# Image Processing and Analysis

## Attention models
## Transformers

olivier.bernard@insa-lyon.fr

nicolas.thome@sorbonne-universite.fr

▶ NLP: Natural Language Processing

> Enable machines to understand, interpret, and generate natural language effectively

▶ Developed for tasks such as
- Automatic translation
- Speech comprehension
- Text generation
- Computer code generation

▶ Known algorithms
- BERT
- GPT 3, 4
- Chat-GPT
- [...]

explain quantum theory to a child in the style of snoop dogg

Yo, little homie, let me break it down for you
Quantum theory, it's a crazy thing, it's true
See, the world is made up of tiny, tiny particles
That are always moving, and changing, and actin' wild

These particles, they don't always behave like we think
Sometimes they're waves, and sometimes they're points, it's a trip
They can be in two places at once, and they can be in two states
It's a mind-blowing concept, but that's the quantum fate

So don't be scared, little homie, it's all good
Quantum theory, it's just the way the world works, it's understood
It's a wild ride, but it's a cool one too
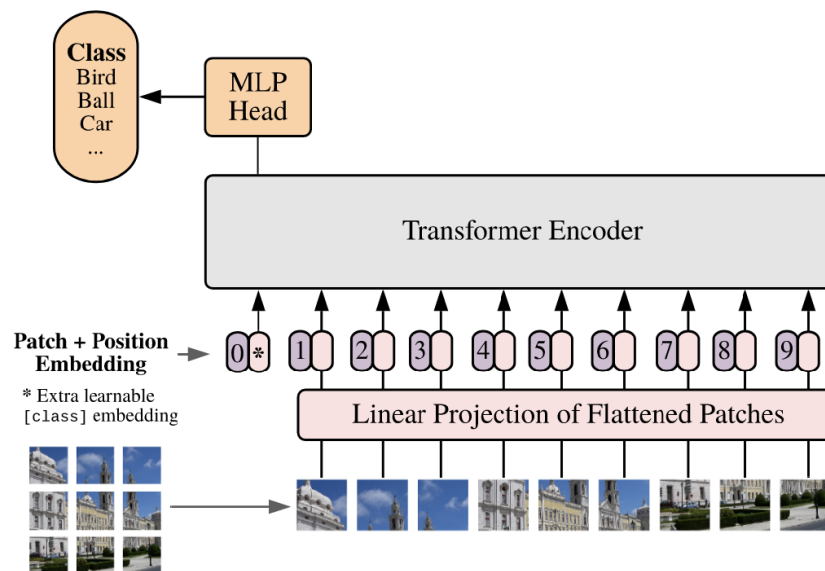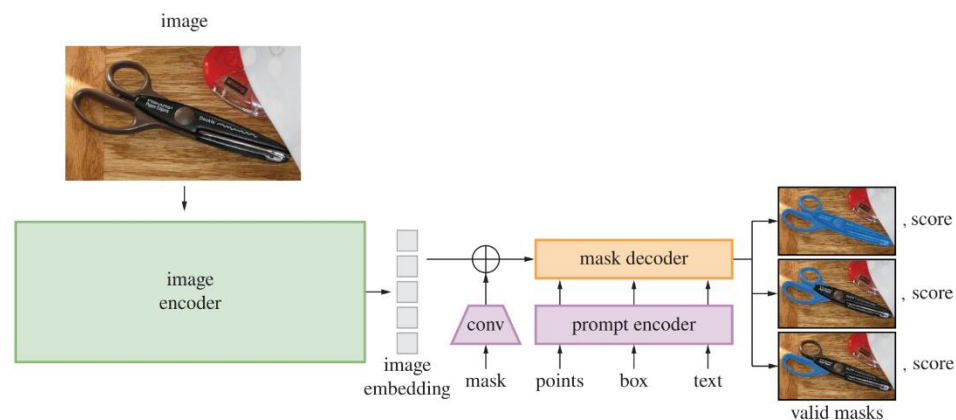Quantum theory, it's the way the world does what it do.

▶ **Also developed in image analysis**

> **Enable machines to understand, analyze, and generate images efficiently**

▶ **State-of-the-art algorithms**

- Classification
- Segmentation
- Segmentation + time
- Tracking
- Image generation

**Vision Transformer - 2020**
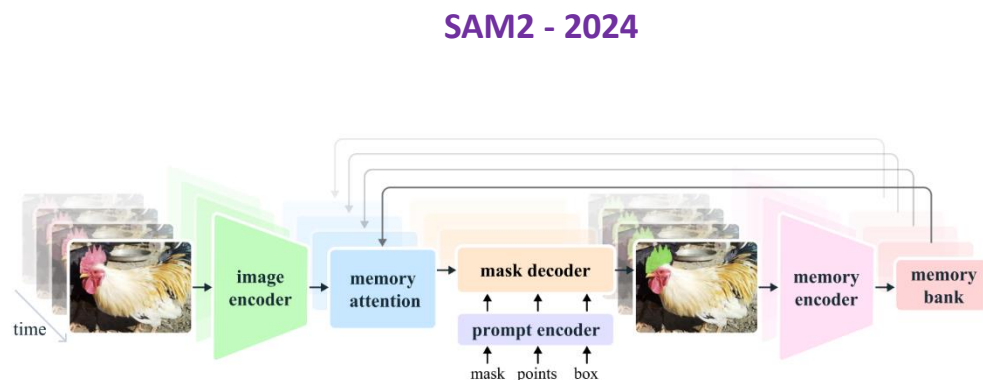
▶ **Also developed in image analysis**

> **Enable machines to understand, analyze, and generate images efficiently**

▶ **State-of-the-art algorithms**

**SAM (Segment Anything) - 2024**

- Classification
- Segmentation
- Segmentation + time
- Tracking
- Image generation

▶ **Also developed in image analysis**

> Enable machines to understand, analyze, and generate images efficiently

▶ **State-of-the-art algorithms**

**SAM2 - 2024**

- Classification
- Segmentation
- **Segmentation + time**
- Tracking
- Image generation

▶ **Also developed in image analysis**

> **Enable machines to understand, analyze, and generate images efficiently**

▶ **State-of-the-art algorithms**

- Classification
- Segmentation
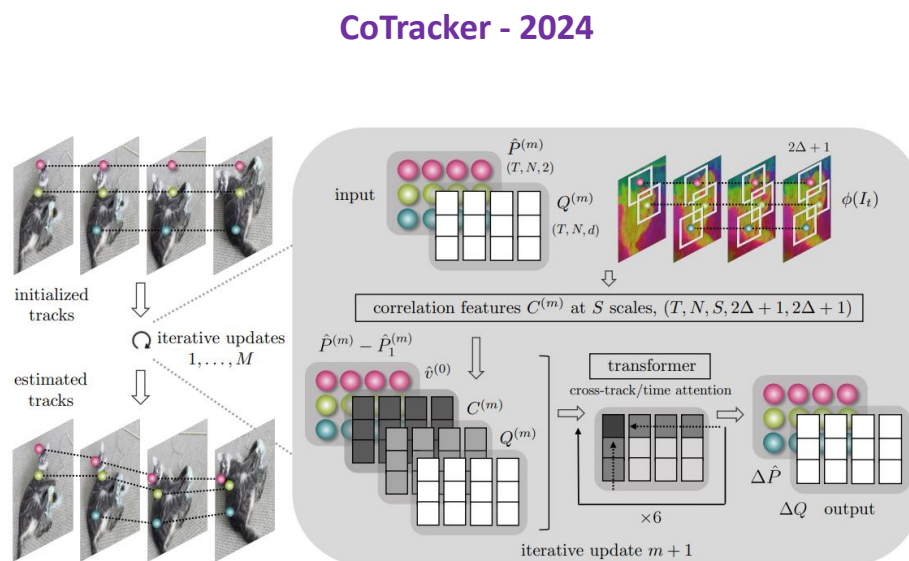- Segmentation + time
- Tracking
- Image generation

**CoTracker - 2024**

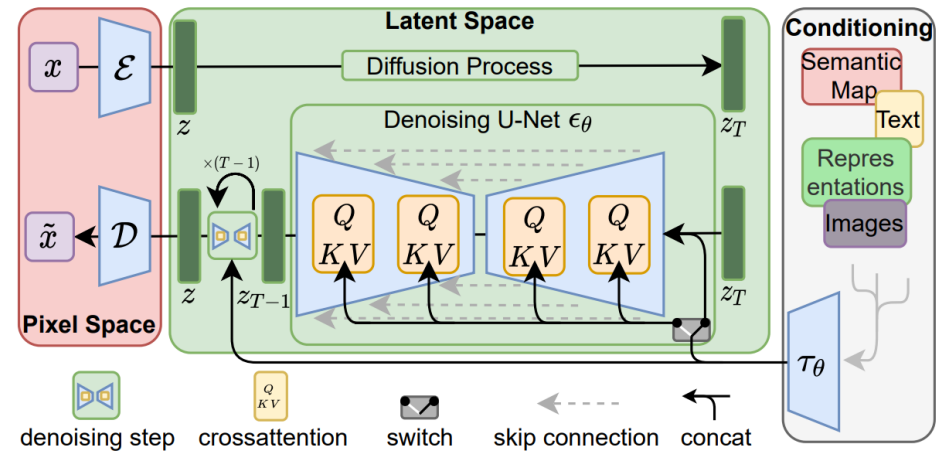# What are the applications of attention models?

▶ **Also developed in image analysis**

> **Enable machines to understand, analyze, and generate images efficiently**

▶ **State-of-the-art algorithms**

- Classification
- Segmentation
- Segmentation + time
- Tracking
- Image generation

**Latent diffusion model - 2023**

# Transformers

*Tokens*

► The input data is structured as *tokens*

➔ Text: token = word of a sentence
➔ Image: token = patch of an image

**Text**

« Hello, I am olivier »

Tokenization procedure

« Hello », « I », « am », « olivier »

« Hello » $\Rightarrow x_i \in \mathbb{R}^t$

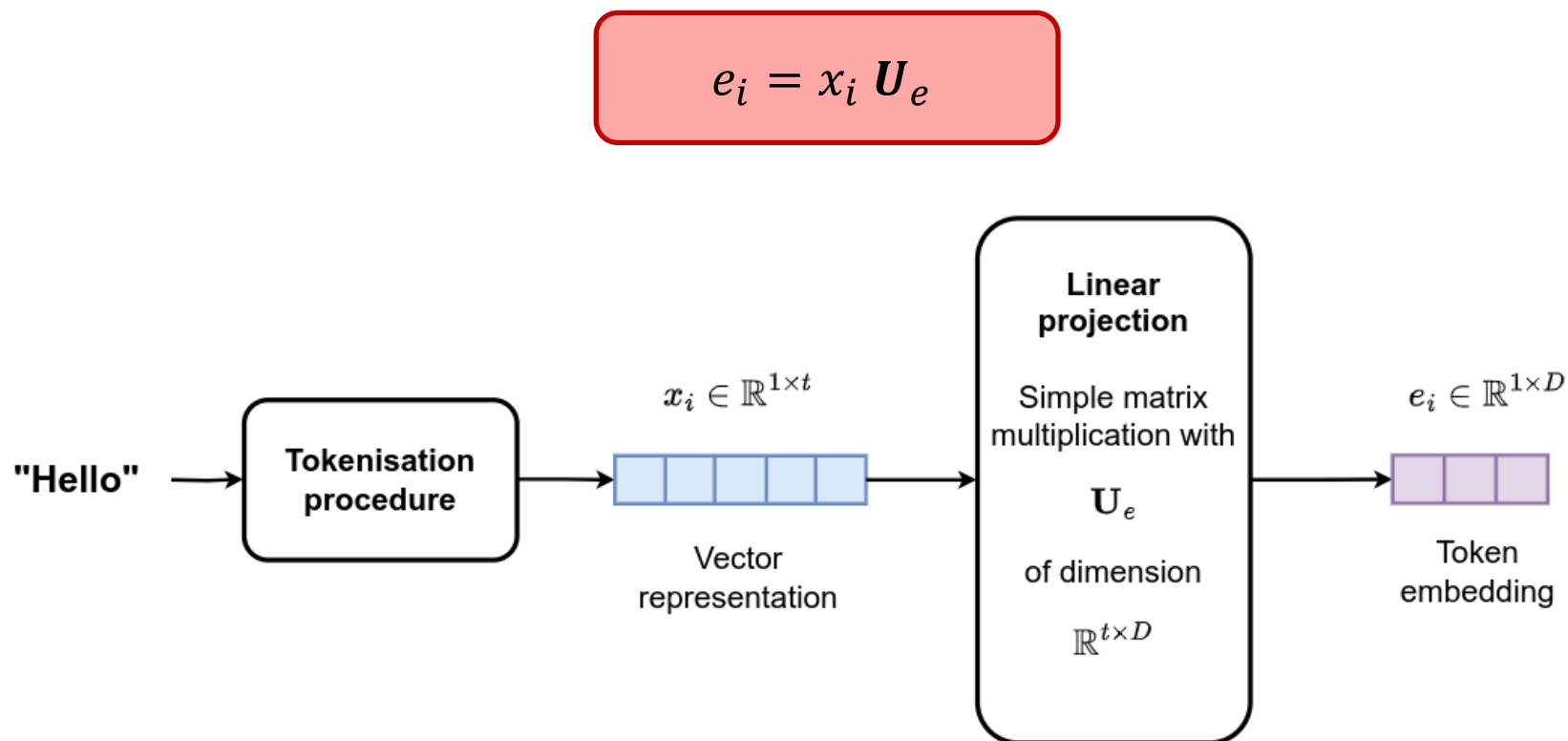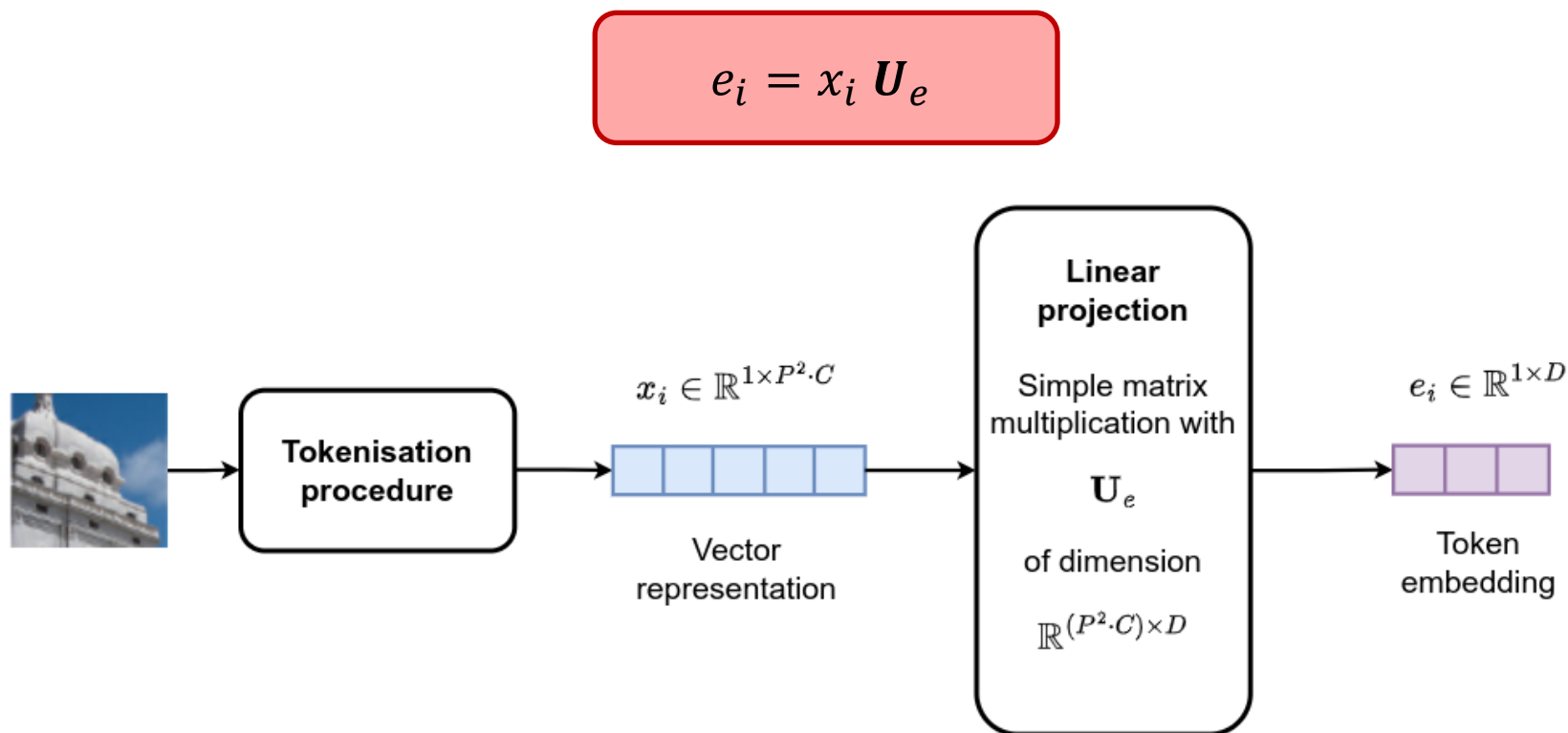| 0 | 1 | 0 | 0 | 0 |
|---|---|---|---|---|

**Image**



Tokenization procedure



$\Rightarrow x_i \in \mathbb{R}^{P^2 C}$

$P$ patch width, $C = 3$ if color image

▶ Creation of a representation (or *embedding*) of the tokens

➔ Simple linear projection

➔ Multiplication by a learnable representation matrix $\boldsymbol{U}_e$

$$e_i = x_i \, \boldsymbol{U}_e$$

"Hello" ⟶ **Tokenisation procedure** ⟶ $x_i \in \mathbb{R}^{1 \times t}$

Vector representation

⟶ **Linear projection**

Simple matrix multiplication with

$\mathbf{U}_e$

of dimension

$\mathbb{R}^{t \times D}$

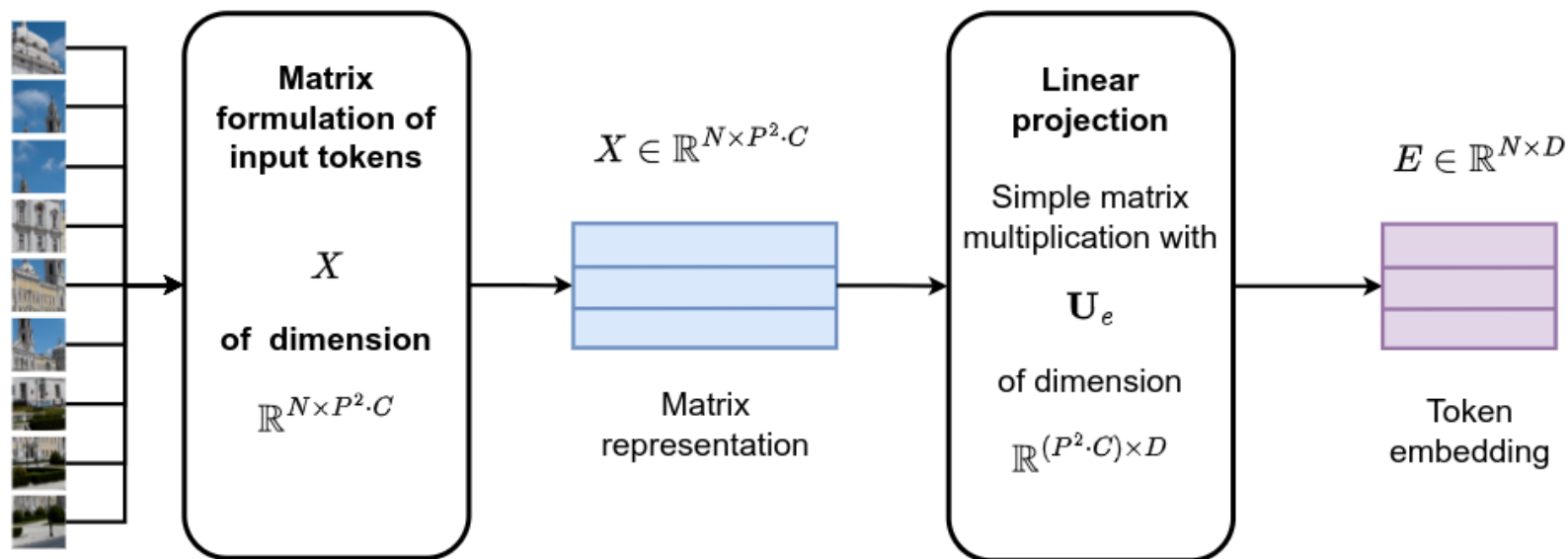⟶ $e_i \in \mathbb{R}^{1 \times D}$

Token embedding

# Tokens generation

► Creation of a representation (or *embedding*) of the tokens

➔ Simple linear projection

➔ Multiplication by a learnable representation matrix $\boldsymbol{U}_e$

$$e_i = x_i \, \boldsymbol{U}_e$$



**Tokenisation procedure**

$x_i \in \mathbb{R}^{1 \times P^2 \cdot C}$

Vector representation

**Linear projection**

Simple matrix multiplication with

$\mathbf{U}_e$

of dimension

$\mathbb{R}^{(P^2 \cdot C) \times D}$

$e_i \in \mathbb{R}^{1 \times D}$

Token embedding

▶ Learning a representation matrix $\boldsymbol{U}_e$ shared by each token

➔ Matrix formulation

$$E = X\,\boldsymbol{U}_e$$



**Matrix formulation of input tokens**

$X$

**of dimension**

$\mathbb{R}^{N \times P^2 \cdot C}$

$X \in \mathbb{R}^{N \times P^2 \cdot C}$

Matrix representation

**Linear projection**

Simple matrix multiplication with

$\mathbf{U}_e$

of dimension

$\mathbb{R}^{(P^2 \cdot C) \times D}$

$E \in \mathbb{R}^{N \times D}$

Token embedding

▶ **Positional encoding**

➔ Sentence / image: a set of independent tokens

➔ Loss of structural information from the input data

▶ **Recovery of structure: positional encoding (PE: positional embedding)**

➔ Correspondence between the position of token $t$ and a vector $p_t \in \mathbb{R}^{1 \times D}$

➔ Classical encoding: sinusoidal function

$$p_t \in \mathbb{R}^{1 \times D}$$

$$p_t = \left[ \sin(\omega_1 t), \cos(\omega_1 t), \cdots, \sin(\omega_{D/2} t), \cos(\omega_{D/2} t) \right]$$

$$\omega_k = \frac{1}{10000^{2k/D}}$$

▶ Sinusoidal positional encoding

➔ Unique vector $p_t$ for each position $t$

➔ $p_t(i) \in [-1,1]$ : Intrinsic normalization of values



Number of tokens = 50, dimension D of each token = 128

▶ Sinusoidal positional encoding

▶ Intrinsic modeling of the relative position of tokens

▶ Position similarity matrix: $K = P \cdot P^t$



Position

► **Final representation**

➜ Final tokens = sum of token and position representations
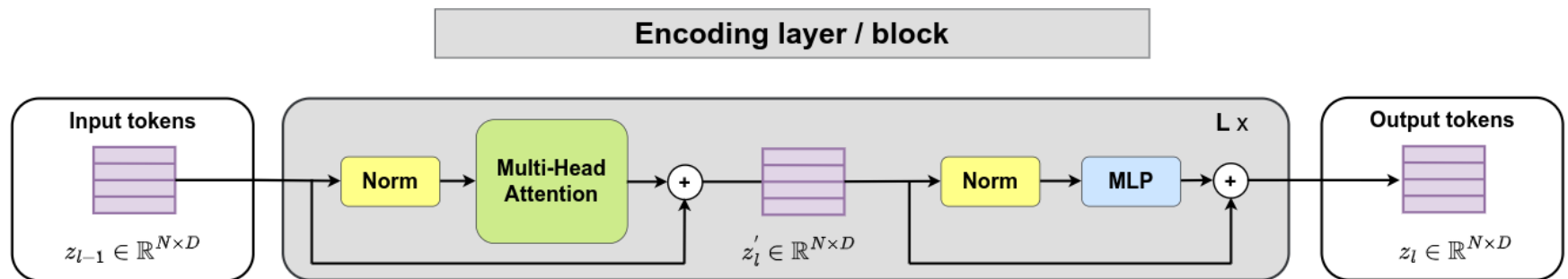
➜ Only the matrix $U_e$ is to be learned for this phase

# Transformers

*Encoding blocks*

► Encoder

➜ Corresponds to N encoding blocks

- Input:       A token representation
- Output:     A new token representation tailored to the target is being optimized



**Encoding layer / block**

Input tokens     $z_{l-1} \in \mathbb{R}^{N \times D}$     Norm → Multi-Head Attention → +     $z'_l \in \mathbb{R}^{N \times D}$     Norm → MLP → +     L x     Output tokens     $z_l \in \mathbb{R}^{N \times D}$

# Transformer: information encoding



**1st step**

➔ Computation of attention maps between tokens

➔ Residual connection   1) Against vanishing gradient

2) Do not forget the positional representation
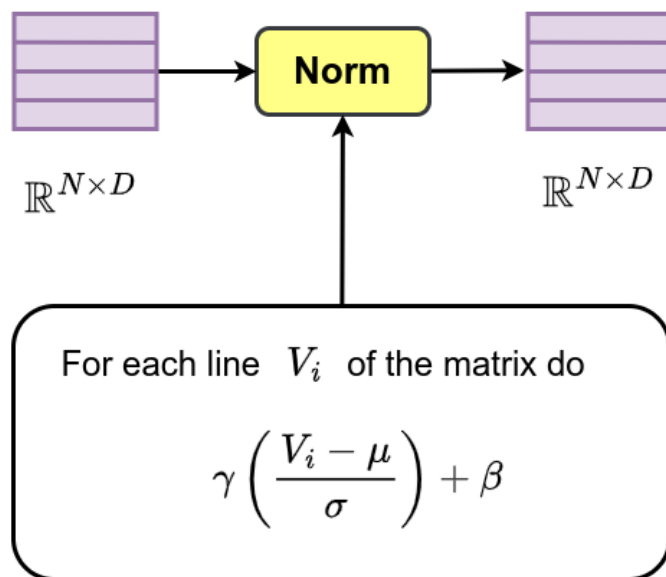
$$z_l' = MHA\big(LN(z_{l-1})\big) + z_{l-1}$$

# Transformer: information encoding



1st step                    2nd step

▶ 2nd step

➔ Introduction of nonlinearities to generate relevant information
➔ Residual connection    1) Against vanishing gradient
                          2) Do not forget the positional representation

$$z_l = MLP\big(LN(z_l')\big) + z_l'$$

► Focus on the 2nd step



**Encoding layer / block**

L x

Norm → MLP → +

$z'_l \in \mathbb{R}^{N \times D}$

**Output tokens**

$z_l \in \mathbb{R}^{N \times D}$

2nd step

▶ **Normalization**

➔ Controls the dynamics of the token values before each key step
- $\mu, \sigma$:     computed over all the tokens corresponding to an image
- $\gamma, \beta$:     parameters to be learned

$$\mathbb{R}^{N \times D} \qquad \boxed{\text{Norm}} \qquad \mathbb{R}^{N \times D}$$

For each line $V_i$ of the matrix do

$$\gamma \left( \frac{V_i - \mu}{\sigma} \right) + \beta$$

▶ MLP

➔ Introduces non-linearity

➔ Enables the generation of relevant information

$$z_l^* = LN(z_l')$$

$$MLP(z_l^*) = max(0, z_l^* W_1 + b_1) W_2 + b_2$$



$\mathbb{R}^{N \times D}$    MLP    $\mathbb{R}^{N \times D}$

For each line $V_i$ of the matrix apply the same fully connected layers



Linear layer    ReLU    Linear layer

$W_1, b_1$      $W_2, b_2$

► Focus on the 1$^{st}$ step



**Encoding layer / block**

Input tokens

$z_{l-1} \in \mathbb{R}^{N \times D}$

Norm

Multi-Head Attention

+

$z_l' \in \mathbb{R}^{N \times D}$

1$^{st}$ step

▶ **Self-attention module**

➔ Management of attention maps $A$ through Q (query), K (key), V (values) matrices

➔ Softmax applied row-wise to the matrix $A$ to normalize the weights that will weight the row vectors of $V$
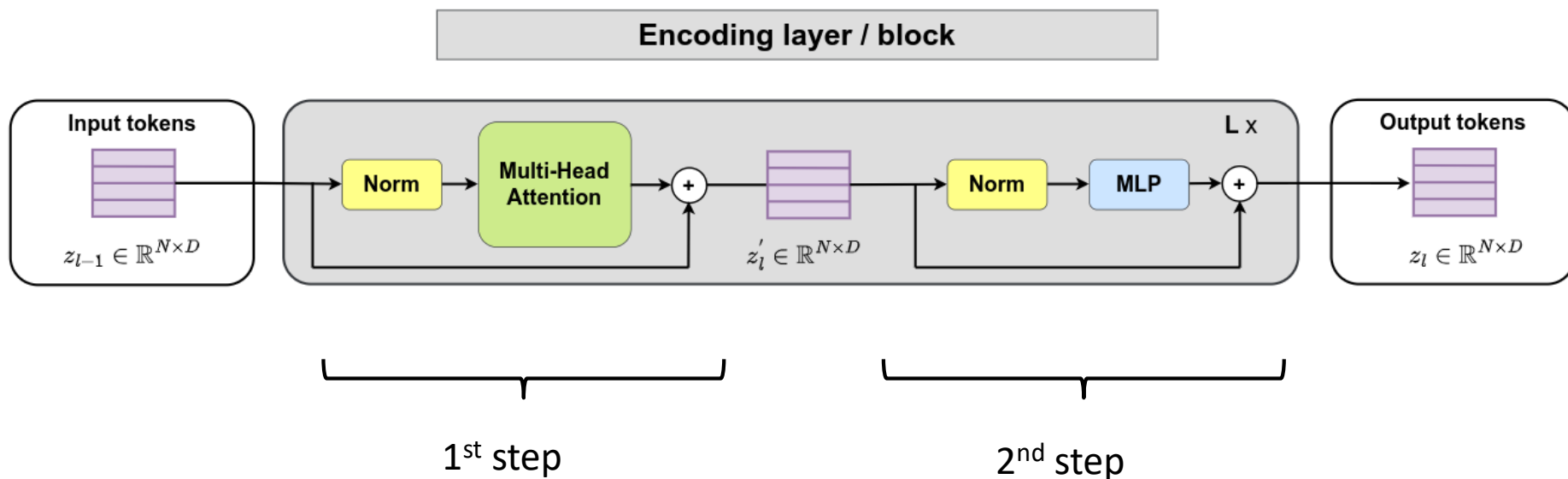
► Multi-Head Attention: Multi-head attention block

➔ Generation of $k$ heads from different self-attention modules

➔ Equivalent to the concept of feature maps in *CNNs*

➔ Linear projection to mix the information from different heads and return to the initial token dimensions

▶ In summary

➔ 1st step: generation of information through attention between tokens

➔ 2nd step: generation of relevant information through non-linearity

# Transformers

*Classification method*

▶ **VIT: reference algorithm**

➔ Trained on JFT (300 million images)

➔ Introduction of the concept of *class token*

> Learning a "pooling" operation with respect to visual tokens



**Vision Transformer (ViT)**

gif animation

► VIT: reference algorithm



Vision Transformer (ViT) · Transformer Encoder

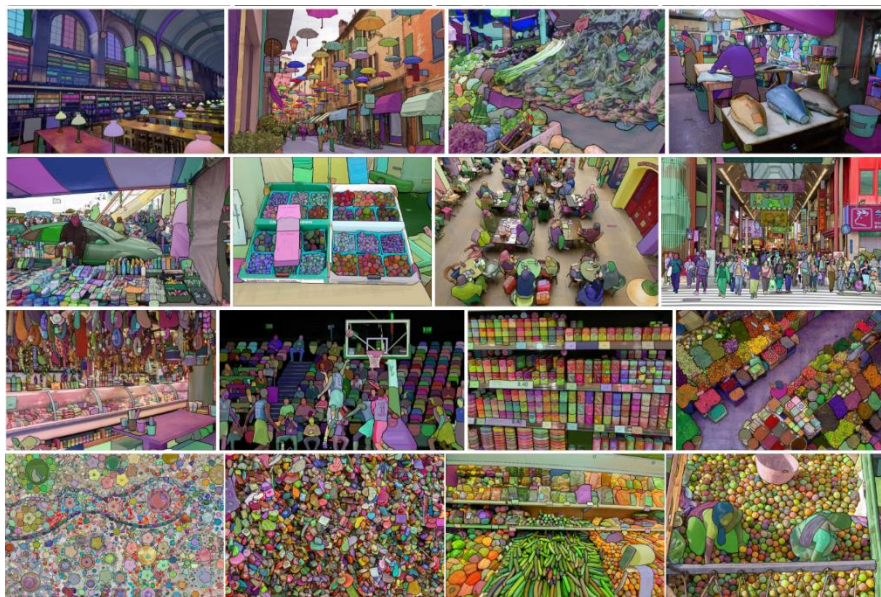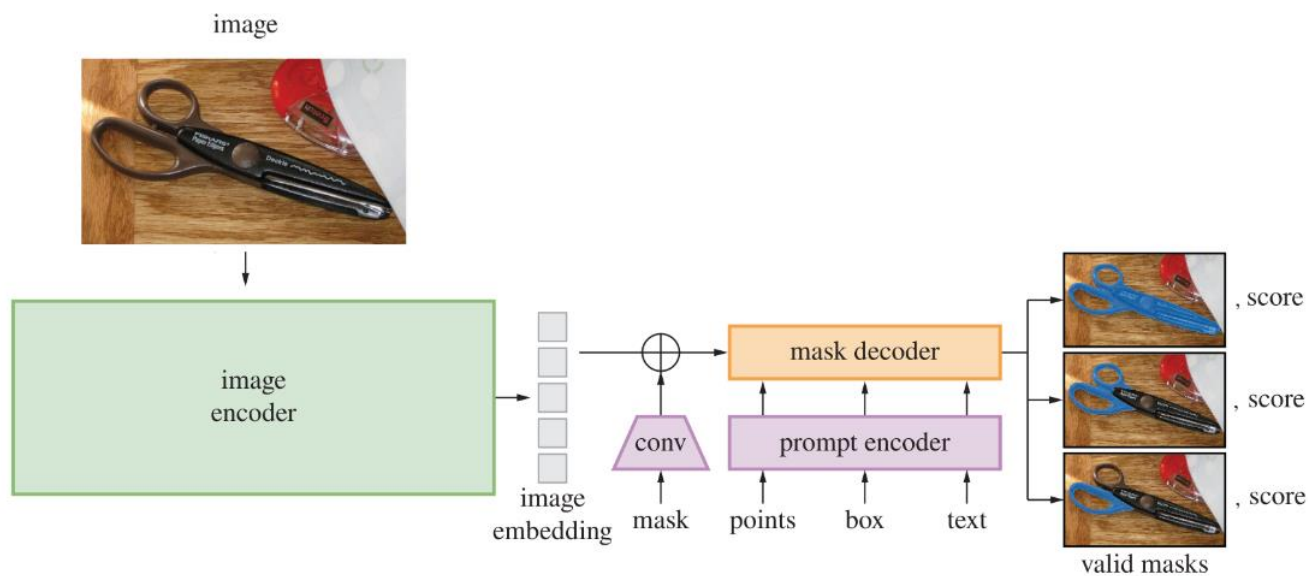| Models | Layers / blocks | Hidden size D | MLP size | Heads | Parameters |
|--------|-----------------|---------------|----------|-------|------------|
| ViT-Base | 12 | 768 | 3072 | 12 | 86 M |
| ViT-Large | 24 | 1024 | 4096 | 16 | 307 M |
| ViT-Huge | 32 | 1280 | 5120 | 16 | 632 M |

# Transformers

*Segmentation model*

► SAM: Segment Anything

➔ Trained on SA-1B (11 million images, 1 billion masks)

➔ 2D images of natural scenes, Minimum side size: 1500px

▶ SAM: Segment Anything

➜ Relatively simple architecture

➜ Interactive segmentation process

➜ Integrates the principle of segmentation ambiguity based on initialization
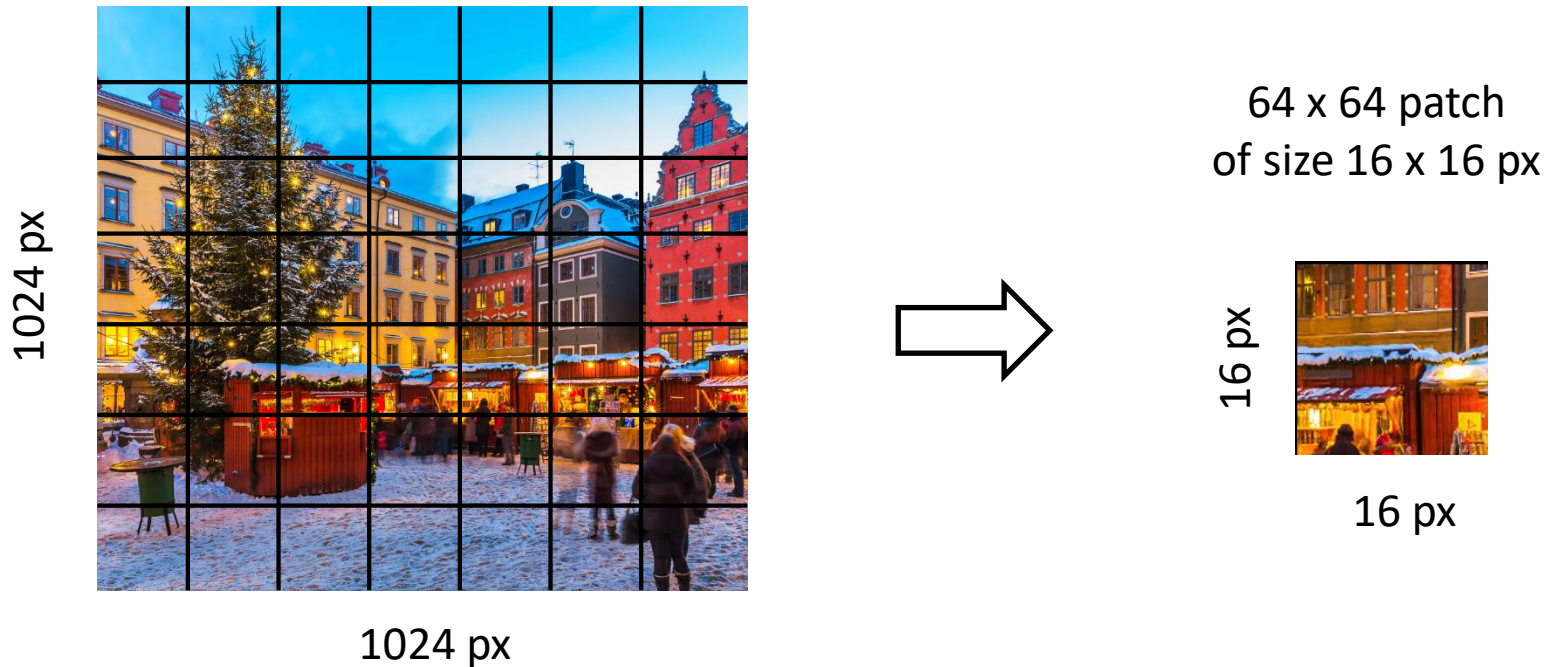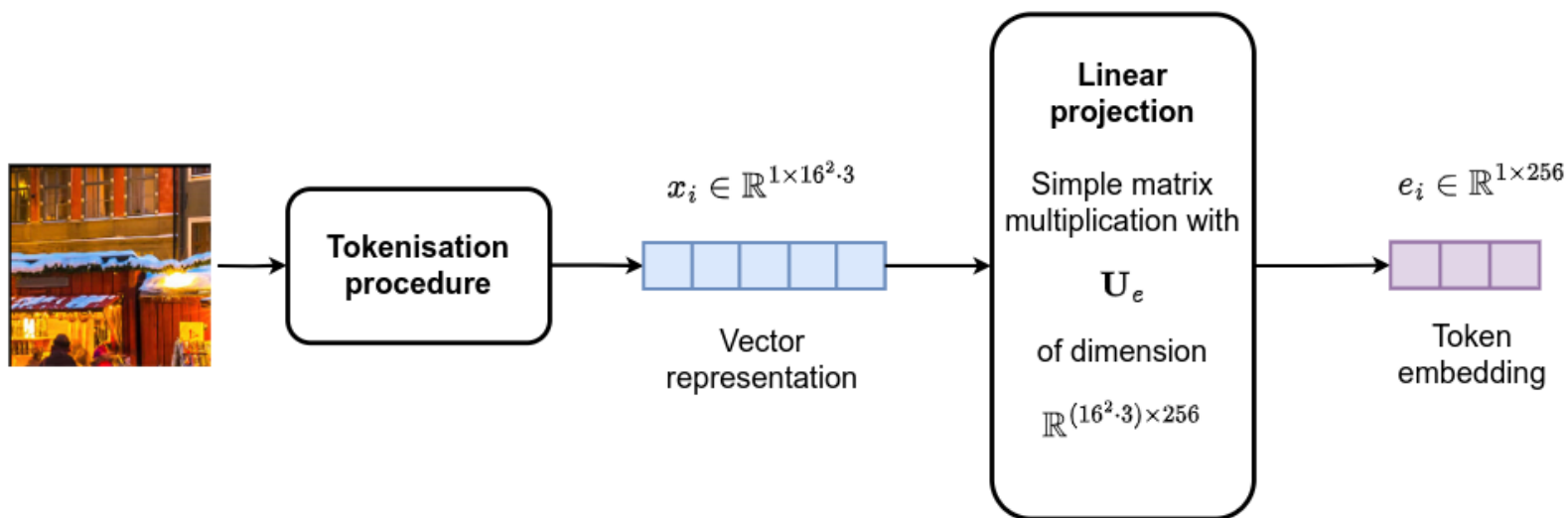
► SAM: Segment Anything

➔ Image encoder

► Input images resized to 1024 x 1014 px

► Architecture: ViT-H/16 (token = patch of size 16 x 16 px)

► Hidden size D: 256

1024 px

1024 px

64 x 64 patch
of size 16 x 16 px
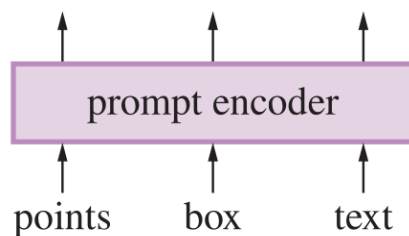
16 px

16 px

▶ **SAM: Segment Anything**

➜ Image encoder

▶ Input images resized to 1024 x 1014 px

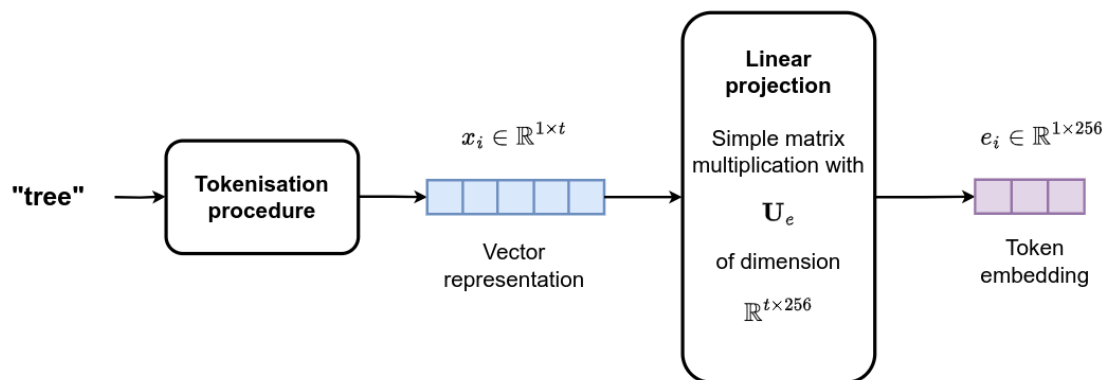▶ Architecture: ViT-H/16 (token = patch of size 16 x 16 px)

▶ Hidden size D: 256



$x_i \in \mathbb{R}^{1 \times 16^2 \cdot 3}$

**Tokenisation procedure**

Vector representation

**Linear projection**

Simple matrix multiplication with

$\mathbf{U}_e$

of dimension

$\mathbb{R}^{(16^2 \cdot 3) \times 256}$

$e_i \in \mathbb{R}^{1 \times 256}$

Token embedding

▶ **SAM: Segment Anything**

➔ Prompt encoder

▶ Several possible prompts: points, bounding boxes, text

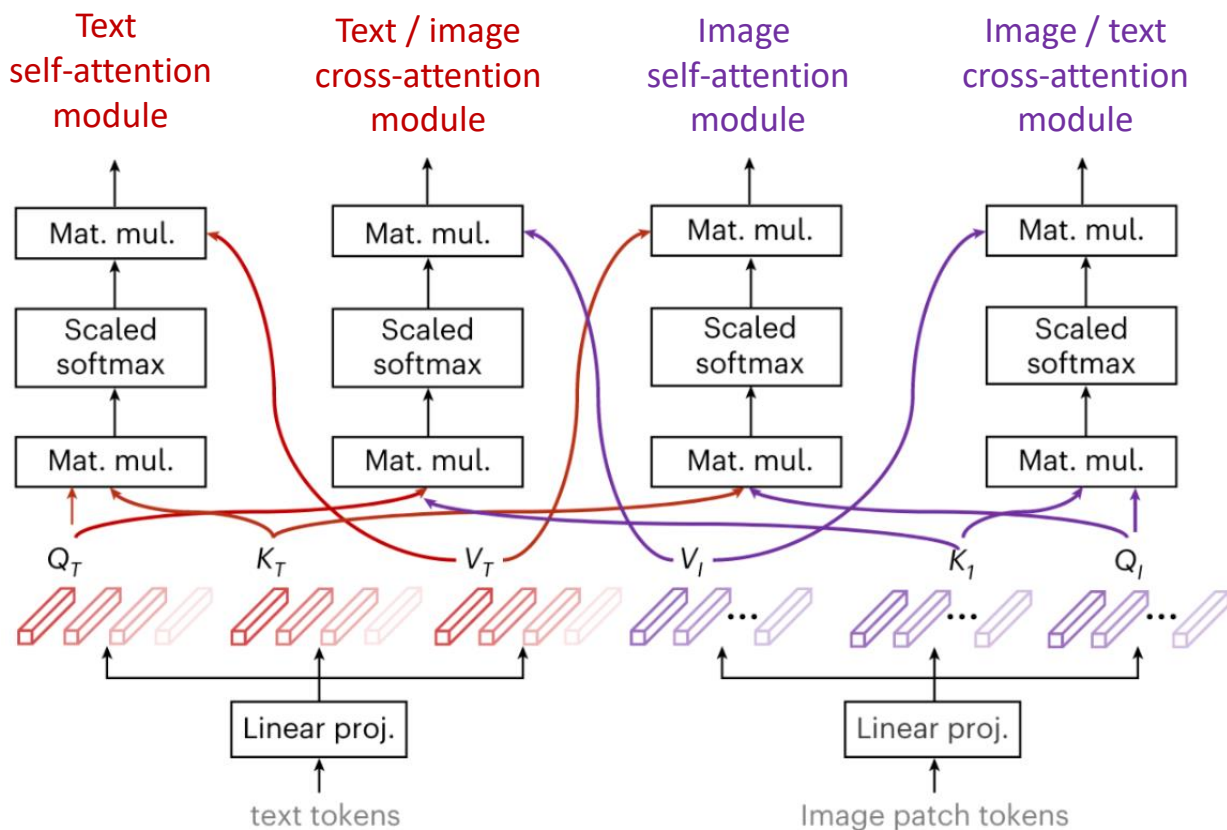▶ Hidden size D: 256



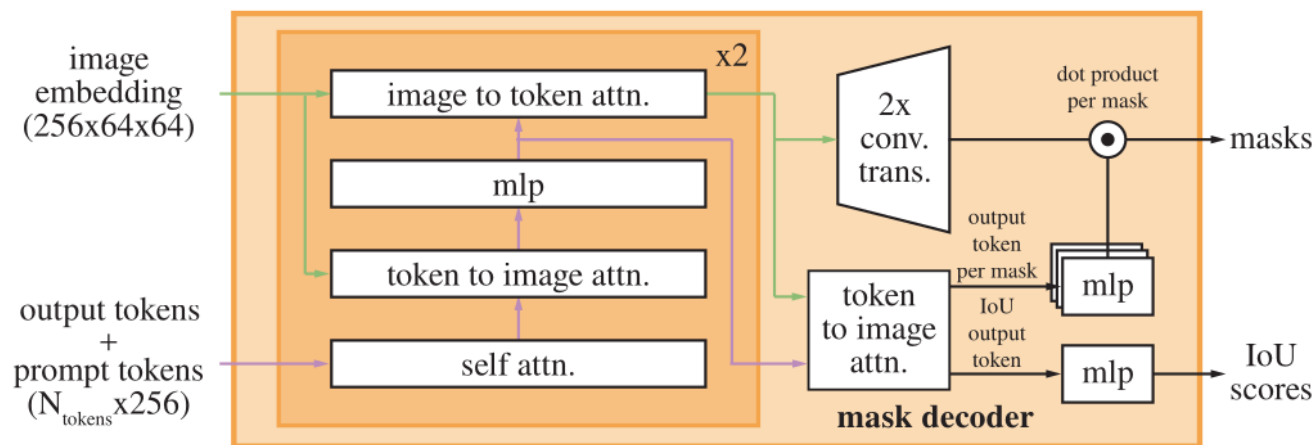▶ Illustration for the text encoder

▶ **SAM: Segment Anything**

➜ Mask decoder

▶ Uses the cross-attention principle

▶ **SAM: Segment Anything**

➔ **Mask decoder**

▶ Integration of a class token (output token)

# Transformer for segmentation

► SAM: Segment Anything

➔ Link to demo

➔ https://segment-anything.com/

# That's all folks