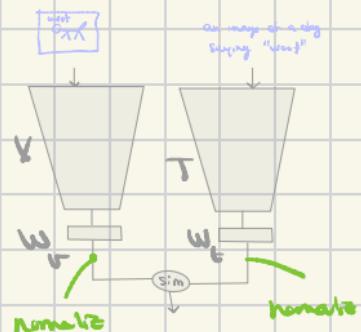




CLIP objective

Let X and Y be an image and text matrix, respectively, such that $\{(x_i, y_i)\}_{i=1}^N$ are image-caption pairs, where

$$x_i \in \mathbb{R}^{(C, H, W)} \quad \text{and} \quad y_i \in \mathbb{R}^L$$



Info NCE / contrastive Loss

$$e_i = w_v V(x) \rightarrow \hat{e}_i = \frac{e_i}{\|e_i\|}$$

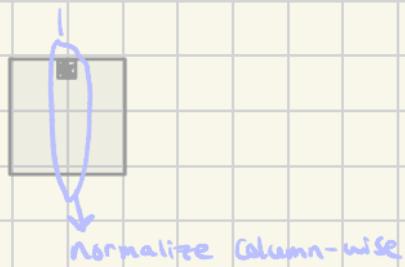
$$w_i = w_t T(y_i) \rightarrow \hat{w}_i = \frac{w_i}{\|w_i\|}$$

Let $S_{ij} = \cos(\hat{e}_i, \hat{w}_j) * \frac{1}{\tau}$, be a similarity matrix for a batch ($B = \text{batch-size}$) of image-caption pairs, where τ is a learnable temperature

$$p(i|j) = \frac{s_{ij}}{\sum_{k=1}^B \exp(s_{ik})}$$



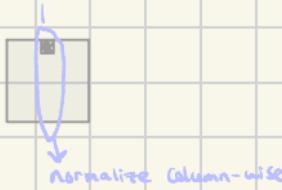
$$q(i|j) = \frac{s_{ij}}{\sum_{k=1}^B \exp(s_{kj})}$$



$$p(j|i) = \frac{s_{ij}}{\sum_{k=1}^B \exp(s_{ik})}$$



$$g(i|j) = \frac{s_{ij}}{\sum_{k=1}^B \exp(s_{kj})}$$



$$\mathcal{L}_{i2t} = \frac{1}{B} \sum_h^B \log p(k|h)$$

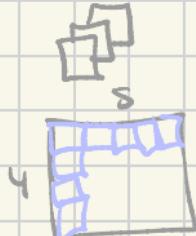


$$\mathcal{L}_{t2i} = \frac{1}{B} \sum_k^B \log g(k|i)$$



$$\mathcal{L}_{clip} = -(\mathcal{L}_{i2t} + \mathcal{L}_{t2i}) / 2$$

Patch Embedding



we want (B, P, D)

Broadcast variable
batch size

batch size

An image of a dog
Bengal cat

clistake $(2, 1, D)$

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

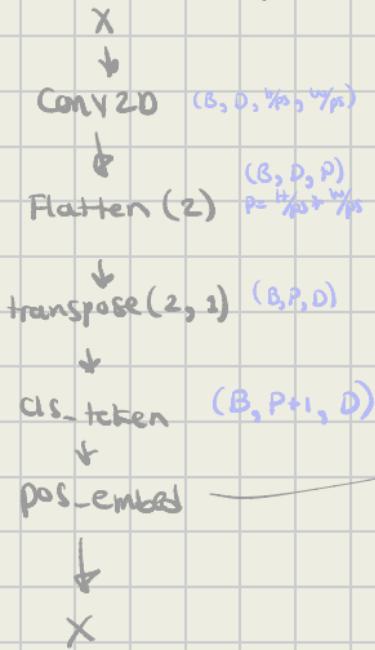
↓

↓

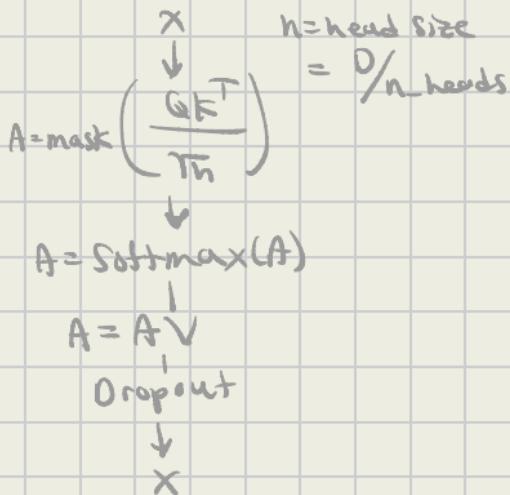
↓

↓

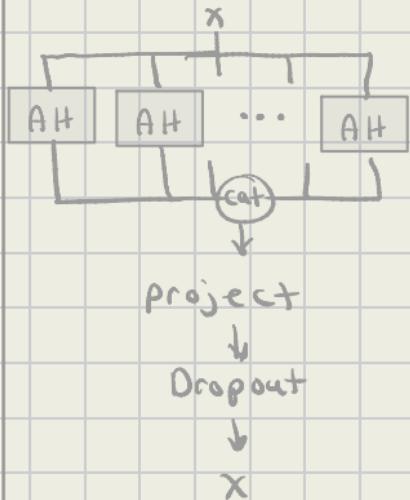
Patch Embedding



Attention Head



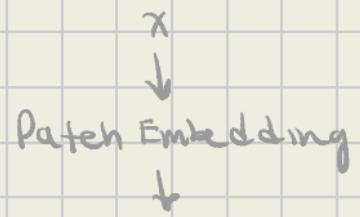
Multi Head Attention



Transformer Block



Vision Transformer



Transformer Block

\vdots : $x \text{ depth}$

Transformer Block



Layer Norm



X

CLIP

