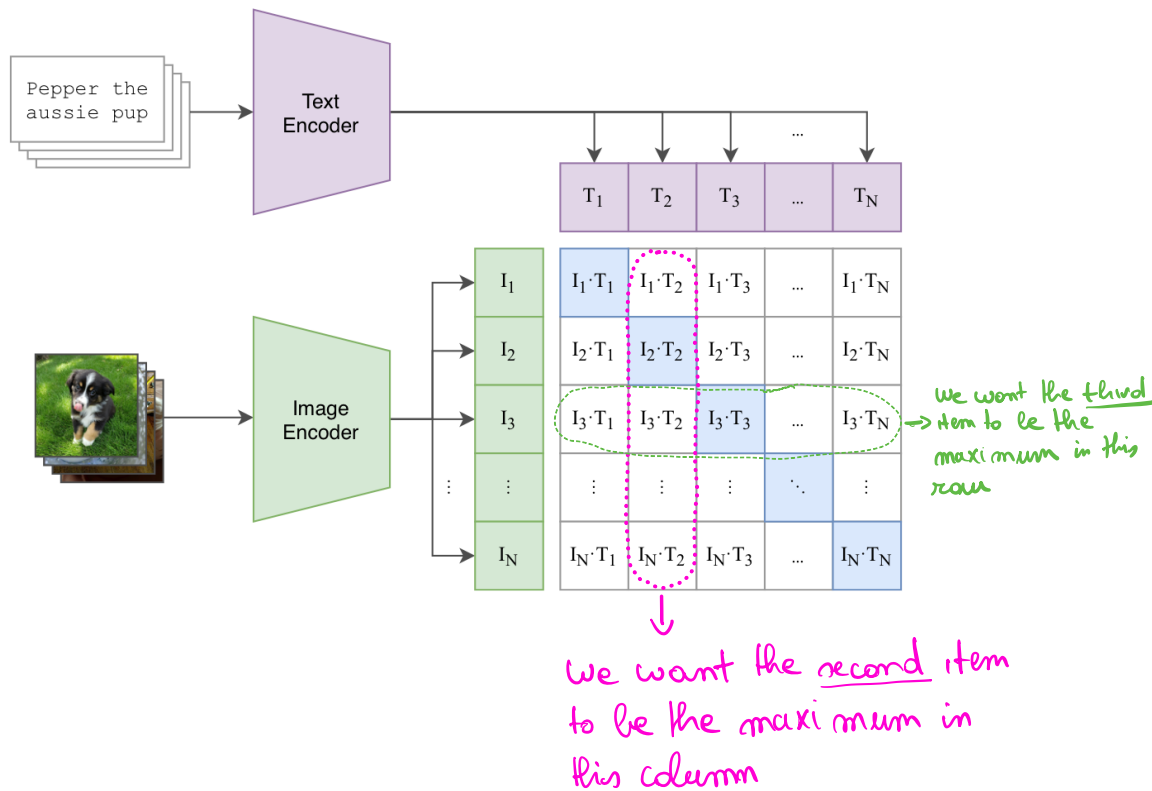


# From CLIP to SigLip

CLIP stands for Contrastive Language-Image Pretraining.

## What is contrastive learning?

(1) Contrastive pre-training



**Problem:** how do we tell the model we want one item in each row/column to be maximized while minimizing all the others?

**Hint:** this is very similar to language modeling in which we want a single token to be the next one given the prompt...

**Solution:** We use the Cross-Entropy Loss!

```

# image_encoder - ResNet or Vision Transformer
# text_encoder - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l] - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) # [n, d_i] → convert a list of images into a list of embeddings
T_f = text_encoder(T) # [n, d_t] → convert a list of prompts into a list of embeddings

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1) } Make sure both image and text
                                              } embeddings have the same number
                                              } of dimensions, and then normalize
                                              } the vectors

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t) → Compute all the possible dot products.

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1) } Teach the model which item in
loss = (loss_i + loss_t)/2 } each row/column needs to be
                           } maximized

```

Figure 3. Numpy-like pseudocode for the core of an implementation of CLIP.

## Numerical stability of the softmax

$\forall i \in 1 \dots N$      $S_i = \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}$     The softmax makes all the elements of a vector in such a way that they're in the real range  $[0, 1]$  and they sum up to 1.

Problem: the softmax is numerically unstable, as the exp function can grow fast and may not fit in a 32 bit floating-point number.

Solution: do not make the exp grow to infinity.

$$S_i = \frac{c \cdot e^{a_i}}{c \cdot \sum_{k=1}^N e^{a_k}} = \frac{e^{\log(c)} e^{a_i}}{e^{\log(c)} \sum_{k=1}^N e^{a_k}} = \frac{e^{a_i + \log(c)}}{\sum_{k=1}^N e^{a_k + \log(c)}}$$

We maximally choose  $\log(c) = -\max_i (a_i)$

This will push the arguments of the exp towards negative numbers and the exp itself towards zero.

# The normalization factor in the softmax

To calculate the normalization factor, we must go through all the elements of each row and each column.

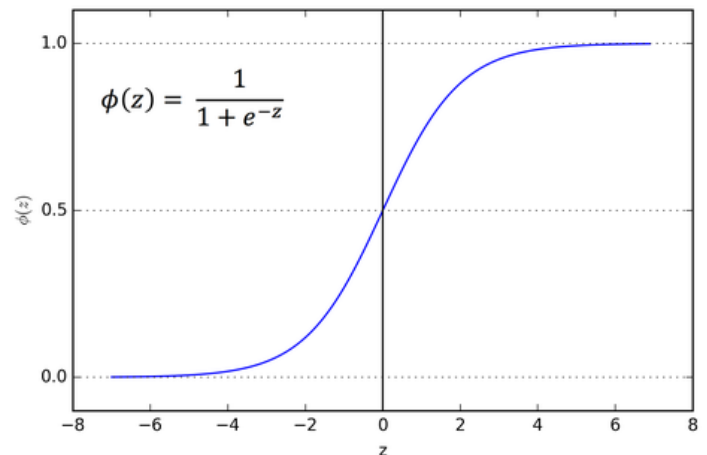
~~the normalization factor is calculated as follows:~~ Note that due to the asymmetry of the softmax loss, the normalization is independently performed two times: across images and across texts [36]. ~~where  $t$  is a global freely learnable parameter~~

$$-\frac{1}{2|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \left( \overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_i \cdot \mathbf{y}_j}}}^{\text{image} \rightarrow \text{text softmax}} + \overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_j \cdot \mathbf{y}_i}}}^{\text{text} \rightarrow \text{image softmax}} \right)$$

The solution is to use ... a Sigmoid!

	I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>	I <sub>4</sub>	I <sub>5</sub>	I <sub>6</sub>	I <sub>7</sub>	I <sub>8</sub>	I <sub>9</sub>	I <sub>10</sub>	I <sub>11</sub>	I <sub>12</sub>
T <sub>1</sub>												
T <sub>2</sub>												
T <sub>3</sub>												
T <sub>4</sub>												
T <sub>5</sub>												
T <sub>6</sub>												
T <sub>7</sub>												
T <sub>8</sub>												
T <sub>9</sub>												
T <sub>10</sub>												
T <sub>11</sub>												
T <sub>12</sub>												

$$-\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \underbrace{\log \frac{1}{1 + e^{z_{ij}(-t\mathbf{x}_i \cdot \mathbf{y}_j + b)}}}_{\mathcal{L}_{ij}}$$



# Parallel computation

