

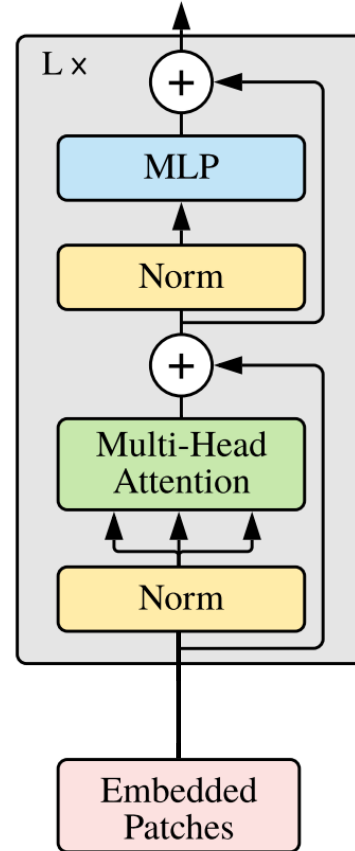
Vision Transformer (ViT)

Arman Afrasiyabi

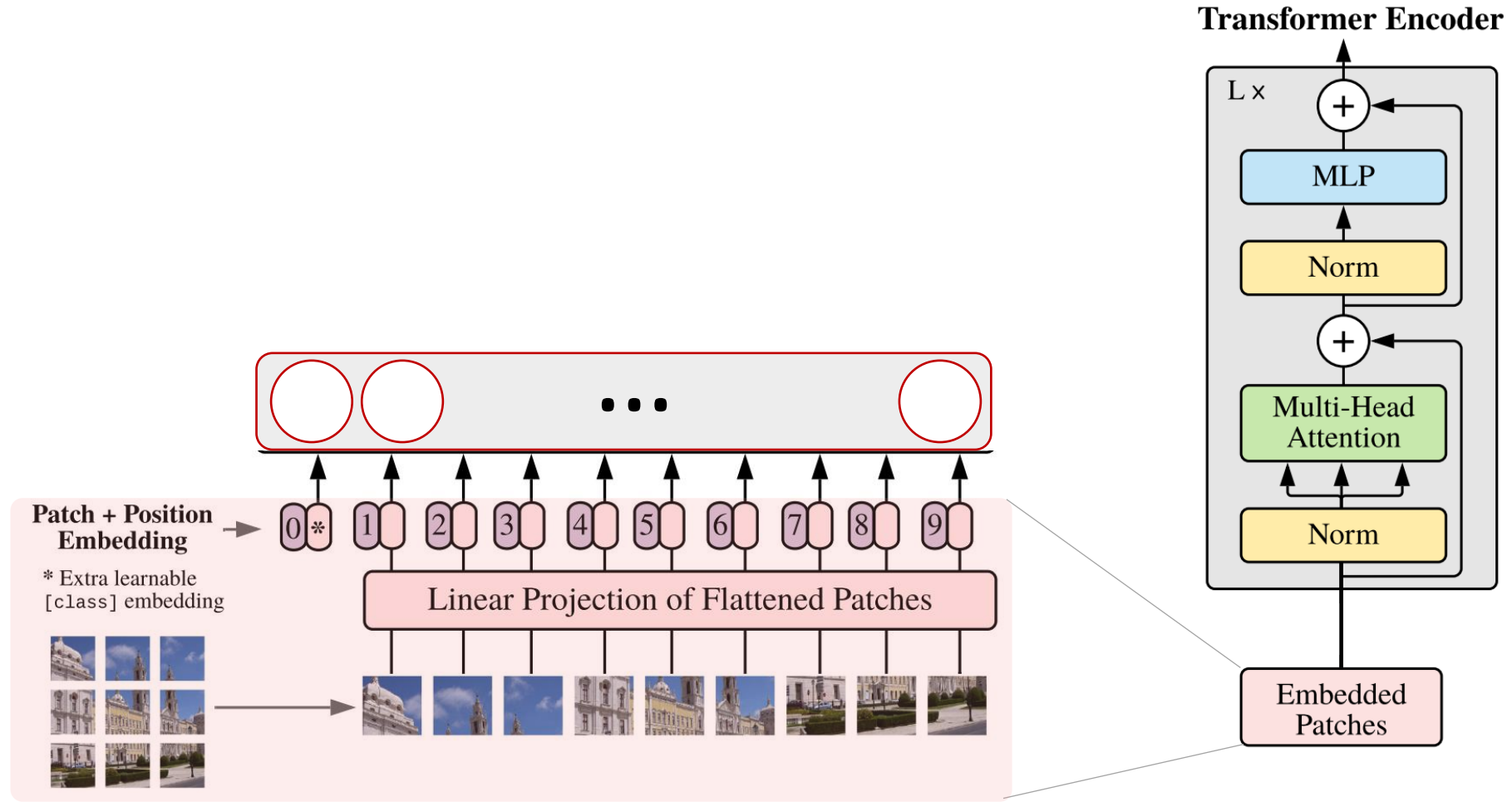
PHD Student at Université Laval

Vision Transformer (ViT)

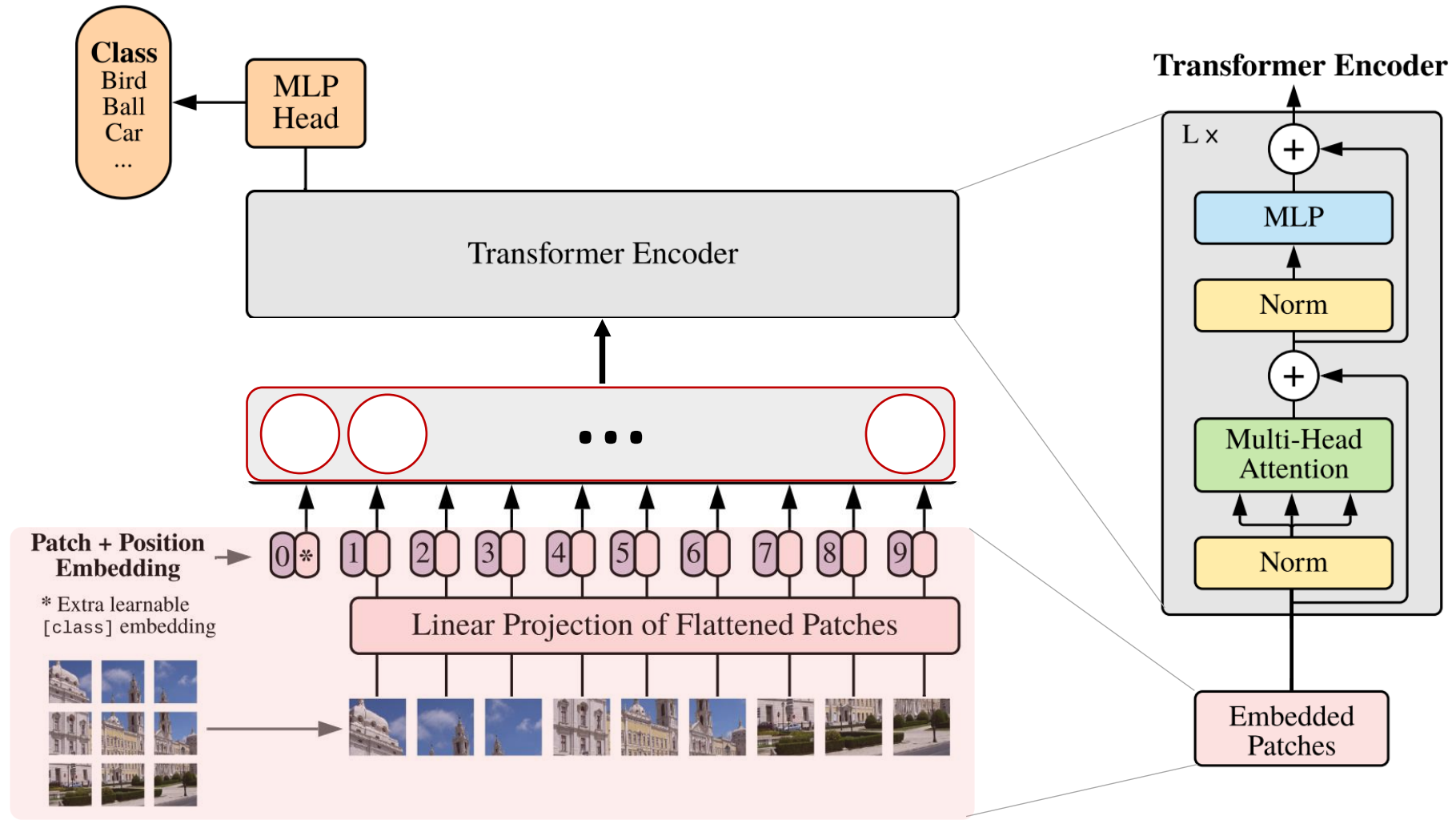
Transformer Encoder



Vision Transformer (ViT)



Vision Transformer (ViT)



Input shape

[b, 3, 84, 84]

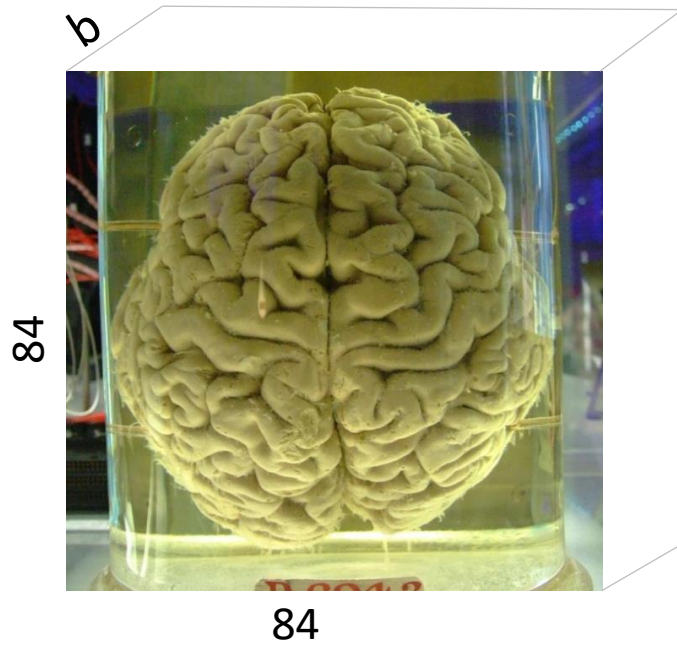
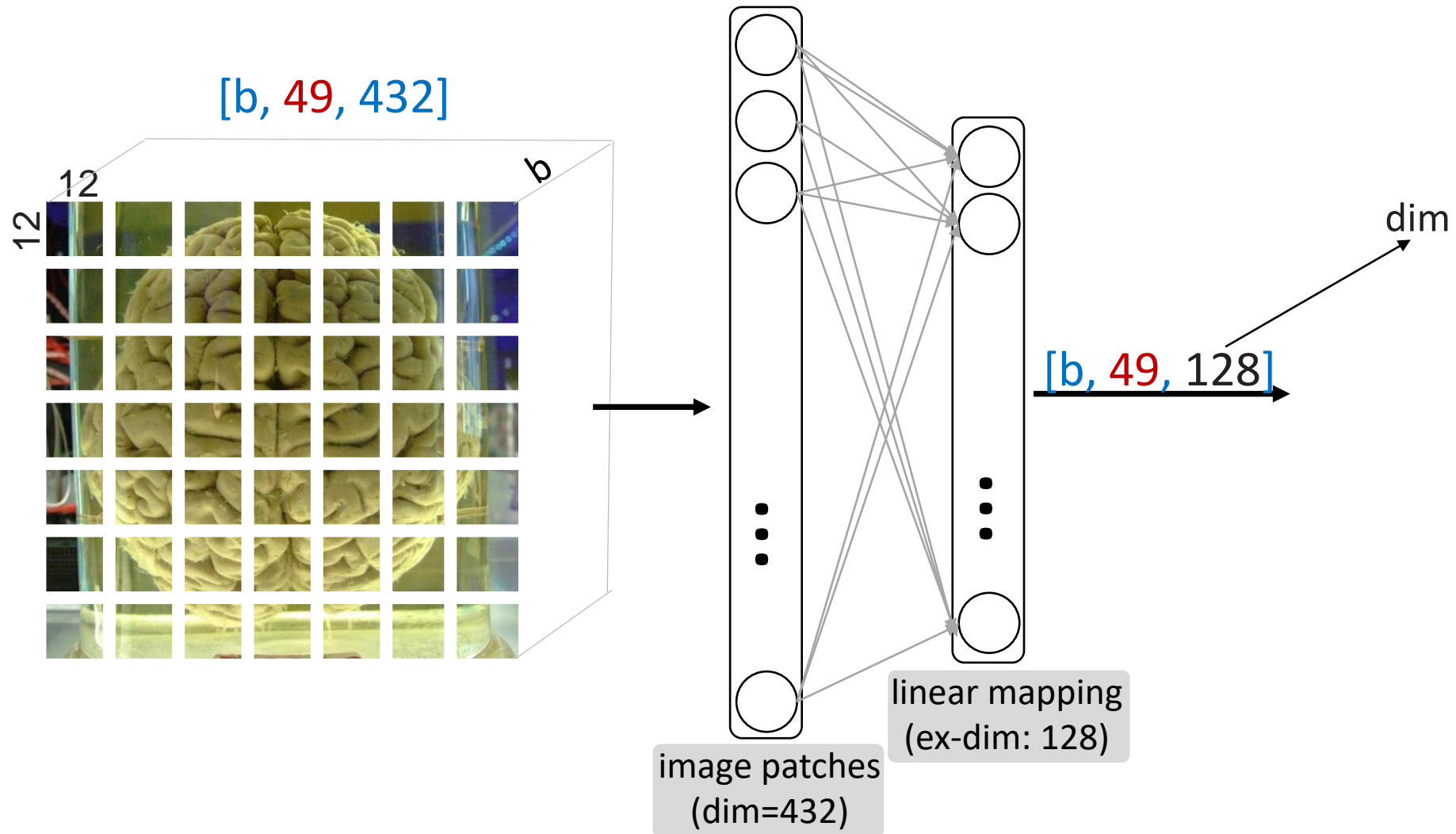


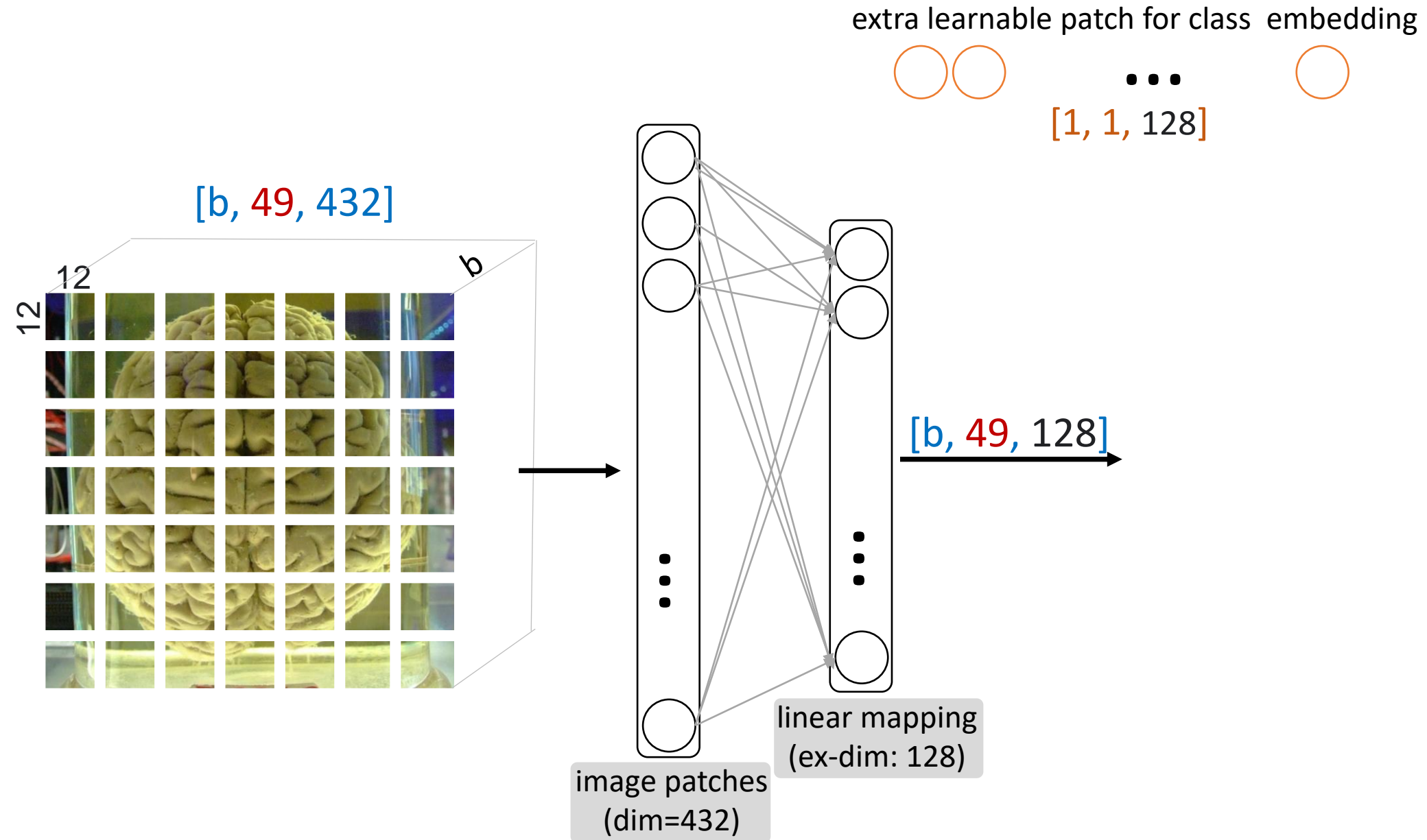
Image to Patches



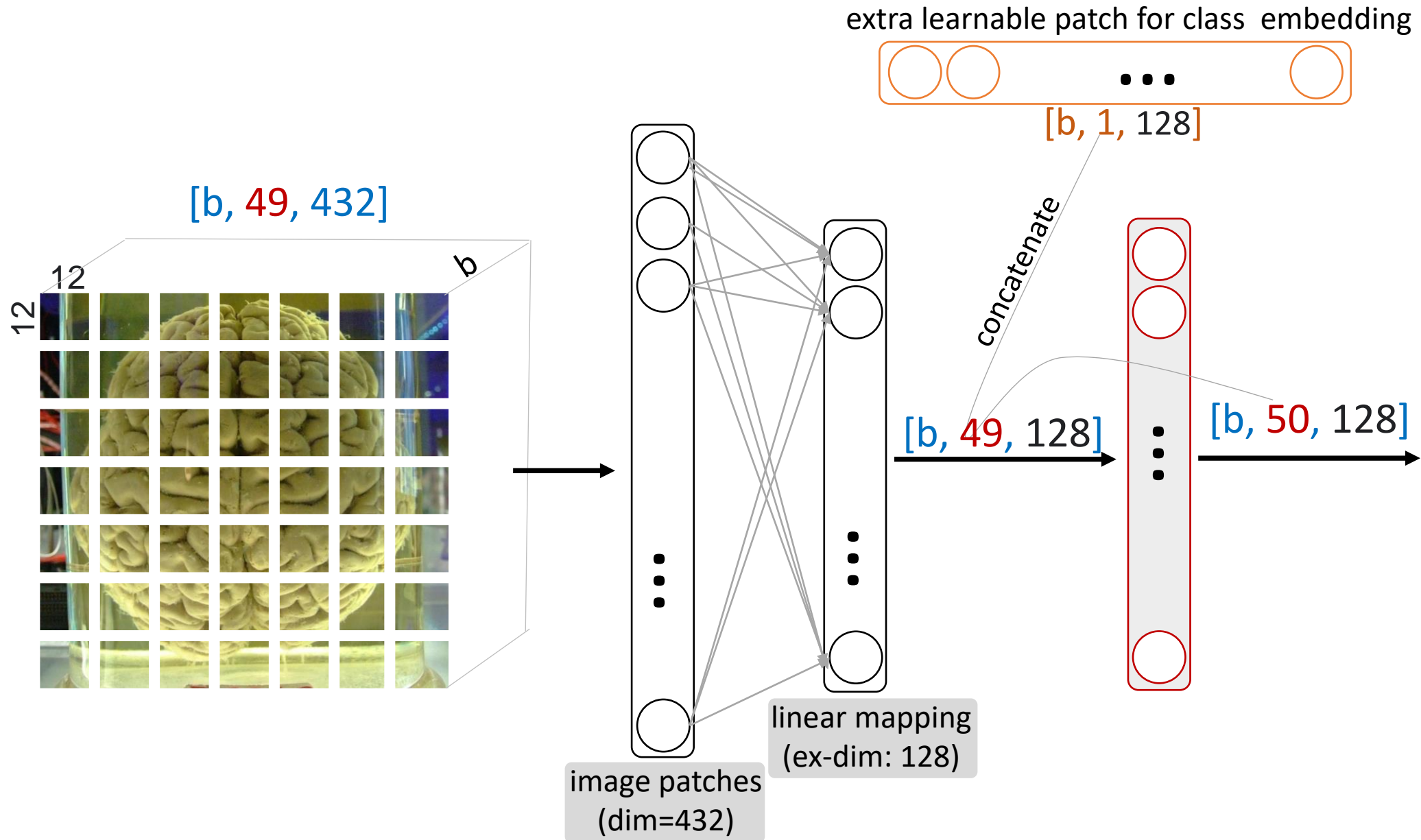
Patch to Token



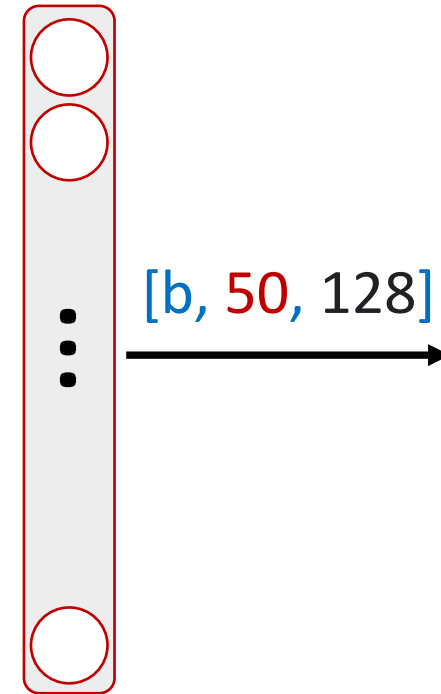
Concat. learnable Class Token



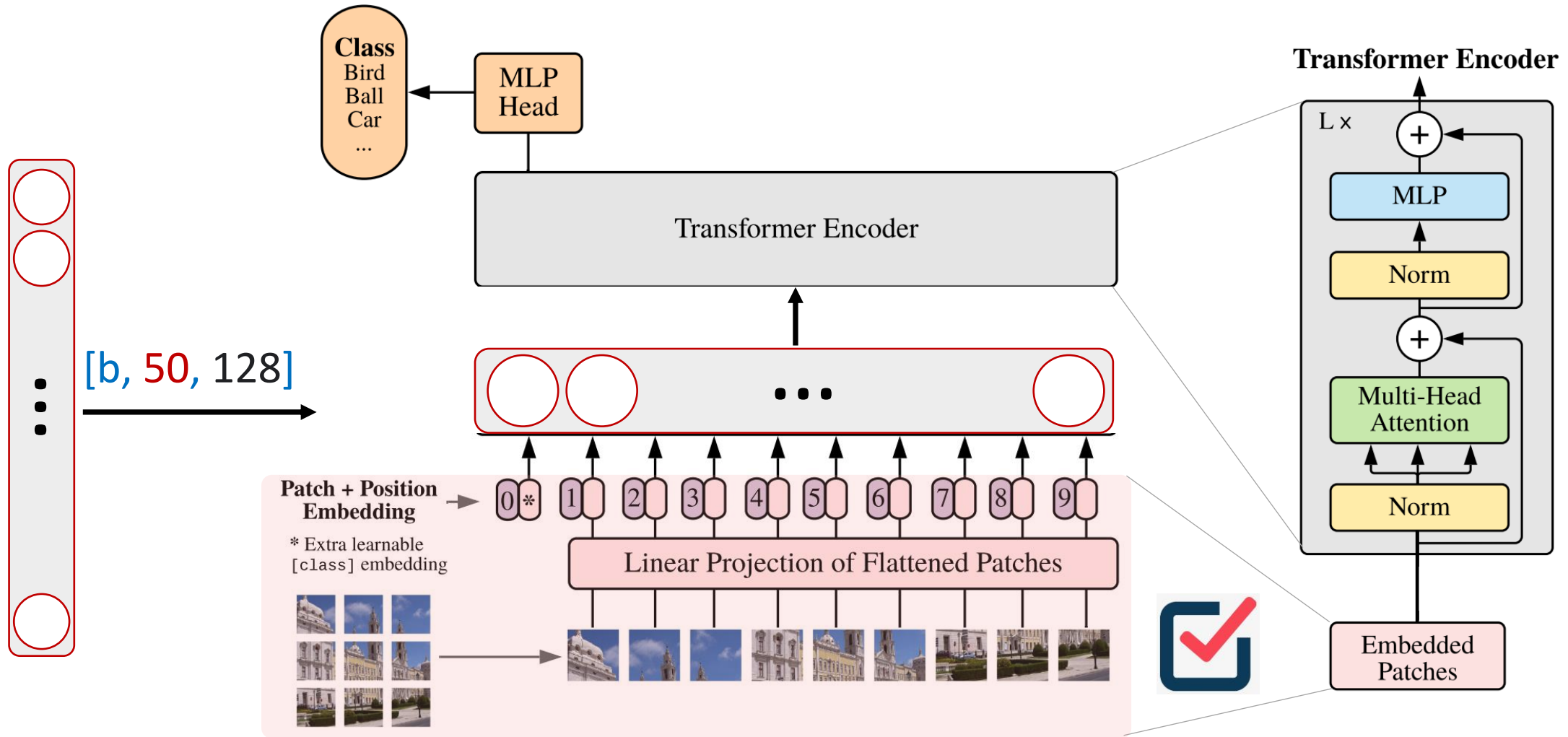
Concat. learnable Class Token



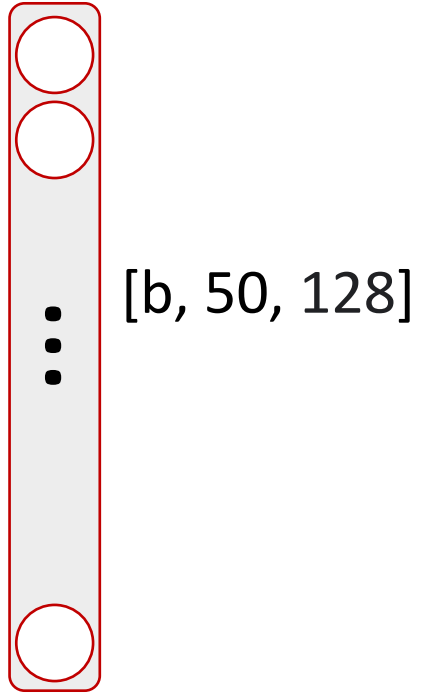
Patch embedding



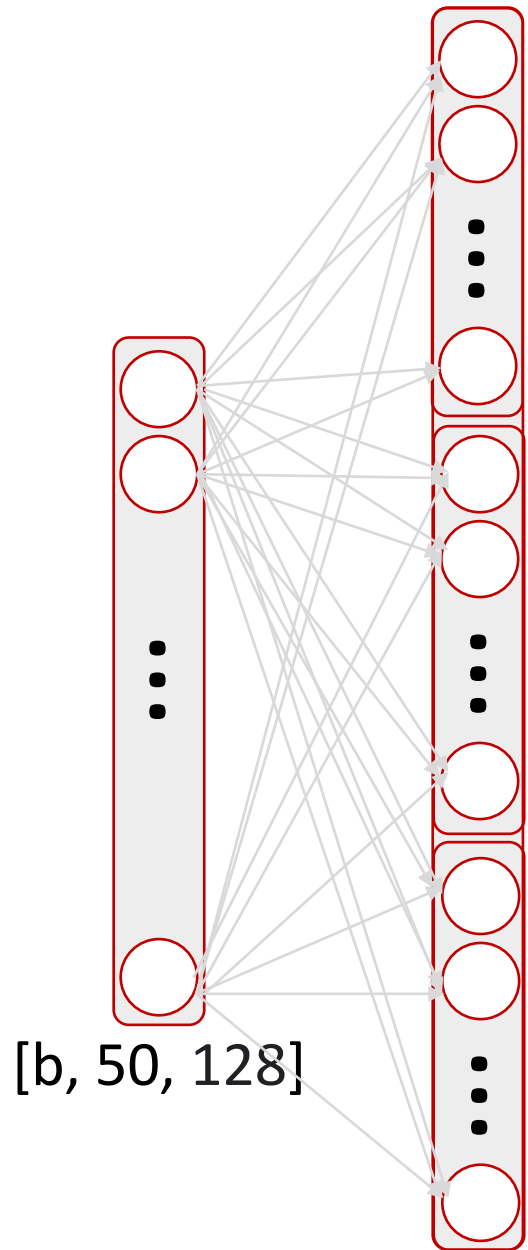
Big picture



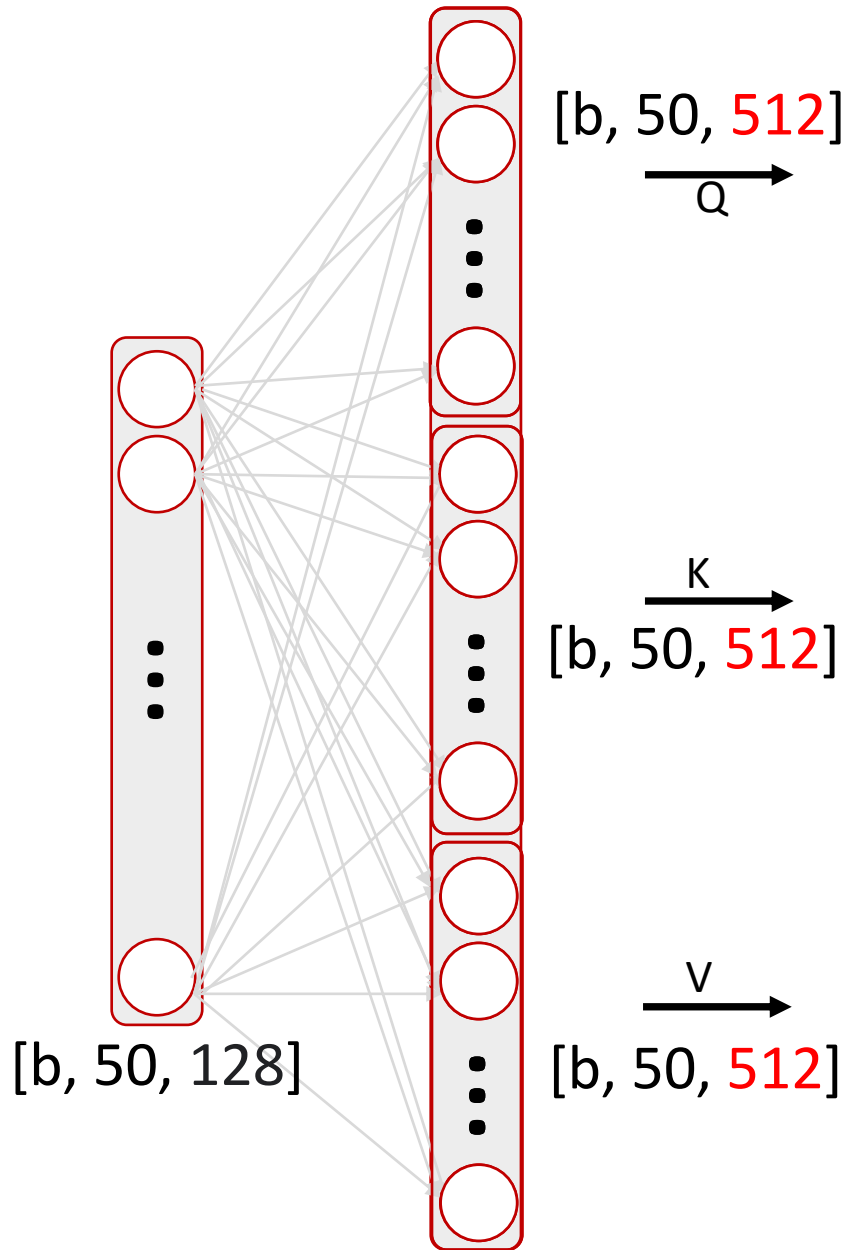
Transformer: MH-Attention



Transformer: MH-Attention



Transformer: MH-Attention

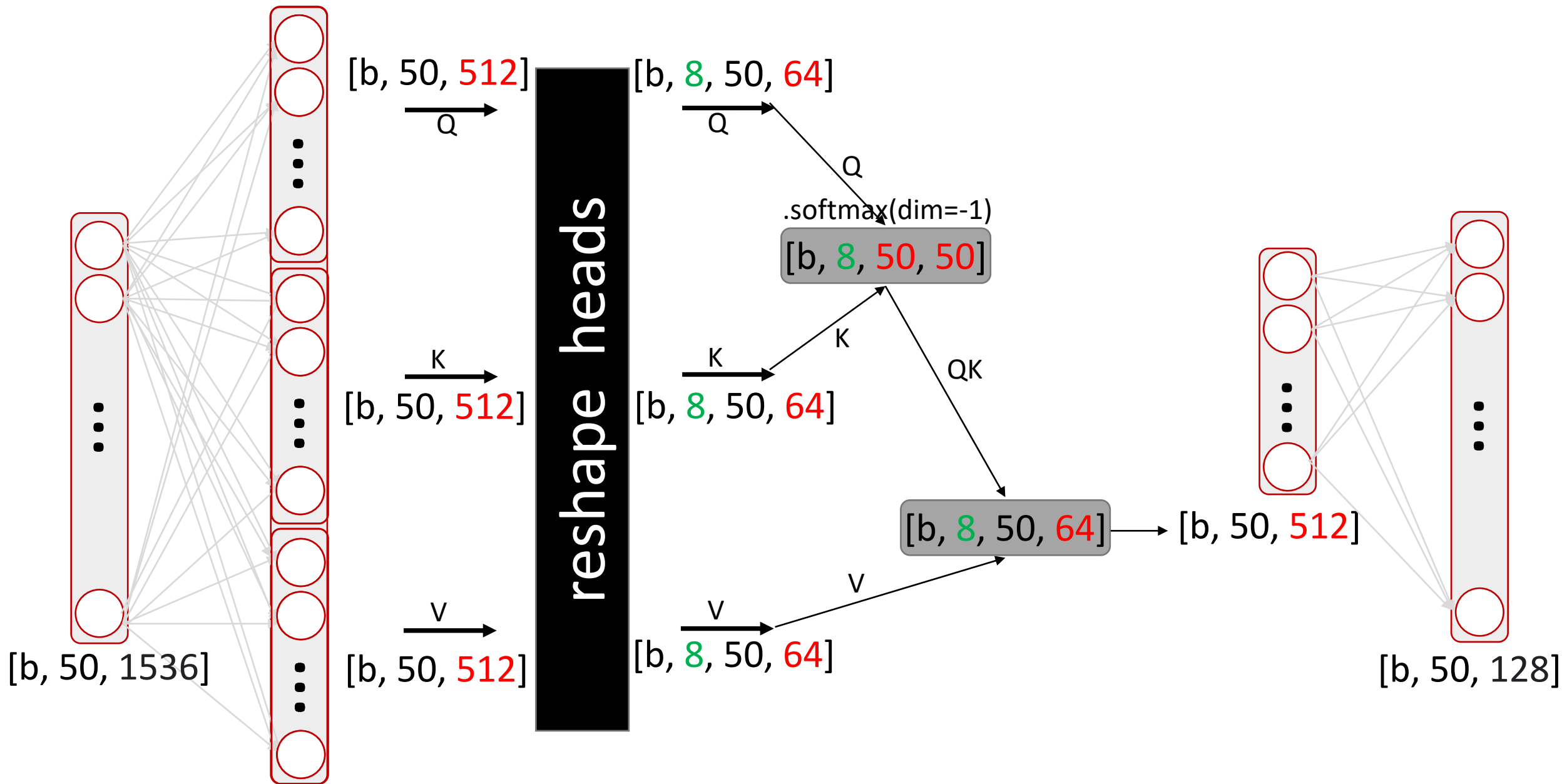


$$512 = 64 * 3$$

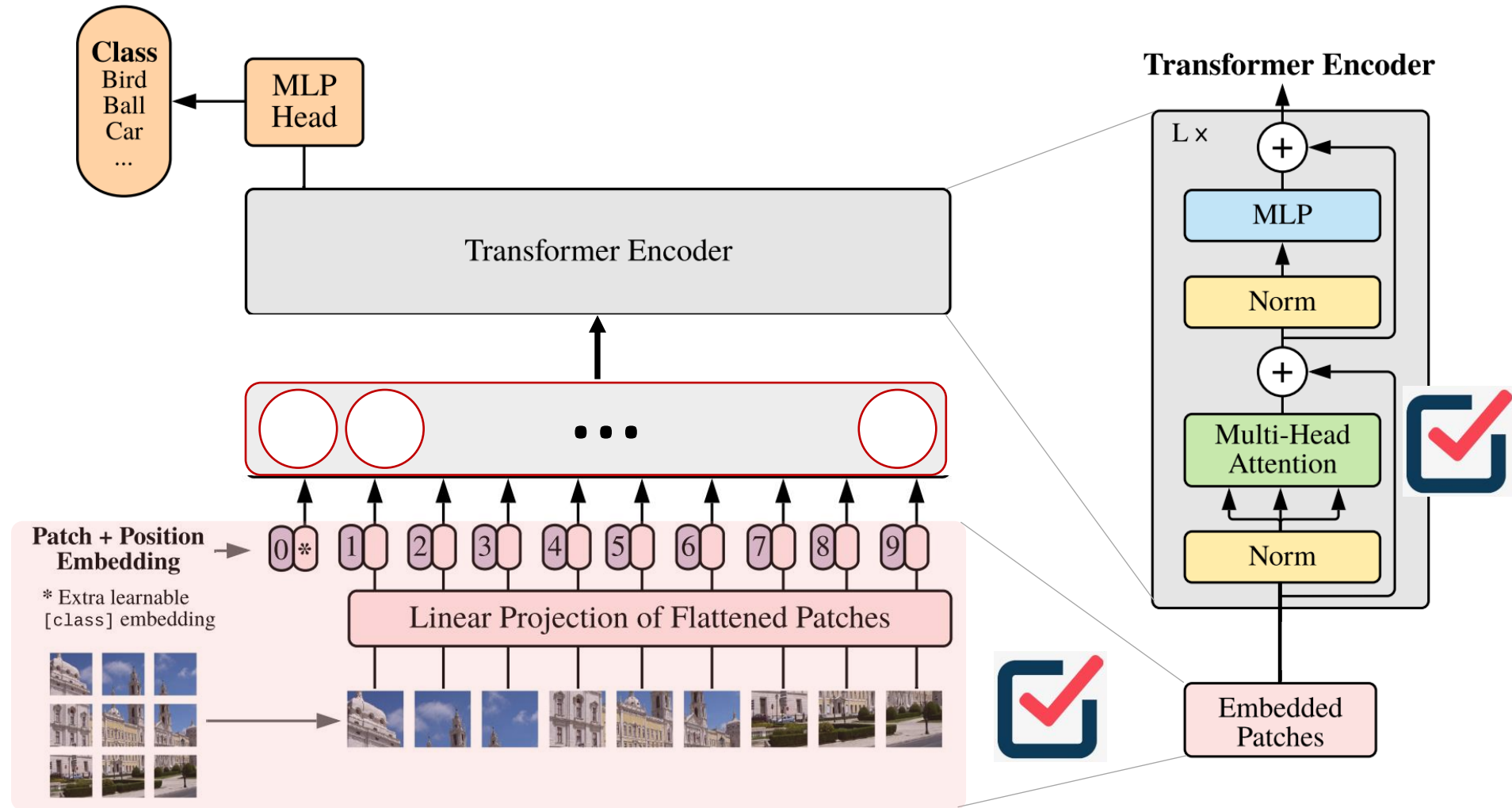
n-heads

head_dim

Transformer: MH-Attention (n_head = 8)



Big picture



Transformer: MLP

