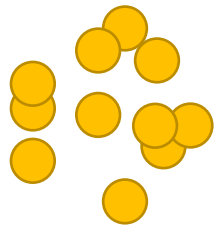
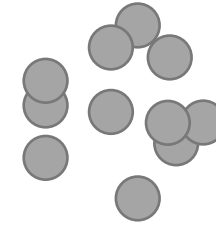


Contrastive Learning

Theory, implementation and a popular example: CLIP

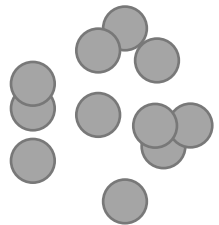


Tim Landgraf - Würzburg - 19 April 2023



Contrastive Learning

Theory, implementation and a popular example: CLIP



Tim Landgraf - Würzburg - 19 April 2023



„Learning“ = finding the best set of parameters

low-level abstraction
high-dimensional inputs



high-level abstraction
low-dimensional outputs

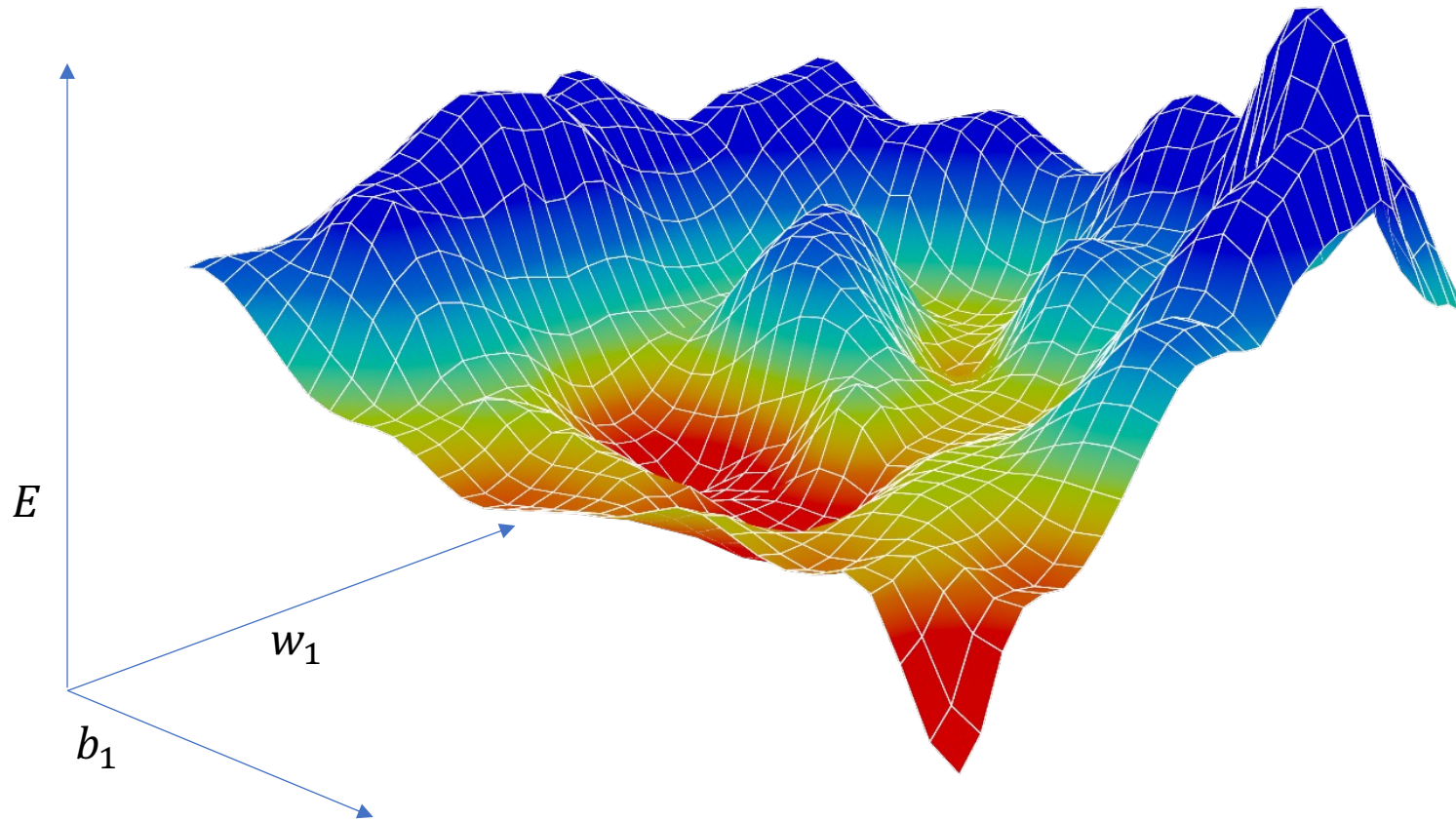


- cat
- dog
- flower
- shoe

```
model.fit(train_images,  
train_labels, epochs=5)
```



Network error as a function of the model parameters (weights and biases)



Goal: Find weight combination that **minimizes network error**

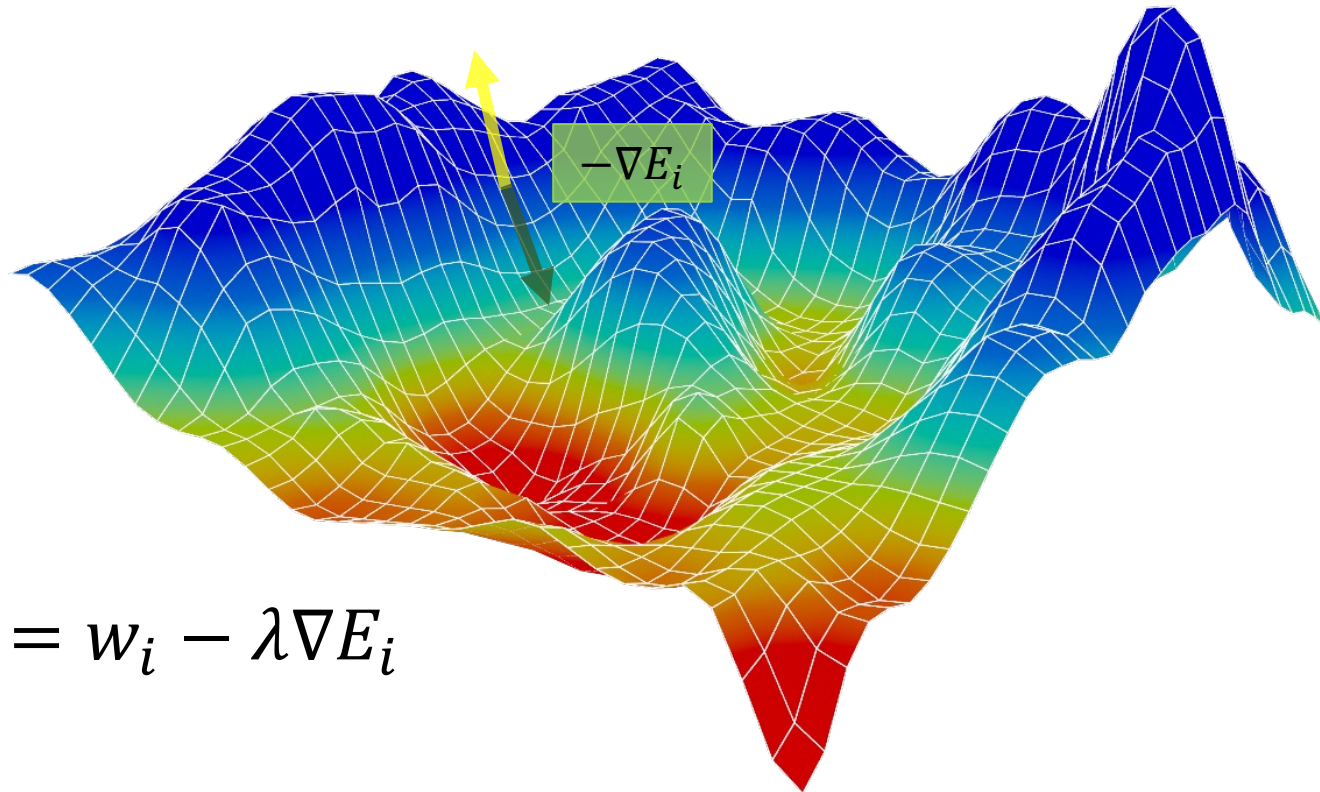


„Learning“ = walking downhill

The **gradient** points **uphill**, so let's walk in **the opposite direction**.

„How much does the error change, for a unit change of weights“

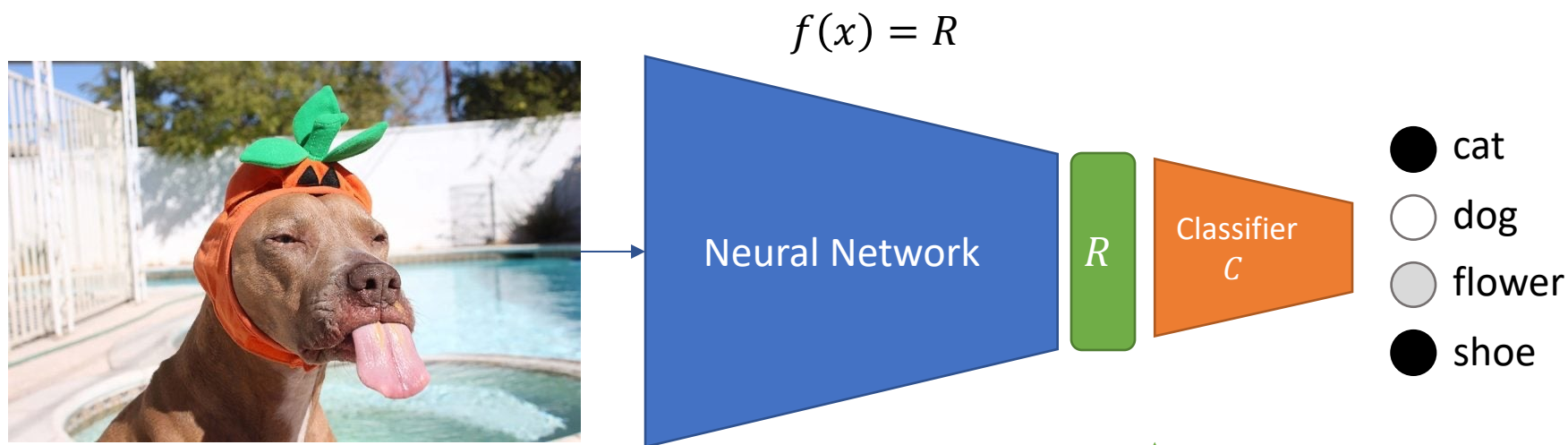
$$\nabla E = \frac{\partial E}{\partial \mathbf{w}} = \begin{pmatrix} \frac{\partial E}{\partial w_1} \\ \vdots \\ \frac{\partial E}{\partial w_N} \end{pmatrix}$$



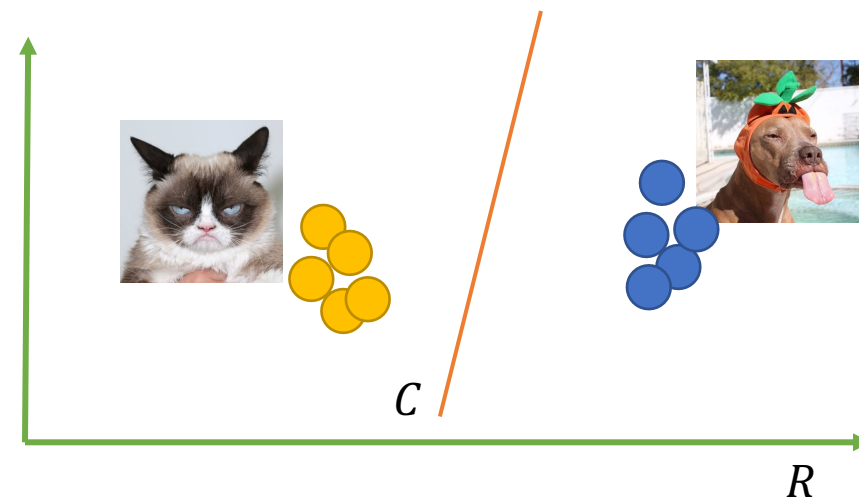
$$w_{i+1} = w_i - \lambda \nabla E_i$$



Learning = finding „good“ *representations*

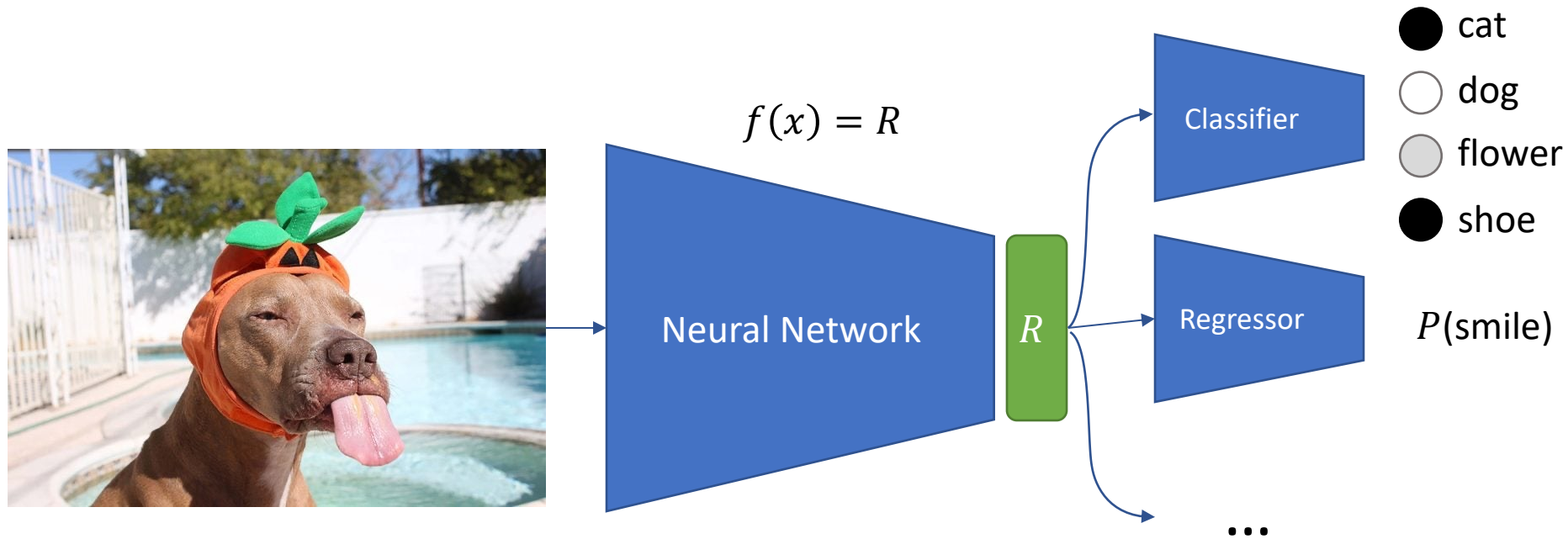


- Neural networks can be seen as simultaneously ...
 - learn **representations** of input data
 - **classify** from those representations



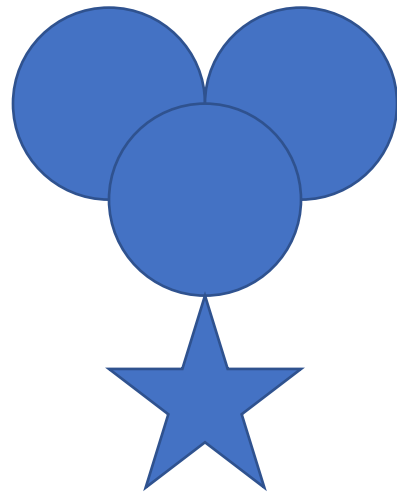


Contrastive Learning used for ...

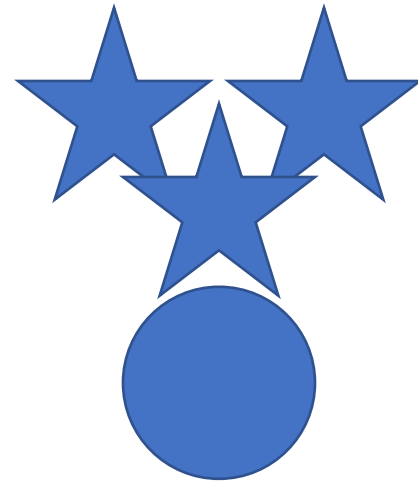


1. Dimensionality *Reduction*
2. Learning useful *representations*
3. *Robust* training objectives
4. *Efficient* use of data

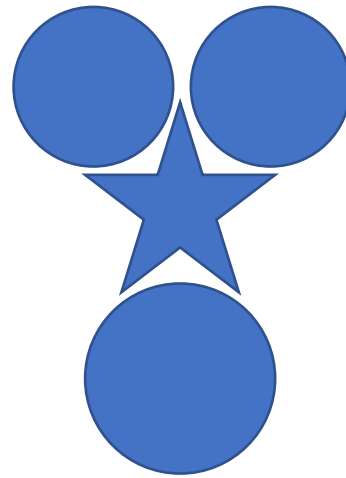
How do networks learn
non-contrastive?



GLORBOXL



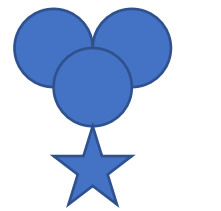
GLIMBIBBLE



GLIMBIBBLE



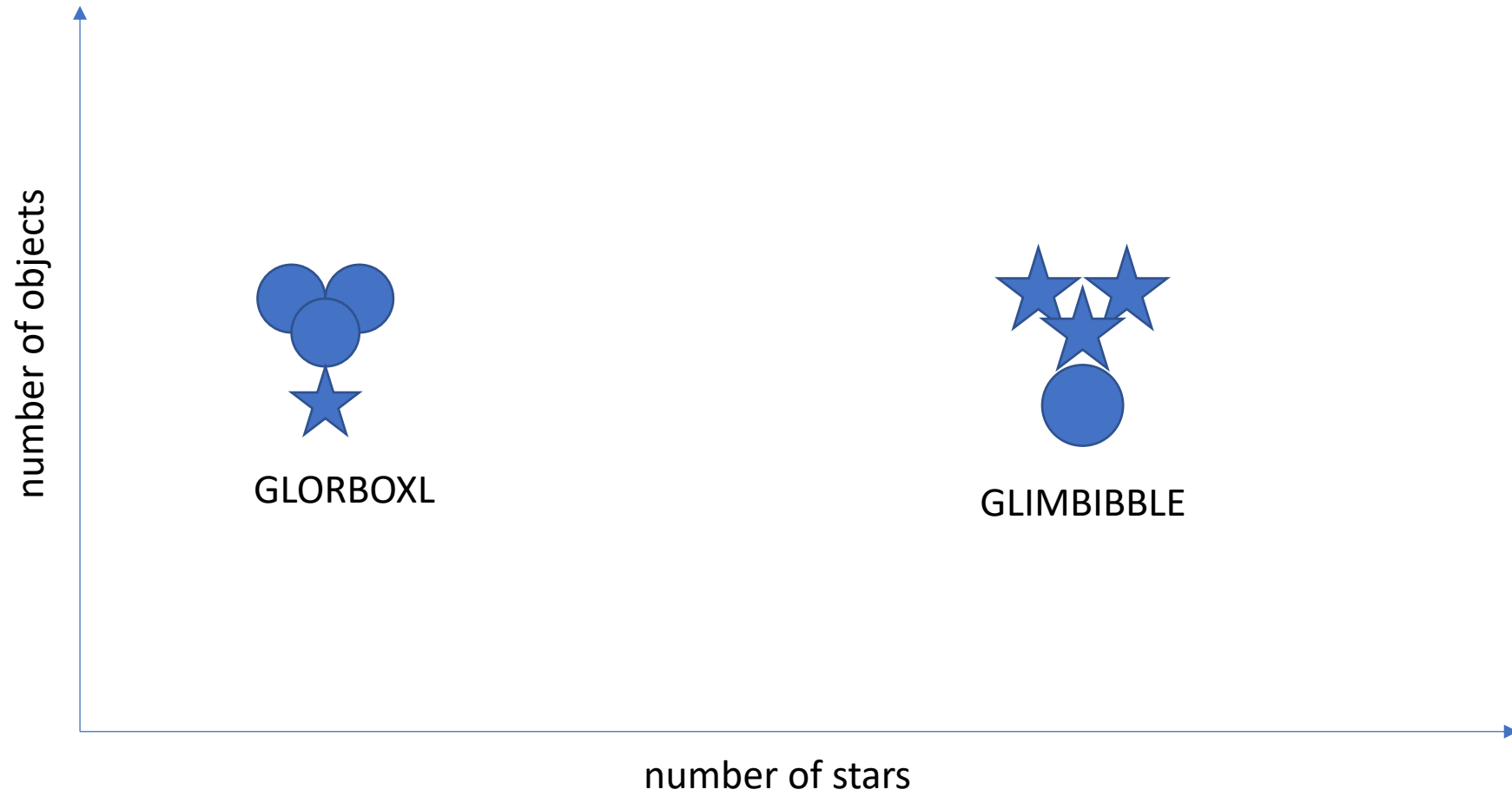
GLORBOXL

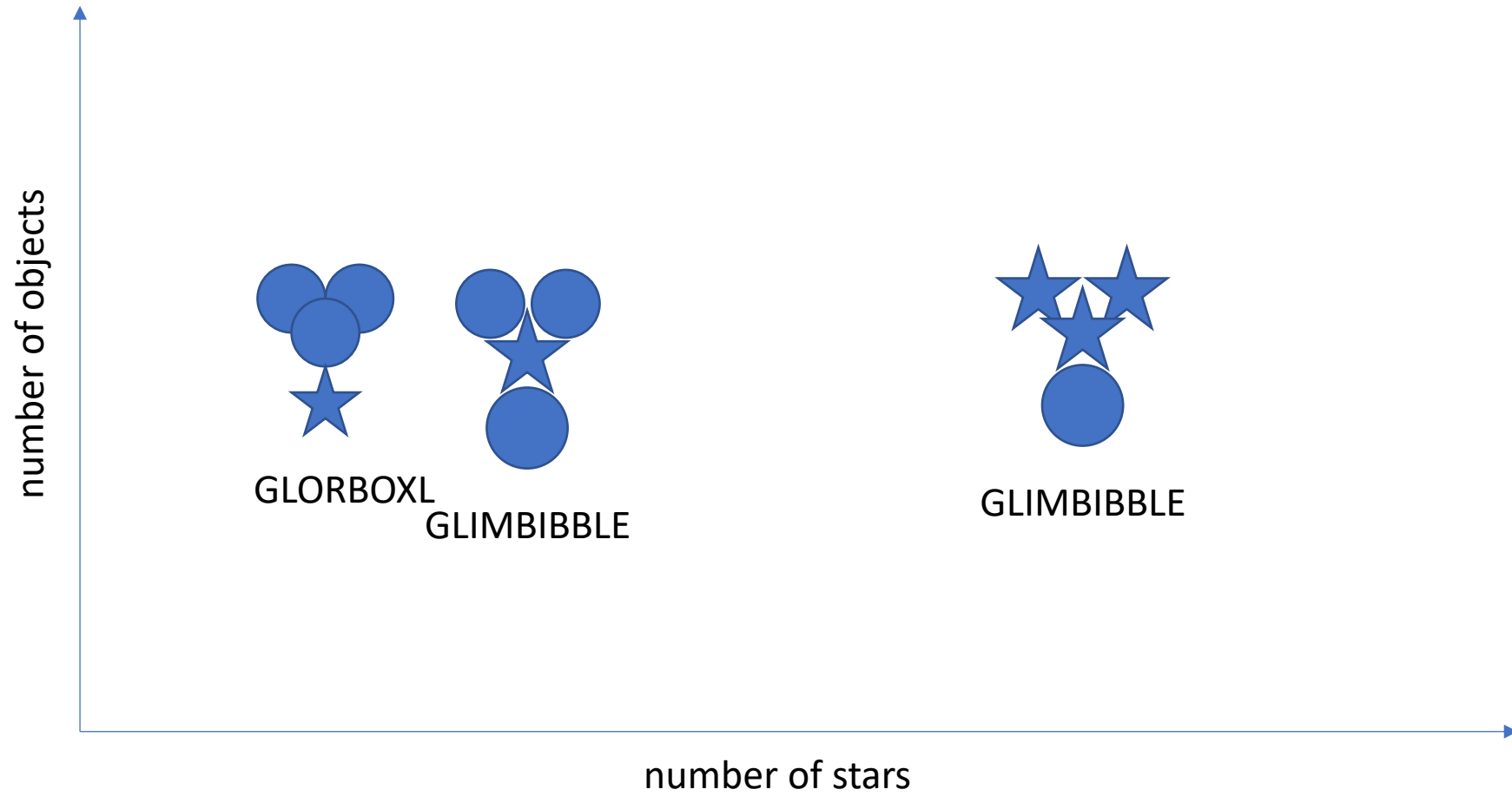


GLORBOXL



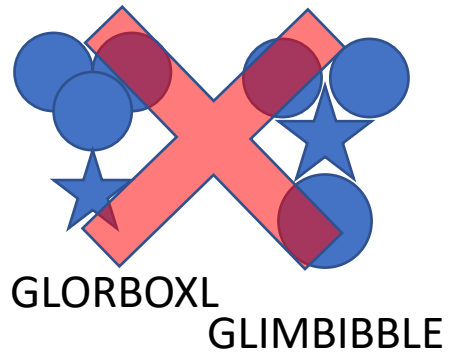
GLIMBIBBLE





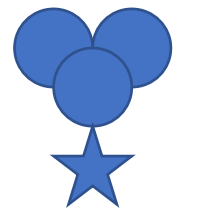


number of objects



number of stars





GLORBOXL



GLORBOXL

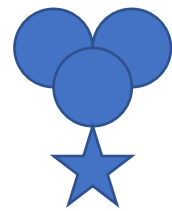


GLIMBIBBLE



GLIMBIBBLE





GLORBOXL



GLORBOXL



GLIMBIBBLE



GLIMBIBBLE

yes

no

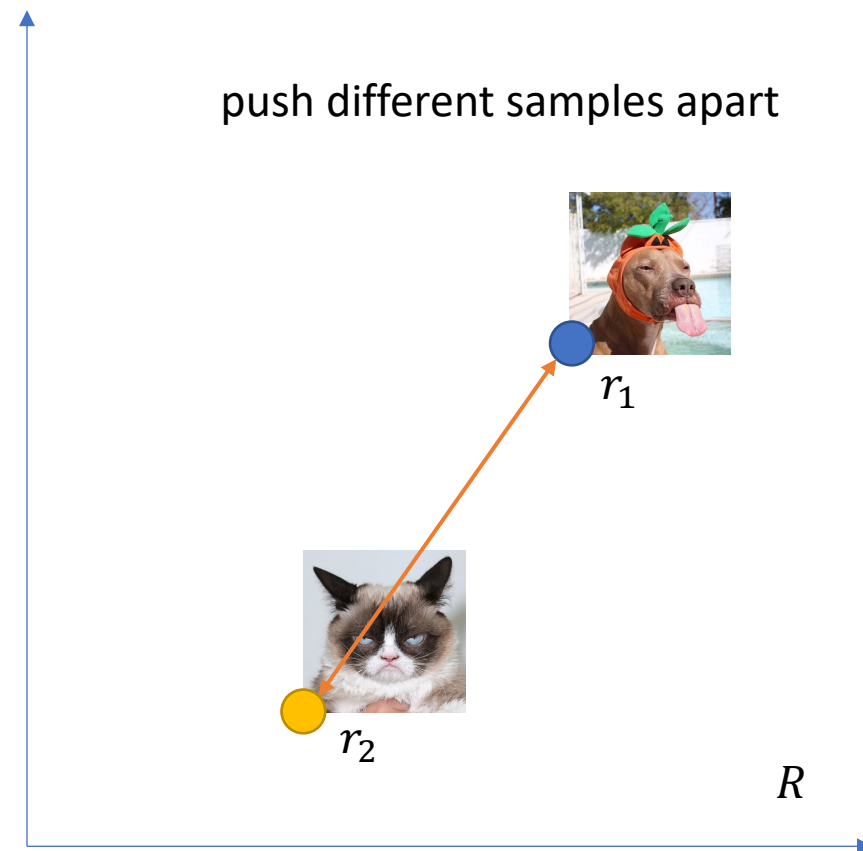
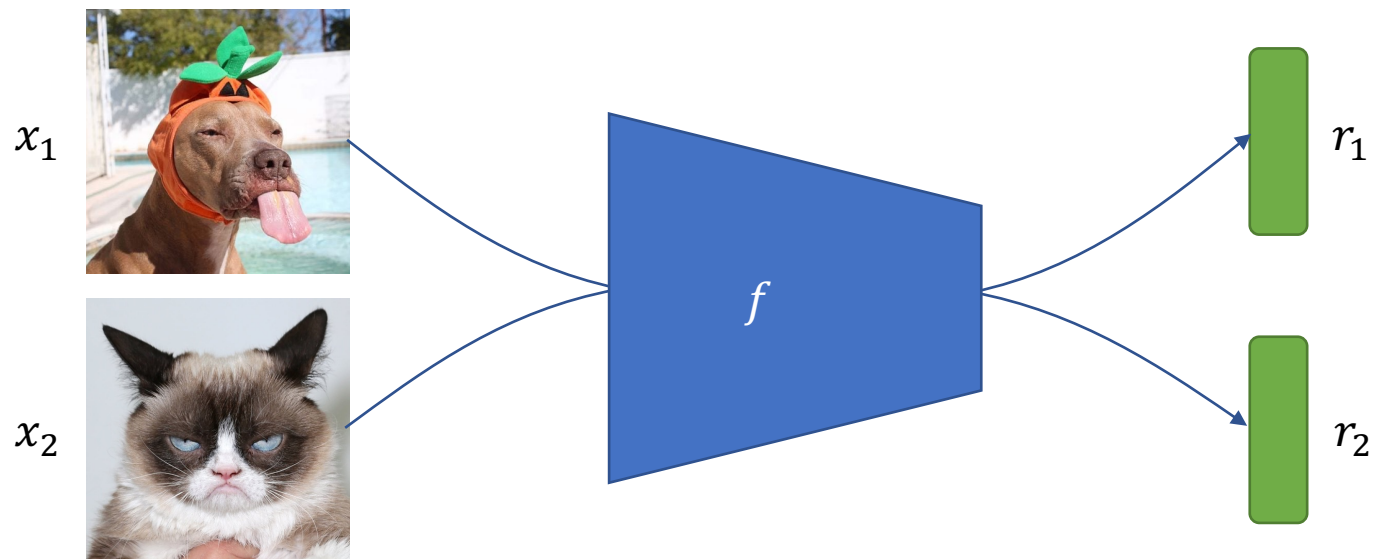
Star in „piercing“ formation?



Sixt, L., Schuessler, M., Popescu, O. I., Weiß, P., & Landgraf, T. Do Users Benefit From Interpretable Vision? A User Study, Baseline, And Dataset. In *International Conference on Learning Representations*.

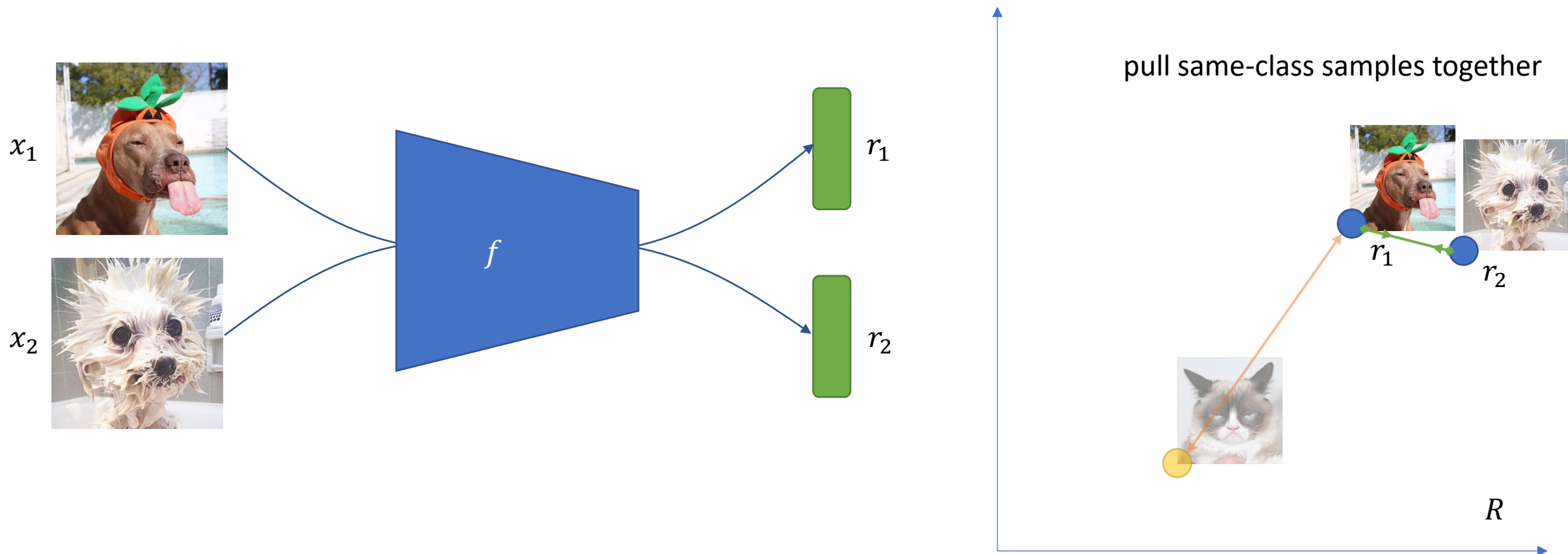


Contrastive Learning Idea (Supervised Case)



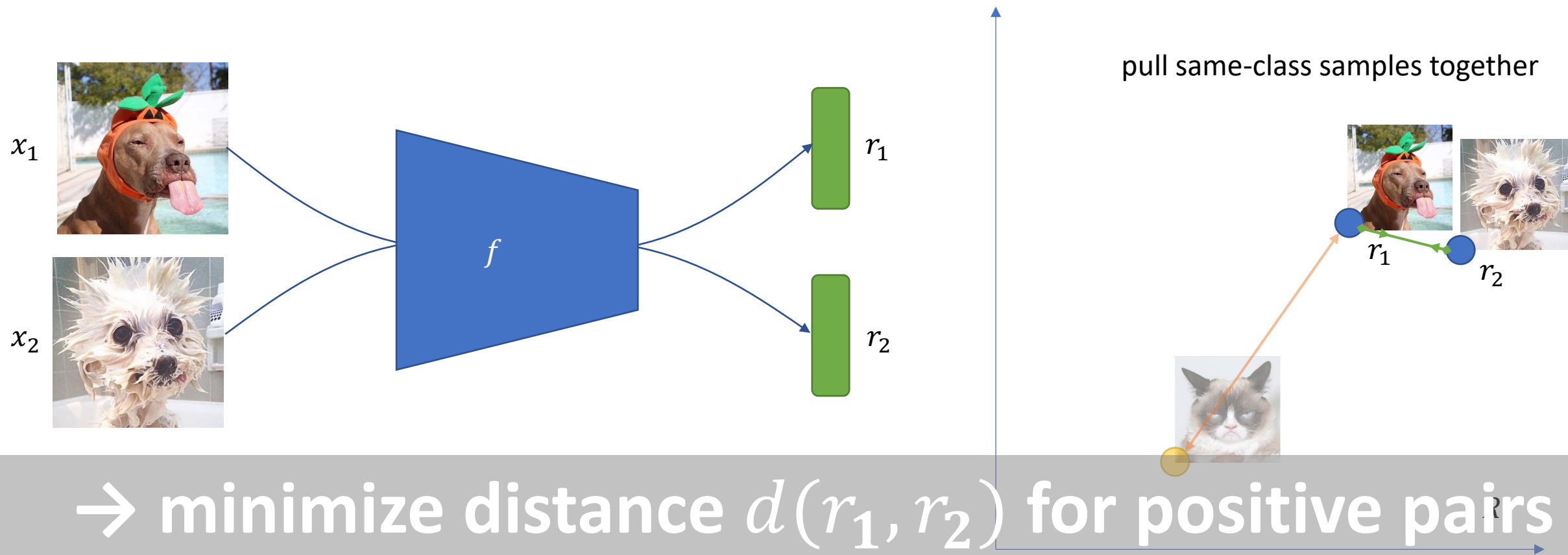


Contrastive Learning Idea (Supervised Case)



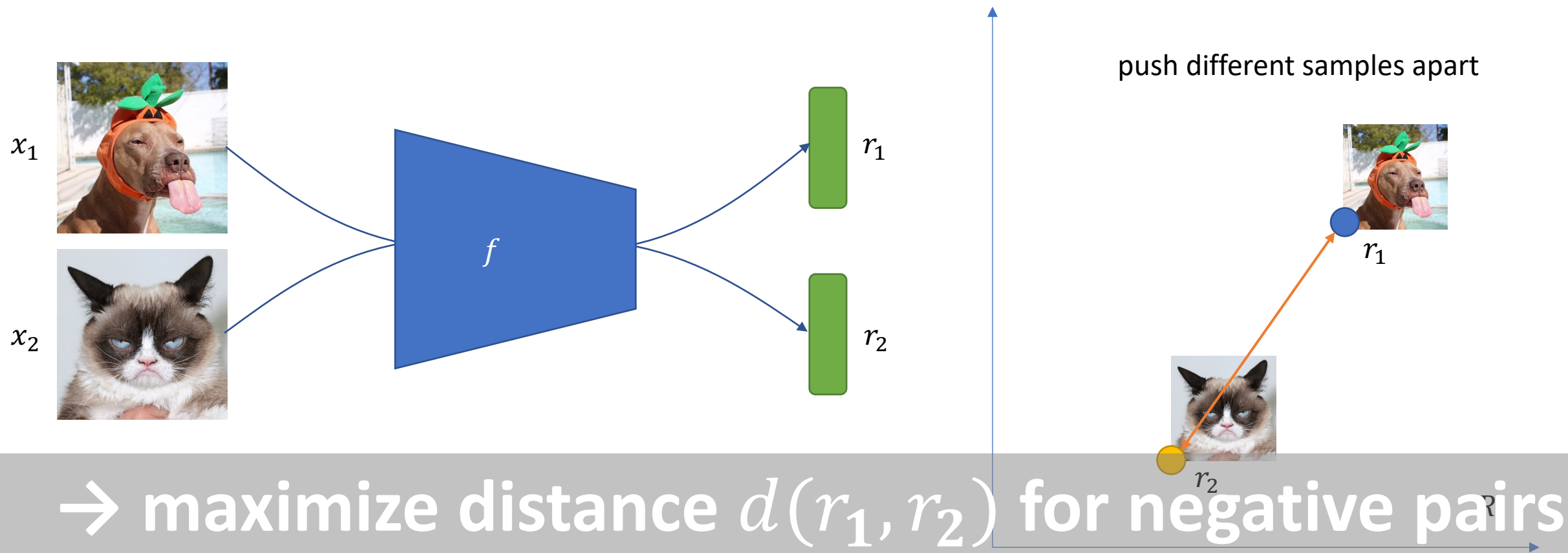


Contrastive Learning Idea (Supervised Case)





Contrastive Learning Idea (Supervised Case)





Simple Contrastive Loss

$$q = \begin{cases} \mathbf{1}, & \mathbf{y}_i = \mathbf{y}_j \quad \text{same label} \\ \mathbf{0}, & \mathbf{y}_i \neq \mathbf{y}_j \quad \text{different label} \end{cases}$$

$$L(\mathbf{x}_i, \mathbf{x}_j) = qd(r_i, r_j) - (1 - q)d(r_i, r_j)$$



Simple Contrastive Loss

$$q = \begin{cases} \mathbf{1}, & \mathbf{y}_i = \mathbf{y}_j \\ \mathbf{0}, & \mathbf{y}_i \neq \mathbf{y}_j \end{cases} \quad \begin{array}{l} \text{same label} \\ \text{different label} \end{array}$$



$$L(\mathbf{x}_i, \mathbf{x}_j) = qd(r_i, r_j) - (1 - q)d(r_i, r_j)$$



Simple Contrastive Loss

$$q = \begin{cases} \mathbf{1}, & \mathbf{y}_i = \mathbf{y}_j \\ \mathbf{0}, & \mathbf{y}_i \neq \mathbf{y}_j \end{cases}$$

same label
different label



$$L(\mathbf{x}_i, \mathbf{x}_j) = \underbrace{q d(\mathbf{r}_i, \mathbf{r}_j) - (1 - q) d(\mathbf{r}_i, \mathbf{r}_j)}$$

is minimized when $d(\mathbf{r}_i, \mathbf{r}_j)$ maximized



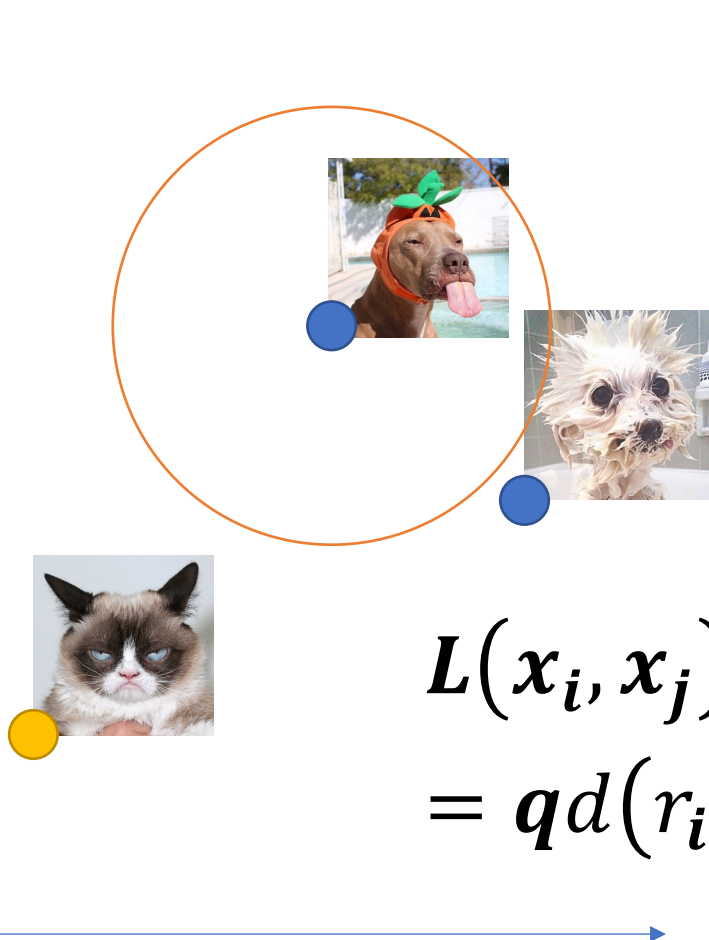
Contrastive Loss with Margin

$$q = \begin{cases} \mathbf{1}, & \mathbf{y}_i = \mathbf{y}_j \\ \mathbf{0}, & \mathbf{y}_i \neq \mathbf{y}_j \end{cases}$$

$$L(\mathbf{x}_i, \mathbf{x}_j) = qd(r_i, r_j) - (\mathbf{1} - q)\max\left(\mathbf{0}, \underset{\substack{\uparrow \\ \text{margin } m}}{m} - d(r_i, r_j)\right)$$



Contrastive Loss with Margin



results in 0, because grumpy cat is too far away

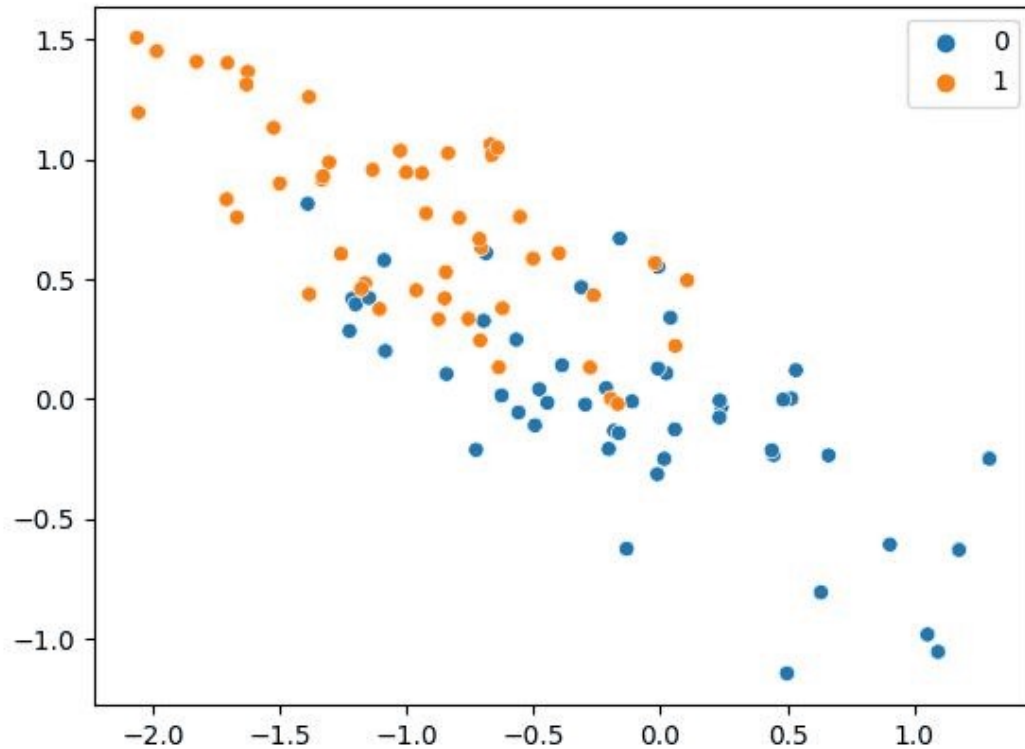
$$L(x_i, x_j) = qd(r_i, r_j) - (1 - q) \max(0, m - d(r_i, r_j))$$



Example with code!

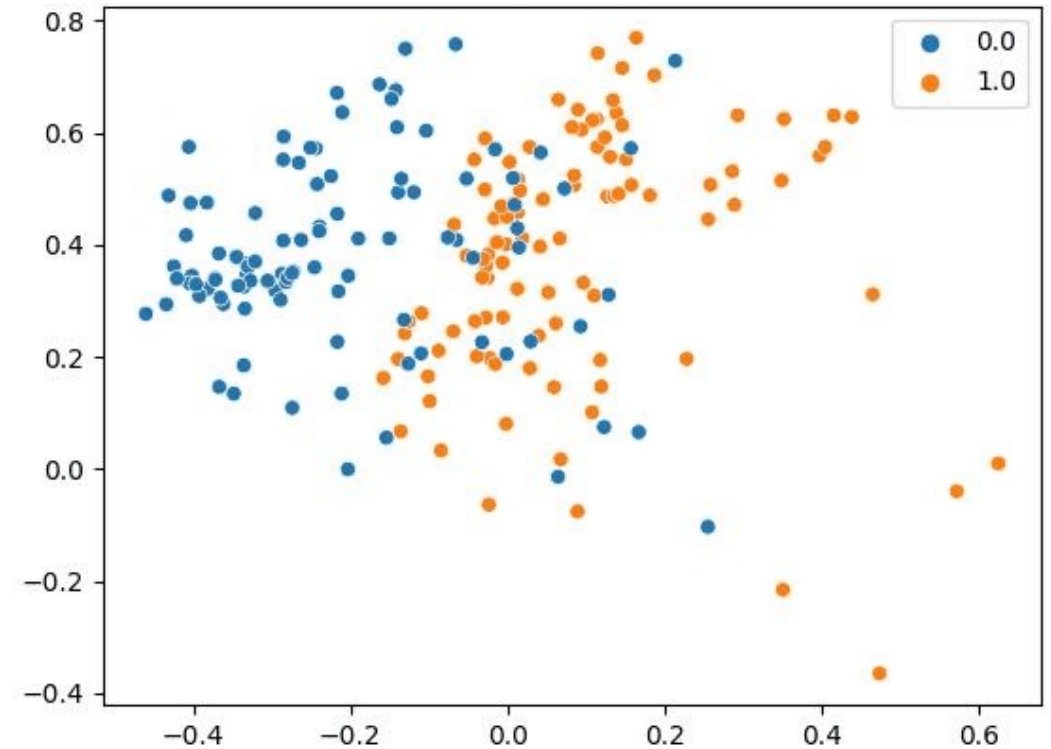
linear mapping

$x \in \mathbb{R}^5, r \in \mathbb{R}^2$



$$Ax = r$$

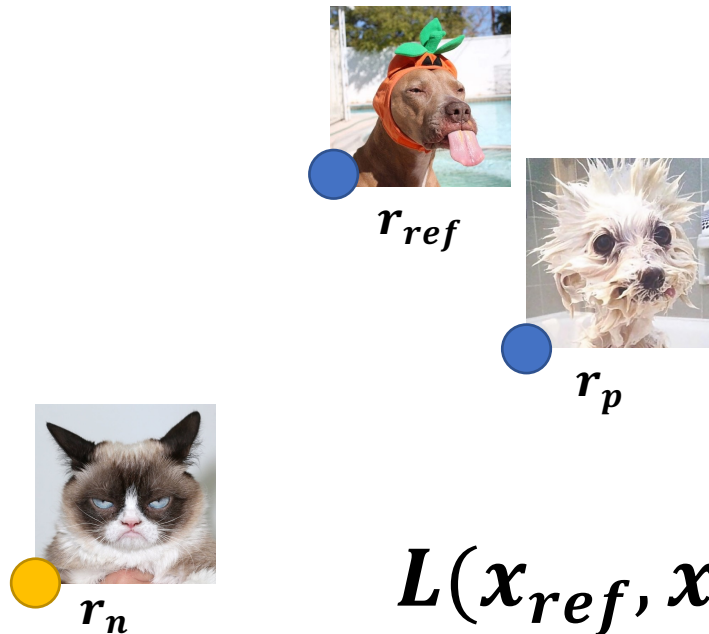
nonlinear mapping



$$NN(x) = r$$



Other Contrastive Losses



„Triplet Loss“

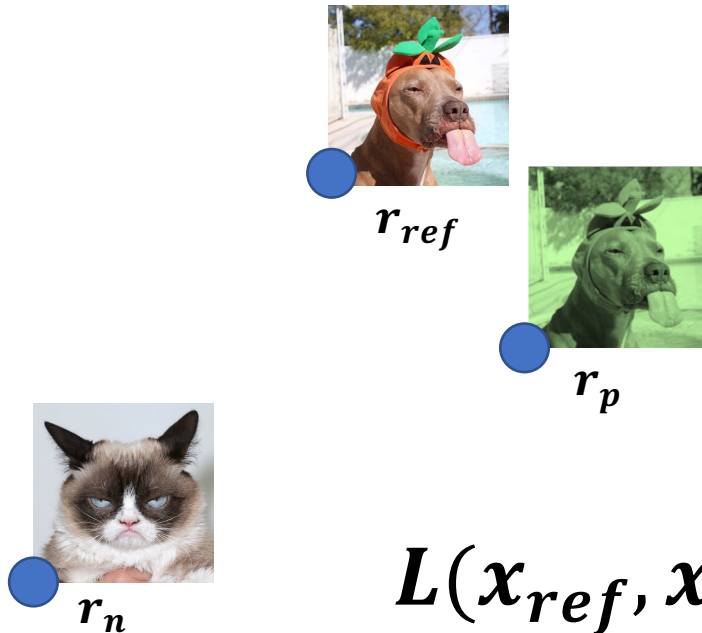
$$L(x_{ref}, x_p, x_n) = d(r_{ref}, r_p) + \alpha - d(r_{ref}, r_n)$$



Contrastive Unsupervised Learning

Positive samples generated by image augmentations

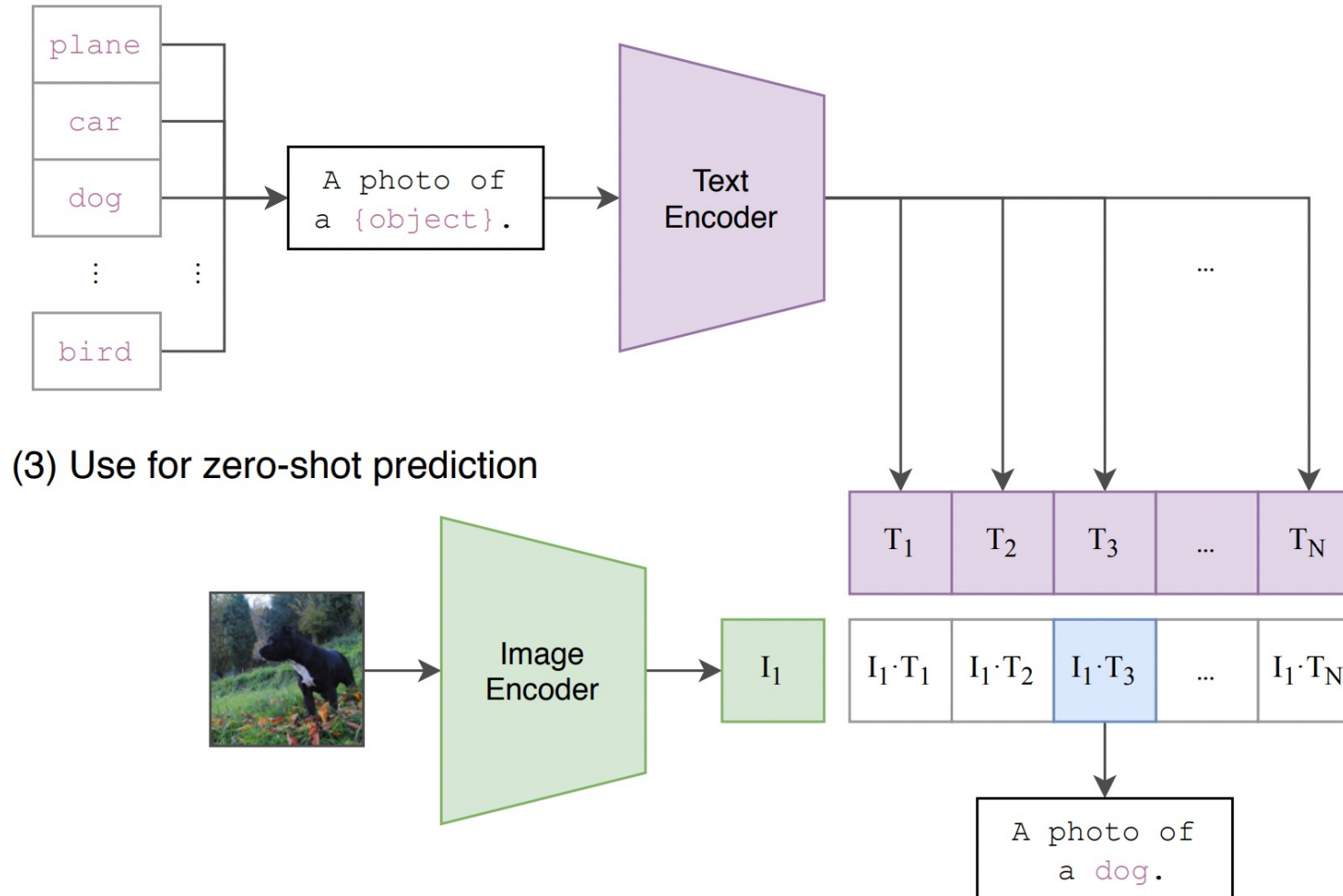
- Color shifts
- Added noise
- Shear, rotations
- Crops
- ...



$$L(x_{ref}, x_p, x_n) = d(r_{ref}, r_p) + \alpha - d(r_{ref}, r_n)$$

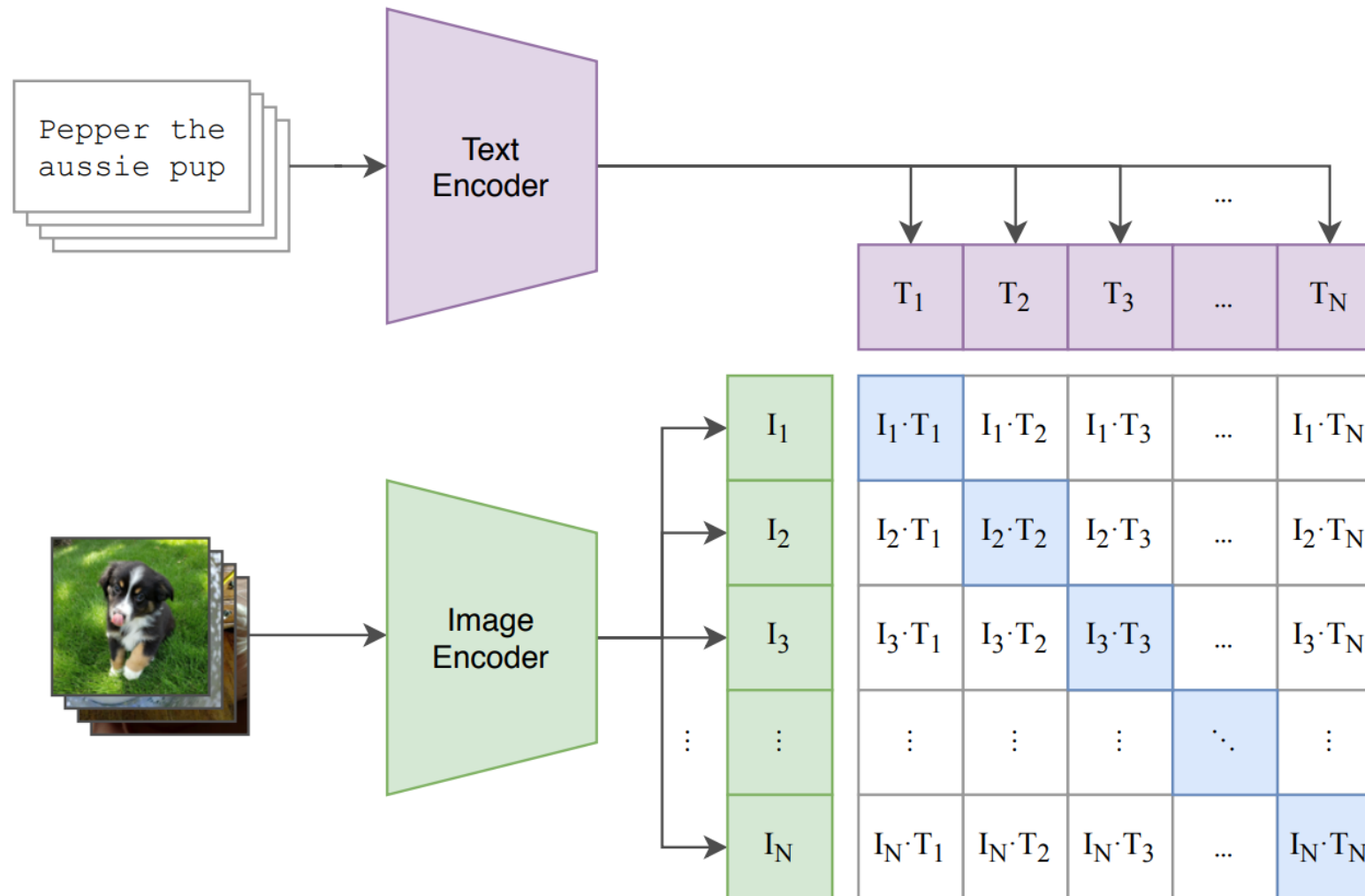


Contrastive Loss for learning multi-modal representations: CLIP!





Contrastive Loss for learning multi-modal representations: CLIP!





Food101

guacamole (90.1%) Ranked 1 out of 101 labels



✓ a photo of **guacamole**, a type of food.

✗ a photo of ceviche, a type of food.

✗ a photo of edamame, a type of food.

✗ a photo of tuna tartare, a type of food.

✗ a photo of hummus, a type of food.

Youtube-BB

airplane, person (89.0%) Ranked 1 out of 23 labels



✓ a photo of a **airplane**.

✗ a photo of a bird.

✗ a photo of a bear.

✗ a photo of a giraffe.

✗ a photo of a car.

SUN397

television studio (90.2%) Ranked 1 out of 397 labels



✓ a photo of a **television studio**.

✗ a photo of a podium indoor.

✗ a photo of a conference room.

✗ a photo of a lecture room.

✗ a photo of a control room.

EuroSAT

annual crop land (46.5%) Ranked 4 out of 10 labels



✗ a centered satellite photo of permanent crop land.

✗ a centered satellite photo of pasture land.

✗ a centered satellite photo of highway or road.

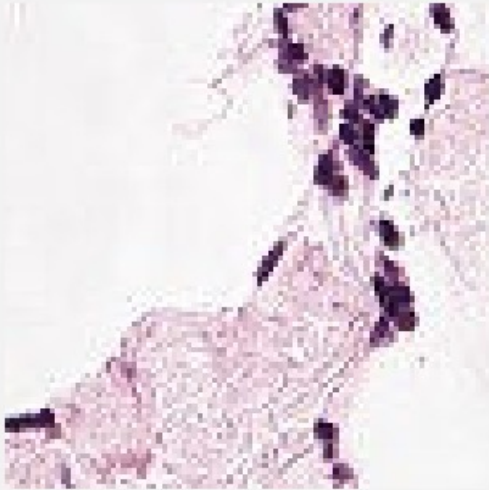
✓ a centered satellite photo of **annual crop land**.

✗ a centered satellite photo of brushland or shrubland.



PatchCamelyon (PCam)

healthy lymph node tissue (77.2%) Ranked 2 out of 2 labels



✗ this is a photo of lymph node tumor tissue

✓ this is a photo of healthy lymph node tissue

ImageNet-A (Adversarial)

lynx (47.9%) Ranked 5 out of 200 labels



Camera Name 30.011n 37°F

01-01-201

✗ a photo of a fox squirrel.

✗ a photo of a mongoose.

✗ a photo of a skunk.

✗ a photo of a red fox.

✓ a photo of a lynx.