# Deep Learning for Computer Vision

## 113-1/Fall 2024

https://cool.ntu.edu.tw/courses/41702 (NTU COOL)
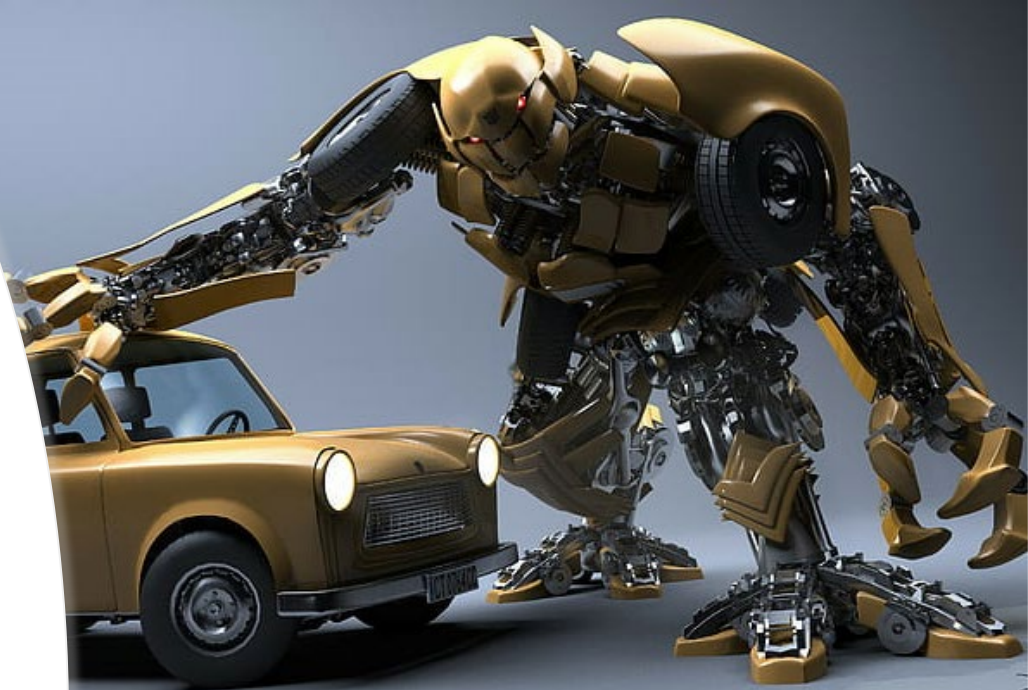
http://vllab.ee.ntu.edu.tw/dlcv.html (Public website)

Yu-Chiang Frank Wang 王鈺強, Professor

Dept. Electrical Engineering, National Taiwan University

2024/10/22

# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - SSL & Beyond
- Vision-Language Model
  - Image2Text
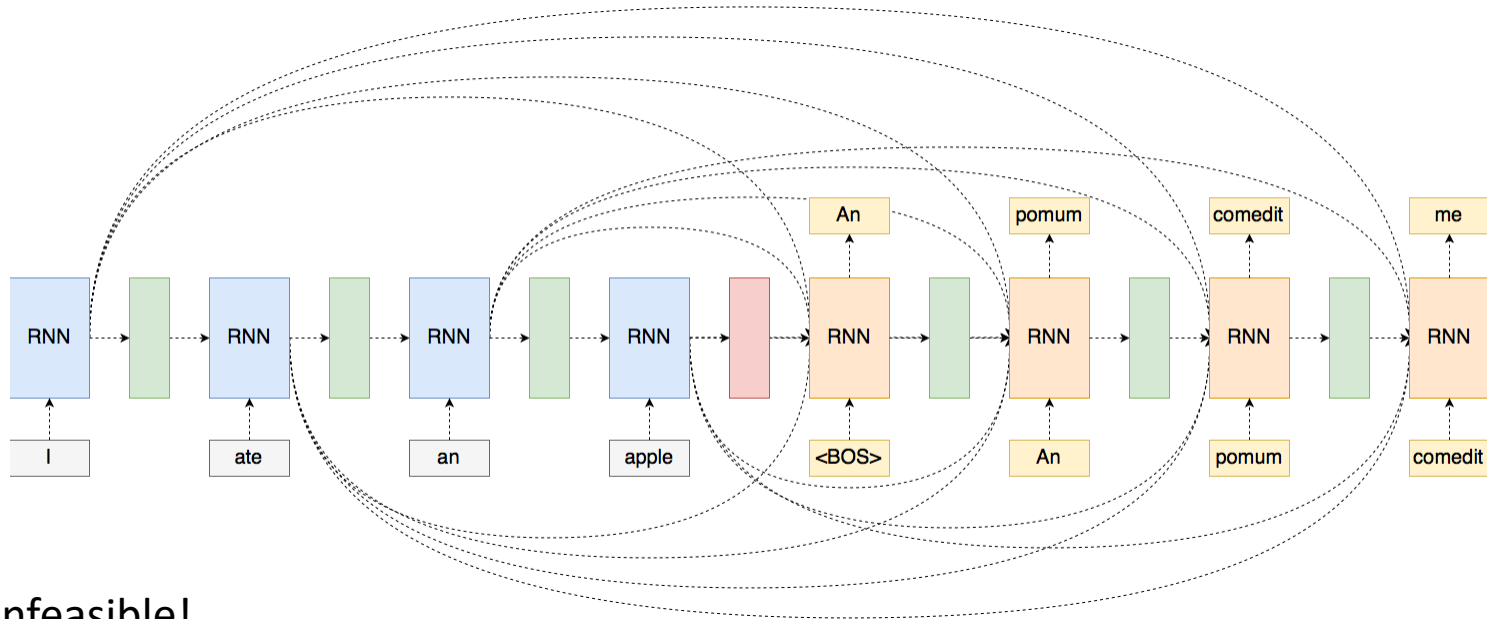  - Text2Image (√)
  - Image-text models

https://medium.com/@navendubrajesh/vision-language-models-use-cases-ee6d54b2c557
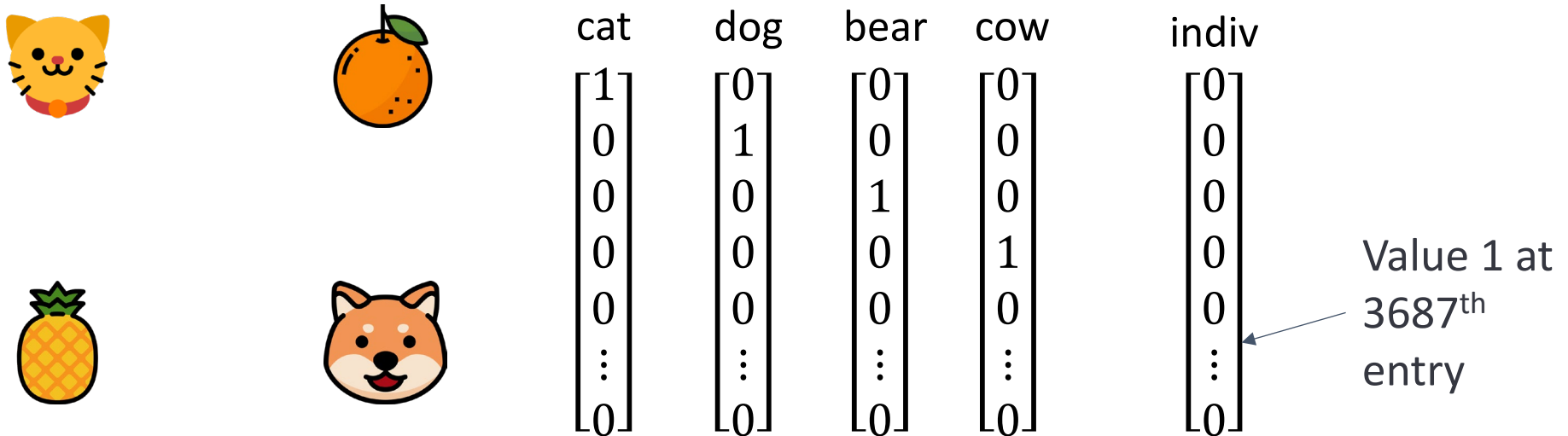
# RNN with Attention is Good, But..

- Attention in a pre-defined sequential order

- Information loss due to long sequences…

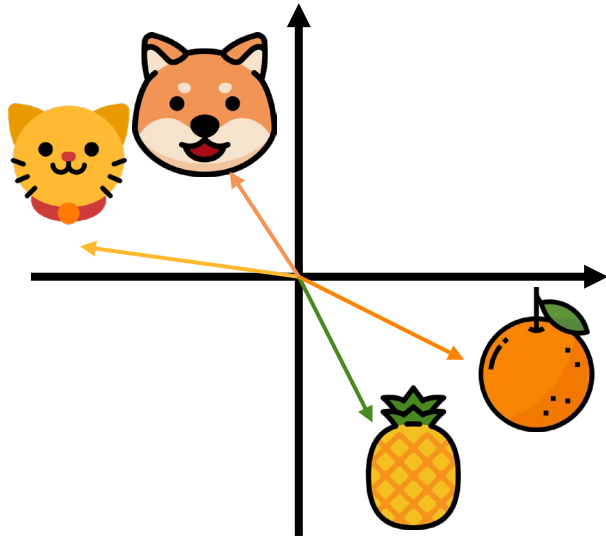- Connecting every hidden state between encoder and decoder?



- Infeasible!
  - Both inputs and outputs are with varying sizes.
  - Overparameterized

# Tokenization

Many words map to one token, but some don't: indivisible.

| 8607 | 4339 | 2472 | 311 | 832 | 4037 | 11 | 719 | 1063 | 1541 | 956 | 25 | 3687 | 23936 | 13 |

**One-hot encoding**

# tokens

$$
\begin{array}{ccccc}
\text{cat} & \text{dog} & \text{bear} & \text{cow} & \text{indiv} \\
\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} &
\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} &
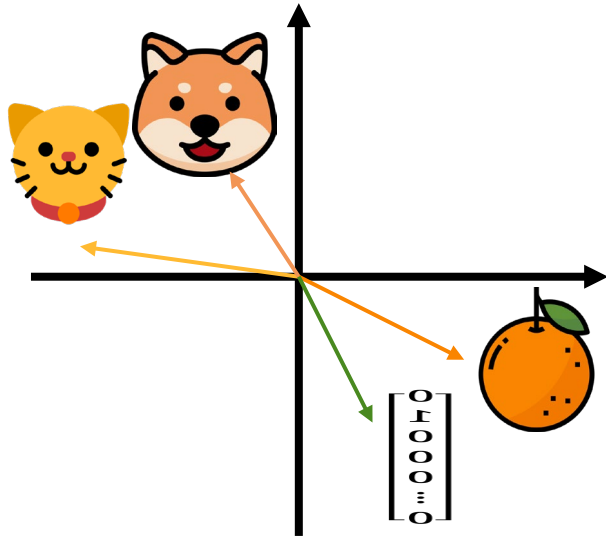\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} &
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} &
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
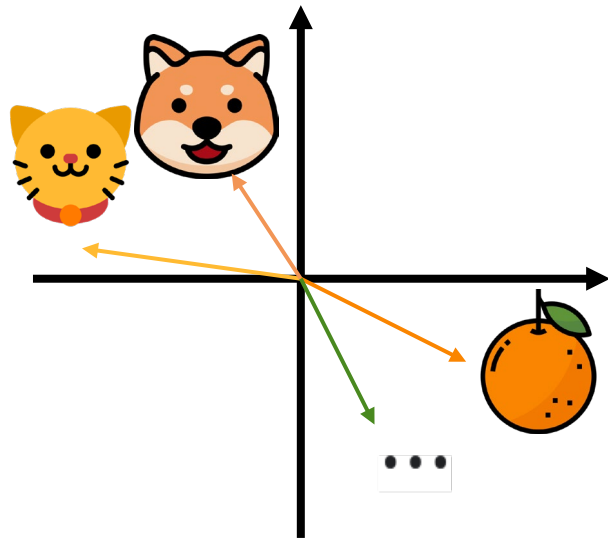\end{array}
$$

Value 1 at 3687th entry

TOKEN EMBEDDING

**One-hot encoding**

|  | cat | dog | bear | cow | indiv |
|---|---|---|---|---|---|
|  | 1 | 0 | 0 | 0 | 0 |
|  | 0 | 1 | 0 | 0 | 0 |
|  | 0 | 0 | 1 | 0 | 0 |
|  | 0 | 0 | 0 | 1 | 0 |
|  | 0 | 0 | 0 | 0 | 0 |
|  | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
|  | 0 | 0 | 0 | 0 | 0 |

Value 1 at $3687^{th}$ entry

# TOKEN EMBEDDING



Embedding Space

cat

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

dog

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

bear

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

cow

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

indiv

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Value 1 at 3687th entry

# TOKEN EMBEDDING



Embedding Space

$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} W_E \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$
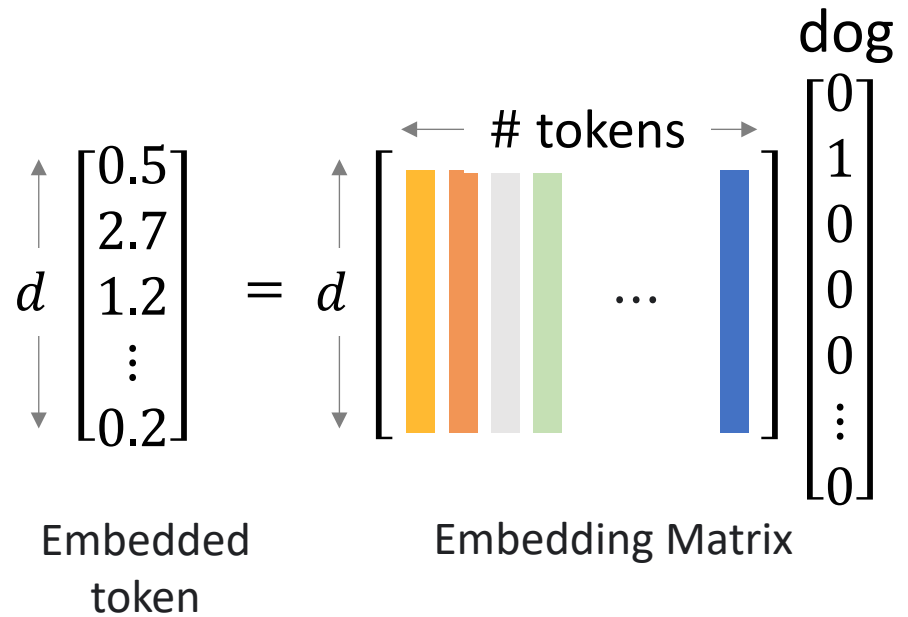
\# tokens

dog

Embedded token

Embedding Matrix

# TOKEN EMBEDDING



Embedding Space

Embedded token

Embedding Matrix

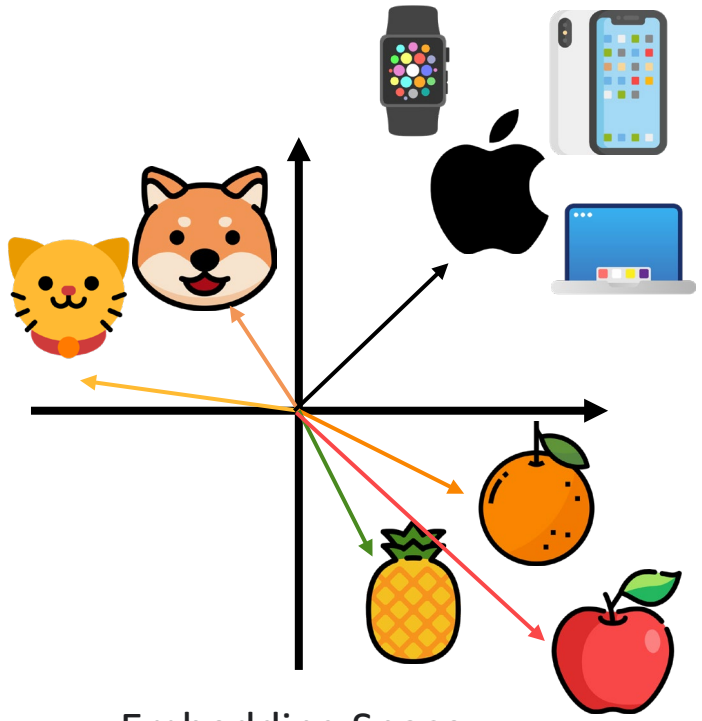$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} & & & & \cdots & \\ & & & & & \end{bmatrix} \overset{\text{dog}}{\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}}$$

\# tokens

# TOKEN EMBEDDING



Apple

dog

$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} & & & & & \\ & & & \cdots & & \\ & & & & & \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

# tokens

Embedding Space

Embedded token

Embedding Matrix

# TOKEN EMBEDDING

Apple

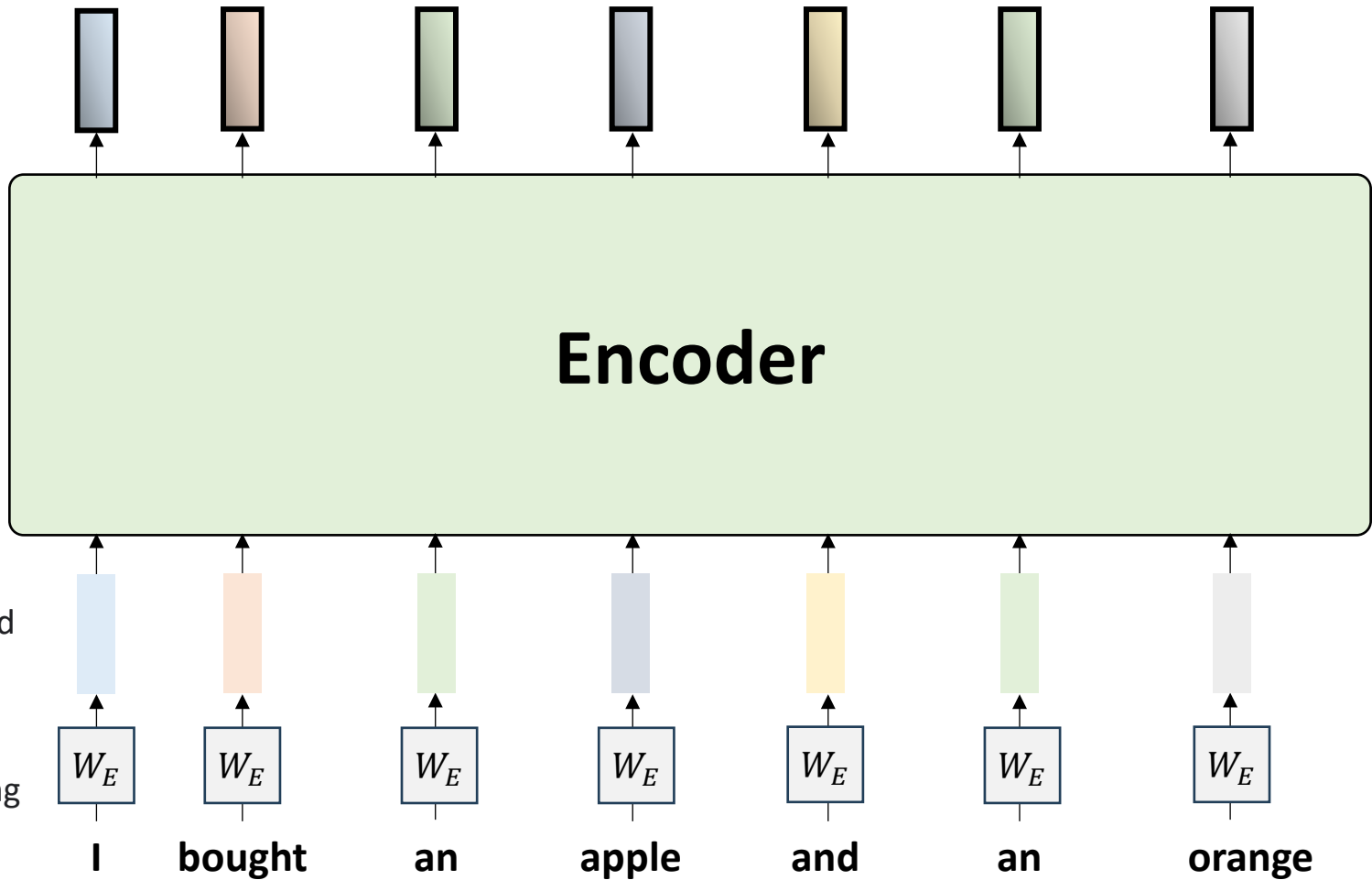I bought an **apple** and an orange.

I bought an **apple** watch.

dog

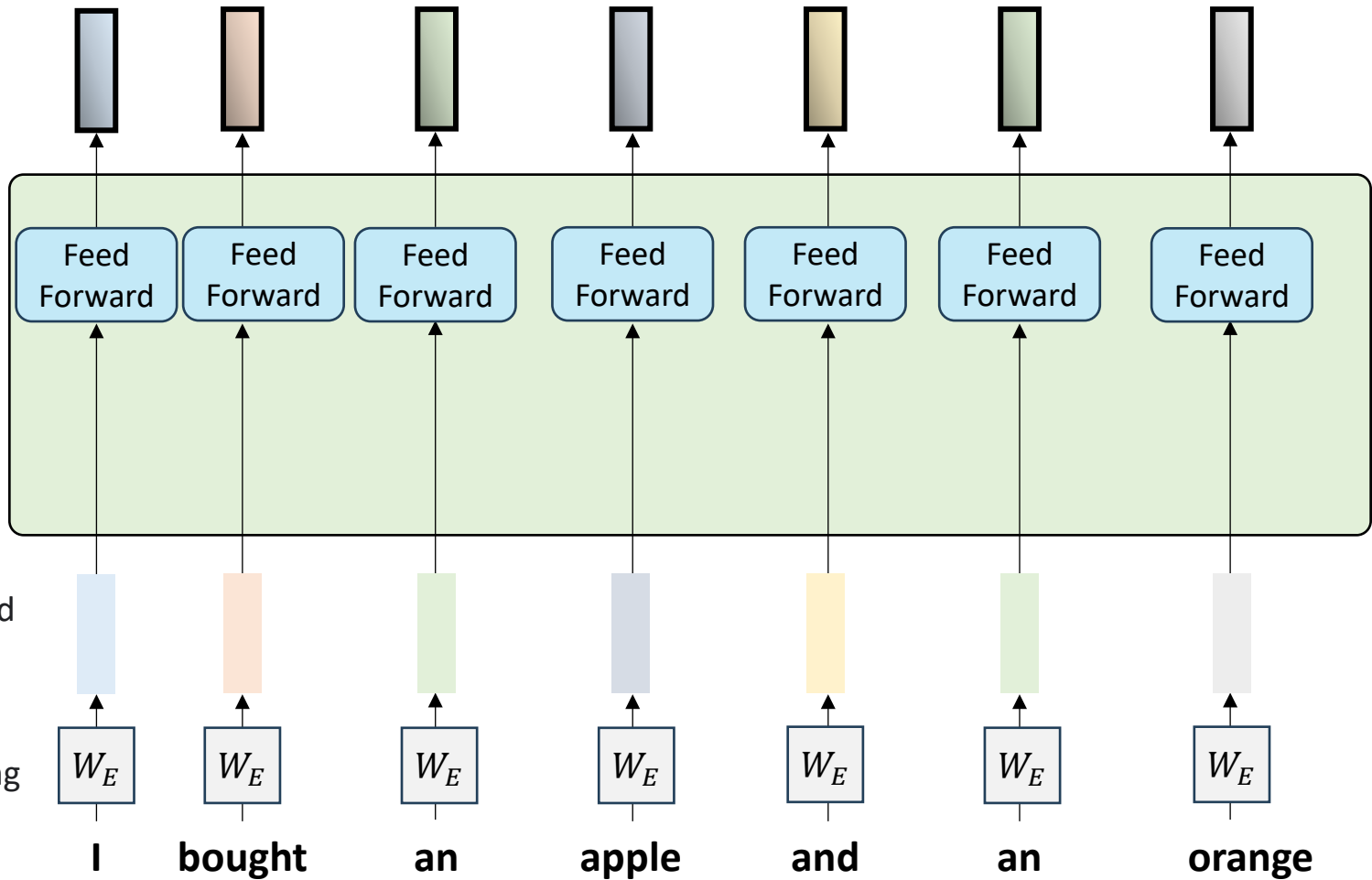$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} & & & & & & \\ & & & & \cdots & & \\ & & & & & & \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

# tokens

Embedding Space

Embedded token

Embedding Matrix

Embedded Tokens

Token Embedding

Tokens

$W_E$ I

$W_E$ bought

$W_E$ an

$W_E$ apple

$W_E$ and

$W_E$ an

$W_E$ orange

Not applicable!

Embedded Tokens

Token Embedding

Tokens: I bought an apple and an orange
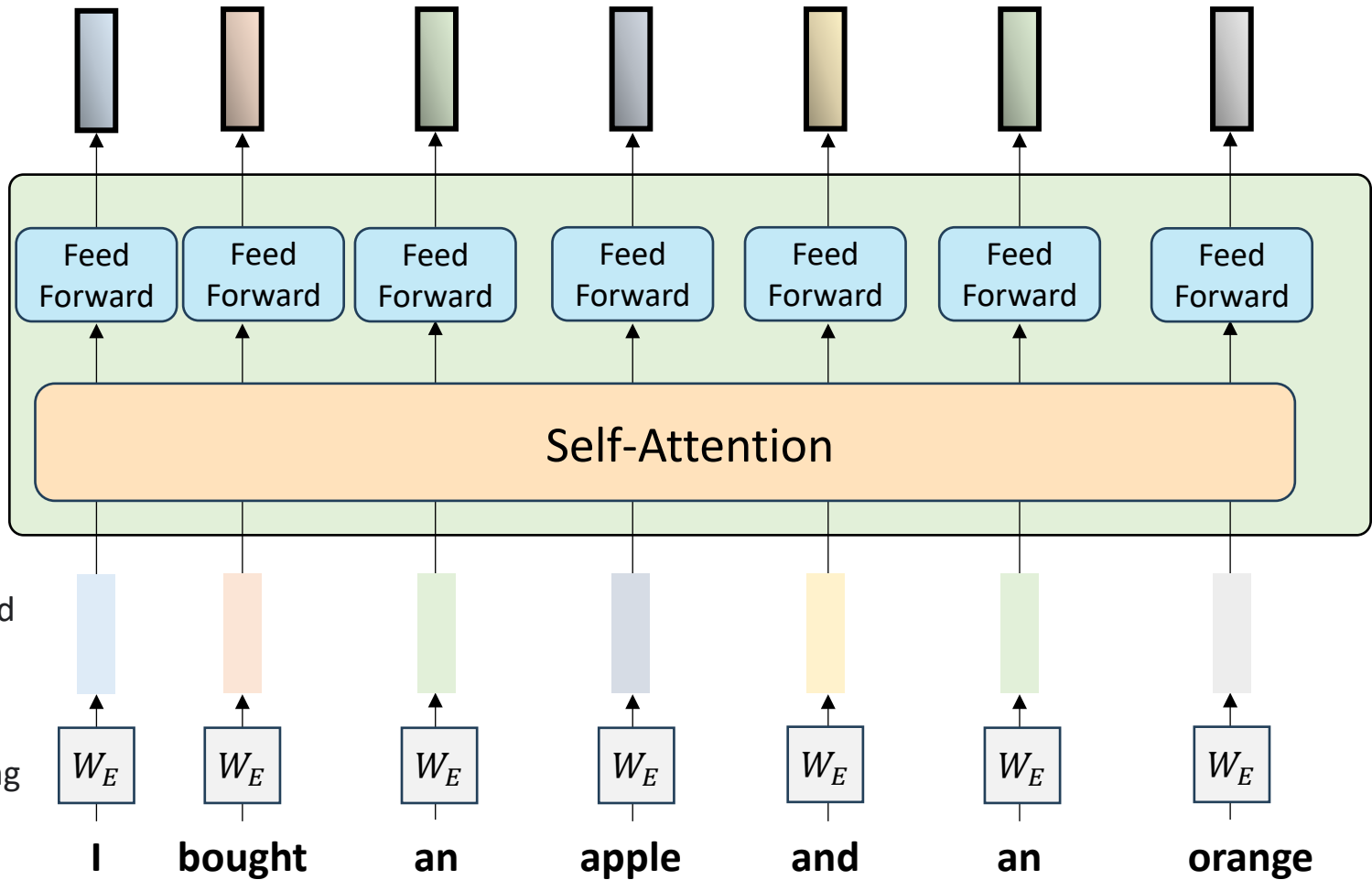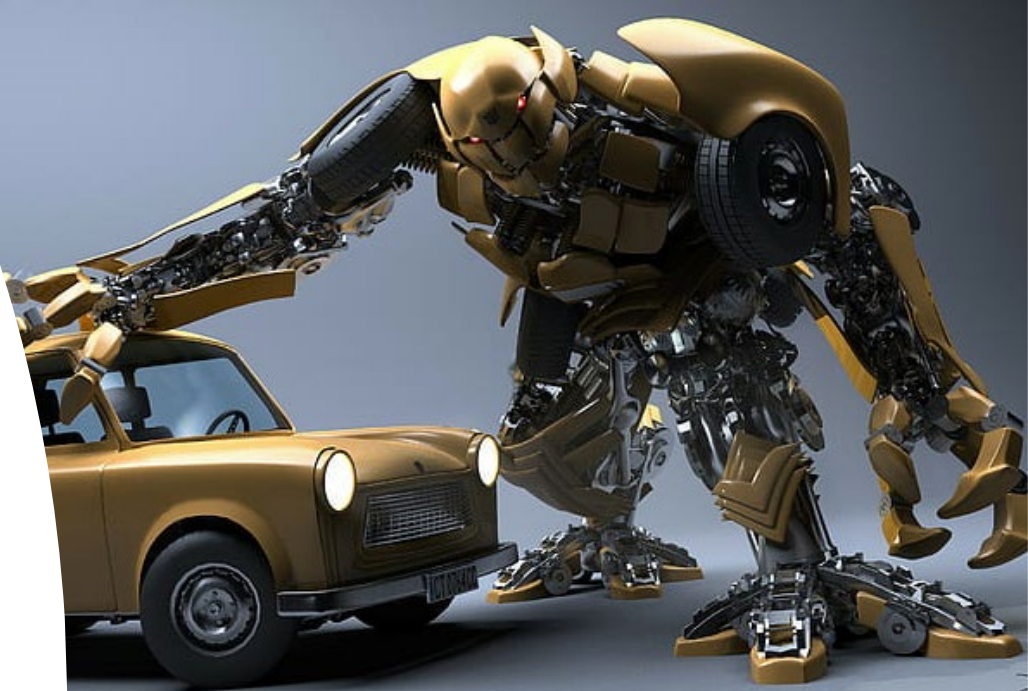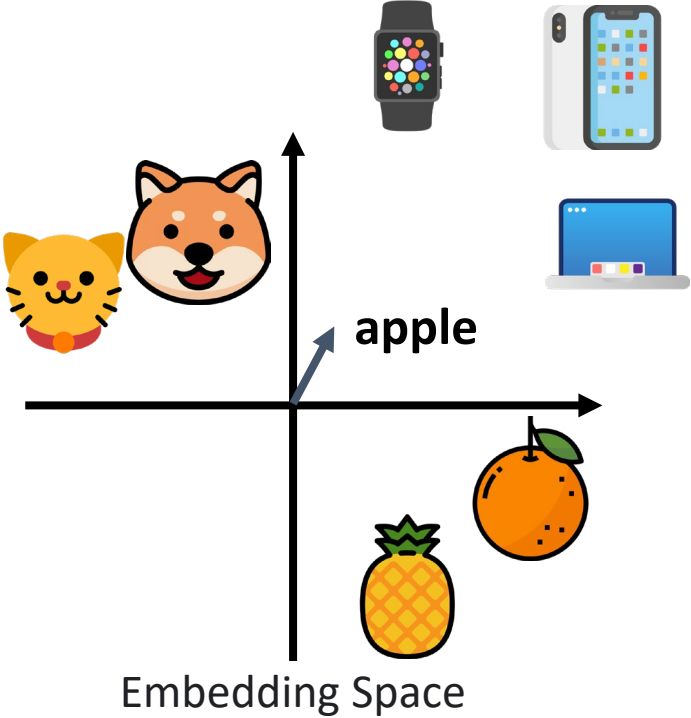
# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - SSL & Beyond
- Vision-Language Model
  - Image2Text
  - Text2Image (√)
  - Image-text models



https://medium.com/@navendubrajesh/vision
-language-models-use-cases-ee6d54b2c557

# Self-Attention



Embedding Space

| Embedded Tokens | | | | | | |
|---|---|---|---|---|---|---|
| Tokens | | | | | | |
| **I** | **bought** | **an** | **apple** | **and** | **an** | **orange** |

# Self-Attention



Embedding Space

| Embedded Tokens | | | | | | |
|---|---|---|---|---|---|---|
| Tokens | | | | | | |
| I | bought | an | apple | and | an | orange |

# Self-Attention

apple

Embedding Space

Embedded Tokens

Tokens | **I** | **bought** | **an** | **apple** | **watch**

# Self-Attention



Embedding Space

| Embedded Tokens | | | | |
|---|---|---|---|---|
| Tokens | | | | |
| **I** | **bought** | **an** | **apple** | **watch** |

# Self-Attention (1/5)

- Query **q**, key **k**, value **v** vectors are learned from each input **x**

$$q_i = W^Q x_i$$
$$k_i = W^K x_i$$
$$v_i = W^V x_i$$

# Self-Attention (2/5)

- Relation between each input is modeled by inner-product of query **q** and key **k**.

$$a_{1,i} = \frac{q_1 \cdot k_i}{\sqrt{d}}, \text{ where } a \in R, q, k \in R^d$$

# Self-Attention (3/5)

- SoftMax is applied:

$$0 \le \hat{a}_i = e^{a_i} / \sum_j^N e^{a_j} \le 1 \text{ , for I =1, ..., N}$$

| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

$\hat{a}_{1,1}$    $\hat{a}_{1,2}$    $\hat{a}_{1,N}$

**SoftMax**

$a_{1,1}$    $a_{1,2}$    $a_{1,N}$

$q_1$ $k_1$ $v_1$    $q_2$ $k_2$ $v_2$    $q_N$ $k_N$ $v_N$

$x_1$    $x_2$    ...    $x_N$

23

# Self-Attention (4/5)

- Value vectors **v** are aggregated
  with attention weight $\hat{a}$ , i.e., $y_1 = \sum_i^N \hat{a}_i \cdot v_i$

# Self-Attention (5/5)

- All $y_i$ can be computed **in parallel**

- Each $y_i$ considers $x_1 \sim x_N$, modeling their **long-distance dependencies**.

- Global feature can be obtained by **average-pooling** over $y_1 \sim y_N$

# Self-Attention: Implementation

- Input sequence can be represented as a N x $d_{in}$ matrix
- * denotes matrix multiplication



$x_i \in R^{d_{in}}$

$x_1$
$x_2$
.
.
.
$x_N$

Input matrix

$= \quad N \quad d_{in}$

$* \quad d_{in} \quad W^Q \quad d$

$* \quad d_{in} \quad W^K \quad d$

$* \quad d_{in} \quad W^V \quad d$

$N \quad Q \quad d$

$N \quad K \quad d$

$N \quad V \quad d$

# Self-Attention: Implementation



- Output matrix **Y**

- All operations are **matrix multiplication**, can be parallelized on GPU.



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

# Multi-Head Self-Attention (1/4)

- Perform self-attention at different subspaces,
  implying performing attention over different input feature types
  (e.g., representations, modalities, positions, etc.)

# Multi-Head Self-Attention (2/4)



- Perform self-attention at different subspaces, implying performing attention over different input feature types

- See example below



Attention weights
of Head 1

Attention weights
of Head 2

# Multi-Head Self-Attention (3/4)

- A two-head example:
  output of two heads are concatenated as the output embedding



$y_{i,1}$

$q_i$  $k_i$  $v_i$    $q_i$  $k_i$  $v_i$    $q_j$  $k_j$  $v_j$    $q_j$  $k_j$  $v_j$

Head 1      Head 2          Head 1      Head 2

$x_i$                    $x_j$

# Multi-Head Self-Attention (4/4)

- A two-head example:
  output of two heads are concatenated as the output embedding

$$y_{i,1} \quad y_{i,2} \quad = \quad y_i$$



$q_i$ $k_i$ $v_i$ $\quad$ $q_i$ $k_i$ $v_i$ $\quad$ $q_j$ $k_j$ $v_j$ $\quad$ $q_j$ $k_j$ $v_j$

Head 1 $\quad$ Head 2 $\quad\quad$ Head 1 $\quad$ Head 2

$x_i$ $\quad\quad\quad\quad$ $x_j$

# The Residuals

- A **residual connection** followed by layer normalization

# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - SSL & Beyond
- Vision-Language Model
  - Image2Text
  - Text2Image (v)
  - Image-text models



https://medium.com/@navendubrajesh/vision-language-models-use-cases-ee6d54b2c557

# Training The Decoder of Transformer

- Encoder-decoder attention
  - Q from <u>self-attn output</u> in decoder, K & V from <u>encoder outputs</u>
- Masked multi-head attention
  - Design similar to that of encoder, except for decoder #1 which takes additional inputs (of GT/predicted word embeddings).
  - Mask unpredicted tokens during softmax: what does this mean & why?

# Training The Decoder of Transformer (cont'd)

- Encoder-decoder attention
  - Q from self-attn in decoder, K & V from encoder outputs
- Masked multi-head attention
  - Design similar to that of encoder, except for decoder #1 which takes additional inputs (of GT/predicted word embeddings).
  - Mask unpredicted tokens during softmax: what does this mean & why?

# Overview of Encoding & Decoding in Transformer

- Encoder/Decoder Self & Cross-Attention

# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT
  - Swin Transformer
- Vision-Language Model
  - Image2Text
  - Text2Image (√)
  - Image-text models

https://medium.com/@navendubrajesh/vision-language-models-use-cases-ee6d54b2c557

| | | | | |
|---|---|---|---|---|
| | | | | |

Embedded Tokens    $\boldsymbol{x}_1$     $\boldsymbol{x}_2$     $\boldsymbol{x}_3$     $\boldsymbol{x}_4$     $\boldsymbol{x}_5$

Tokens     **I**     **bought**     **an**     **apple**     **watch**

Position     $k = 1$     $k = 2$     $k = 3$     $k = 4$     $k = 5$

Position $k$

Angular frequency

$w_i = N^{-2i/d}$

$N = 100{,}000$

$$d\begin{bmatrix} \sin(w_0 k) \\ \cos(w_0 k) \\ \sin(w_1 k) \\ \cos(w_1 k) \\ \vdots \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} k\right) \\ \cos\left(w_{\frac{d}{2}-1} k\right) \end{bmatrix}$$

←— Fast oscillating

←— Slow oscillating

Embedded Tokens $\boldsymbol{x}_1$ $\boldsymbol{x}_2$ $\boldsymbol{x}_3$ $\boldsymbol{x}_4$ $\boldsymbol{x}_5$

Tokens **I** **bought** **an** **apple** **watch**

Position $k = 1$ $k = 2$ $k = 3$ $k = 4$ $k = 5$

$$d \begin{bmatrix} \sin(w_0 k) \\ \cos(w_0 k) \\ \sin(w_1 k) \\ \cos(w_1 k) \\ \vdots \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} k\right) \\ \cos\left(w_{\frac{d}{2}-1} k\right) \end{bmatrix}$$

Embedded Tokens $d$ $\boldsymbol{x}_1$ $\boldsymbol{x}_2$ $\boldsymbol{x}_3$ $\boldsymbol{x}_4$ $\boldsymbol{x}_5$

Tokens **I** **bought** **an** **apple** **watch**

Position $k = 1$ $k = 2$ $k = 3$ $k = 4$ $k = 5$

$d$ $\boldsymbol{P}_1$ $\boldsymbol{P}_2$ $\boldsymbol{P}_3$ $\boldsymbol{P}_4$ $\boldsymbol{P}_5$

$$d \begin{bmatrix} \sin(w_0 k) \\ \cos(w_0 k) \\ \sin(w_1 k) \\ \cos(w_1 k) \\ \vdots \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} k\right) \\ \cos\left(w_{\frac{d}{2}-1} k\right) \end{bmatrix}$$

$\boldsymbol{x}_i$ [concat] $\boldsymbol{P}_i$

$\boldsymbol{x}_i$ [MLP] $\boldsymbol{P}_i$

$\boldsymbol{x}_i$ [ + ] $\boldsymbol{P}_i$

| Position | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| | I | walk | my | dog | every | day |

| Position | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| | every | day | I | walk | my | dog |

From *absolute* to *relative* positional embedding
"RoFormer: Enhanced Transformer with Rotary Position Embedding", arxiv 2021

# Extension:
# BERT - Bidirectional Encoder Representation from Transformers

- Proposed by Google AI Language

- Two additional objectives
  - Masked language model (MLM)
  - Next sentence prediction (NSP)



Pre-training                     Fine-Tuning

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, arxiv. 2018                43

# Transformer is promising, but…

- Concerns of Transformer?
  - Computation
  - Space/memory
- Potential solutions
  - Sparse Transformer
  - Linformer
  - Linearized attention, etc.
  - Ever heard of *Mamba*?

# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - SSL & Beyond
- Vision-Language Model
  - Image2Text
  - Text2Image (√)
  - Image-text models

https://medium.com/@navendubrajesh/vision
-language-models-use-cases-ee6d54b2c557

Vision
Language
Model

# Vision Transformer

- "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR, 2021. (Google Research)

- Only the encoder part is utilized

# Vision Transformer (cont'd)

- Partition the input image into a **patch sequence**

- An additional **token** (**\***) is appended to perform attention on patches

- Both the "**\***" token and positional embeddings (denoted by 0, 1, 2 …) are **trainable vectors**.



**Vision Transformer (ViT)**

# Vision Transformer (cont'd)

- "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR, 2021. (Google Research)

$$y = LN(z_L^0)$$

MLP Head

Transformer Encoder

$$z_0 = \left[ x_{Class}; x_p^1 E; ...; x_p^N E \right] + E_{pos}$$

Trainable Vector [CLS token]

Convert integer to D-dimension vector — Position Embedding

Linear Projection — Linear transform from P*P*C to D

Flatten

Split the image to 16*16 Patches

$$\mathbf{z}_0 = [\mathbf{x}_{class}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \cdots ; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{pos}, \qquad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \; \mathbf{E}_{pos} \in \mathbb{R}^{(N+1) \times D}$$

$$\mathbf{z'}_\ell = \mathrm{MSA}(\mathrm{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \qquad \ell = 1 \ldots L \qquad \text{Multiheaded self-attention (MSA)}$$

$$\mathbf{z}_\ell = \mathrm{MLP}(\mathrm{LN}(\mathbf{z'}_\ell)) + \mathbf{z'}_\ell, \qquad \ell = 1 \ldots L \qquad \text{Layer norm (LN)}$$

$$\mathbf{y} = \mathrm{LN}(\mathbf{z}_L^0)$$

Figure credit: CS886 Univ. Waterloo

# Query-Key-Value Attention in ViT

- E.g., An input image is partitioned into 4 patches, with feature dimension = 3 (i.e., P=4 and D=3).

- Note that there are (P+1) rows since we have an additional token "*".

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$



$P$: Patch Number
$D$: Dimension
$W_{Q,K,V}$: Learnable Matrices
$Q, K, V$: Query, Key, Value
$A$: Attention Matrix

# Query-Key-Value Attention in ViT (cont'd)

- In the standard vision transformer, we only take the **first output token** of the output sequence (the **first row** of Y) for classification purposes

- This corresponds to the output when **token "0"** serves as query



$P$: Patch Number
$D$: Dimension
$W_{Q,K,V}$: Learnable Matrices
$Q, K, V$: Query, Key, Value
$A$: Attention Matrix

50

# Visualization of ViT

- To visualize the **attention maps**, we take the attention scores from the **first row** of A (when token "0" serves as query)

- Note the first element is excluded, and thus there are **P scores** corresponding to the P image patches



$P$: Patch Number
$D$: Dimension
$W_{Q,K,V}$: Learnable Matrices
$Q, K, V$: Query, Key, Value
$A$: Attention Matrix

# ViT Results

- ViT outperforms CNN-based models
  - Pretrained on JFT & ImageNet
  - Can be trained using TPUv3 w/ 8 cores in 30 days; faster than CNN
  - ViT for visualization/interpretability?

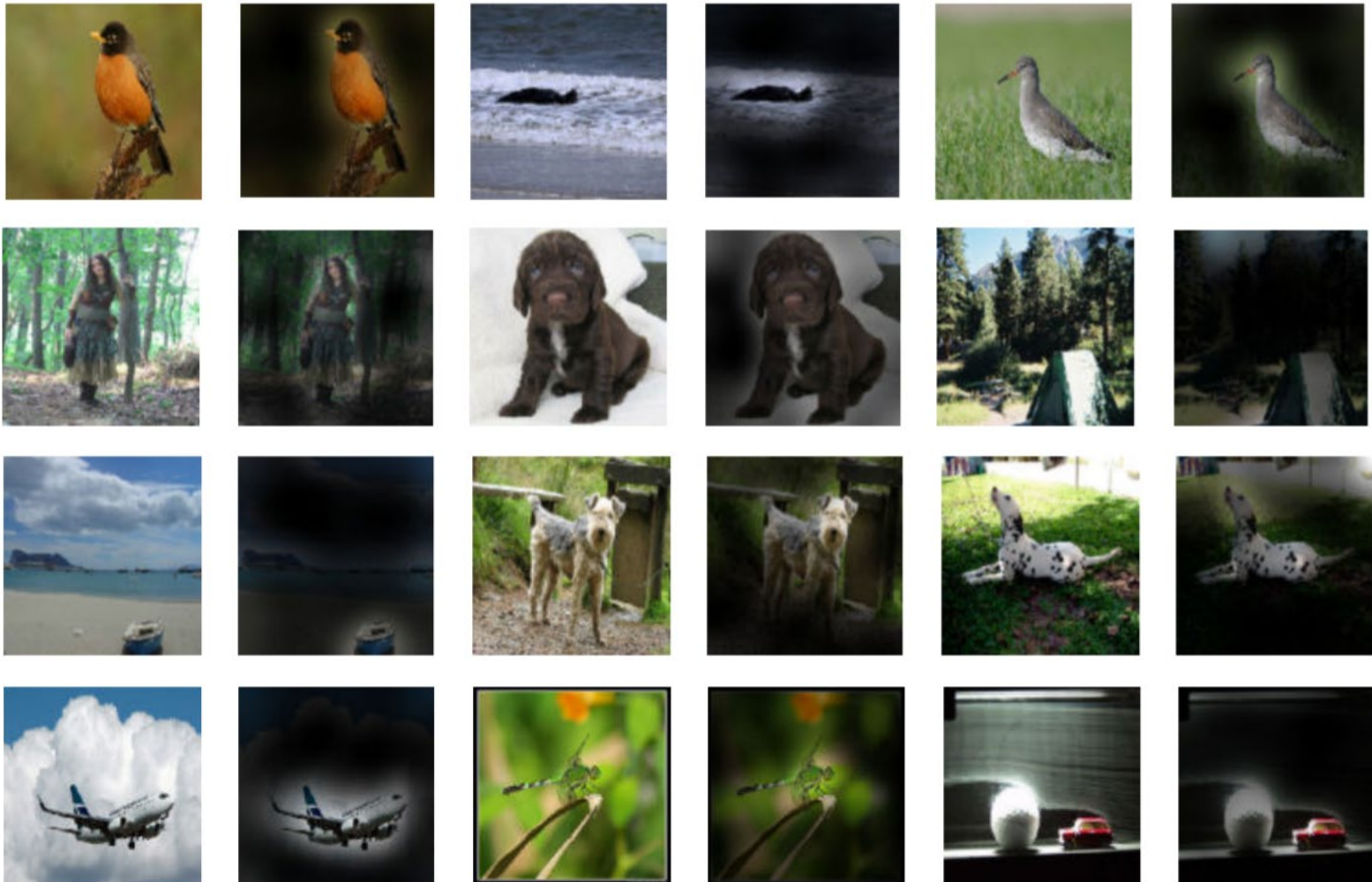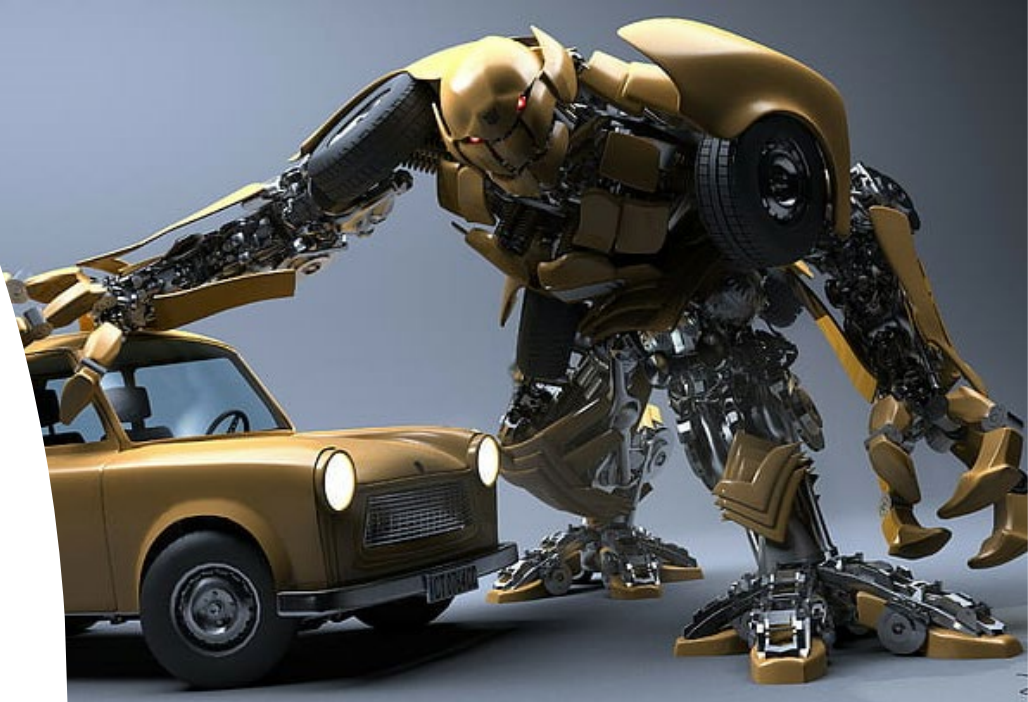| | Ours-JFT (ViT-H/14) | Ours-JFT (ViT-L/16) | Ours-I21k (ViT-L/16) | BiT-L (ResNet152x4) | Noisy Student (EfficientNet-L2) |
|---|---|---|---|---|---|
| ImageNet | $88.55 \pm 0.04$ | $87.76 \pm 0.03$ | $85.30 \pm 0.02$ | $87.54 \pm 0.02$ | 88.4/88.5* |
| ImageNet ReaL | $90.72 \pm 0.05$ | $90.54 \pm 0.03$ | $88.62 \pm 0.05$ | 90.54 | 90.55 |
| CIFAR-10 | $99.50 \pm 0.06$ | $99.42 \pm 0.03$ | $99.15 \pm 0.03$ | $99.37 \pm 0.06$ | – |
| CIFAR-100 | $94.55 \pm 0.04$ | $93.90 \pm 0.05$ | $93.25 \pm 0.05$ | $93.51 \pm 0.08$ | – |
| Oxford-IIIT Pets | $97.56 \pm 0.03$ | $97.32 \pm 0.11$ | $94.67 \pm 0.15$ | $96.62 \pm 0.23$ | – |
| Oxford Flowers-102 | $99.68 \pm 0.02$ | $99.74 \pm 0.00$ | $99.61 \pm 0.02$ | $99.63 \pm 0.03$ | – |
| VTAB (19 tasks) | $77.63 \pm 0.23$ | $76.28 \pm 0.46$ | $72.72 \pm 0.21$ | $76.29 \pm 1.70$ | – |
| TPUv3-core-days | 2.5k | 0.68k | 0.23k | 9.9k | 12.3k |

# Example Visualization for Object Recognition

# ViT Results

- ViT outperforms CNN-based models
  - Pretrained on JFT & ImageNet
  - Can be trained using TPUv3 w/ 8 cores in 30 days; faster than CNN
  - However, JFT is not publicly available…

| | Ours-JFT (ViT-H/14) | Ours-JFT (ViT-L/16) | Ours-I21k (ViT-L/16) | BiT-L (ResNet152x4) | Noisy Student (EfficientNet-L2) |
|---|---|---|---|---|---|
| ImageNet | $88.55 \pm 0.04$ | $87.76 \pm 0.03$ | $85.30 \pm 0.02$ | $87.54 \pm 0.02$ | $88.4/88.5^*$ |
| ImageNet ReaL | $90.72 \pm 0.05$ | $90.54 \pm 0.03$ | $88.62 \pm 0.05$ | $90.54$ | $90.55$ |
| CIFAR-10 | $99.50 \pm 0.06$ | $99.42 \pm 0.03$ | $99.15 \pm 0.03$ | $99.37 \pm 0.06$ | — |
| CIFAR-100 | $94.55 \pm 0.04$ | $93.90 \pm 0.05$ | $93.25 \pm 0.05$ | $93.51 \pm 0.08$ | — |
| Oxford-IIIT Pets | $97.56 \pm 0.03$ | $97.32 \pm 0.11$ | $94.67 \pm 0.15$ | $96.62 \pm 0.23$ | — |
| Oxford Flowers-102 | $99.68 \pm 0.02$ | $99.74 \pm 0.00$ | $99.61 \pm 0.02$ | $99.63 \pm 0.03$ | — |
| VTAB (19 tasks) | $77.63 \pm 0.23$ | $76.28 \pm 0.46$ | $72.72 \pm 0.21$ | $76.29 \pm 1.70$ | — |
| TPUv3-core-days | 2.5k | 0.68k | 0.23k | 9.9k | 12.3k |

# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- **Transformer for Visual Analysis**
  - Vision Transformer (ViT)
  - **DeiT & Swin Transformer**
  - SSL & Beyond
- Vision-Language Model
  - Image2Text
  - Text2Image (√)
  - Image-text models

https://medium.com/@navendubrajesh/vision-language-models-use-cases-ee6d54b2c557

# DeiT

- "Training data-efficient image transformers & distillation through attention", ICML 2021. (Facebook AI)

- ViT outperforms CNN, but the dataset is not publicly available.

- By distillation, DeiT only requires ImageNet for pretraining (10 times smaller).



Accuracy vs. throughput on ImageNet (Ours = DeiT)

# DeiT: Distill through Attention



viewed as a student

Can be a convnet, or a mixture of classifier

$\psi$: Softmax
$\lambda$: Weight of loss
$\tau$: Softmax temperature
$L_{CE}$: Cross-entropy loss

Soft Distillation: $L_t = (1 - \lambda)L_{CE}(\psi(Z_s), y) + \lambda\tau^2 KL(\psi(Z_s/\tau), \psi(Z_t/\tau))$

Hard Distillation: $L_t = \frac{1}{2}L_{CE}(\psi(Z_s), y) + \frac{1}{2}L_{CE}(\psi(Z_s), y_t)$

Figure credit: CS886 Univ. Waterloo

# DeiT Results

- Variants of DeiT architectures (adopted from the ViT backbone)

| Model | embedding dimension | #heads | #layers | #params | training resolution | throughput (im/sec) |
|---|---|---|---|---|---|---|
| DeiT-Ti | 192 | 3 | 12 | 5M | 224 | 2536 |
| DeiT-S | 384 | 6 | 12 | 22M | 224 | 940 |
| DeiT-B | 768 | 12 | 12 | 86M | 224 | 292 |

Same as the original ViT model →

- Choices of different teacher models & ablation studies

| Teacher Models | acc. | Student: DeiT-B pretrain | ↑384 |
|---|---|---|---|
| DeiT-B | 81.8 | 81.9 | 83.1 |
| RegNetY-4GF | 80.0 | 82.7 | 83.6 |
| RegNetY-8GF | 81.7 | 82.7 | 83.8 |
| RegNetY-12GF | 82.4 | 83.0 | 83.9 |
| RegNetY-16GF | 82.9 | 83.0 | 84.0 |

| DeiT: method ↓ | supervision label | teacher | ImageNet top-1 (%) Ti 224 | S 224 | B 224 | B↑384 |
|---|---|---|---|---|---|---|
| no distillation | ✓ | ✗ | 72.2 | 79.8 | 81.8 | 83.1 |
| usual distillation | ✗ | soft | 72.2 | 79.8 | 81.8 | 83.2 |
| hard distillation | ✗ | hard | 74.3 | 80.9 | 83.0 | 84.0 |
| class embedding | ✓ | hard | 73.9 | 80.9 | 83.0 | 84.2 |
| distil. embedding | ✓ | hard | 74.6 | 81.1 | 83.1 | 84.4 |
| DeiT⚗: class+distil. | ✓ | hard | 74.5 | 81.2 | 83.4 | 84.5 |

# Swin Transformer

- "Swin Transformer: Hierarchical Vision Transformer Shifted Windows", ICCV 2021. (MSRA)

- ViT's computation complexity vs. image size ->

- Propose to perform patch merging & Swin Transformer architecture



(a) Swin Transformer (ours)     (b) ViT
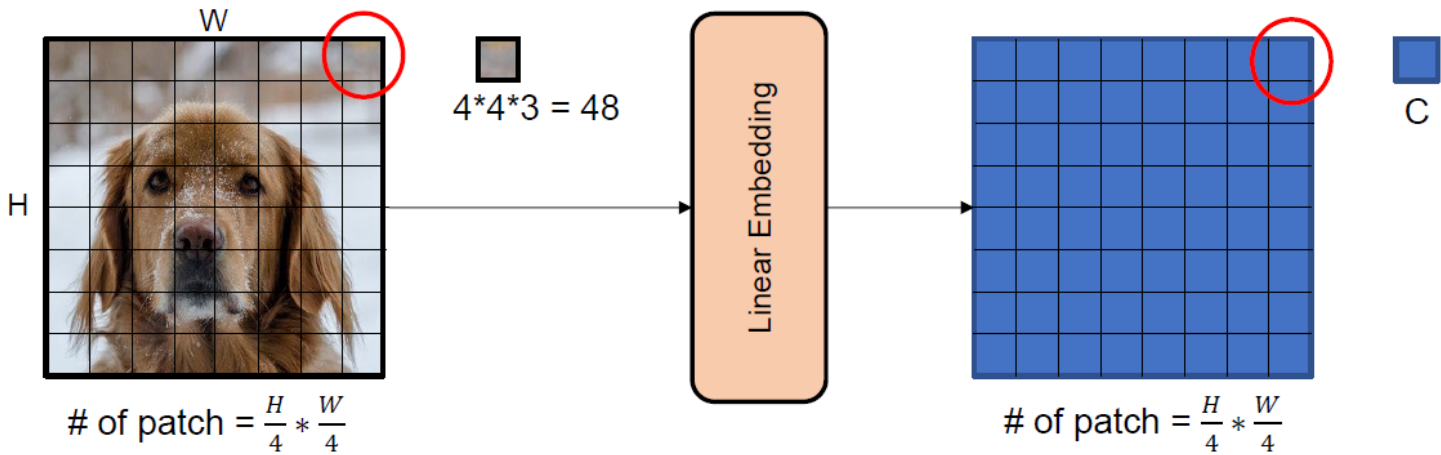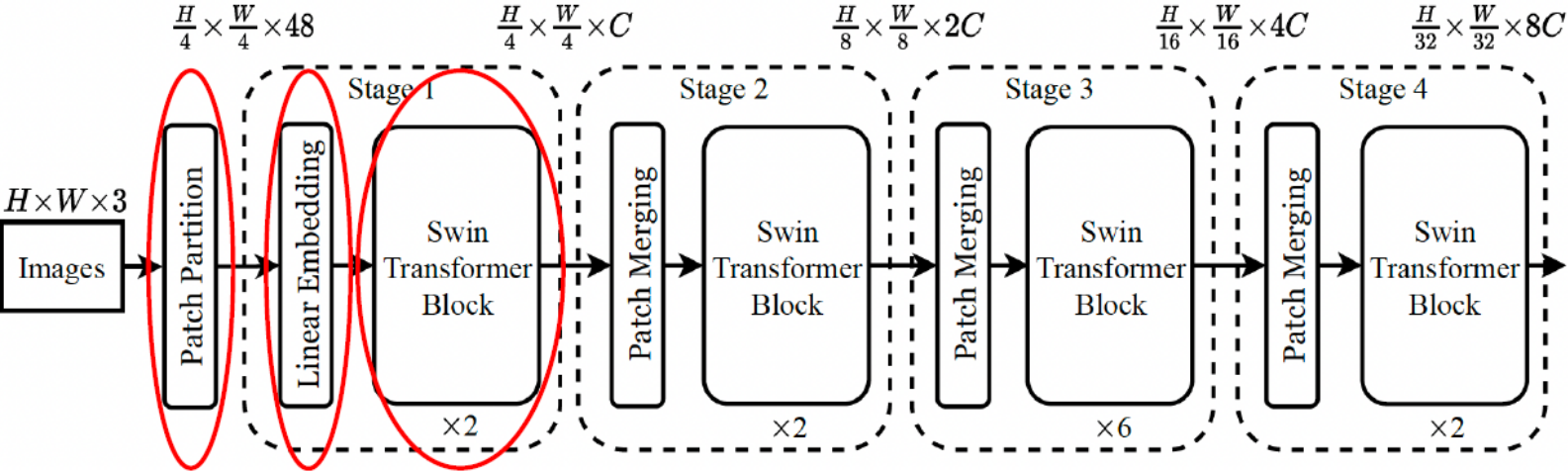
# Swin Architecture (1/6)



Figure credit: CS886 Univ. Waterloo

# Swin Architecture (2/6)

Identical to Transformer but replaced the standard multi-head self-attention (MSA) with:

- **Window MSA (W-MSA)**
- **Shifted Window MSA (SW-MSA)**

# Swin Architecture (3/6)

- Window MSA (W-MSA)
    - Compute attention only within each window
    - Linear complexity (wrt the # of patches) due to the fixed window size
    - What about attention across different windows?



Figure credit: CS886 Univ. Waterloo

# Swin Architecture (4/6)

- Shifted Window MSA (SW-MSA)
    - How to perform across different windows?
    - Shift the window by half the window size (M/2)
    - Additional problem?
      9 instead of 4 windows, plus padding?

# Swin Architecture (5/6)

- Shifted Window MSA (SW-MSA)
  - How to perform across different windows?
  - Shift the window by half the window size (M/2)
  - Additional problem?
    9 instead of 4 windows, plus padding?
  - Introduce *cycle shift*

# Swin Architecture (6/6)

# Swin Transformer Output



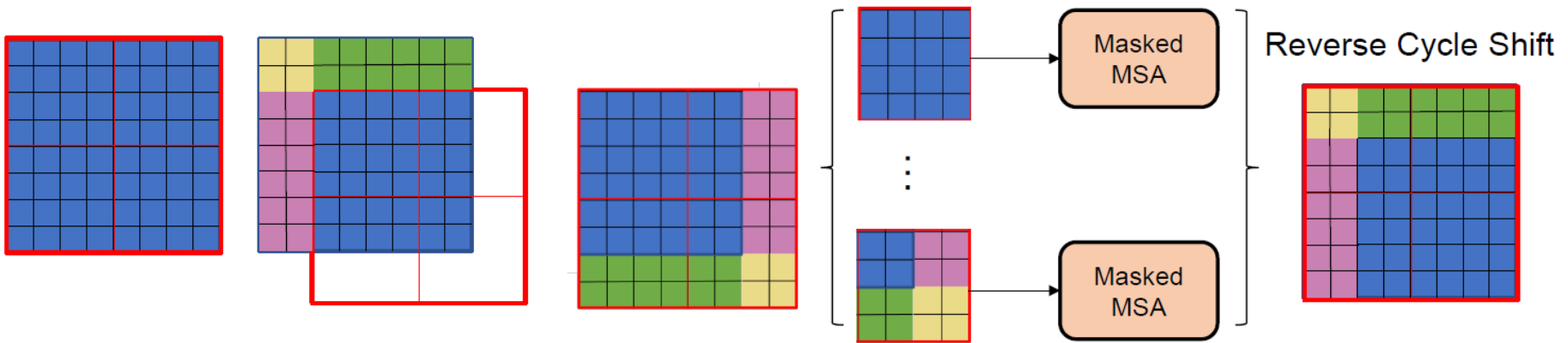$\frac{H}{4} \times \frac{W}{4} \times 48$      $\frac{H}{4} \times \frac{W}{4} \times C$      $\frac{H}{8} \times \frac{W}{8} \times 2C$      $\frac{H}{16} \times \frac{W}{16} \times 4C$      $\frac{H}{32} \times \frac{W}{32} \times 8C$

$H \times W \times 3$

Images → Patch Partition → **Stage 1**: Linear Embedding → Swin Transformer Block ×2 → **Stage 2**: Patch Merging → Swin Transformer Block ×2 → **Stage 3**: Patch Merging → Swin Transformer Block ×6 → **Stage 4**: Patch Merging → Swin Transformer Block ×2

- **Image classification:** Use the last output
- **Object detection and Image segmentation:** Use the output of all the stages

Figure credit: CS886 Univ. Waterloo

# Swin Transformer Results

- Image classification

## (a) Regular ImageNet-1K trained models

| method | image size | #param. | FLOPs | throughput (image / s) | ImageNet top-1 acc. |
|---|---|---|---|---|---|
| RegNetY-4G [48] | $224^2$ | 21M | 4.0G | 1156.7 | 80.0 |
| RegNetY-8G [48] | $224^2$ | 39M | 8.0G | 591.6 | 81.7 |
| RegNetY-16G [48] | $224^2$ | 84M | 16.0G | 334.7 | 82.9 |
| EffNet-B3 [58] | $300^2$ | 12M | 1.8G | 732.1 | 81.6 |
| EffNet-B4 [58] | $380^2$ | 19M | 4.2G | 349.4 | 82.9 |
| EffNet-B5 [58] | $456^2$ | 30M | 9.9G | 169.1 | 83.6 |
| EffNet-B6 [58] | $528^2$ | 43M | 19.0G | 96.9 | 84.0 |
| EffNet-B7 [58] | $600^2$ | 66M | 37.0G | 55.1 | 84.3 |
| ViT-B/16 [20] | $384^2$ | 86M | 55.4G | 85.9 | 77.9 |
| ViT-L/16 [20] | $384^2$ | 307M | 190.7G | 27.3 | 76.5 |
| DeiT-S [63] | $224^2$ | 22M | 4.6G | 940.4 | 79.8 |
| DeiT-B [63] | $224^2$ | 86M | 17.5G | 292.3 | 81.8 |
| DeiT-B [63] | $384^2$ | 86M | 55.4G | 85.9 | 83.1 |
| Swin-T | $224^2$ | 29M | 4.5G | 755.2 | 81.3 |
| Swin-S | $224^2$ | 50M | 8.7G | 436.9 | 83.0 |
| Swin-B | $224^2$ | 88M | 15.4G | 278.1 | 83.5 |
| Swin-B | $384^2$ | 88M | 47.0G | 84.7 | 84.5 |

## (b) ImageNet-22K pre-trained models

| method | image size | #param. | FLOPs | throughput (image / s) | ImageNet top-1 acc. |
|---|---|---|---|---|---|
| R-101x3 [38] | $384^2$ | 388M | 204.6G | - | 84.4 |
| R-152x4 [38] | $480^2$ | 937M | 840.5G | - | 85.4 |
| ViT-B/16 [20] | $384^2$ | 86M | 55.4G | 85.9 | 84.0 |
| ViT-L/16 [20] | $384^2$ | 307M | 190.7G | 27.3 | 85.2 |
| Swin-B | $224^2$ | 88M | 15.4G | 278.1 | 85.2 |
| Swin-B | $384^2$ | 88M | 47.0G | 84.7 | 86.4 |
| Swin-L | $384^2$ | 197M | 103.9G | 42.1 | 87.3 |

Figure credit: CS886 Univ. Waterloo

# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- **Transformer for Visual Analysis**
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - **SSL & Beyond**
- Vision-Language Model
  - Image2Text
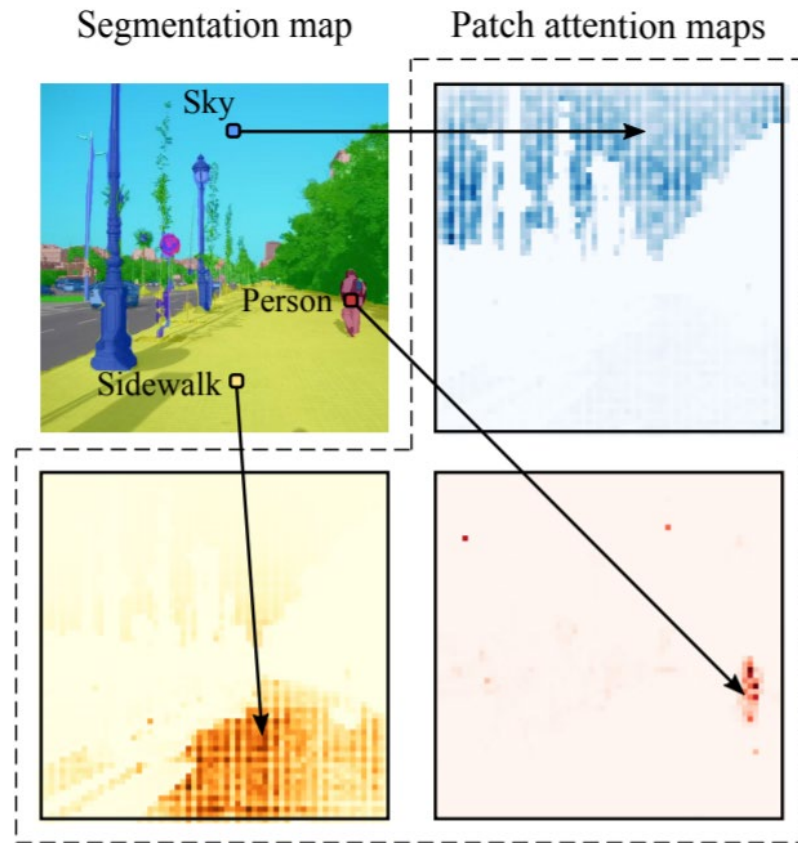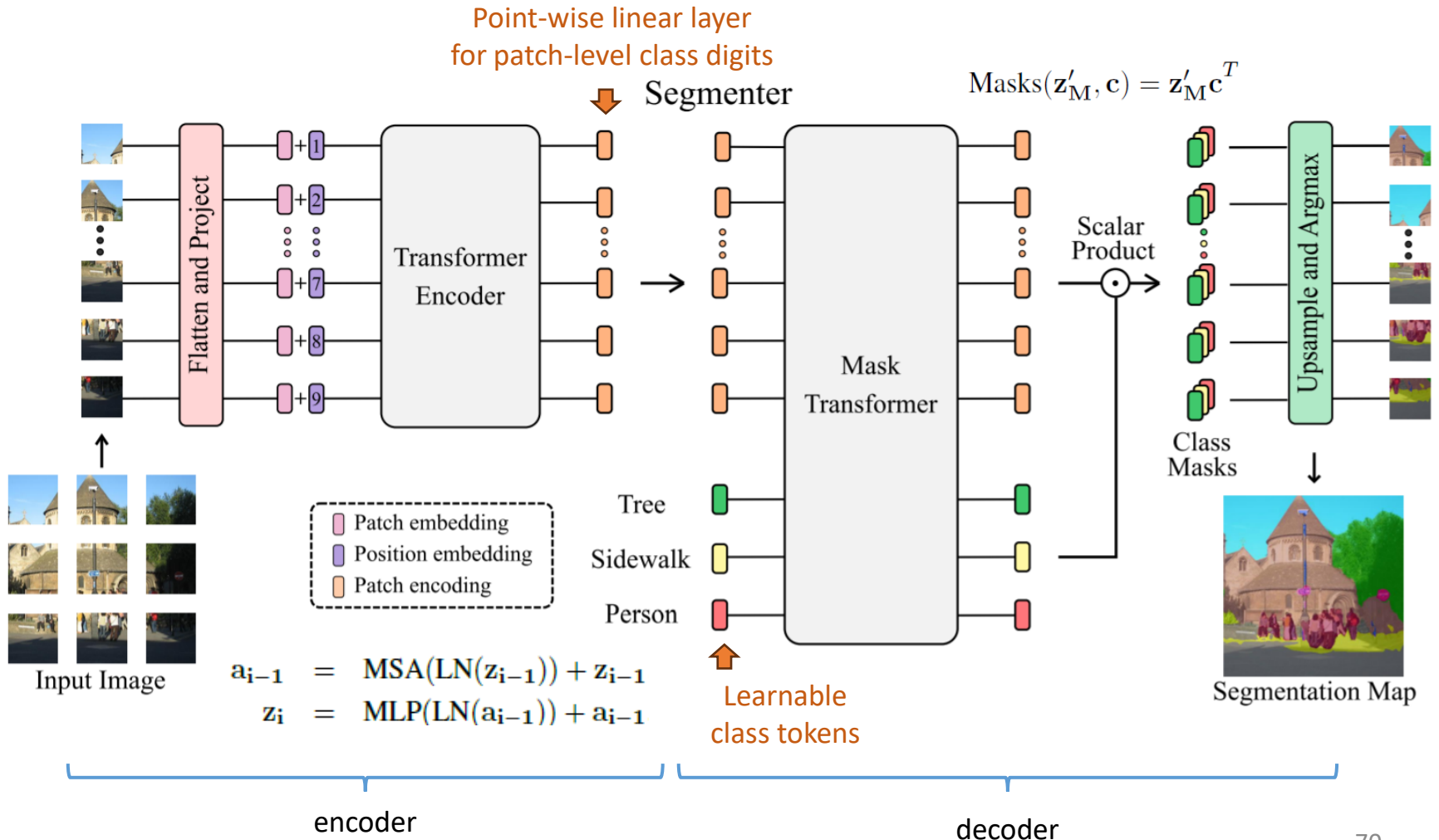  - Text2Image (V)
  - Image-text models

https://medium.com/@navendubrajesh/vision-language-models-use-cases-ee6d54b2c557

# Transformer for Semantic Segmentation

- Segmentation via attention



Strudel et al. "Segmenter: Transformer for Semantic Segmentation." ICCV 2021

# Transformer for Semantic Segmentation (cont'd)

- Inspired by object detection models of DETR (ECCV'20), etc.



Point-wise linear layer for patch-level class digits

$$\text{Masks}(\mathbf{z}'_M, \mathbf{c}) = \mathbf{z}'_M \mathbf{c}^T$$

$$a_{i-1} = \text{MSA}(\text{LN}(z_{i-1})) + z_{i-1}$$
$$z_i = \text{MLP}(\text{LN}(a_{i-1})) + a_{i-1}$$

Learnable class tokens

encoder

decoder

# Example Visualization



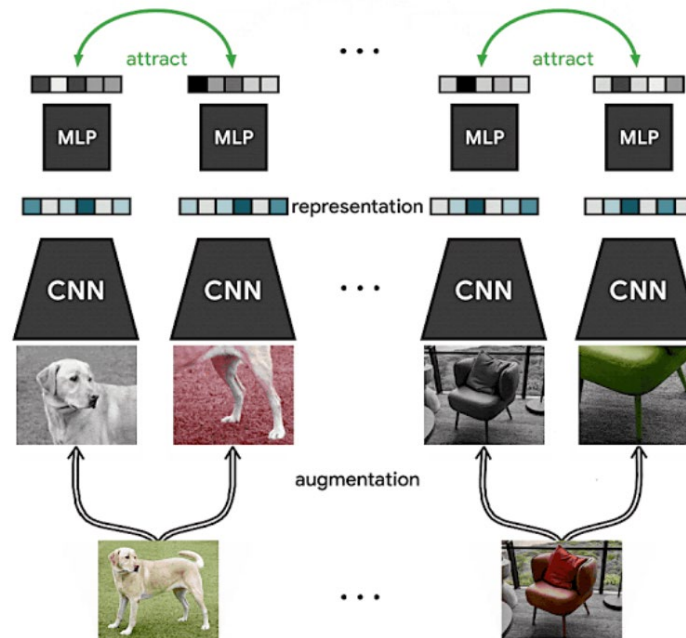(a) Patch size $32 \times 32$    (b) Patch size $16 \times 16$    (c) Patch size $8 \times 8$    (d) Ground Truth
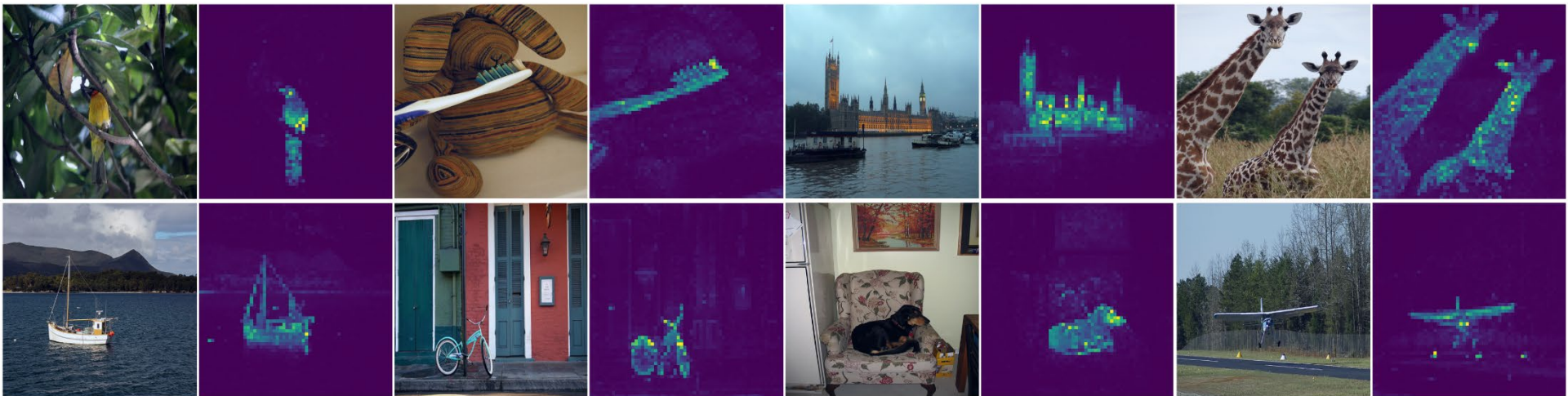
# Self-Supervised Learning (SSL)

- Learning (somewhat) discriminative representations from **unlabeled** data

- Create self-supervised tasks via **data augmentation**

- Recall: SSL for CNN using image data:

Chen et al. "A simple framework for contrastive learning of visual representations." ICML 2020

# Self-Supervised Learning (SSL) for Transformer (cont'd)

- SSL ViT/DeiT features contain info about semantic segmentation

- The above features are excellent k-NN classifiers

- By visualizing self-attention of the CLS token on different heads of the last layer:
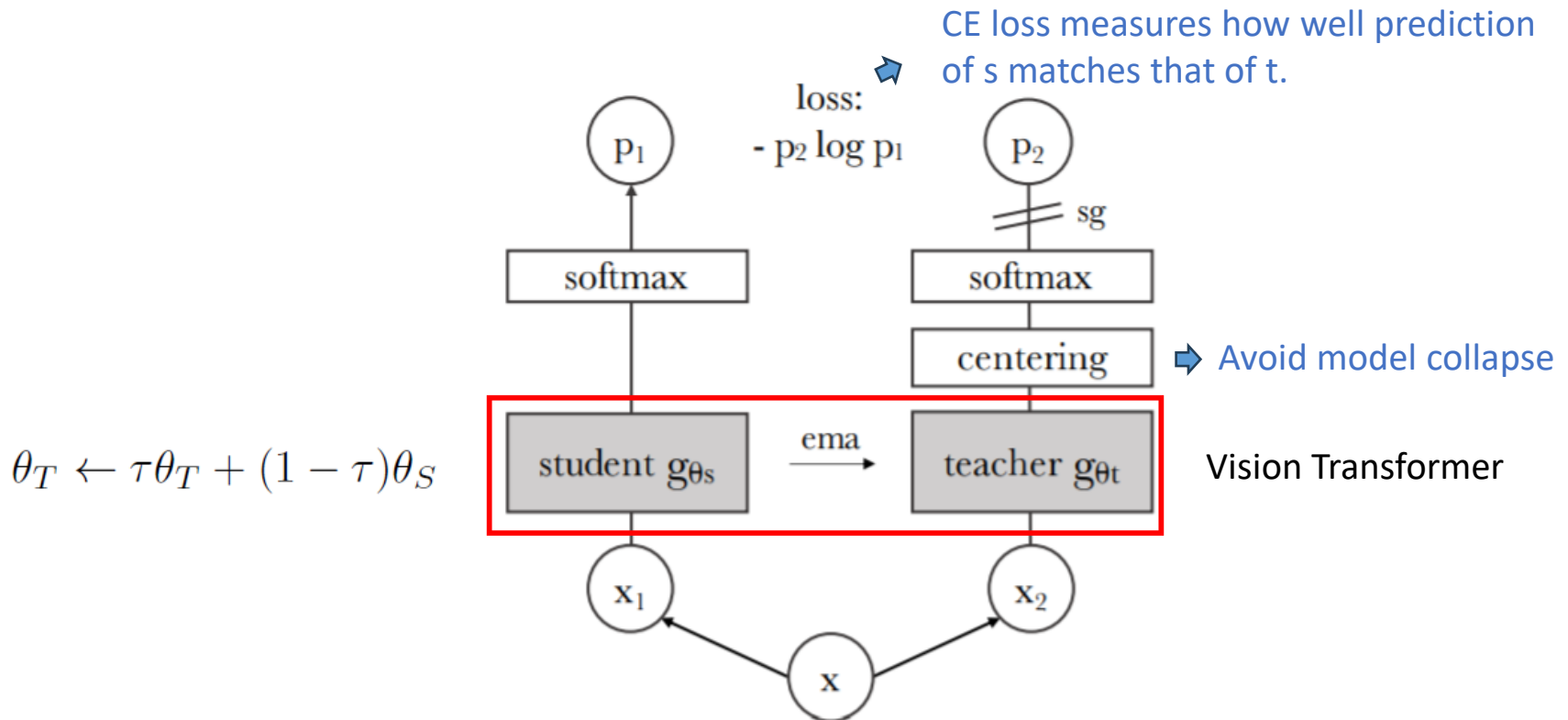
Facebook AI Research & INRIA, "Emerging Properties in Self-Supervised Vision Transformers." ICCV 2021

# Self-Supervised Learning (SSL) for Transformer (cont'd)

- Illustration of the proposed idea:

Facebook AI Research & INRIA, "Emerging Properties in Self-Supervised Vision Transformers." ICCV 2021
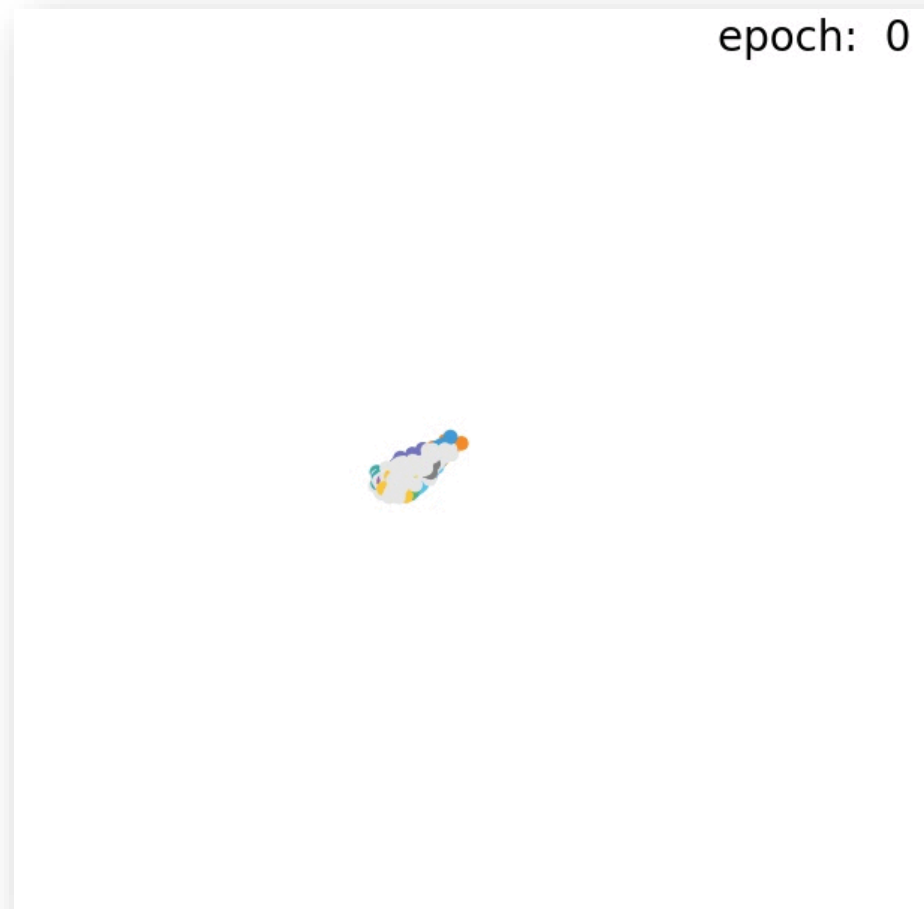
# Self-Distillation with No Labels (DINO)

- Vision Transformer + **SSL**

- Maximize the prediction similarity btw input & its augmented version

- Idea: a **teacher-student** network (i.e., knowledge distillation)
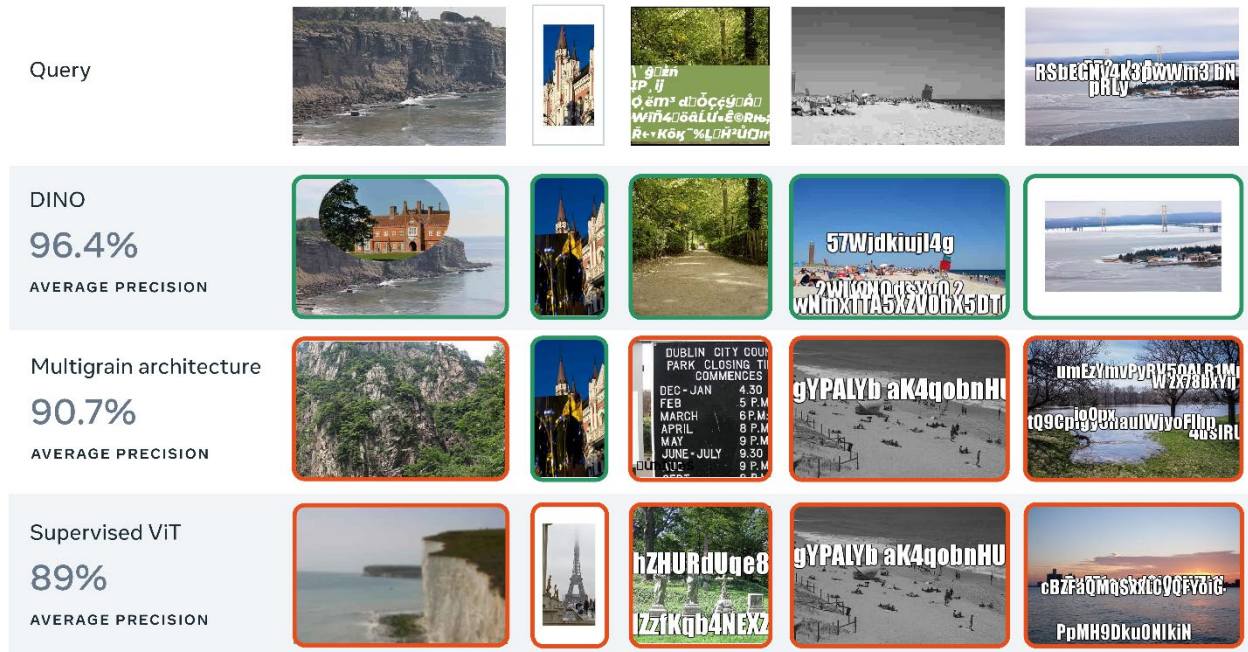  & **EMA** (exponential moving average)

CE loss measures how well prediction of s matches that of t.

loss:
$-p_2 \log p_1$

$$\theta_T \leftarrow \tau\theta_T + (1 - \tau)\theta_S$$

Avoid model collapse

Vision Transformer

Caron et al. "Emerging properties in self-supervised vision transformers." ICCV 2021

75

# Self-Supervised Learning (SSL) for Transformer (cont'd)

- Illustration of the learned features:

epoch: 0

Facebook AI Research & INRIA, "Emerging Properties in Self-Supervised Vision Transformers." ICCV 2021
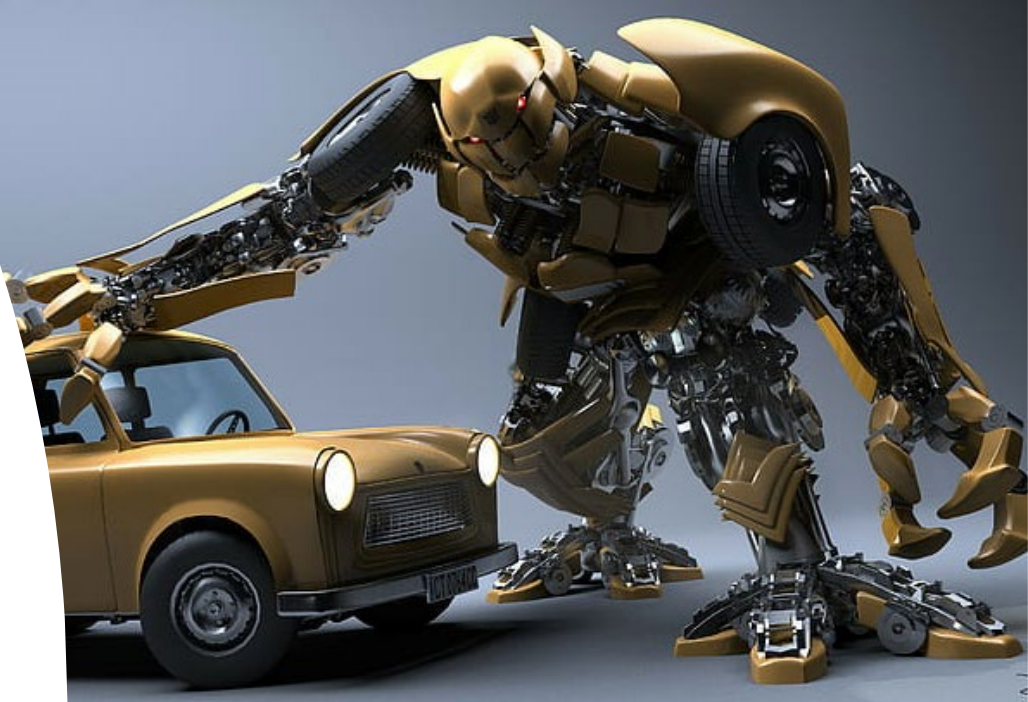
# Highlights & Example Results of DINO

- Learning discriminative representations from **unlabeled** data
- Create self-supervised tasks via **data augmentation**

Facebook AI Research & INRIA, "Emerging Properties in Self-Supervised Vision Transformers." ICCV 2021

# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - SSL & Beyond
- Vision-Language Model
  - Image2Text
  - Text2Image (√)
  - Image-text models

https://medium.com/@navendubrajesh/vision-language-models-use-cases-ee6d54b2c557

Vision Language Model

# *A picture is worth a thousand words...*



- Thing
- Airplane
- Flying airplane in blue sky
- A Lufthansa MD-11 cargo plane in blue sky flying over mountainous terrain

# Vision and Language

- Image-to-Text: Image Captioning

- Text-to-Image: Image Manipulation

- Composed Image Retrieval

- Visual Question Answering (VQA)
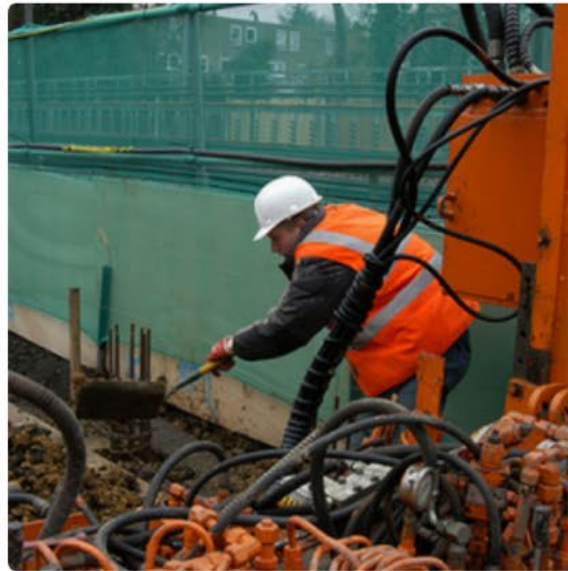
- VQA + Natural Language Explanation
and many more…

- *E.g.,* Question: Is the xiaolongbao fresh at Din Tai Fung?
Answer: Yes.
Explanation: Because the xiaolongbao is made to order at the restaurant.

**Visual Input**

# Image Captioning



Applications: semantics understanding, image-text retrieval,  medical AI, etc.
How does it help GenAI? (e.g., text-to-image generation)

# Image Captioning (cont'd)

- *Image Captioning: Transforming Objects into Words*, Yahoo Research, NeurIPS 2019

- Motivation: mid-level image understanding for captioning
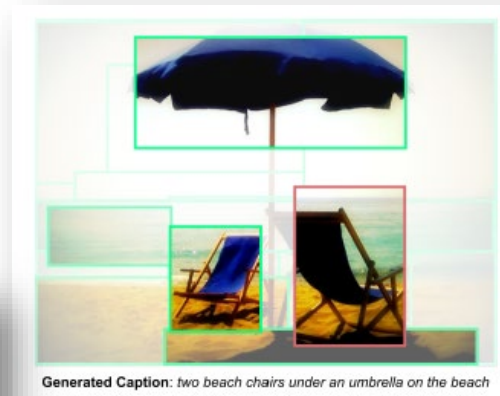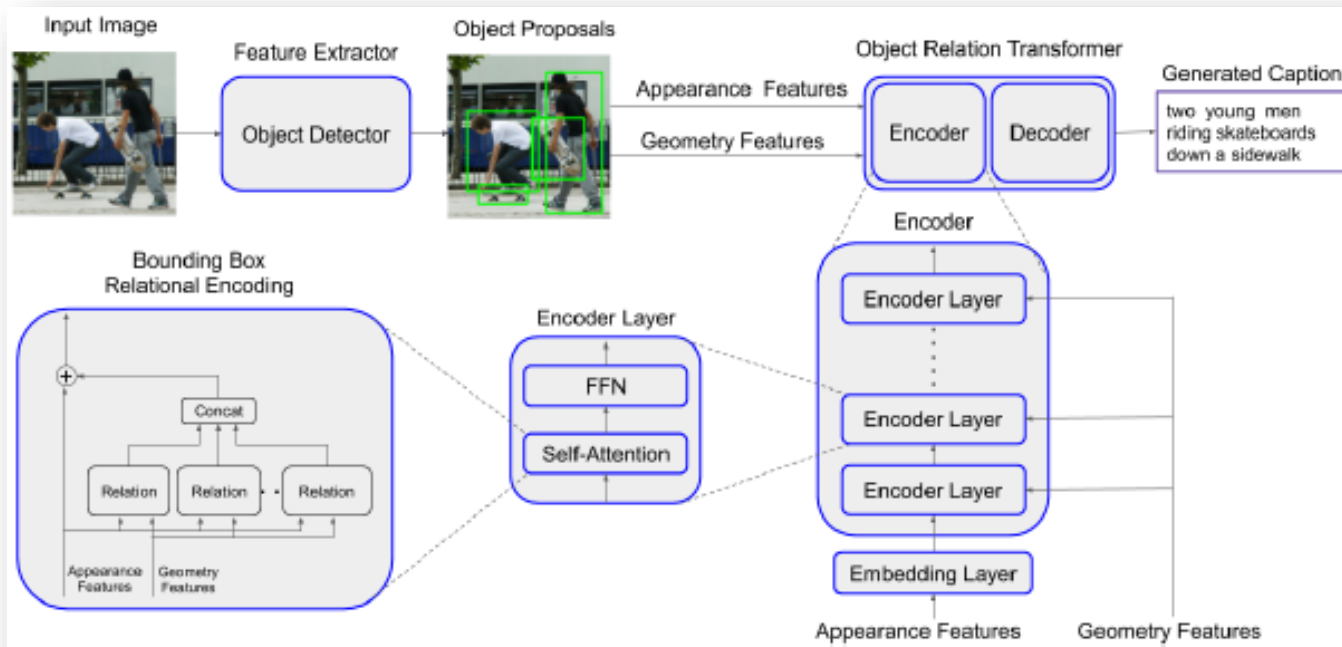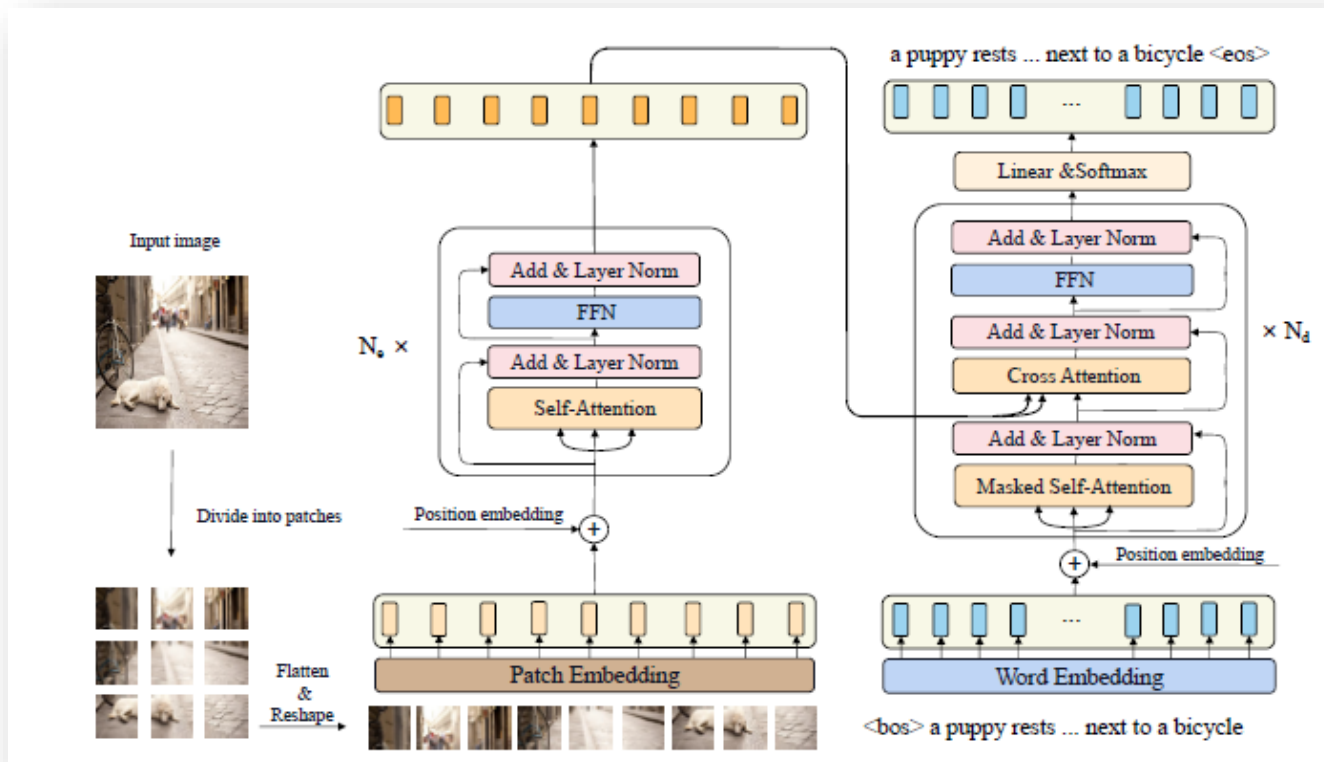
- Framework:



Generated Caption: two beach chairs under an umbrella on the beach

# Image Captioning (cont'd)

- **C**ap**t**ion **T**ransforme**r** (CPTR) -
  *CPTR: Full Transformer Network for Image Captioning,* CAS, arxiv 2021

- Motivation: patch translation for image captioning

# Remarks & Extension:

- Training a captioning model requires a large amount of image-caption data pairs

- Image captioning in the wild:
    - Describing images with novel content during inference
    - For example, COCO dataset has 80 object categories.
      How to generalize captioning models to Open Image (w/ 600 classes)?

- Domain-specific image captioning:
    - From general-purpose captioning to task-oriented captioning -> **finetuning?**



COCO (80 classes)

Two pug **dogs** sitting on a **bench** at the beach.

A **child** is sitting on a **couch** and holding an **umbrella**.



Open Images (600 classes)

goat          artichoke     accordion

dolphin       waffle        balloon
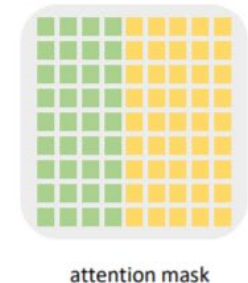
# Image Captioning *in the Wild*

- **Novel Object Captioning (NOC)**
  - Training with captioned and uncaptioned data
    captioned data: labeled image data with captions (e.g., COCO)
    uncaptioned data: only labels of novel classes available (e.g., Open Images)
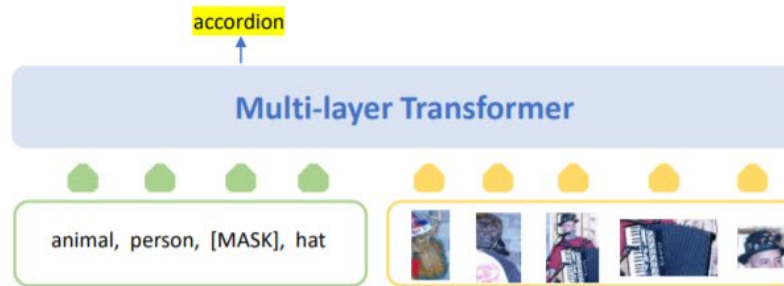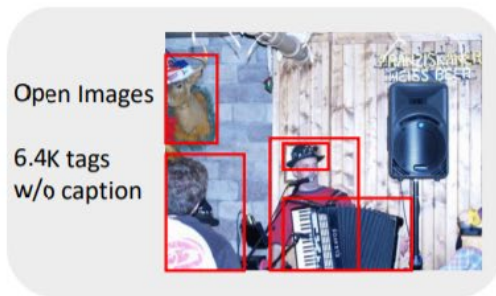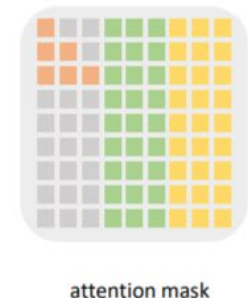

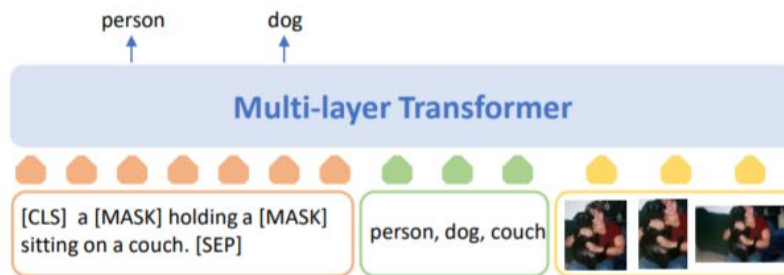
**We have image-caption data**



**Data with labels for <u>novel objects</u>
but w/o captions**

# Novel Object Captioning (cont'd)

- "**Paraphrasing Is All You Need for Novel Object Captioning**", NTU, NeurIPS'22

- **VIVO**: **Vi**sual **Vo**cabulary Pre-Training for Novel Object Caption Captioning, Microsoft, AAAI'21
    - Pre-training: uncaptioned image data containing novel class labels
    - Fine-tuning: (a limited amount of) image data with class labels & descriptions



(a) Pre-training: learn visual vocabulary

(b) Fine-tuning: learn sentence description

# Novel Object Captioning (cont'd)

- **VIVO**: **Vi**sual **Vo**cabulary Pre-Training for Novel Object Caption Captioning (AAAI'21)
  - Pre-training: uncaptioned image data containing novel class labels
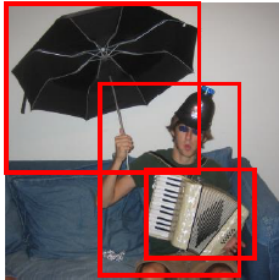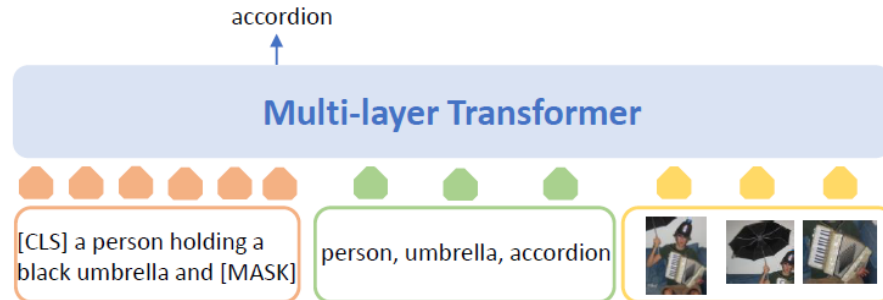  - Fine-tuning: (a limited amount of) image data with class labels & descriptions
  - Inference:
    - Inputs: image (with region features & tags) & [CLS]
    - Output: caption


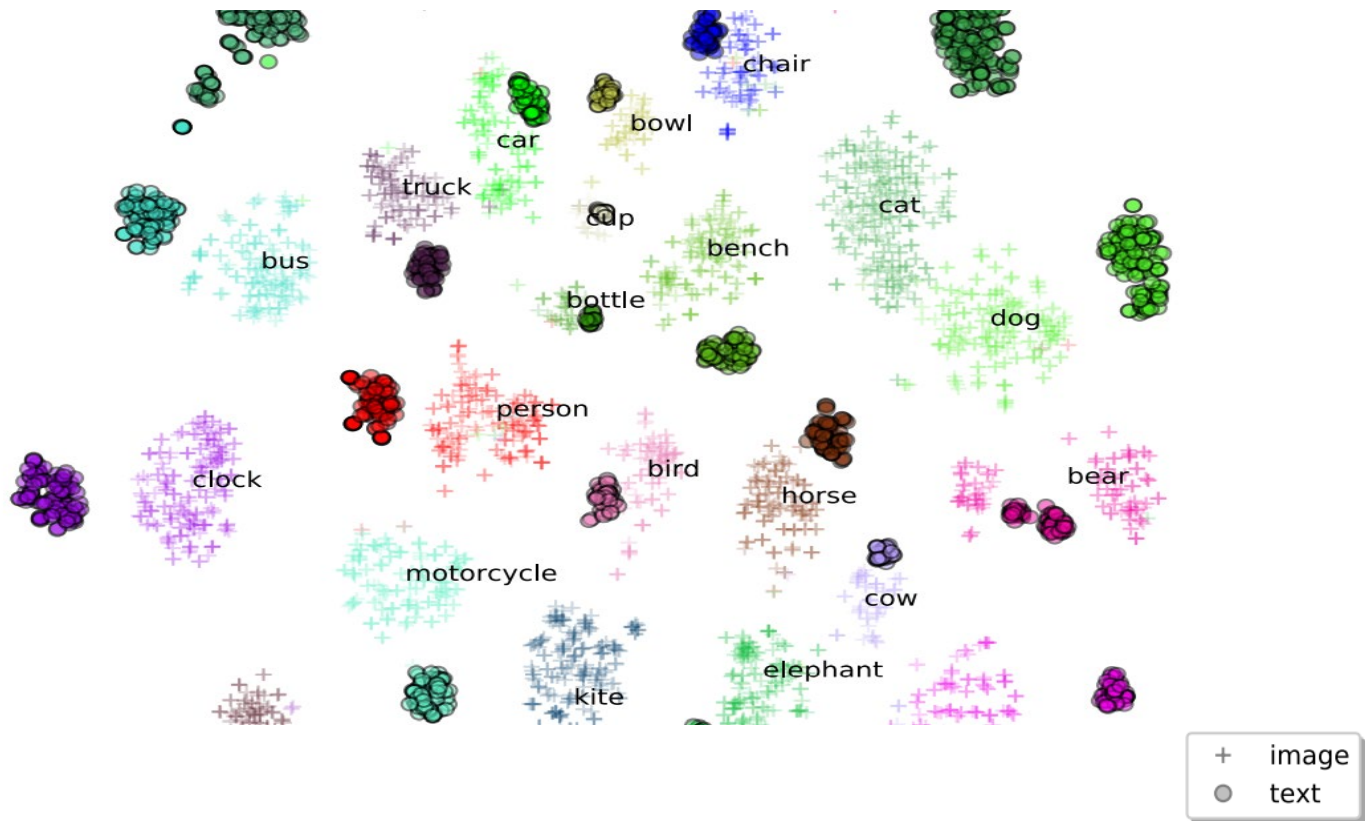
**(c) Inference: novel object captioning**

accordion

**Multi-layer Transformer**

[CLS] a person holding a black umbrella and [MASK]

person, umbrella, accordion

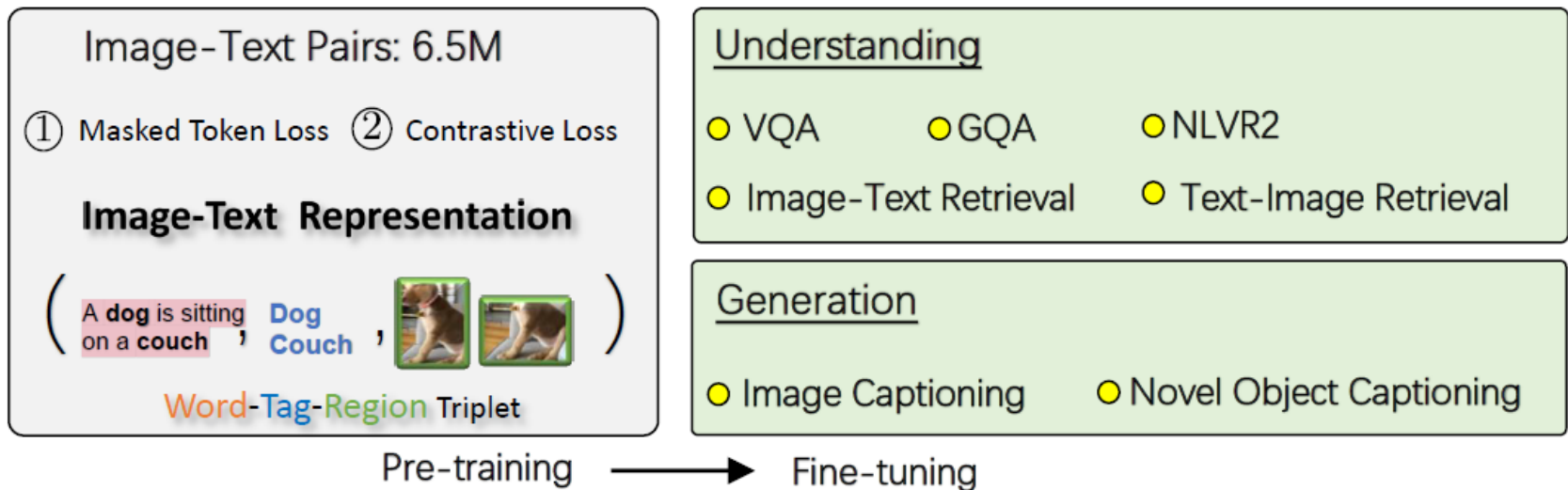A person holding a black umbrella and accordion.

# Novel Object Captioning (cont'd)

- VIVO: Visual Vocabulary Pre-Training for Novel Object Caption Captioning
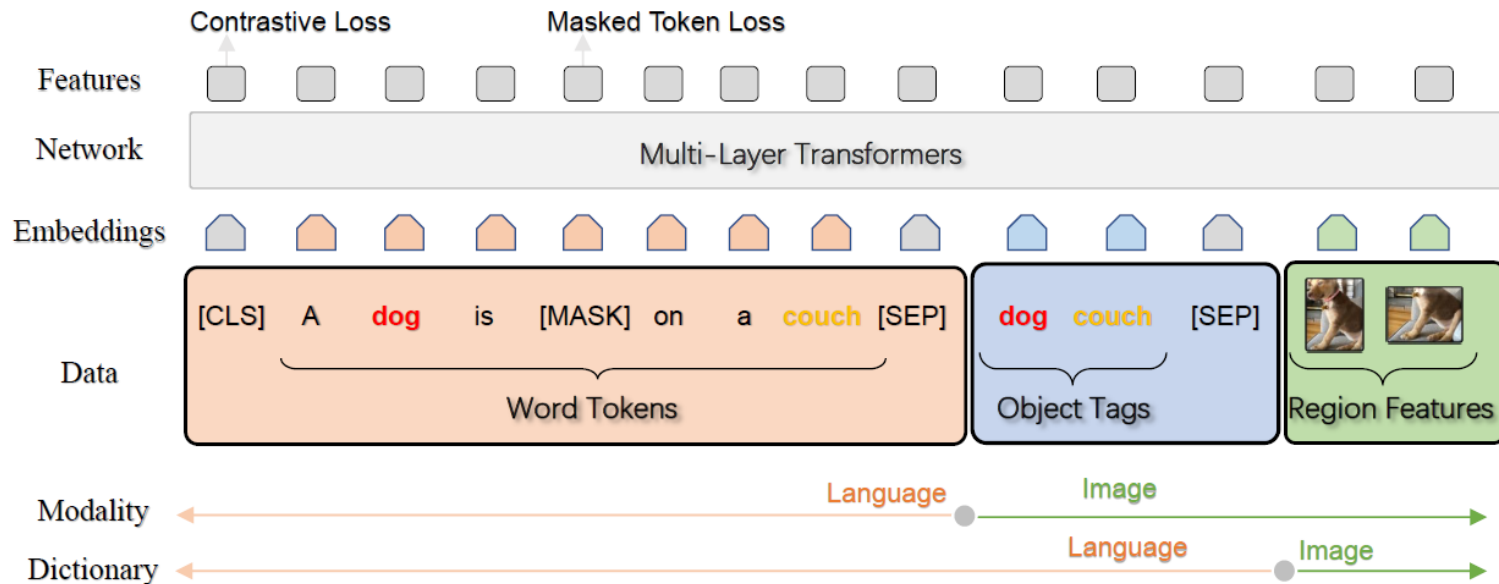  - Visualization image-text alignment

# Beyond Image Captioning:
# Unified Vision & Language Model

- **Oscar**: **O**bject-**S**emanti**c**s **A**ligned **P**re-training for Vision-Language Tasks, Microsoft, ECCV'20
  - Training data:
    triplets of caption-tag-region
  - Objectives:
    1. Masked token loss for words & tags
    2. Contrastive loss tags and others
  - Fine-tuning:
    5 vision & language tasks (VQA, image-text retrieval, image captioning, NOC, etc.)
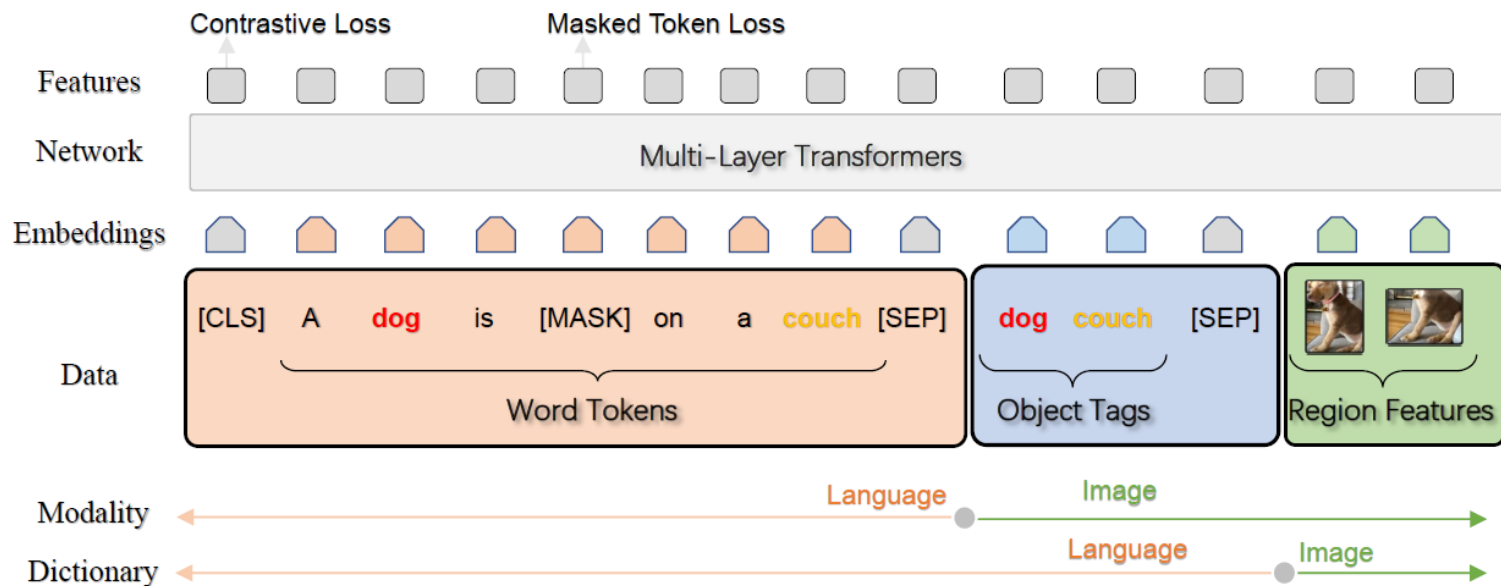
# Semantics-Aligned Pre-training for V+L Tasks

- **Oscar**: **O**bject-**S**emantics **A**ligned **P**re-training for Vision-Language Tasks
  - Training:
    - Inputs: triplets of caption-tag-region
    - Objectives: Masked token loss for words & tags + Contrastive loss tags and others
  - Fine-tuning:
    5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)

# Semantics-Aligned Pre-training for V+L Tasks (cont'd)

- **Oscar**: **O**bject-**S**emanti**c**s **A**ligned **P**re-training for Vision-Language Tasks (ECCV'20)
  - Training:
    - Inputs: triplets of word-tag-region
    - Objectives: Masked token loss for words & tags + Contrastive loss tags and others
  - Fine-tuning:
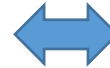    - 5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)
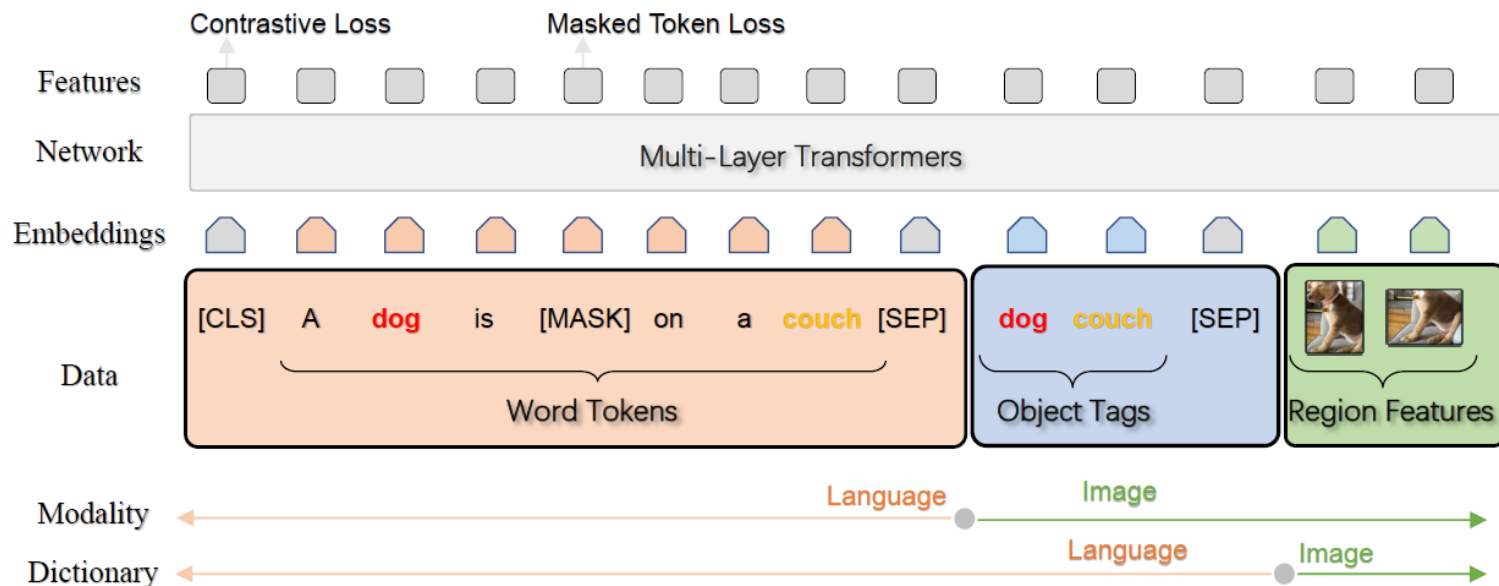
*Holding an apple* ⟷  or 

- **Oscar** (cont'd)
  - Fine-tuning:
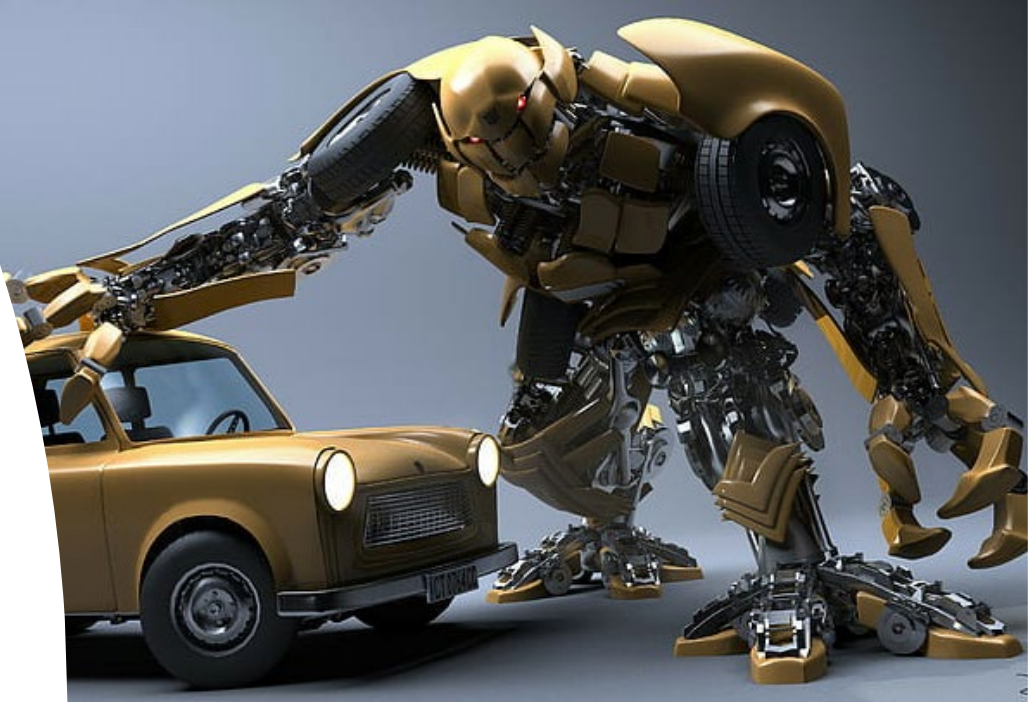    5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)
  - Take **image-text retrieval** as an example
    - Training: aligned/mis-aligned image-text pairs as positive/negative input pairs,
      with **[CLS]** for binary classification (1/0)
    - Inference: for either image or text retrieval,
      calculate <u>classification score</u> of **[CLS]** for the input query
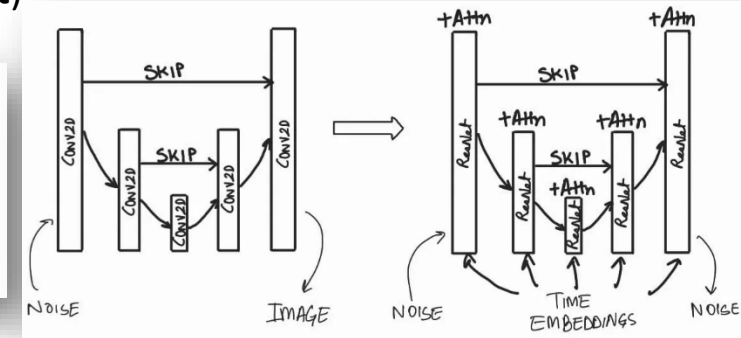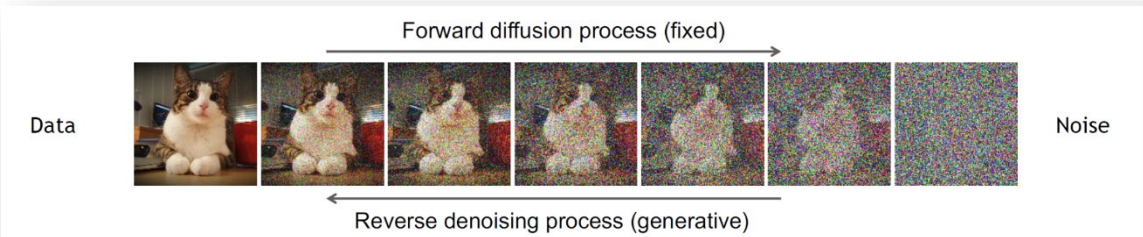
# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - SSL & Beyond
- **Vision-Language Model**
  - Image2Text
  - **Text2Image (√)**
  - Image-text models

https://medium.com/@navendubrajesh/vision-language-models-use-cases-ee6d54b2c557

# Recap: **D**enoising **D**iffusion **P**robabilistic **M**odels (DDPM)

- Training:
  - Forward/reverse diffusion & denoising process
    - learns to generate/restore data by denoising
    - typically implemented via a **conditional U-net**)



Forward diffusion process (fixed)

Data — Noise

Reverse denoising process (generative)

- Pseudo Code for Training/Inference (Sampling):

**Algorithm 1** Training

1: **repeat**
2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3: $t \sim \text{Uniform}(\{1, \ldots, T\})$
4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
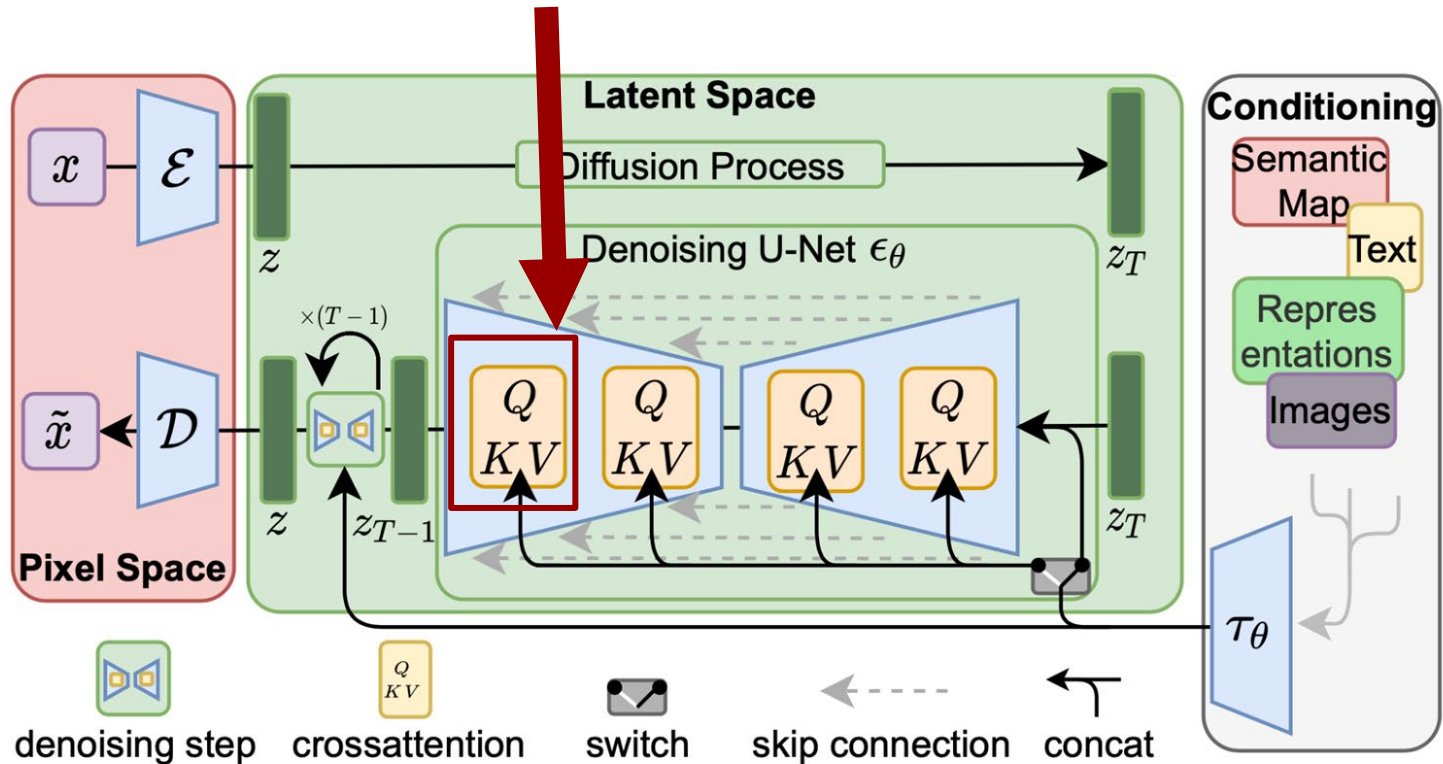5: Take gradient descent step on
$$\nabla_\theta \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon}, t) \right\|^2$$
6: **until** converged

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\right) + \sigma_t \mathbf{z}$
5: **end for**
6: **return** $\mathbf{x}_0$

Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020
Song et al., Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021
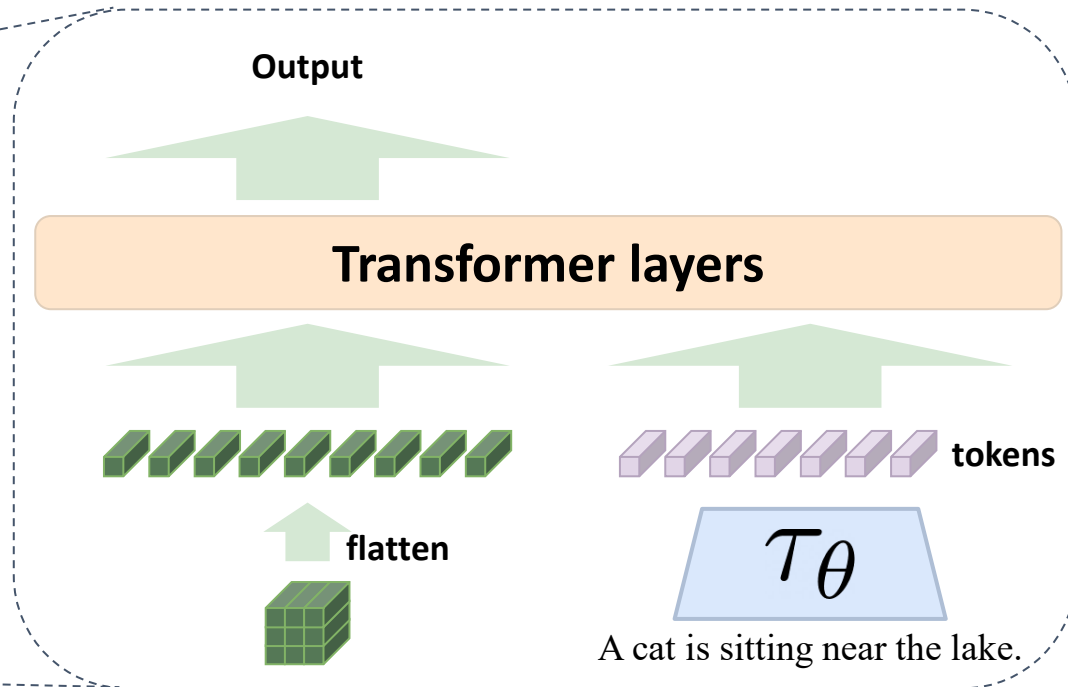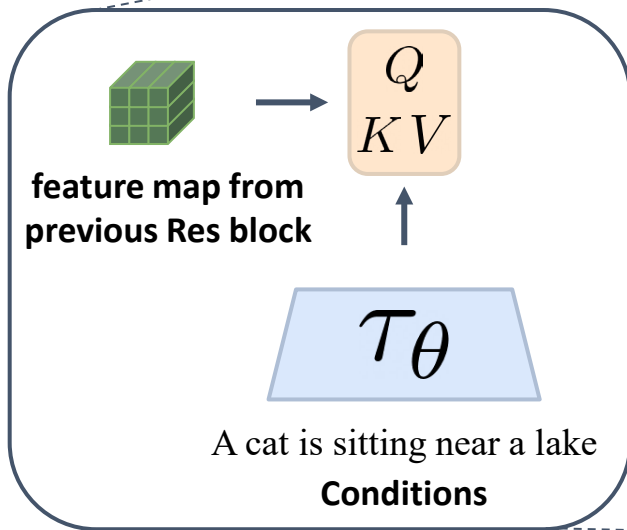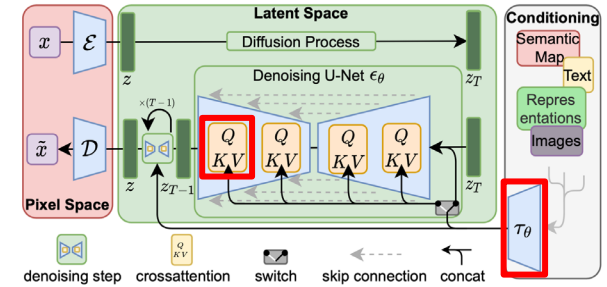
# From Unconditional to Conditional Latent Diffusion Model

- Latent diffusion model (LDM), CVPR'22: DDPM in latent space
- Condition mechanism: cross-attention at transformer layers

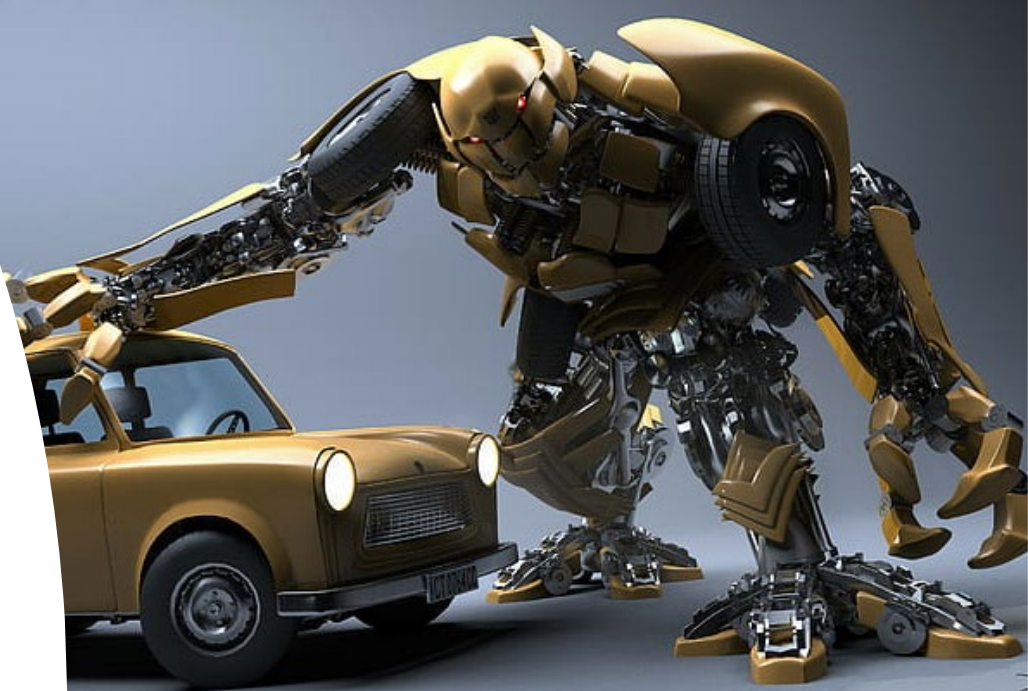# From Unconditional to Conditional Latent Diffusion Model (cont'd)



- **Condition mechanism:** using transformer layers
- $\tau_\theta$ is the embedding module for conditions
  e.g., BERT, CLIP text embedding, etc.



feature map from previous Res block

$Q$

$KV$

$\tau_\theta$

A cat is sitting near a lake

**Conditions**

Output

**Transformer layers**

flatten

tokens

$\tau_\theta$

A cat is sitting near the lake.

# What to Be Covered?

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - SSL & Beyond
- **Vision-Language Model**
  - Image2Text
  - Text2Image (√)
  - **Image-text models**



https://medium.com/@navendubrajesh/vision-language-models-use-cases-ee6d54b2c557
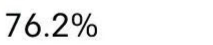
# CLIP: Contrastive Language-Image Pretraining

- OpenAI, *Learning Transferable Visual Models From Natural Language Supervision*, NeurIPS WS 2021 (w/ 9000+ citations)

- Why DL/CNN not good enough?
    - Require annotated data for training image classification
    - Domain gap between closed-world and open-world domain data
    - Lack of ability for zero-shot classification



...Net — 76.2%

geNet V2 — 64.3%

...nageNet Rendition — 37.7%

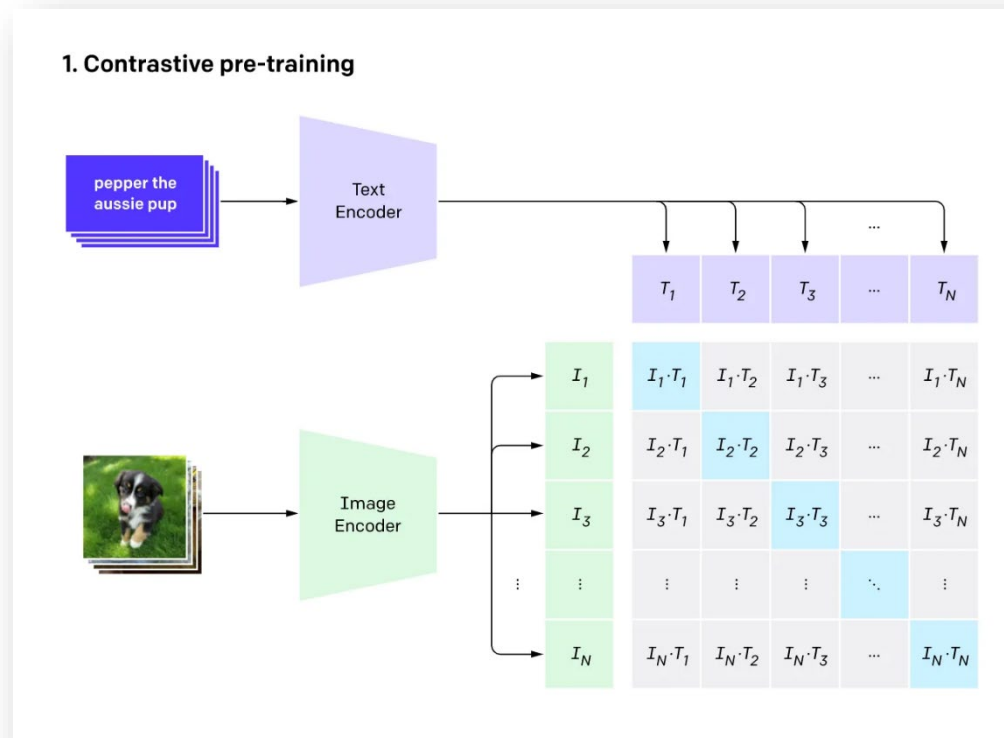...jectNet — 32.6%

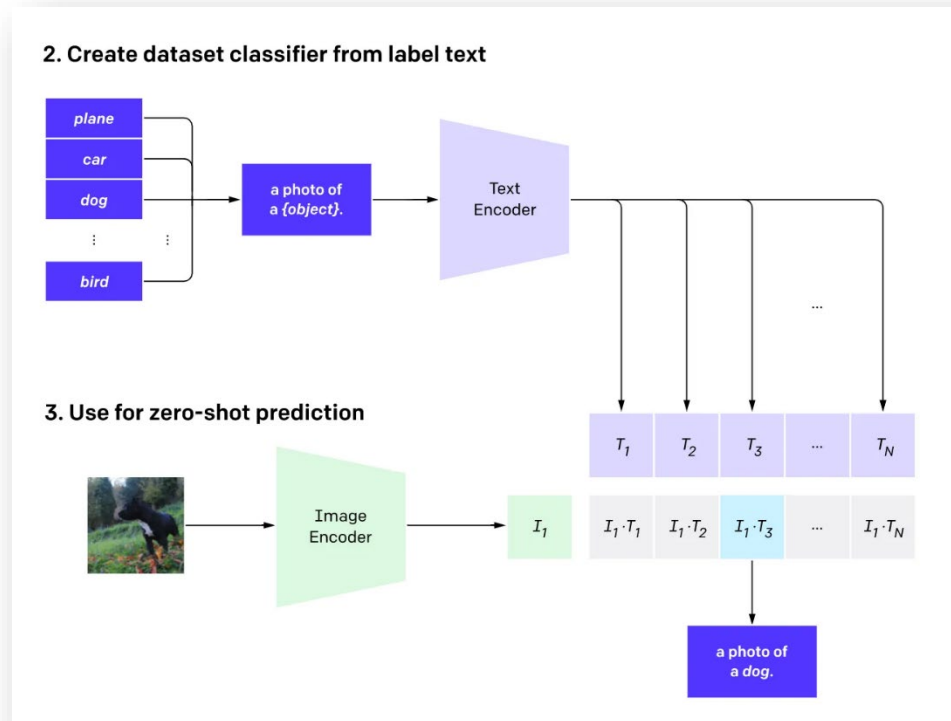...geNet Sketch — 25.2%

...Net Adversarial — 2.7%

# CLIP (cont'd)

- Why DL/CNN not good enough?
    - Require annotated data for training image classification
    - Domain gap between closed-world and open-world domain data
    - Lack of ability for zero-shot classification
- Motivation/Objectives
    - Cross-domain contrastive learning from large-scale image-language data
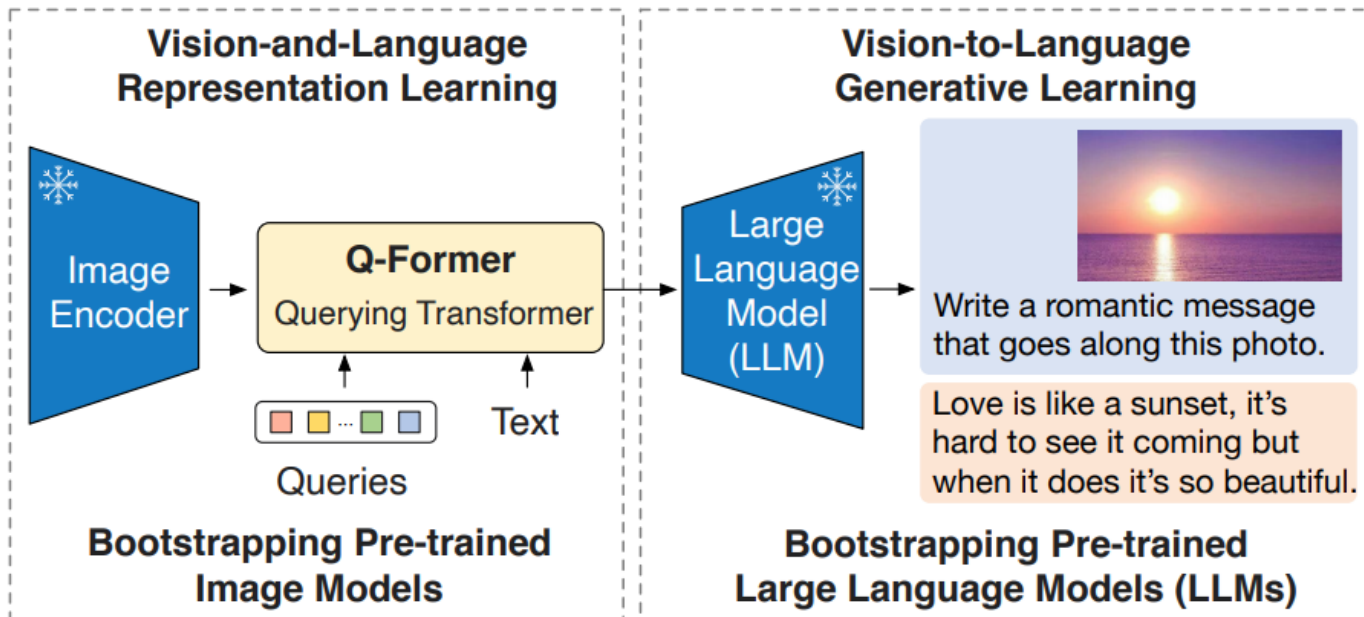
# CLIP (cont'd)

- (Zero-shot) Inference:



- Limitation
  - Fine-grained description ability
  - Any examples?

# BLIP-2 (ICML'23)

- **BLIP:**
  Bootstrapping Language-Image Pre-training for Unified Vision-Language
  Understanding and Generation, Salesforce Research, NeurIPS 2021
- **Goal:**
  Bridge the modality gap between off-the-shelf **frozen pre-trained image encoders** and
  **frozen large language models** with a lightweight Querying Transformer (Q-Former).
- **Advantages:**
  1. No need to train from scratch
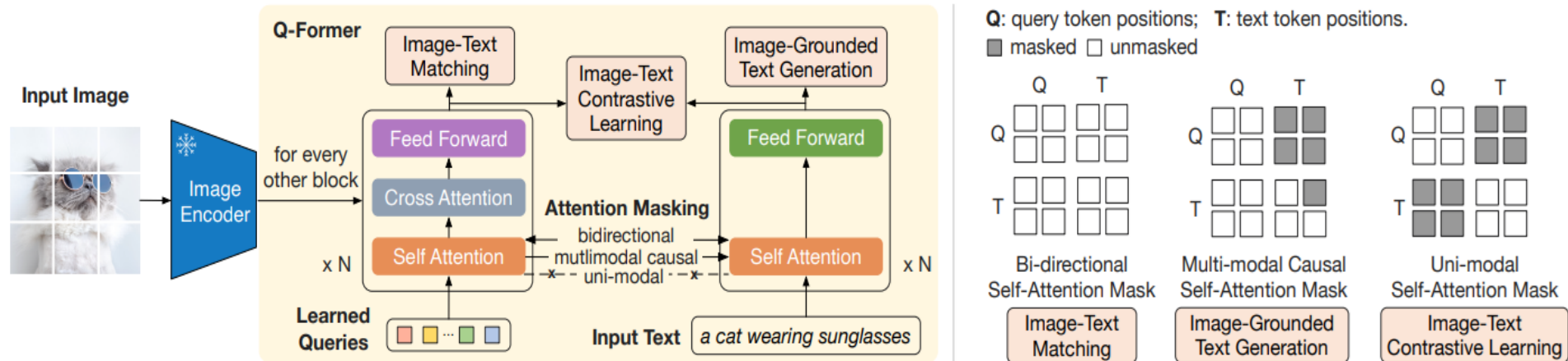  2. Avoid catastrophic forgetting (w/ fixed VLM & LLM)

# Pre-training

- A **two-stage** pre-training strategy
  - **Stage 1**: Representation Learning
    - enforce **Q-Former** to learn **visual representation**
      that is most relevant to the text description
  - **Stage 2**: Generative Learning
    - make the output representation of **Q-Former** to be understood by **LLM**

# Pre-training Stage 1 - VL Representation Learning

- **Goal:** enforce **Q-Former** to extract visual representation relevant to text
- **Method:** three pre-training tasks
  - **Image-Text Contrastive Learning (ITC)**:
    self-attn in Q/T, followed by max (sim(Q, T)) -> can be viewed as CLIP training
  - **Image-Text Matching (ITM)**:
    for each learnable query -> linear classifier for binary decision
  - **Image-grounded Text Generation (ITG)**:
    self-attn in Q for encoder training; T->Q for image-to-text generation

# Pre-training Stage 2 - VL Generative Learning

- **Goal:**
  Learning with LLM guidance
  i.e., make the output representation of **Q-Former** to be understood by **LLMs**.
- **Method:**
  pre-training with Image-grounded Text Generation (ITG)

# Quantitative Results

- Comparison on zero-shot visual question answering (VQA)

| Models | #Trainable Params | #Total Params | VQAv2 | | OK-VQA | GQA |
|---|---|---|---|---|---|---|
| | | | val | test-dev | test | test-dev |
| VL-T5$_{no-vqa}$ | 224M | 269M | 13.5 | - | 5.8 | 6.3 |
| FewVLM (Jin et al., 2022) | 740M | 785M | 47.7 | - | 16.5 | 29.3 |
| Frozen (Tsimpoukelli et al., 2021) | 40M | 7.1B | 29.6 | - | 5.9 | - |
| VLKD (Dai et al., 2022) | 406M | 832M | 42.6 | 44.5 | 13.3 | - |
| Flamingo3B (Alayrac et al., 2022) | 1.4B | 3.2B | - | 49.2 | 41.2 | - |
| Flamingo9B (Alayrac et al., 2022) | 1.8B | 9.3B | - | 51.8 | 44.7 | - |
| Flamingo80B (Alayrac et al., 2022) | 10.2B | 80B | - | 56.3 | **50.6** | - |
| BLIP-2 ViT-L OPT$_{2.7B}$ | 104M | 3.1B | 50.1 | 49.7 | 30.2 | 33.9 |
| BLIP-2 ViT-g OPT$_{2.7B}$ | 107M | 3.8B | 53.5 | 52.3 | 31.7 | 34.6 |
| BLIP-2 ViT-g OPT$_{6.7B}$ | 108M | 7.8B | 54.3 | 52.6 | 36.4 | 36.4 |
| BLIP-2 ViT-L FlanT5$_{XL}$ | 103M | 3.4B | 62.6 | 62.3 | 39.4 | 44.4 |
| BLIP-2 ViT-g FlanT5$_{XL}$ | 107M | 4.1B | 63.1 | 63.0 | 40.7 | 44.2 |
| BLIP-2 ViT-g FlanT5$_{XXL}$ | 108M | 12.1B | **65.2** | **65.0** | 45.9 | **44.7** |

# Quantitative Results (cont'd)

- Comparison on image captioning

| Models | #Trainable Params | NoCaps Zero-shot (validation set) | | | | | | | | COCO Fine-tuned Karpathy test | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | in-domain | | near-domain | | out-domain | | overall | | | |
| | | C | S | C | S | C | S | C | S | B@4 | C |
| OSCAR (Li et al., 2020) | 345M | - | - | - | - | - | - | 80.9 | 11.3 | 37.4 | 127.8 |
| VinVL (Zhang et al., 2021) | 345M | 103.1 | 14.2 | 96.1 | 13.8 | 88.3 | 12.1 | 95.5 | 13.5 | 38.2 | 129.3 |
| BLIP (Li et al., 2022) | 446M | 114.9 | 15.2 | 112.1 | 14.9 | 115.3 | 14.4 | 113.2 | 14.8 | 40.4 | 136.7 |
| OFA (Wang et al., 2022a) | 930M | - | - | - | - | - | - | - | - | **43.9** | <u>145.3</u> |
| Flamingo (Alayrac et al., 2022) | 10.6B | - | - | - | - | - | - | - | - | - | 138.1 |
| SimVLM (Wang et al., 2021b) | ~1.4B | 113.7 | - | 110.9 | - | 115.2 | - | 112.2 | - | 40.6 | 143.3 |
| BLIP-2 ViT-g OPT$_{2.7B}$ | 1.1B | <u>123.0</u> | <u>15.8</u> | 117.8 | <u>15.4</u> | 123.4 | **15.1** | 119.7 | <u>15.4</u> | <u>43.7</u> | **145.8** |
| BLIP-2 ViT-g OPT$_{6.7B}$ | 1.1B | **123.7** | <u>15.8</u> | <u>119.2</u> | 15.3 | <u>124.4</u> | 14.8 | <u>121.0</u> | 15.3 | 43.5 | 145.2 |
| BLIP-2 ViT-g FlanT5$_{XL}$ | 1.1B | **123.7** | **16.3** | **120.2** | **15.9** | **124.8** | **15.1** | **121.6** | **15.8** | 42.4 | 144.5 |

C: CIDEr    S: SPICE    B@4:BLEU@4

# Visualization

- Instructed zero-shot image-to-text generation examples



**Explain the advantages of this product.**

The audi e-tron quattro concept is a plug-in hybrid electric sports car that has a range of 310 miles and can accelerate from 0 to 60 mph in just 4 seconds.

**Tell me something about the history of this place.**

The great wall of china was built in 221 bc by emperor qin shihuang to protect the emperor's capital from invaders from the north.

**Write down the facts that you know about this flower.**

It is a genus of flowering plants in the orchid family, orchidaceae. It is native to eastern asia and south east asia. It is a perennial flowering plant that grows up to a height of 1 m.

**Is this photo unusual?**

Yes, it's a house that looks like it's upside down.

**How could someone get out of the house?**

It has a slide on the side of the house.

**What are shown in the photo?**

A man and a chicken.

**What does the man feel and why?**

He is scared of the chicken because it is flying at him.

**What are the ingredients I need to make this?**

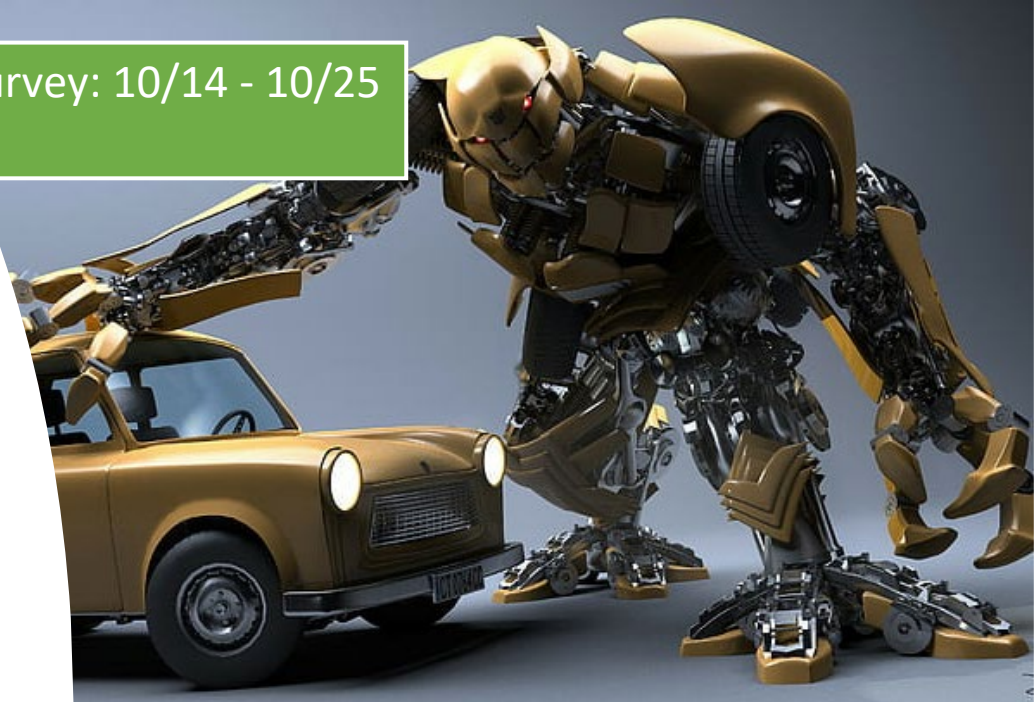Pepperoni, mozzarella cheese, pizza sauce, olive oil, salt, pepper, basil.

**What is the first step?**

Place the pizza dough on a baking sheet, brush with olive oil, sprinkle with salt, pepper, and basil.

# What We have Covered Today

- Transformer
  - Self-Attention
  - Cross-Attention
  - Positional Embedding
- Transformer for Visual Analysis
  - Vision Transformer (ViT)
  - DeiT & Swin Transformer
  - SSL & Beyond
- Vision-Language Model
  - Image2Text
  - Text2Image (√)
  - Image-text models

https://medium.com/@navendubrajesh/vision
-language-models-use-cases-ee6d54b2c557

Vision
Language
Model