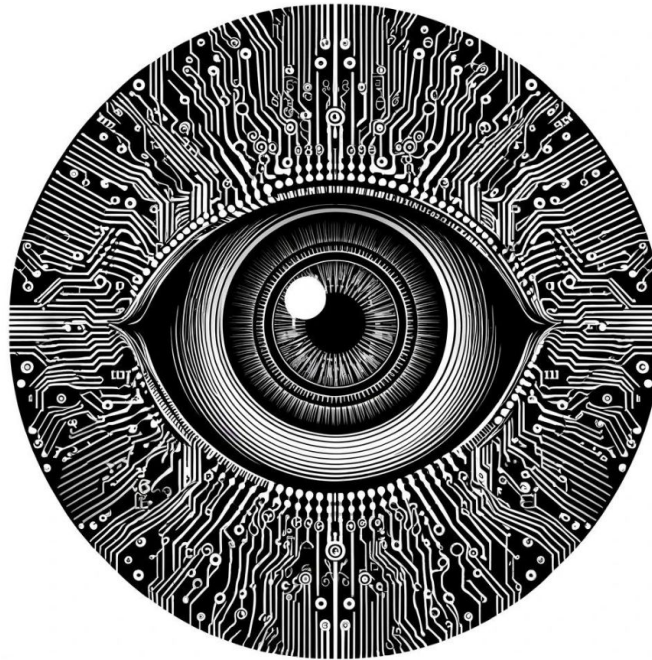


Image Embeddings



Antonio Rueda-Toicen

SPONSORED BY THE



Federal Ministry
of Education
and Research

Learning goals

- Understand dense vector embeddings as learned representations
- Describe how image embeddings are produced in neural networks
- Extract image embeddings out of a pretrained Resnet50
- Inspect image neighborhoods using the cosine similarity metric

Sparse vs dense vector representations

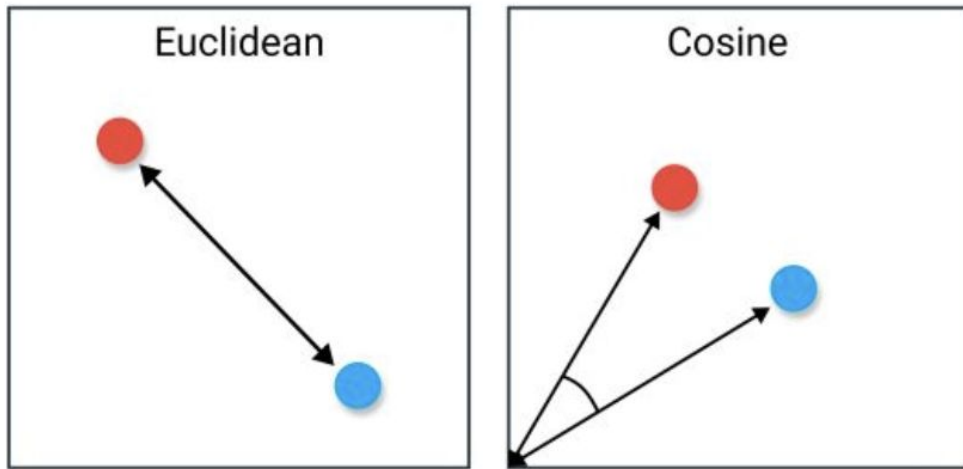
$$\mathbf{x}_1 = \text{[Image of handwritten 3]} \quad \mathbf{x}_2 = \text{[Image of handwritten 3]} \quad 28 \times 28 \text{ pixels} = 784 \text{ values}$$

$$\mathbf{x}_1 \text{ sparse} = \mathbf{x}_2 \text{ sparse} = [0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0]^\top$$

$$\mathbf{x}_2 \text{ dense} = [0.10 \quad -0.75 \quad 0.35 \quad 0.41 \quad -0.20 \quad 0.89 \quad -0.60 \quad 0.03 \quad 0.59 \quad -0.50]^\top$$

$$\mathbf{x}_1 \text{ dense} = [0.12 \quad -0.85 \quad 0.37 \quad 0.44 \quad -0.22 \quad 0.90 \quad -0.63 \quad 0.01 \quad 0.58 \quad -0.49]^\top$$

Common metrics to compare embeddings



$$d_{\text{euclid}}(x_1, x_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}$$

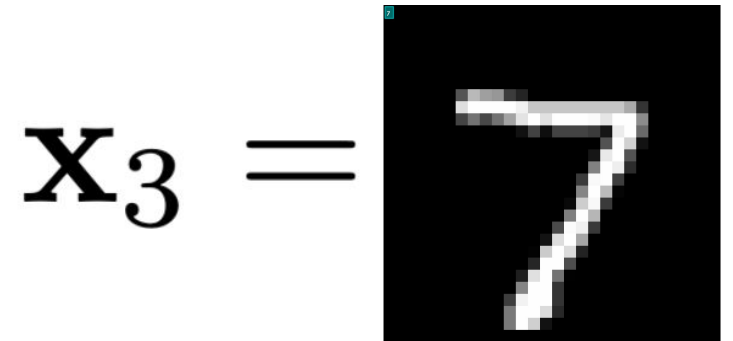
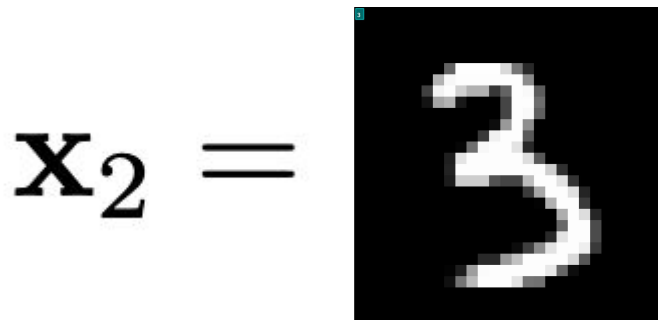
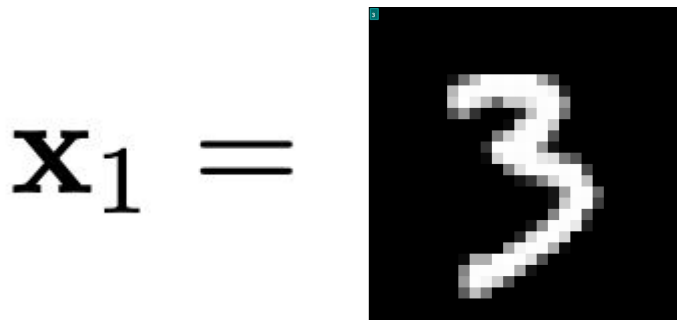
$$\text{cosine similarity}(x_1, x_2) = \frac{x_1 \cdot x_2}{\|x_1\| \|x_2\|}$$

Embeddings as representations learned during training

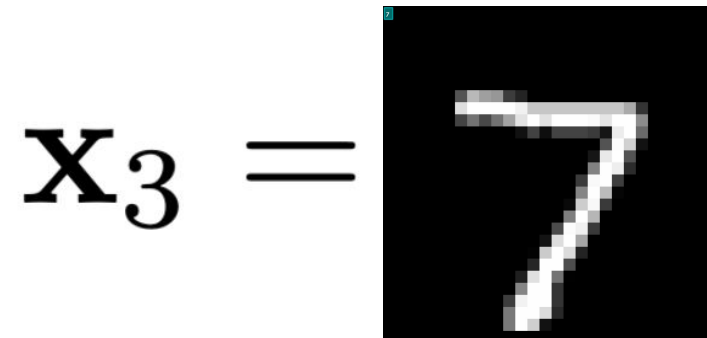
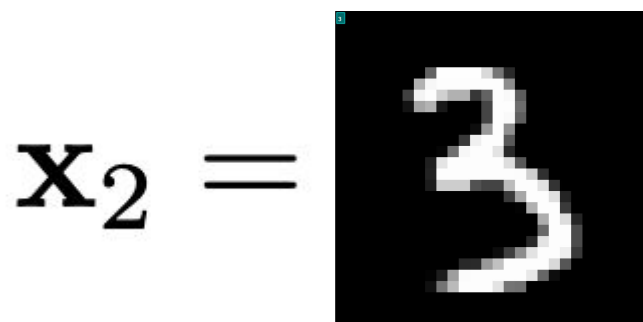
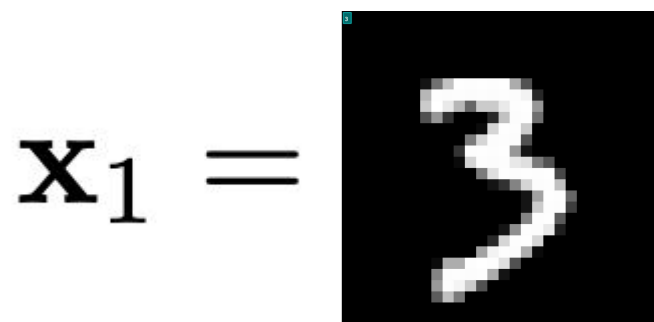
```
import torch.nn as nn
```

```
conv_linear_embedder = nn.Sequential(  
    nn.Conv2d(1, 32, kernel_size=3), # 26x26x32  
    nn.ReLU(),  
    nn.Conv2d(32, 64, kernel_size=3), # 24x24x64  
    nn.ReLU(),  
    nn.Flatten(), # 24*24*64  
    # This layer can be extracted as 64-dimensional  
    # embedding  
    nn.Linear(24*24*64, 64),  
    nn.ReLU(),  
    # This layer can be a 10-dimensional embedding  
    nn.Linear(in_features=64, out_features=10)  
)
```

$$H(p, q) = -\frac{1}{n} \sum_{i=1}^n p(i) \log q(i)$$



Similar objects have similar embeddings

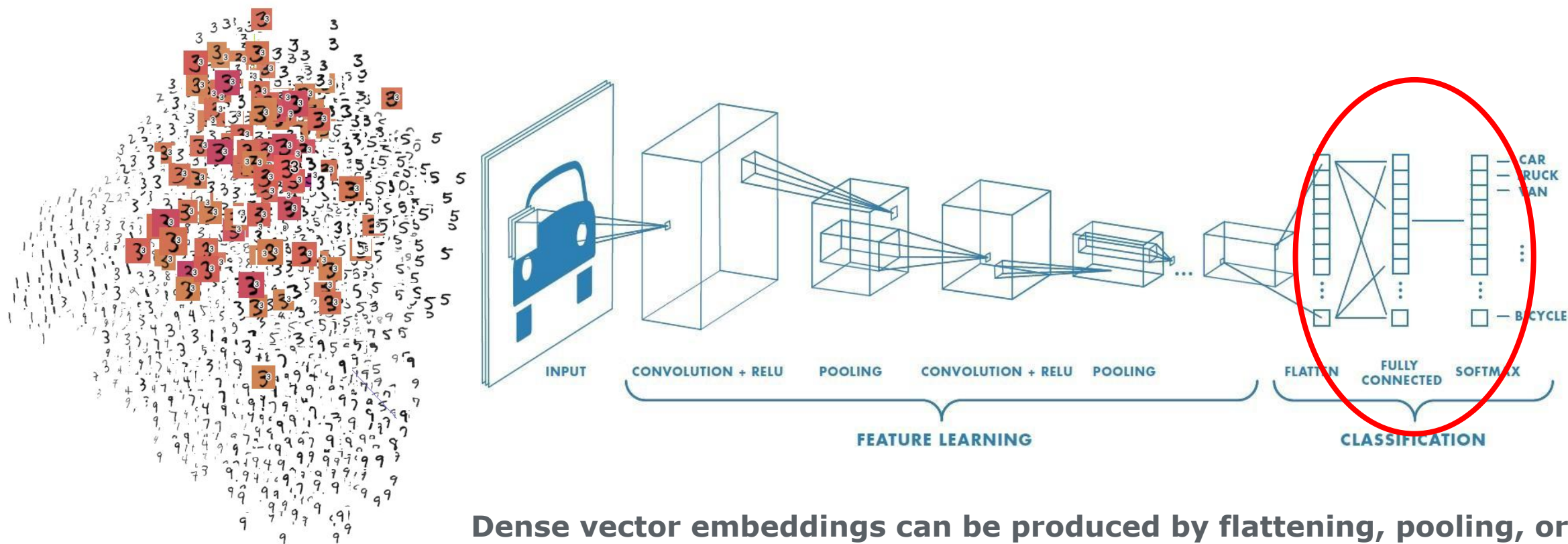


Similar digits (X_1, X_2): High cos similarity (≈ 0.98), Low Euclidean distance (≈ 0.3)

Different digits (X_1, X_3): Low cos similarity (≈ 0.4), High Euclidean distance (≈ 2.1)

Image embeddings as a byproduct of training

When we train a neural network, we make it **embed** inputs that are semantically similar, close to each other



Dense vector embeddings can be produced by flattening, pooling, or feeding (learned) features to linear layers

[Source of image](#) (try exploring interactively)

Approaches to produce embeddings with neural networks

```
import torch.nn as nn
```

```
# Approach 1: Direct embedding
```

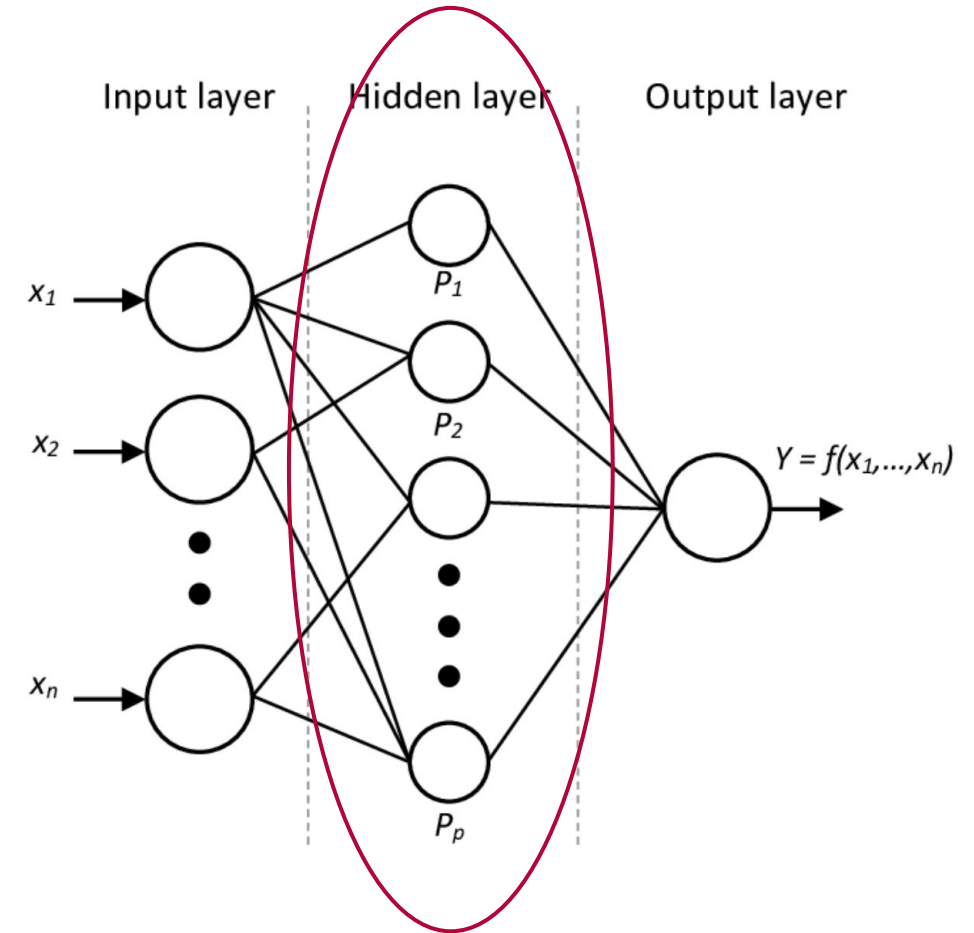
```
linear_embedder = nn.Sequential(  
    nn.Flatten(), # 784  
    nn.Linear(784, 10)  
)
```

```
# Approach 2: Convolution and then fully connected layer
```

```
conv_linear_embedder = nn.Sequential(  
    nn.Conv2d(1, 32, kernel_size=3), # 26x26x32  
    nn.ReLU(),  
    nn.Conv2d(32, 64, kernel_size=3), # 24x24x64  
    nn.ReLU(),  
    nn.Flatten(), # 24*24*64  
    nn.Linear(24*24*64, 10)  
)
```

```
# Approach 3: Global pooling after convolutions
```

```
conv_pool_embedder = nn.Sequential(  
    nn.Conv2d(1, 32, kernel_size=3), # 26x26x32  
    nn.ReLU(),  
    nn.Conv2d(32, 10, kernel_size=3), # 24x24x10  
    nn.ReLU(),  
    nn.AdaptiveAvgPool2d((1, 1)), # 1x1x10  
    nn.Flatten() # 10  
)
```



$$L_{L1} = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

Extract embeddings from a pretrained Resnet50

```
import torch
```

```
from torchvision import models
```

```
# Load pre-trained ResNet50 weights and transforms
```

```
weights = models.ResNet50_Weights.IMAGENET1K_V2
```

```
model = models.resnet50(weights=weights)
```

```
# Preprocess the image using the transformations from the imagenet dataset
```

```
preprocess = weights.transforms()
```

```
# Remove final Imagenet classification layer, keep 2048 dense vector
```

```
# Take time to unpack this syntax in the practical notebook
```

```
embedder = torch.nn.Sequential(*list(model.children())[:-1])
```

```
# Disable BatchNormalization
```

```
embedder.eval()
```

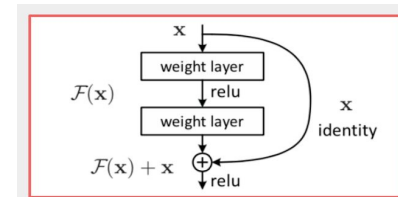
```
# Disable gradient computation
```

```
with torch.inference_mode():
```

```
    # unsqueeze(0) adds the batch size for inference
```

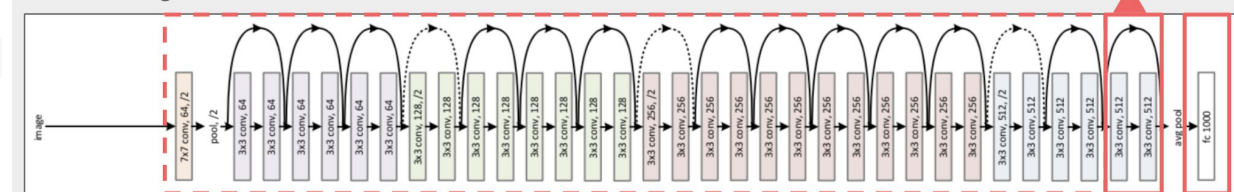
```
    # squeeze() removes the singleton dimensions from the output vector
```

```
    embedding = embedder(preprocess(image).unsqueeze(0)).squeeze()
```

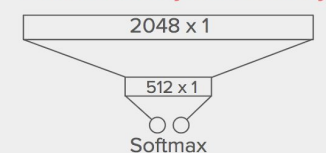


Residual Learning Block

ResNet50 Diagram



Re-architect fully-connected layers



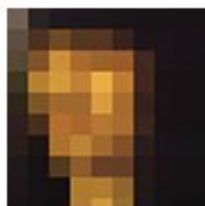
README.md

Art Recommendation system

This is the repository of a portfolio project at [DSR](#). This project aims to identify similar images using pre-trained computer vision networks. For an explanation of the technology see the [technology section](#).

Contributors

- [Catarina Ferreira](#)
- [Gargi Maheshwari](#)

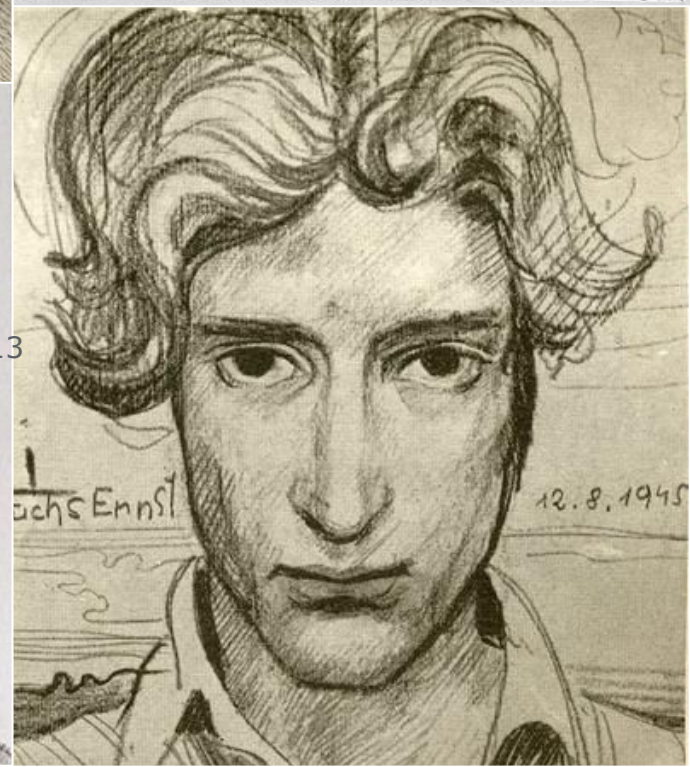
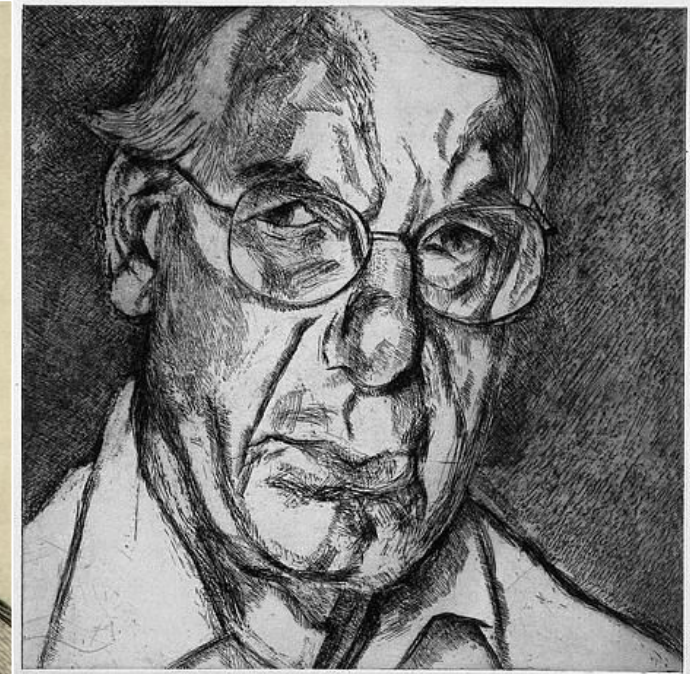


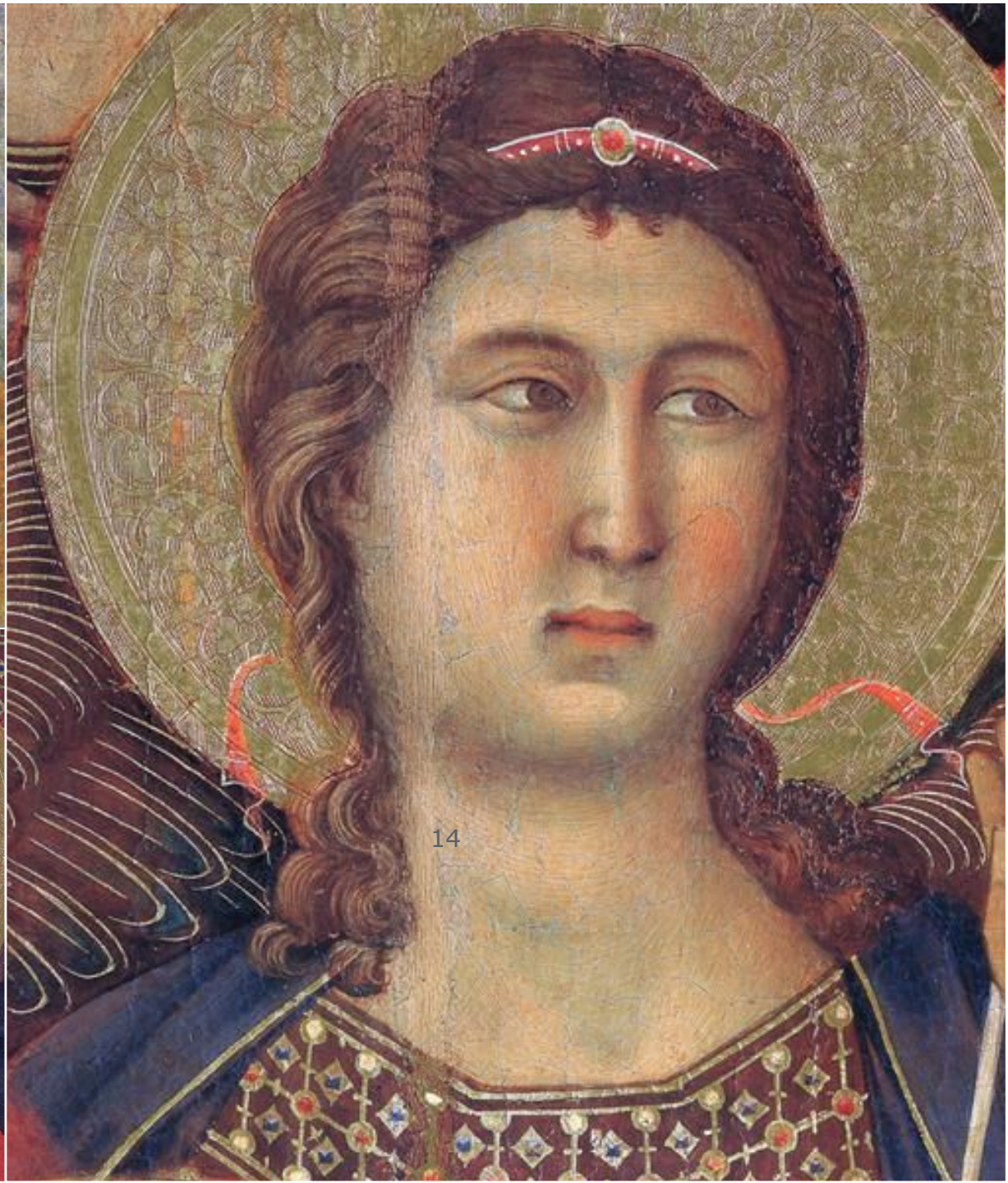
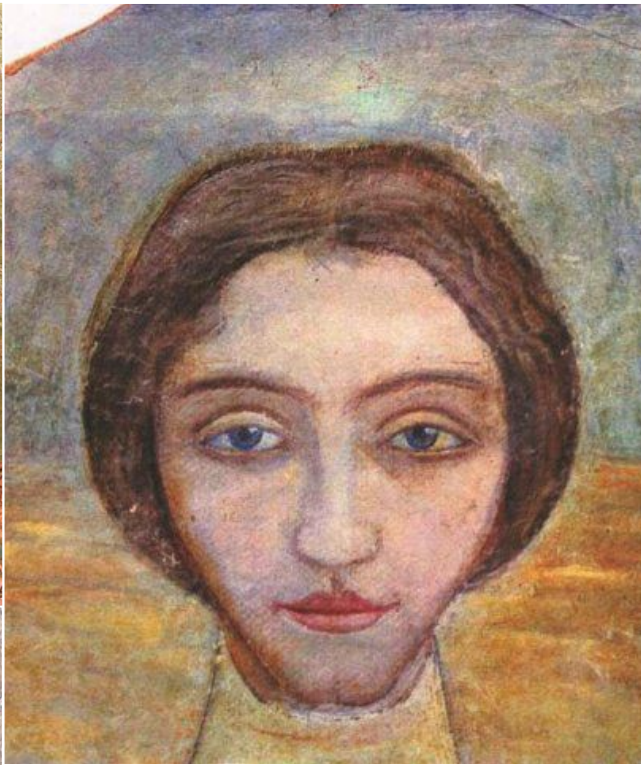
WIKIART
VISUAL ART ENCYCLOPEDIA

<https://github.com/gargimaheshwari/Wikiart-similar-art>









Fine-grained embeddings with triplets

```
import torch
import torch.nn.functional as F
def triplet_loss(query, positive,
                 negative, margin=1.0):
    pos_dist = F.pairwise_distance(query, positive)
    neg_dist = F.pairwise_distance(query, negative)
    loss = F.relu(pos_dist - neg_dist + margin)
    return loss.mean()
```

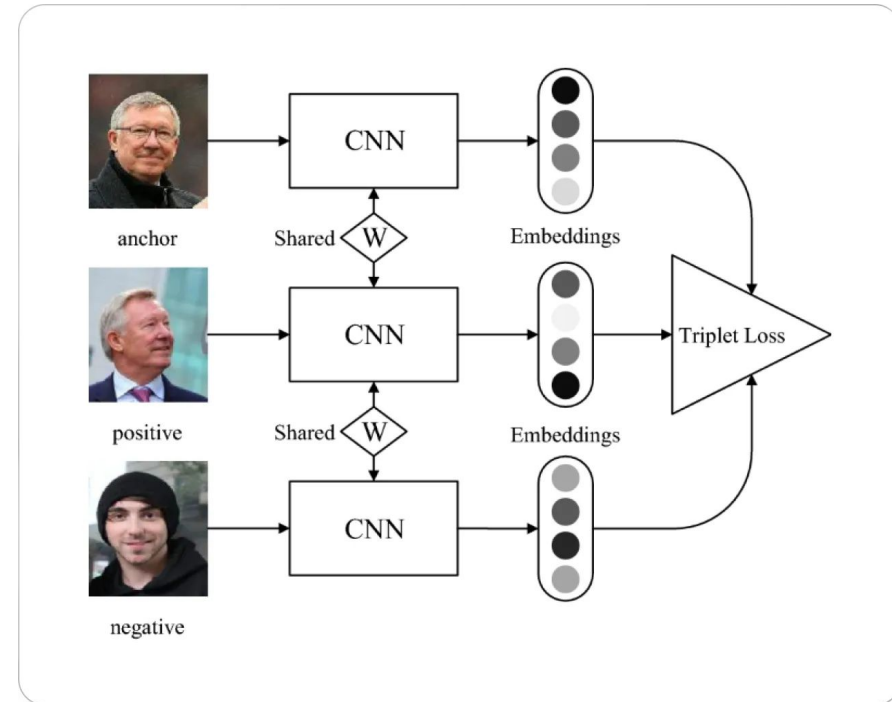


Image from [Triplet Loss: Intro, Implementation, Use Cases](#)

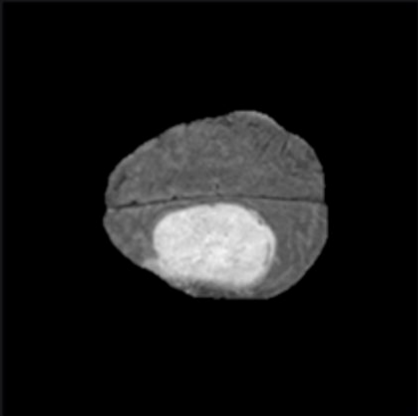
$$L = \max \left(\|f(a) - f(p)\|_2^2 - \|f(a) - f(n)\|_2^2 + \alpha, 0 \right)$$

Application: differential diagnosis

PIXEL

DIAGNOSE

Select Image



Cancer Type:

☒ Gliomas
☒ Meningiomas
☒ Metastasis

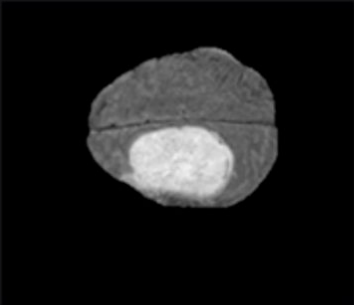
Image Type:

☐ t1
☐ t1c
☐ t2
☒ flair
☐ seg_t1
☐ seg_t1c
☐ seg_t2
☐ seg_flair

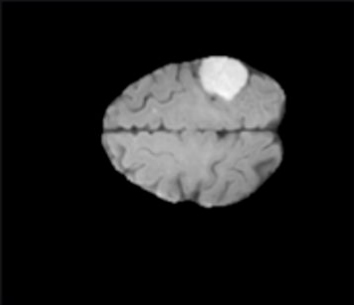
Number of Results per Type:

9

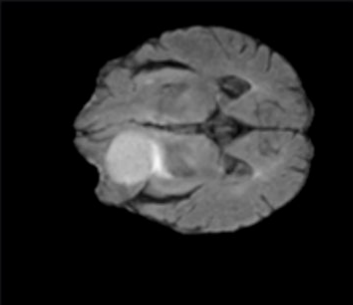
Start Search



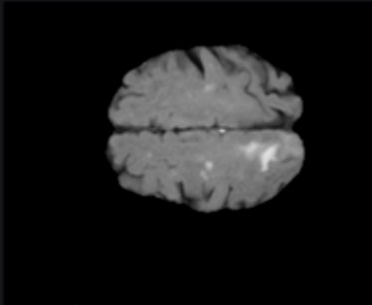
Slice no.: 112 / 150
Similarity Score: 0.9999999




Slice no.: 114 / 150
Similarity Score: 0.8529618



Slice no.: 62 / 150
Similarity Score: 0.82390577



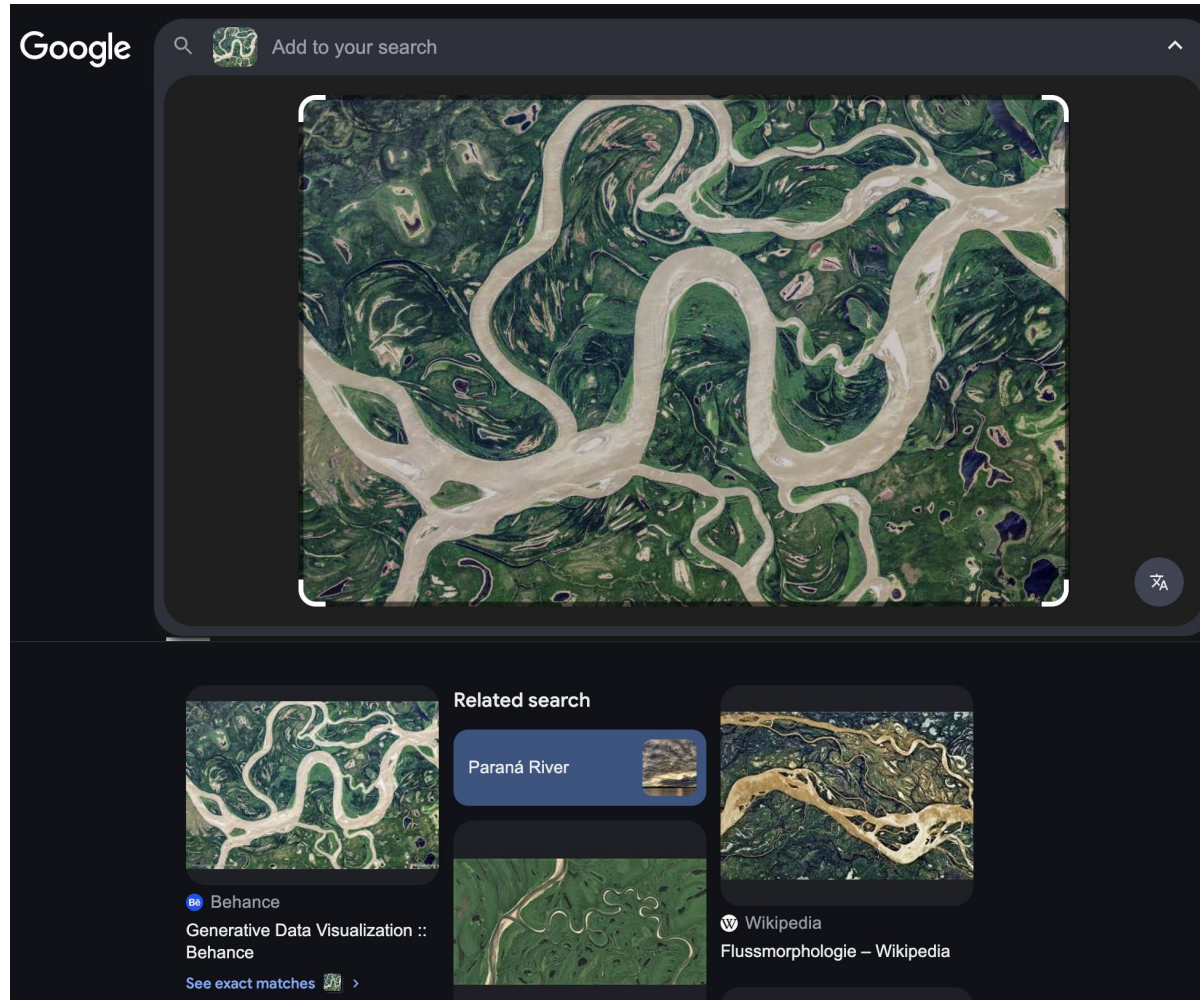
Slice no.: 109 / 150
Similarity Score: 0.8092099



Slice no.: 119 / 150
Similarity Score: 0.8026103

Image from <https://pixel-diagnose.github.io/>

Application: search by image



Search for similar images in images.google.com

Summary

Embeddings are learned vector representations

- Neural networks learn to map similar images to similar vectors through training
- We measure the similarity of embeddings using metrics like cosine similarity and Euclidean distance

There are different architectural approaches to create embeddings

- We can use linear layers, or flattened and pooled convolutions to produce embeddings
- Imagenet-pretrained models produce rich representations in their embeddings
- Embeddings can also be fine tuned through the triplet loss

Embeddings have multiple applications in semantic search and dataset curation

- They can be used for reverse image search, recommendation systems, and deduplication

SPONSORED BY THE