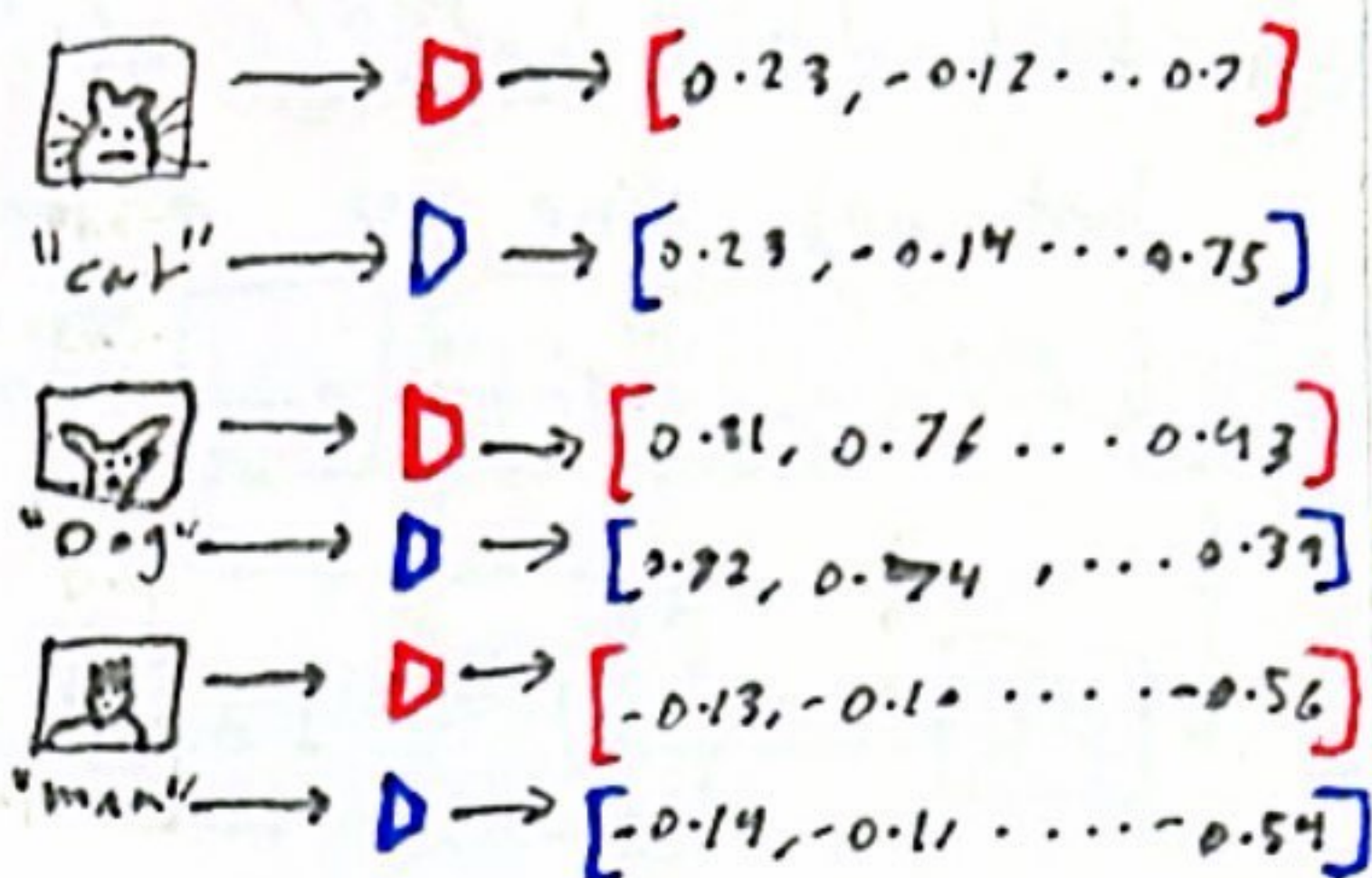


- in Feb 2021 OpenAI released a new model architecture called CLIP trained on 400 million image and caption pairs from the web.

Image Encoder Text Encoder

- CLIP is composed of two models one that processes text and one that processes images. The output of each is a 512 vector. The idea is that both vectors should be similar.



- from here the idea is to make use of the similarity not just between the corresponding images and captions but all image caption pairs in the batch when training the model.

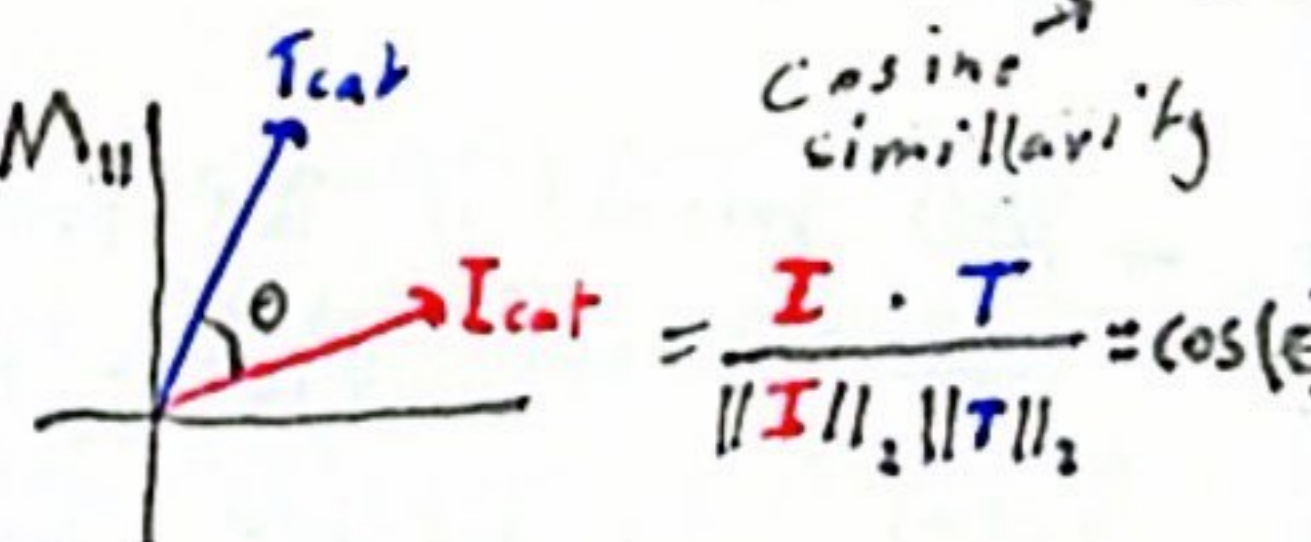
- If we arrange our img vectors as a col and text vector as rows in a matrix the pairs of vectors along the diagonal $= \dots$ correspond to matching imgs and caption

M

"cat"	$\begin{bmatrix} 0.23 \\ \vdots \end{bmatrix}$ +	$\begin{bmatrix} \dots \end{bmatrix}$ -	$\begin{bmatrix} \dots \end{bmatrix}$ -
"dog"	$\begin{bmatrix} \dots \end{bmatrix}$ -	$\begin{bmatrix} \dots \end{bmatrix}$ +	$\begin{bmatrix} \dots \end{bmatrix}$ -
"man"	$\begin{bmatrix} \dots \end{bmatrix}$ -	$\begin{bmatrix} \dots \end{bmatrix}$ -	$\begin{bmatrix} \dots \end{bmatrix}$ +

- all other pairs are non matching images and captions. The CLIP training objective is to maximize the similarity between matching pairs $+$ while minimizing the similarity between non matching pairs $-$

- The algorithm measures similarity between vectors using cosine similarity as the metric. The vectors point in some high dim space.



- The cosine measures the cosine of the angle between our vectors " θ " in this space. So if the angle between the vectors is 0 the cosine score will be the max $= 1$. So the img and text models that make up CLIP are trained to maximize the alignment of related img caption pairs while minimizing alignment for non related pairs.
- This shared vector space is 78 called the embedding space.