# Vision Transformer (ViT)

Arman Afrasiyabi

PHD Student at Université Laval

# AN IMAGE IS WORTH 16x16 WORDS:
# TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

**Alexey Dosovitskiy**[*,†]**, Lucas Beyer**[*]**, Alexander Kolesnikov**[*]**, Dirk Weissenborn**[*]**,**
**Xiaohua Zhai**[*]**, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,**
**Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby**[*,†]

[*]equal technical contribution, [†]equal advising
Google Research, Brain Team
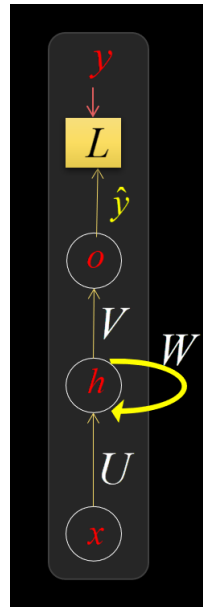`{adosovitskiy, neilhoulsby}@google.com`

# Introduction

First, they project the input image to an embedding, then they apply self-transformer.

When trained on mid-sized datasets such as ImageNet, such models yield modest accuracies of a few percentage points below ResNets of comparable size. This seemingly discouraging outcome may be expected: Transformers lack some of the inductive biases inherent to CNNs, such as translation equivariance and locality, and therefore do not generalize well when trained on insufficient amounts of data.
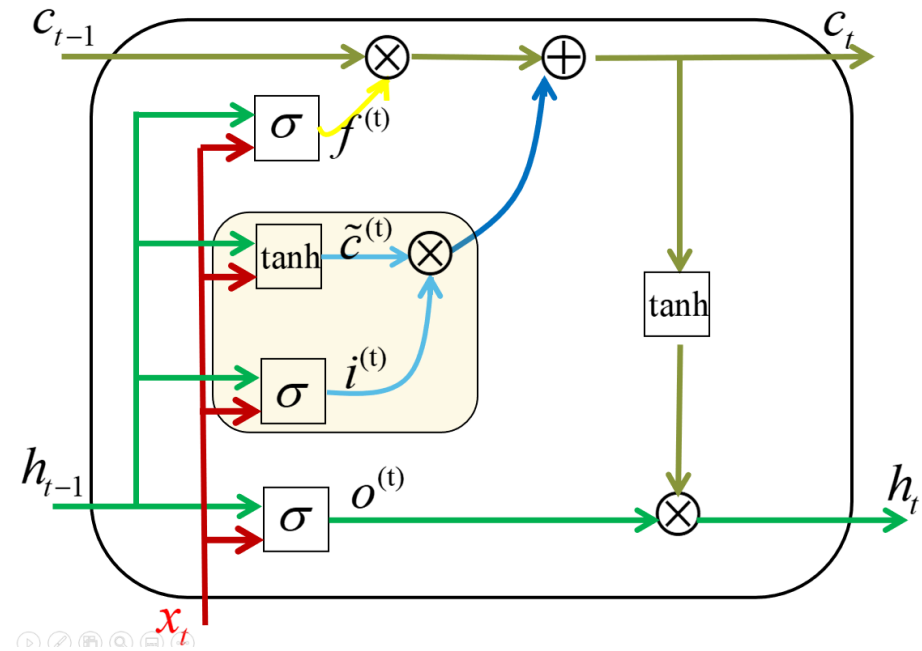
Dosovitskiy  et al. "An image is worth 16x16 words: Transformers for image recognition at scale", ICLR-2020.

Why?

# RNNs and LSTMs

RNN
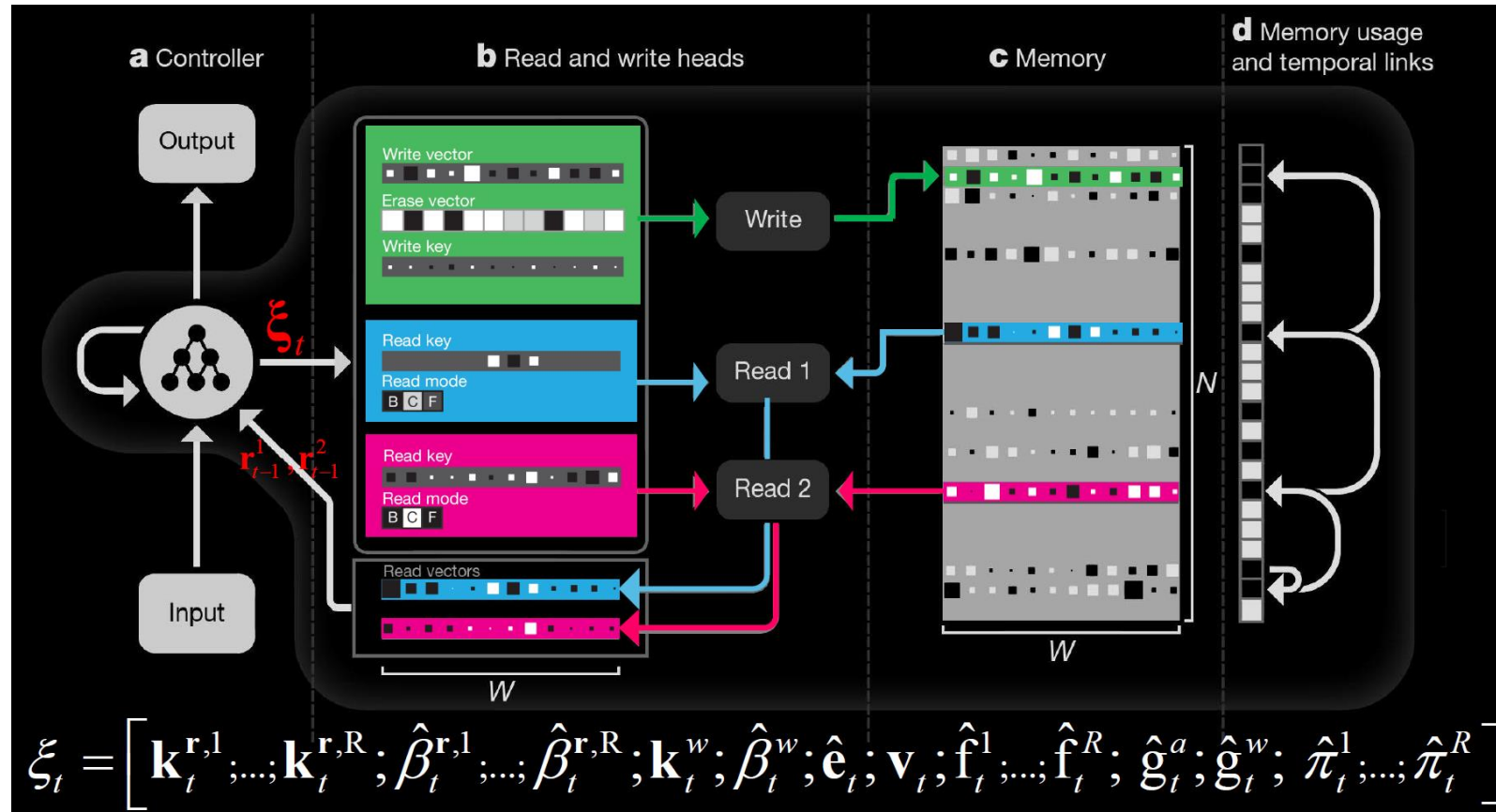


LSTM



See this for Arman's slides on RNN and LSTM.

# Aim at moving forward with NTM

Differentiable Neural Machines (Neural Turing Machines)



$$\xi_t = \left[ \mathbf{k}_t^{\mathbf{r},1};...;\mathbf{k}_t^{\mathbf{r},R}; \hat{\beta}_t^{\mathbf{r},1};...;\hat{\beta}_t^{\mathbf{r},R}; \mathbf{k}_t^w; \hat{\beta}_t^w; \hat{\mathbf{e}}_t; \mathbf{v}_t; \hat{\mathbf{f}}_t^1;...;\hat{\mathbf{f}}_t^R; \hat{\mathbf{g}}_t^a; \hat{\mathbf{g}}_t^w; \hat{\pi}_t^1;...;\hat{\pi}_t^R \right]$$

See this for Arman's slides on NTM.

# Keeping to move forward with Transformers



Vaswani  et al. "Attention is all you need", NeurIPS-2017.

How?

# Vision Transformer (ViT)

**Transformer Encoder**

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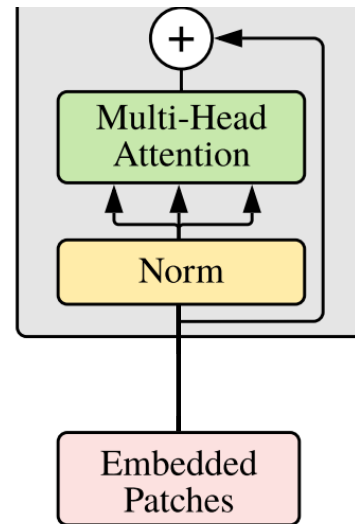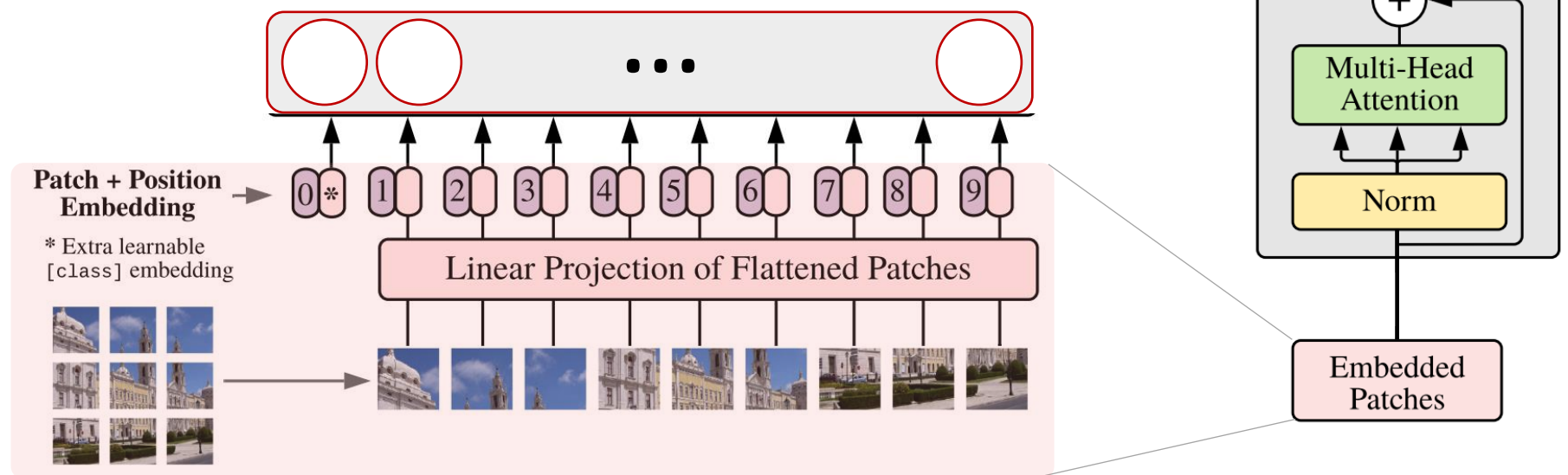
$$\mathbf{z}'_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \qquad \ell = 1 \ldots L \qquad (2)$$

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, \qquad \ell = 1 \ldots L \qquad (3)$$

$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0) \qquad (4)$$



Dosovitskiy  et al. "An image is worth 16x16 words: Transformers for image recognition at scale", ICLR-2020.

# Vision Transformer (ViT)

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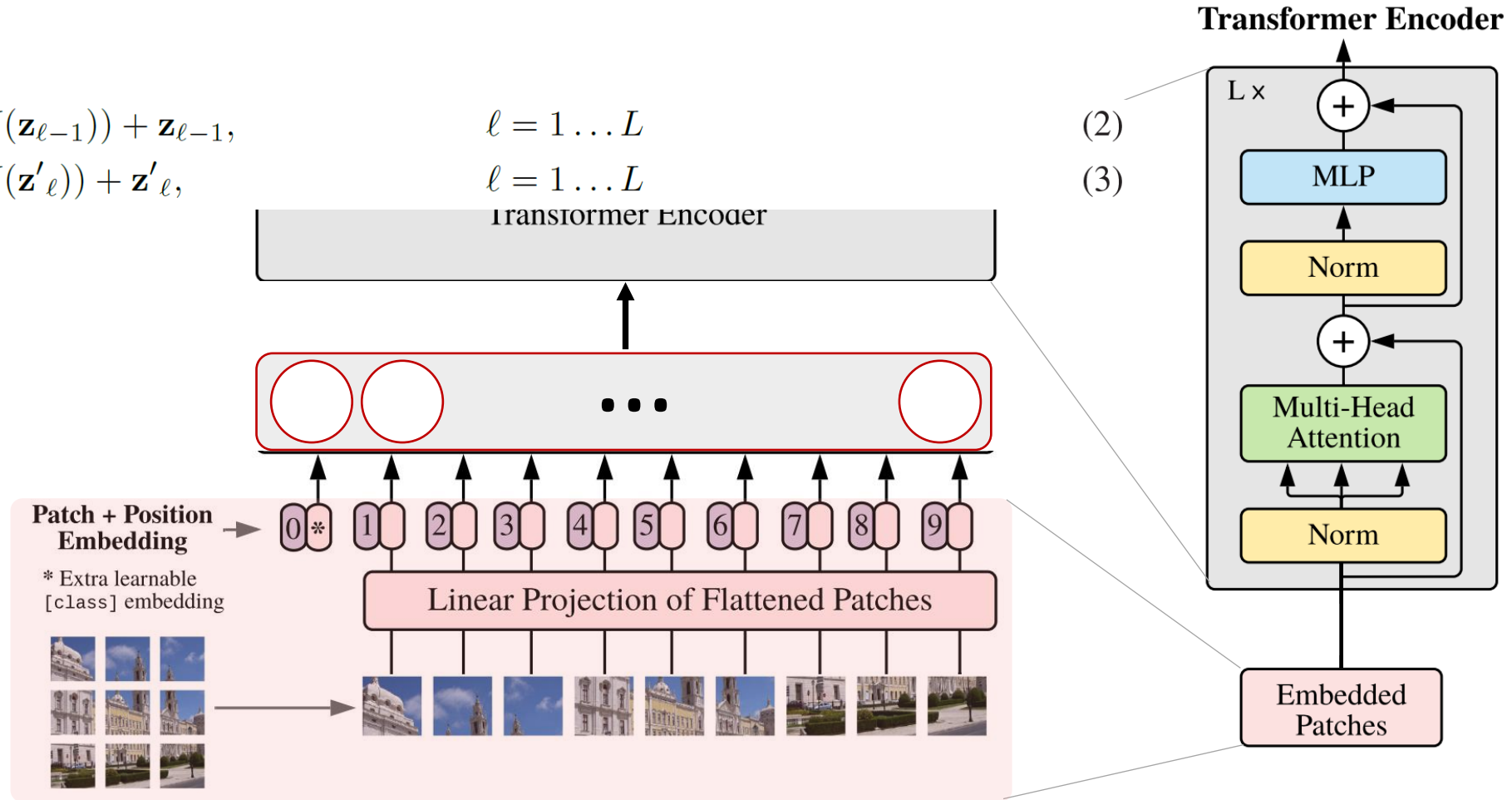


Dosovitskiy  et al. "An image is worth 16x16 words: Transformers for image recognition at scale", ICLR-2020.

# Vision Transformer (ViT)

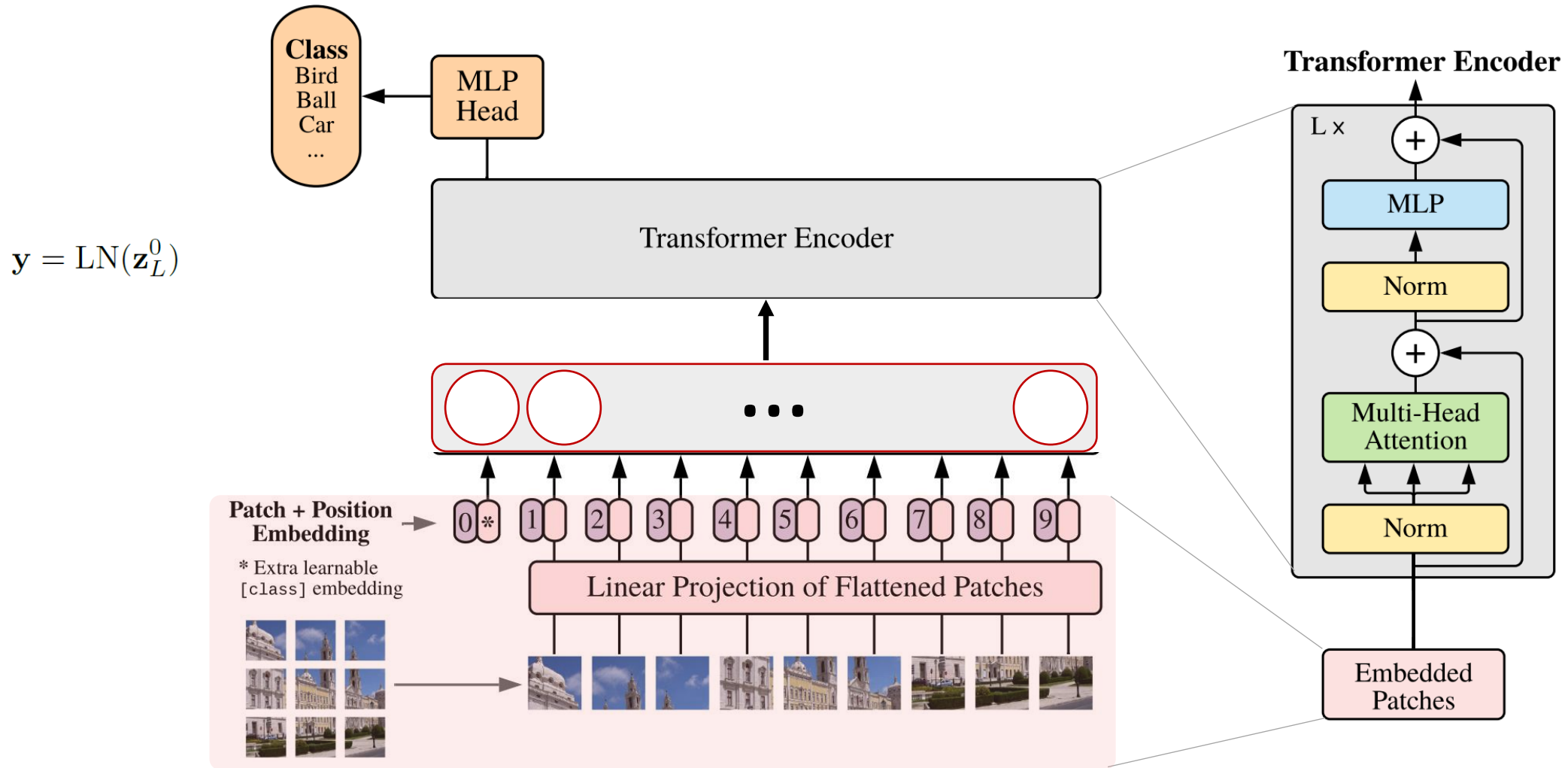$$\mathbf{z}'_{\ell} = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \qquad \ell = 1 \ldots L \qquad (2)$$

$$\mathbf{z}_{\ell} = \text{MLP}(\text{LN}(\mathbf{z}'_{\ell})) + \mathbf{z}'_{\ell}, \qquad \ell = 1 \ldots L \qquad (3)$$



Dosovitskiy et al. "An image is worth 16x16 words: Transformers for image recognition at scale", ICLR-2020.

# Vision Transformer (ViT)



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

Dosovitskiy et al. "An image is worth 16x16 words: Transformers for image recognition at scale", ICLR-2020.

# Input image format
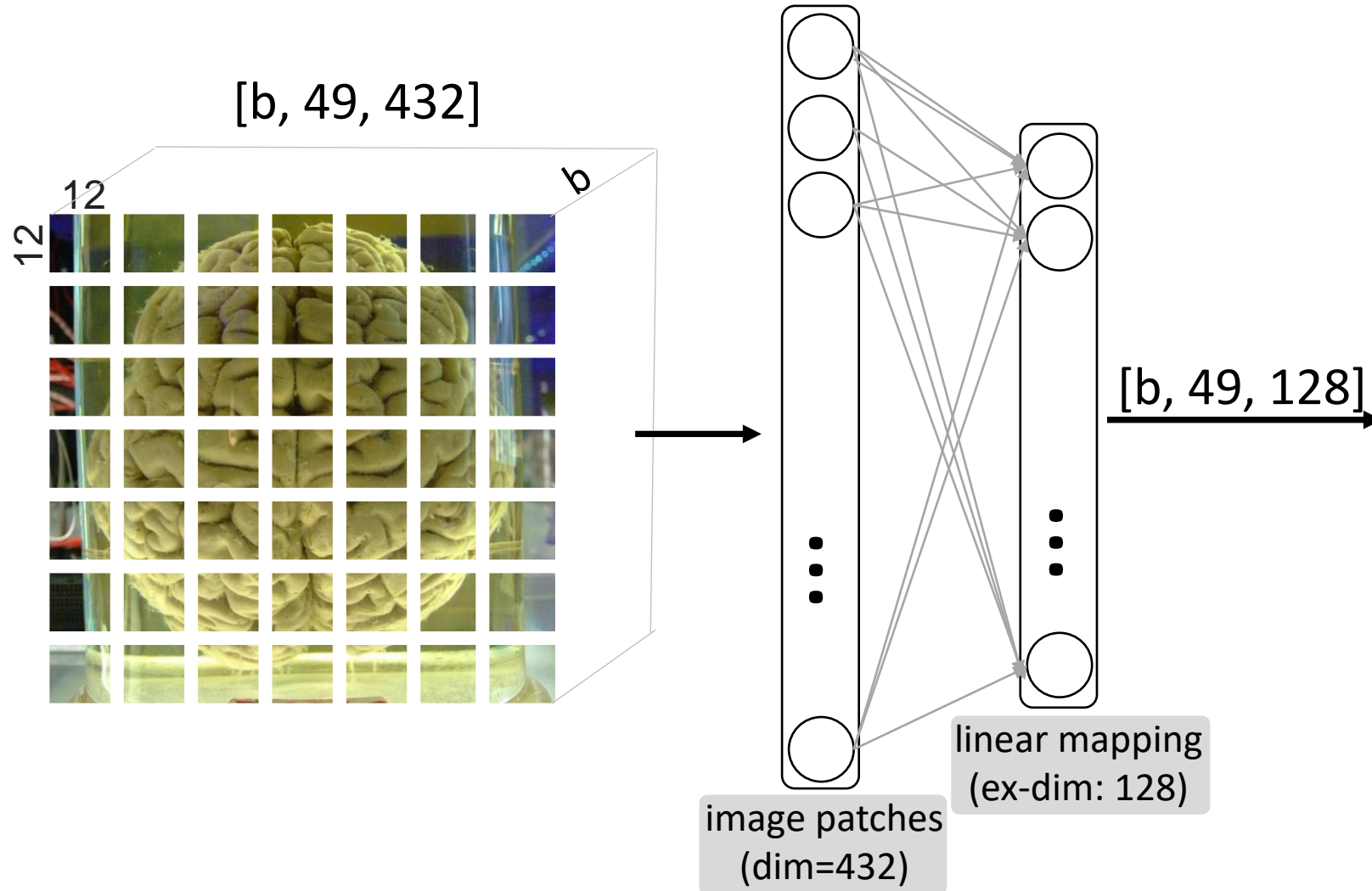
[b, 3, 84, 84]



$$\mathbf{x} \in \mathbb{R}^{H \times W \times C}$$

# Image to Patches

[b, 49, 432]



12

12
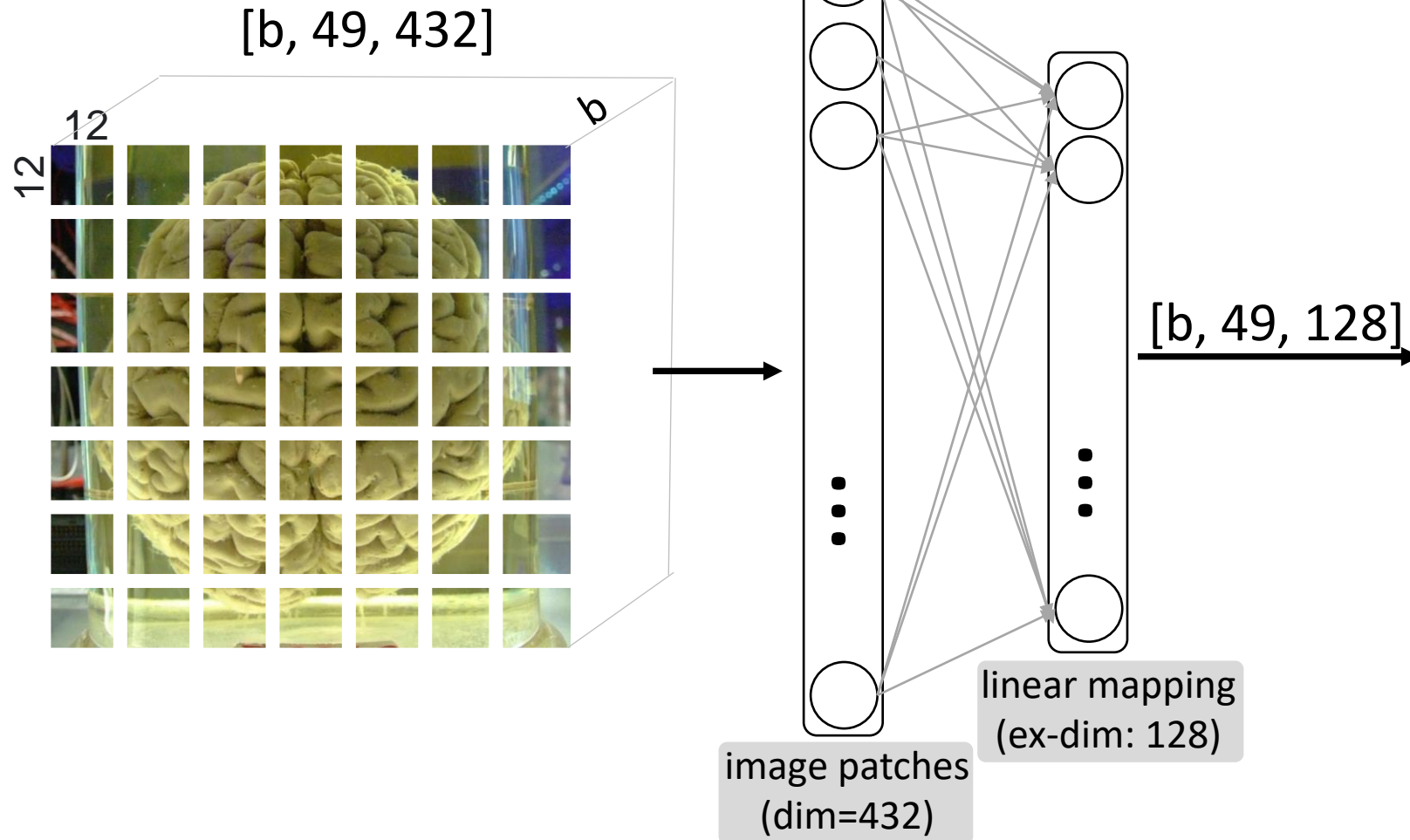
b

$$\mathbf{x}_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$$

# Patch to Token



[b, 49, 432]

12

12

b
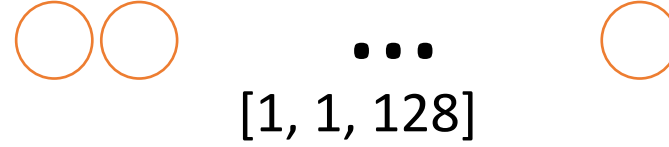
[b, 49, 128]

linear mapping
(ex-dim: 128)

image patches
(dim=432)

# Concat. of the Learnable Class Token

```python
self.cls_token = nn.Parameter(torch.randn(1, 1, dim))
```
extra learnable patch for class embedding

[1, 1, 128]

[b, 49, 432]

12

12

b

[b, 49, 128]

linear mapping
(ex-dim: 128)

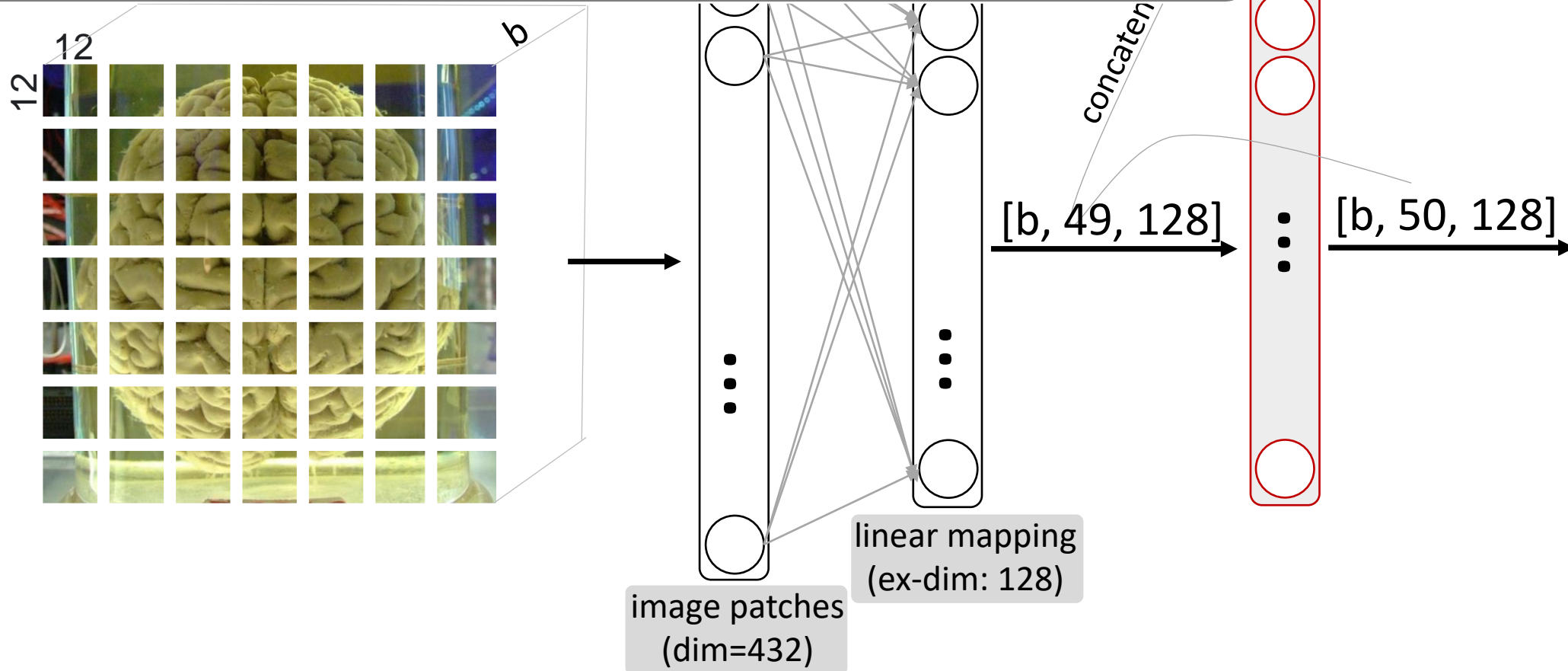image patches
(dim=432)

# Concat. learnable Class Token

# Concat. learnable Class Token

```
self.cls_token = nn.Parameter(torch.randn(1, 1, dim))
```
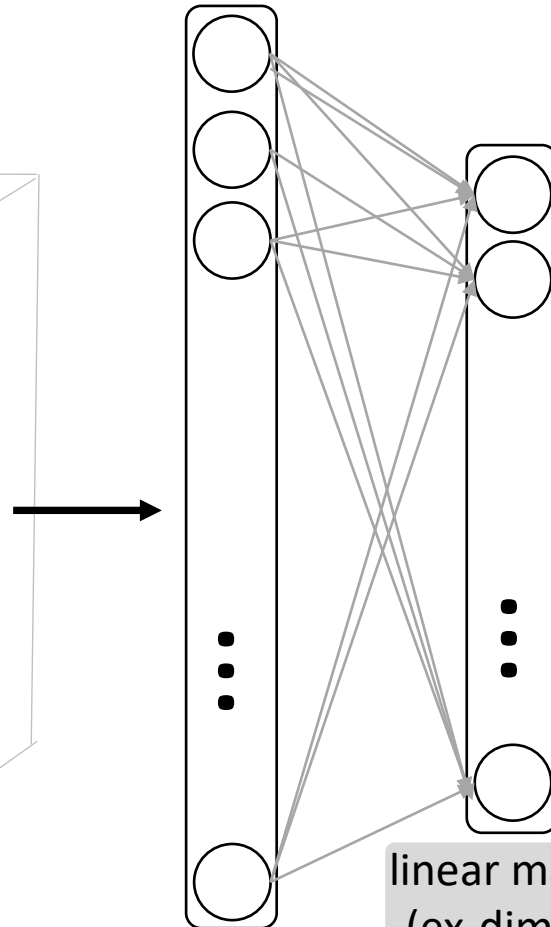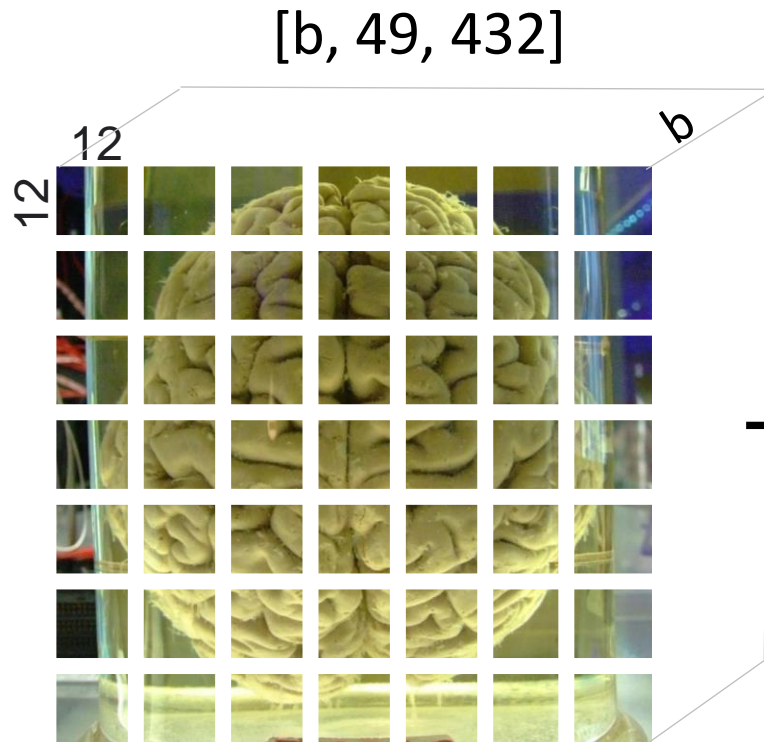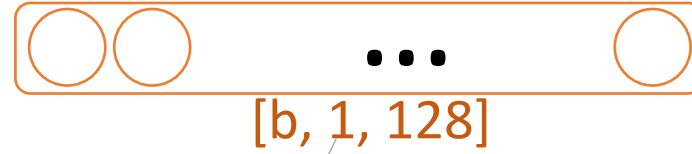extra learnable patch for class embedding

$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \cdots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{pos}, \qquad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{pos} \in \mathbb{R}^{(N+1) \times D} \qquad (1)$$
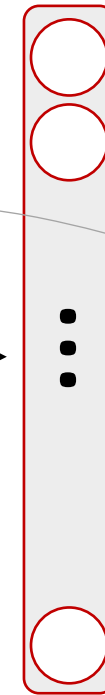
[128]

12

12

b

concaten

[b, 49, 128]

[b, 50, 128]

image patches
(dim=432)

linear mapping
(ex-dim: 128)

# Adding position embeddings



```
cls_token = nn.Parameter(torch.randn(1, 1, dim))
```
extra learnable patch for class embedding

[b, 1, 128]

[b, 49, 432]

12

12

b

[b, 49, 128]

concatenate

[b, 50, 128]

linear mapping
(ex-dim: 128)

image patches
(dim=432)

[1, 50, 128]

```
pos_embedding = nn.Parameter(torch.randn(1, num_patches + 1, dim))
```

# Adding position embeddings



cls_token = nn.Parameter(torch.randn(1, 1, dim))
extra learnable patch for class embedding
[b, 1, 128]

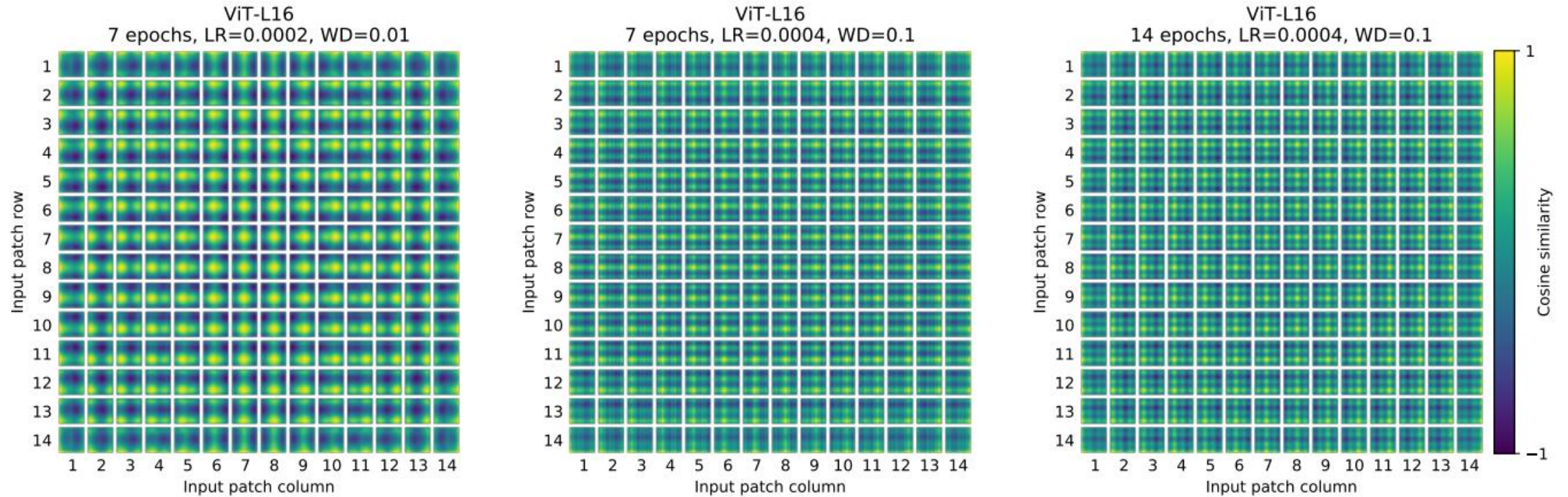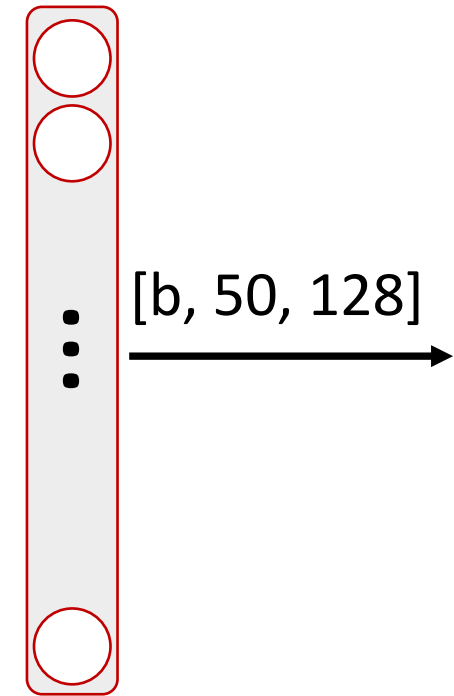D.3 POSITIONAL EMBEDDING (ablation of position embedding)
- Providing no positional information
- 1-dimensional positional embedding (paper)
- 2-dimensional positional embedding
- Relative positional embeddings

Not much difference

[b, 49, 128]

[b, 50, 128]

Concatenate

linear mapping
(ex-dim: 128)

image patches
(dim=432)

[1, 50, 128]

pos_embedding = nn.Parameter(torch.randn(1, num_patches + 1, dim))

# Position embeddings depends on hyperparameters



Figure 9: Position embeddings of models trained with different hyperparameters.

# Patch embedding

[b, 50, 128]

# Big picture



Class
Bird
Ball
Car
...

MLP
Head

[b, 50, 128]

Transformer Encoder

Patch + Position
Embedding
* Extra learnable
[class] embedding

0 * 1 2 3 4 5 6 7 8 9

of Flattened Patches

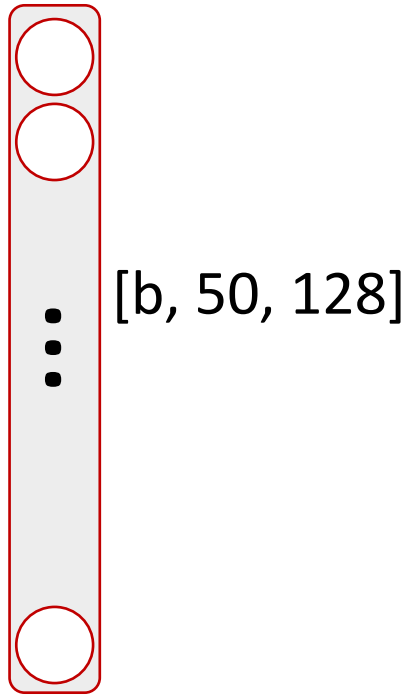Transformer Encoder

L ×

MLP

Norm
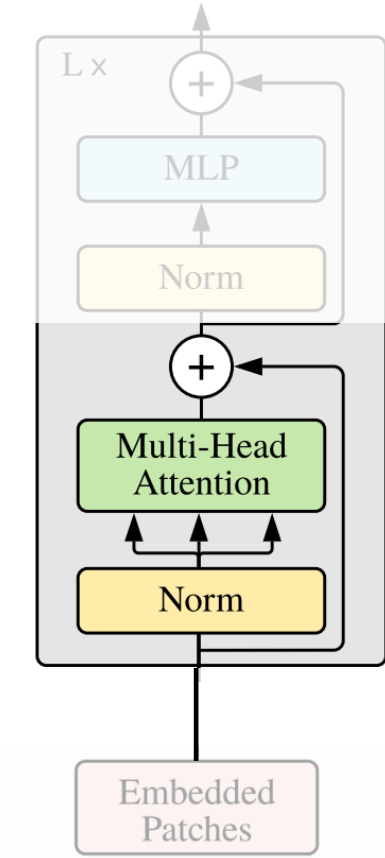
Multi-Head
Attention

Norm

Embedded
Patches

# Transformer: MH-Self Attention (MSA)

$$\mathbf{z}'_\ell = \mathrm{MSA}(\mathrm{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \qquad \ell = 1 \ldots L$$
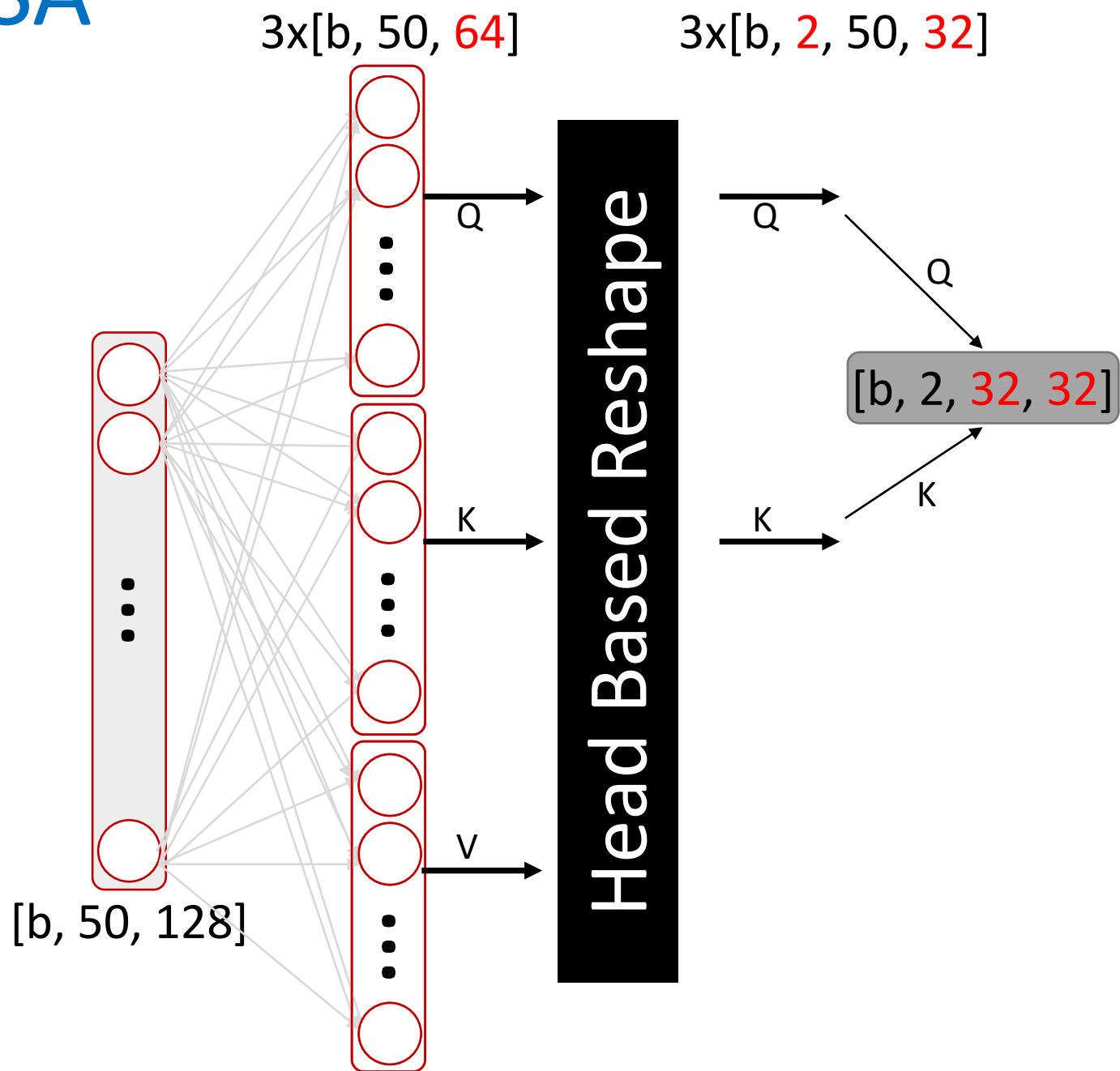
[b, 50, 128]

# MSA

3x[b, 50, 64]

3x[b, 2, 50, 32]

$$\text{softmax}(\frac{QK^T}{\sqrt{d_k}})$$

[b, 50, 128]

Head Based Reshape

Q

K

V

Q

K

Q

K

[b, 2, 32, 32]

# MSA



3x[b, 50, 64]    3x[b, 2, 50, 32]

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

Q

Q

.softmax(dim=-1)

[b, 2, 32, 32]

K    K

$$\text{softmax}(\frac{QK^T}{\sqrt{d_k}})$$

Head Based Reshape

K

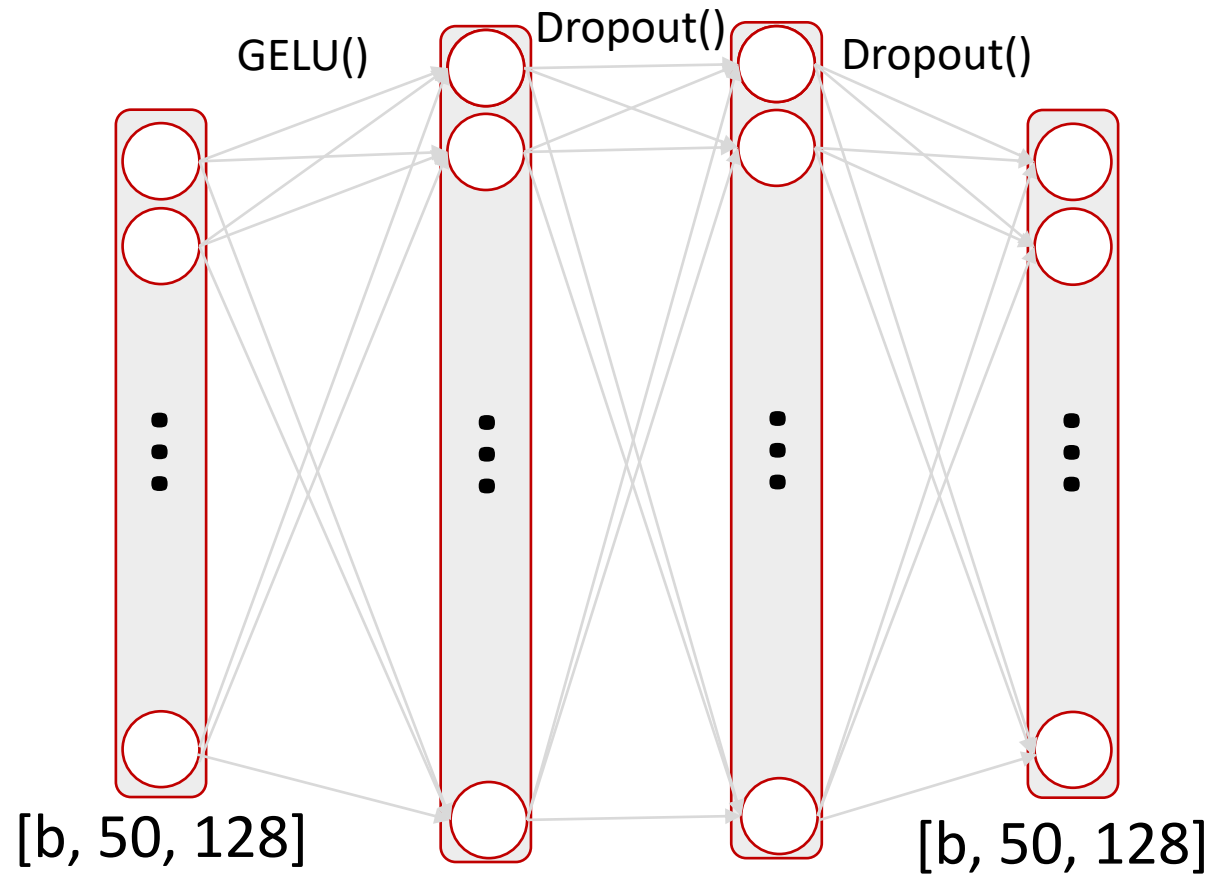V    V    V

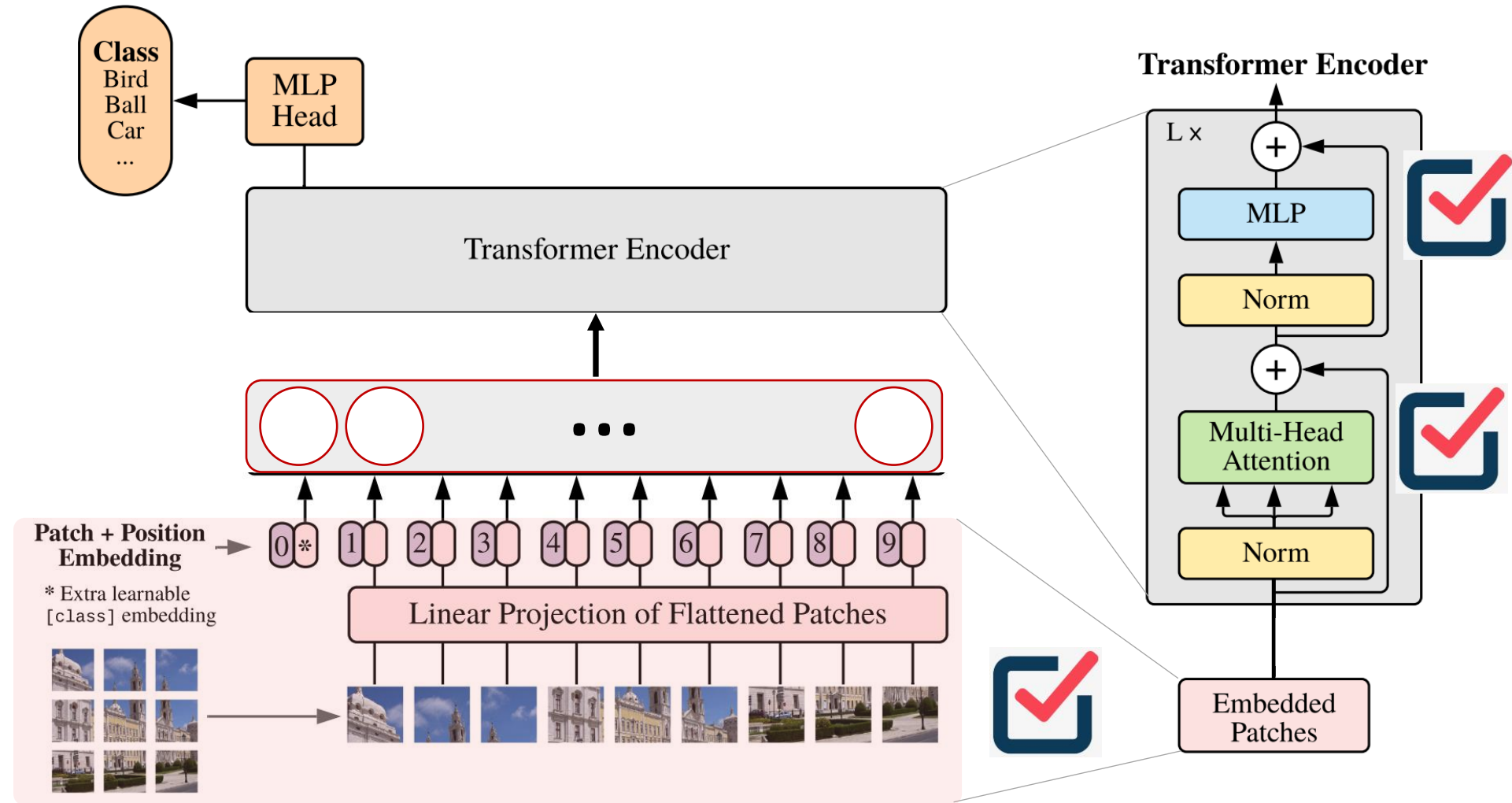[b, 2, 50, 32]    →    [b, 50, 64]

[b, 50, 128]    [b, 50, 128]

# Big picture

# Transformer: MLP

# Big picture

# Transformers



Vaswani et al. "Attention is all you need", NeurIPS-2017.