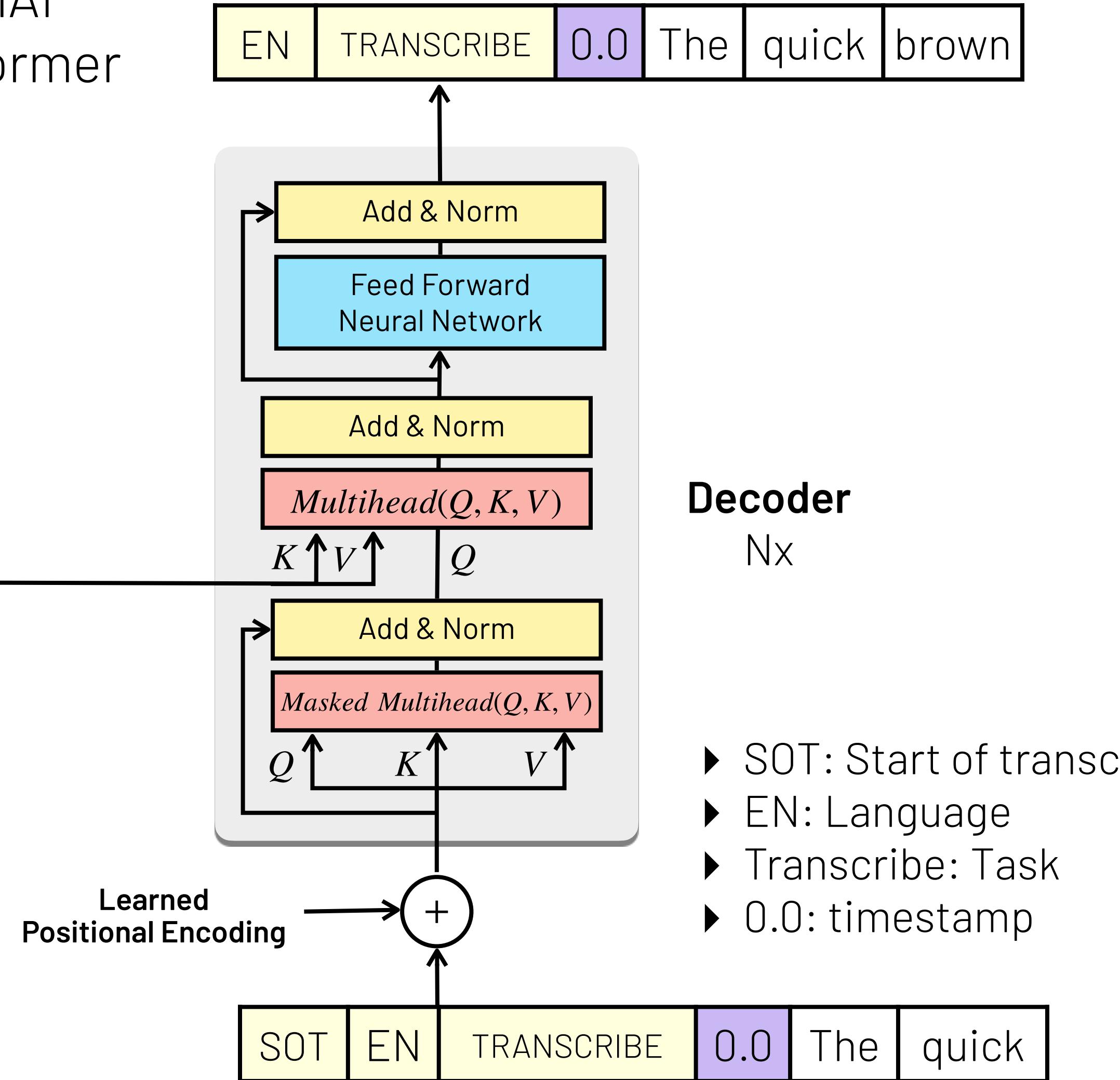
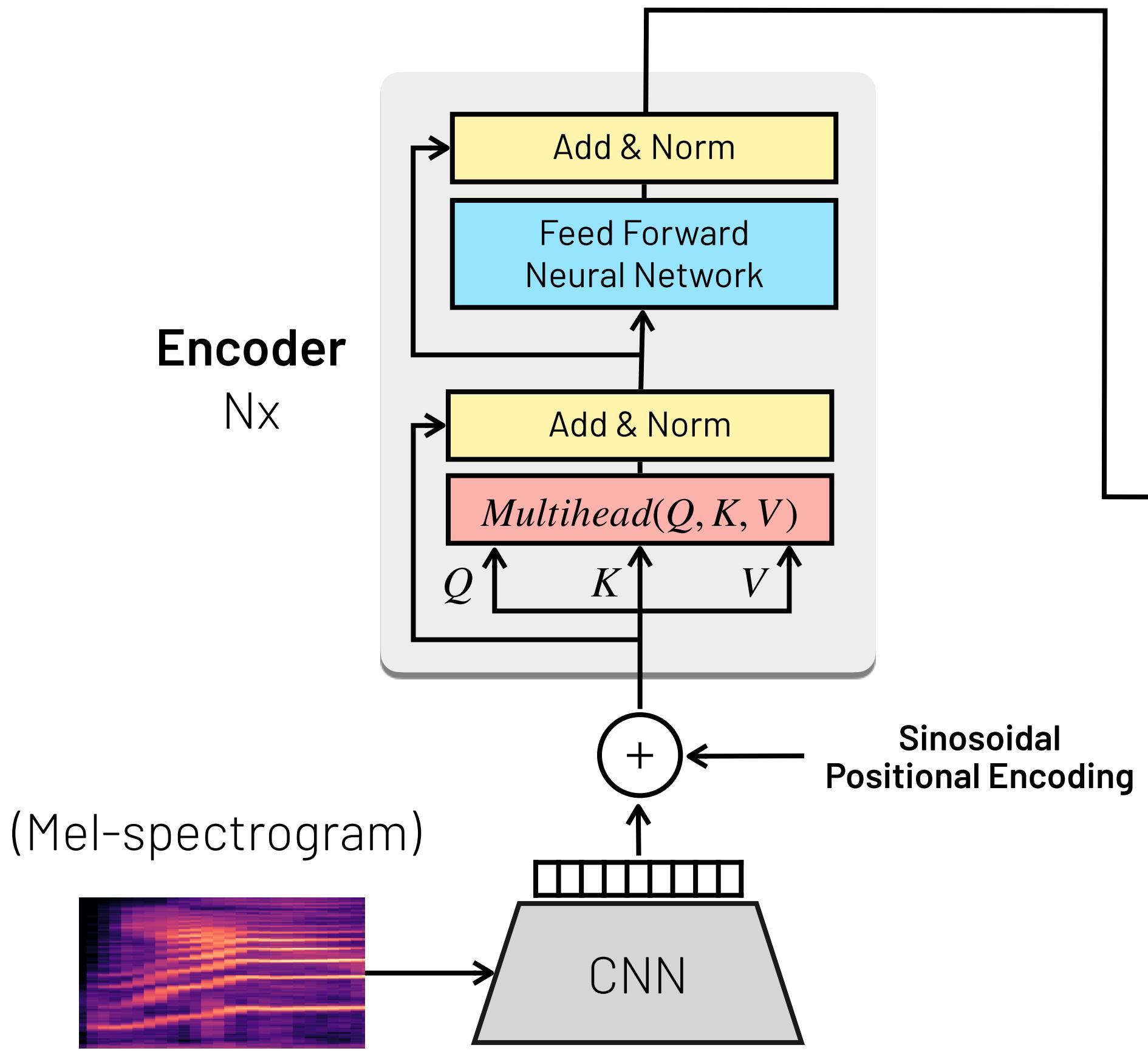
Examples:

- ▶ Audio-to-Text (Speech Recognition)
- ▶ Text-to-Image



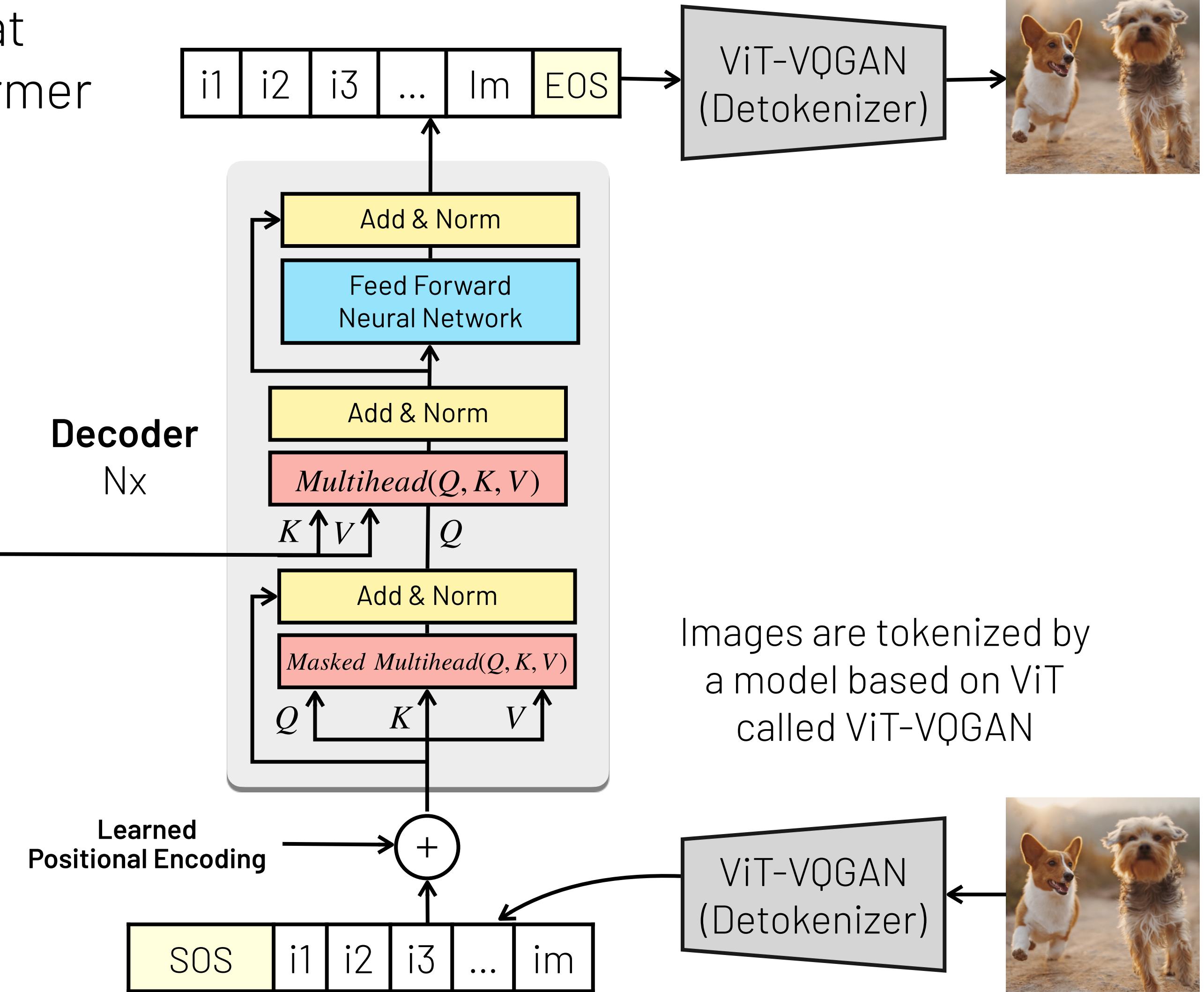
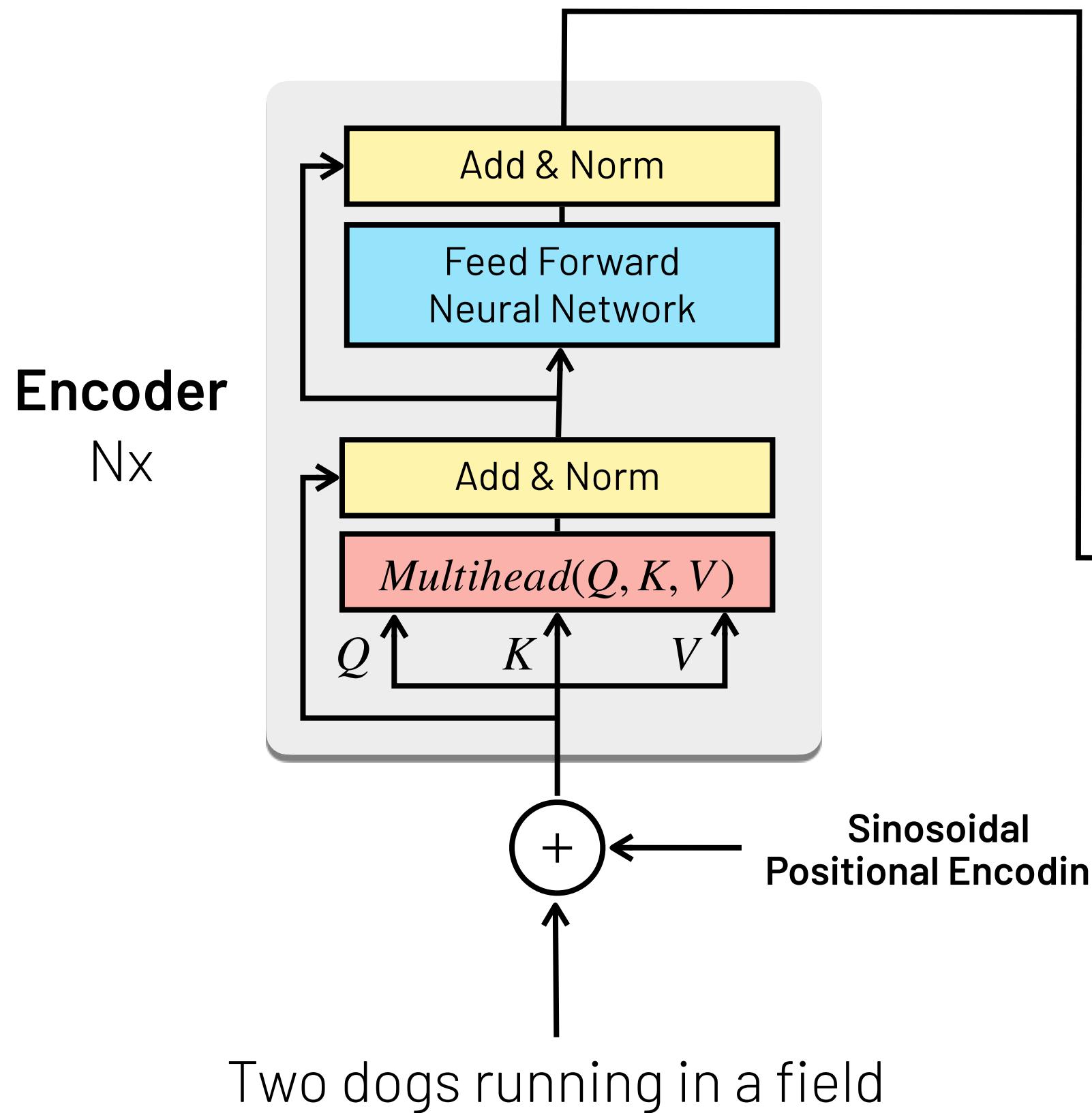
Whisper: Audio-to-Text (Speech Recognition)

Whisper is a speech recognition model by OpenAI that uses a traditional encoder-decoder transformer



Parti (Pathways Autoregressive Text-to-Image)

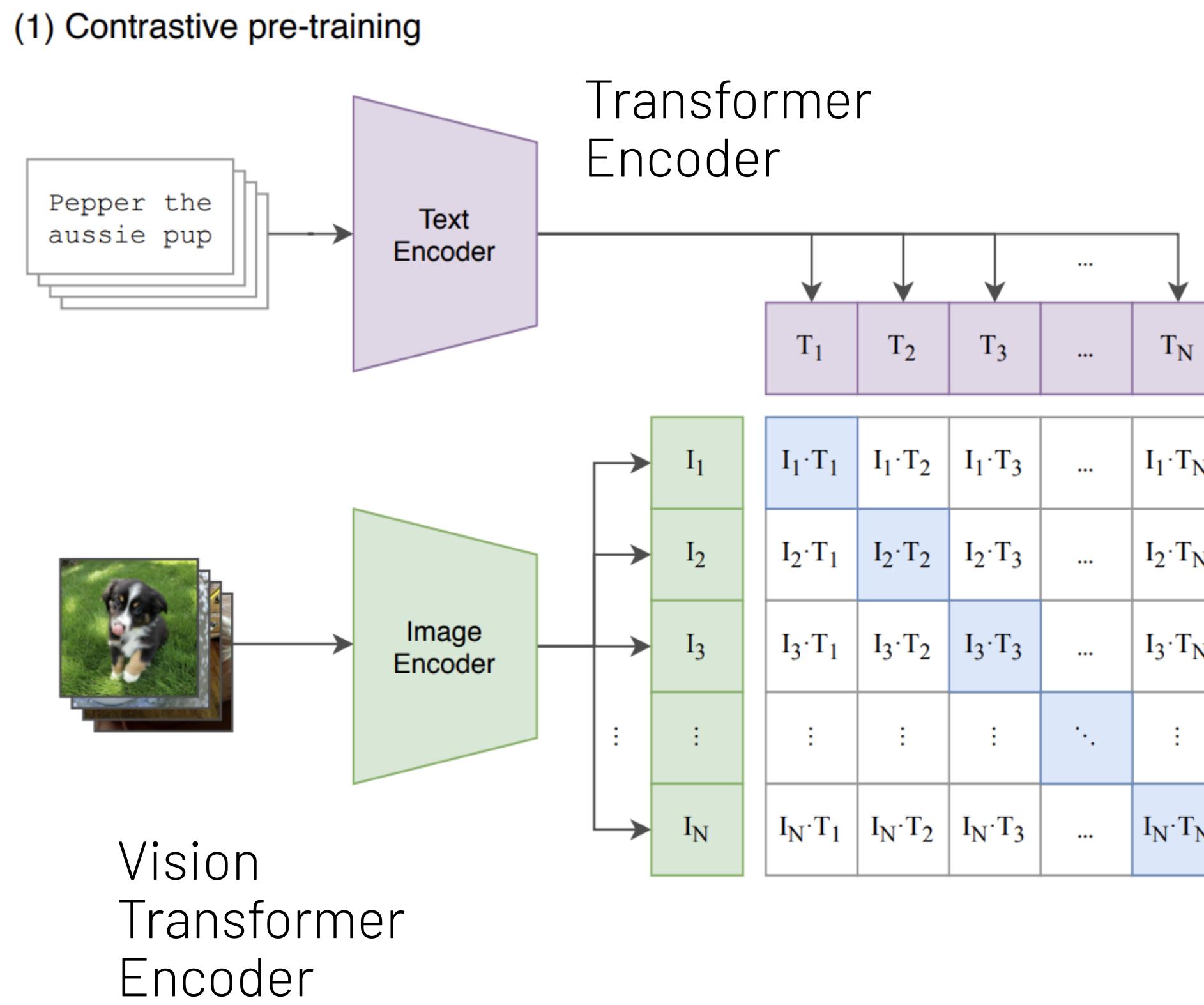
Parti is a text-to-image model by Google that uses a traditional encoder-decoder transformer



Images are tokenized by a model based on ViT called ViT-VQGAN

CLIP (Contrastive Language-Image Pre-training)

CLIP is a model developed by OpenAI to learn embeddings for Image and text in the same vector space and can be directly compared



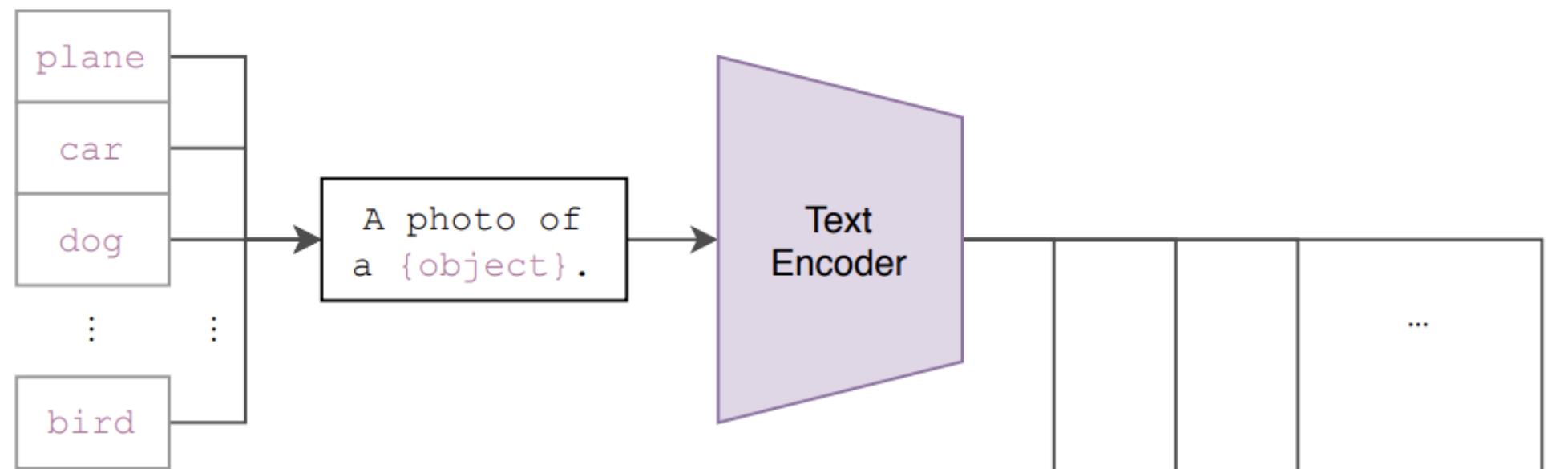
Training

1. Takes batch of (image, text) pairs
2. Passes images through vision encoder
3. Passes captions through text encoder
4. Maximizes similarity between correct image-text pairs
5. Minimizes similarity between incorrect pairs using contrastive loss

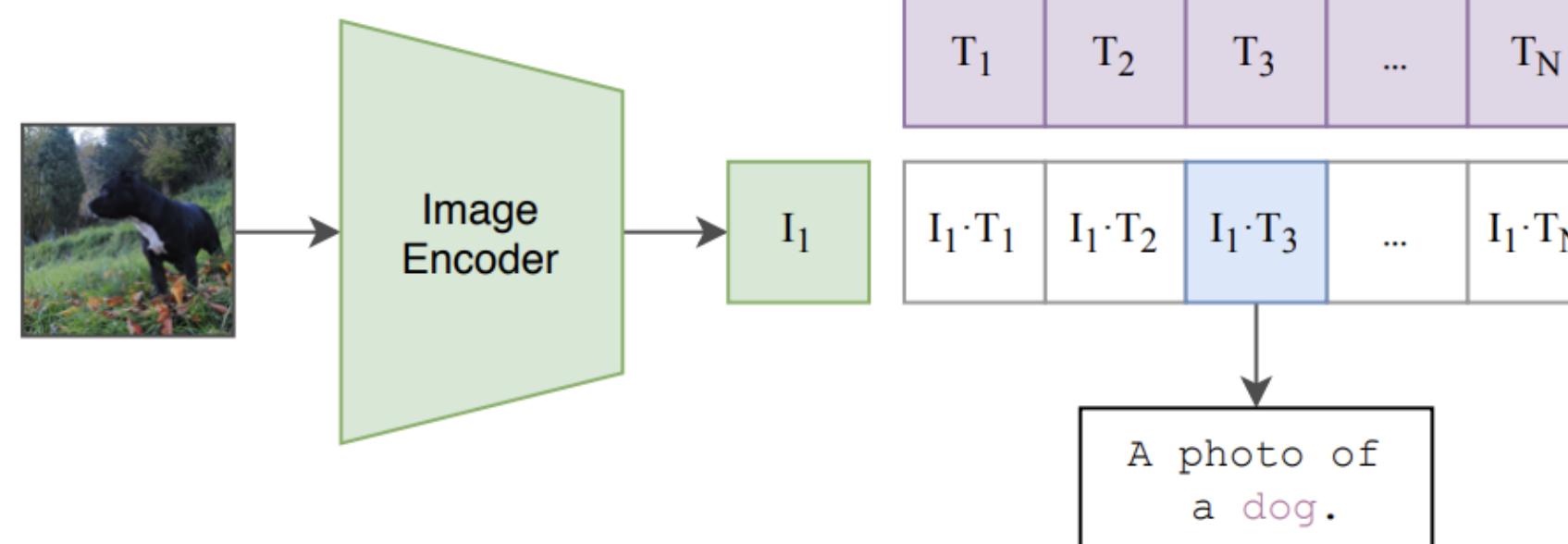
CLIP (Contrastive Language-Image Pre-training)

CLIP is a model developed by OpenAI to learn embeddings for Image and text in the same vector space and can be directly compared

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



Inference

1. Can measure similarity between any image and text embedding
2. Enables zero-shot classification by comparing image embeddings to text embeddings of class names
3. Supports open-ended image-text matching tasks

Next Lecture

L20: GANs

Generating images with Generative Adversarial Networks (GANs)