🌲 Pinecone                                                    **Sign Up Free**            ☰

Learn | Article

# Text-to-Image and Image-to-Image Search Using CLIP

**Jump to section** ▾

# Introduction

Industries today deal with ever increasing amounts of data. Especially in retail, fashion, and other industries where the image representatio...

In such a situation, we can often describ... perform accurate and least time-consu...

*Could I take advantage of state-of-the-art artificial intelligence solutions to tackle such a challenge?*

This is where OpenAI's CLIP comes in handy. A deep learning algorithm that makes it easy to connect text and images.

After completing this conceptual blog, you will understand: (1) what CLIP is, (2) how it works and why you should adopt it, and finally, (3) how to implement it for your own use case using both local and cloud-based vector indexes.

# What is CLIP?

Contrastive Language-Image Pre-training (CLIP for short) is a state-of-the-art model introduced by OpenAI in February 2021 [1].

CLIP is a neural network trained on about 400 million (text and image) pairs. Training uses a contrastive learning approach that aims to unify text and images, allowing tasks like image classification to be done with text-image similarity.

This means that CLIP can find whether a given image and textual description match without being trained for a specific domain. Making CLIP powerful for out-of-the-box text and image search, which is the main focus of this article.

Besides text and image search, we can apply CLIP to image classification, image generation, image similarity search, image ranking, object tracking, robotics control, image captioning, and more.

# Why should you adopt the CLIP models?

Below are some reasons that increased the adoption of the CLIP models by the AI community.
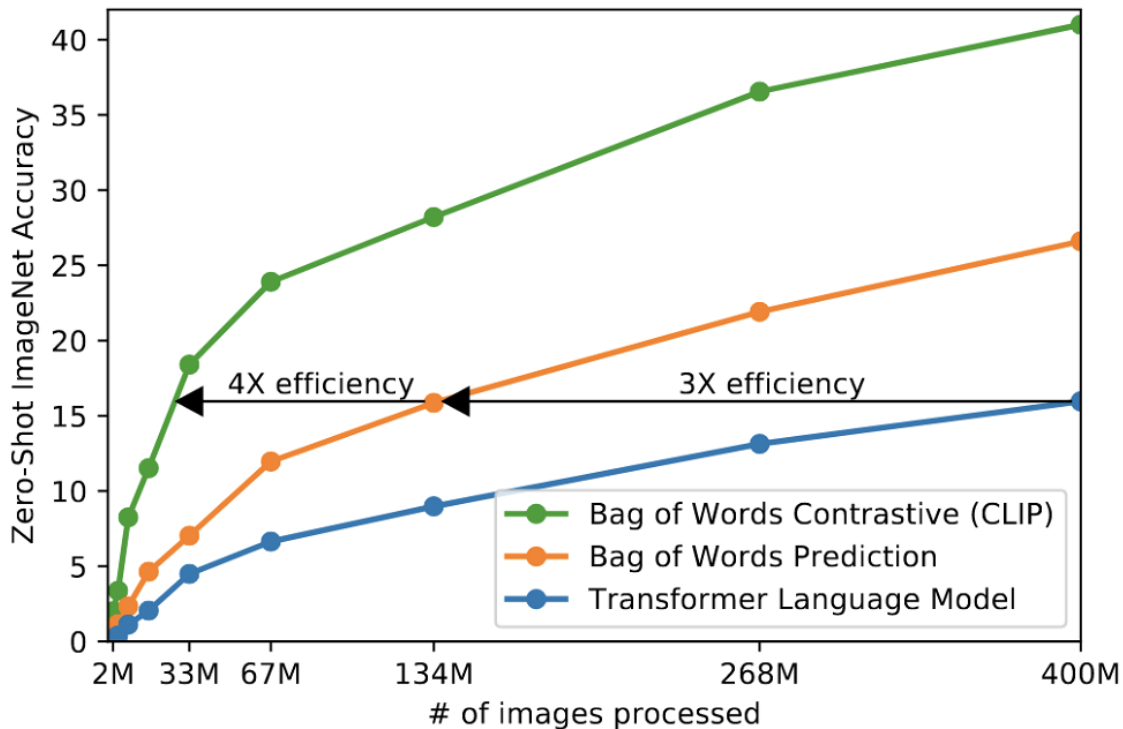
## Efficiency

Don't miss the next one...

Get an email the next time we publish an article about machine learning and similarity search.

The use of the contrastive objective increased the efficiency of the CLIP model by 4-to-10x more at zero-shot ImageNet classification.

Also, the adoption of the Vision Transformer created an additional 3x gain in compute efficiency compared to the standard ResNet.



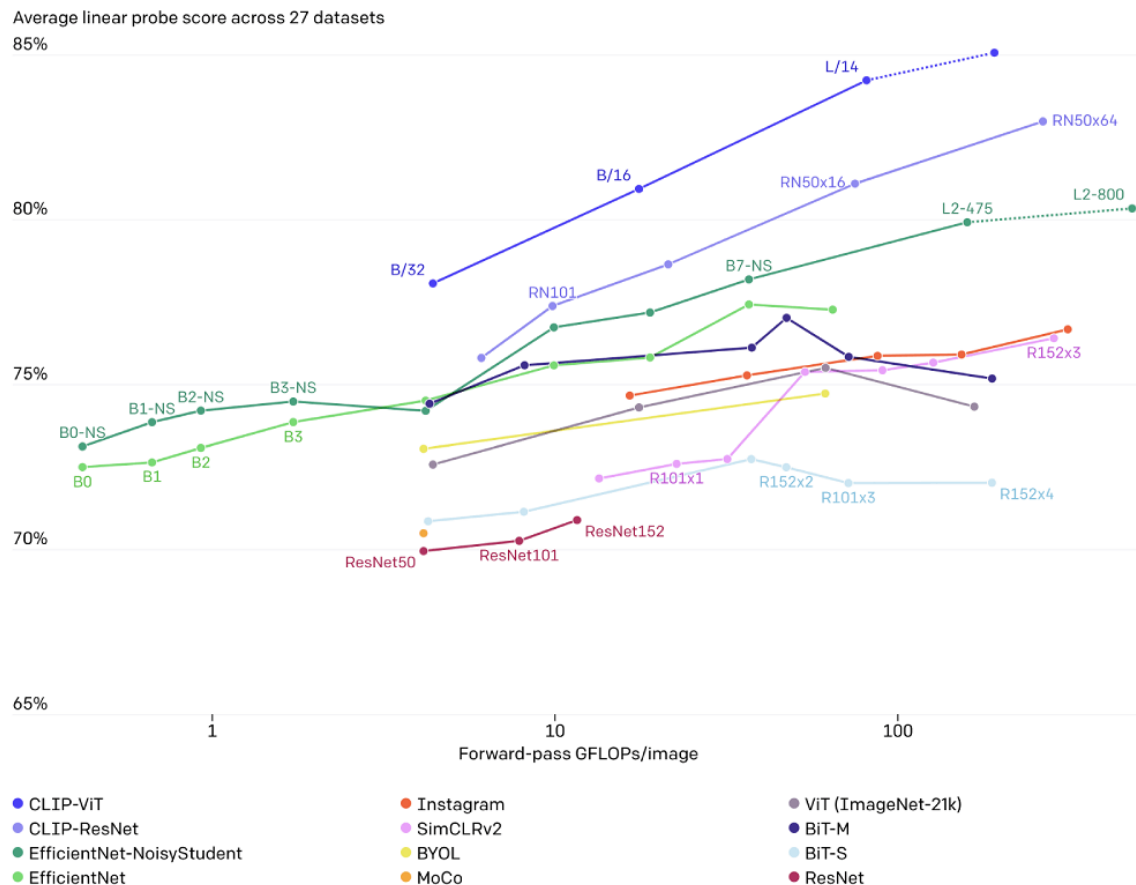The efficiency of CLIP at zero-shot transfer
Source: arxiv

# More general & flexible

CLIP outperforms existing ImageNet models in new domains because of its ability to learn a wide range of visual representations directly from natural language.

The following graphic highlights CLIP zero-shot performance compared to ResNet models few-shot linear probe performance on fine-grained object detection, geo-localization, action recognition, and optical character recognition tasks.

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

Average linear probe score across 27 datasets
Source: OpenAI

# CLIP Architecture

CLIP architecture consists of two main components: (1) a text encoder, and (2) an Image encoder. These two encoders are jointly trained to predict the correct pairings of a batch of training (image, text) examples.

- The *text encoder's* backbone is a transformer model [2], and the base size uses 63 millions-parameters, 12 layers, and a 512-wide model containing 8 attention heads.

- The *image encoder*, on the other hand, uses both a Vision Transformer (ViT) and a ResNet50 as its backbone, responsible for generating the feature representation of the image.
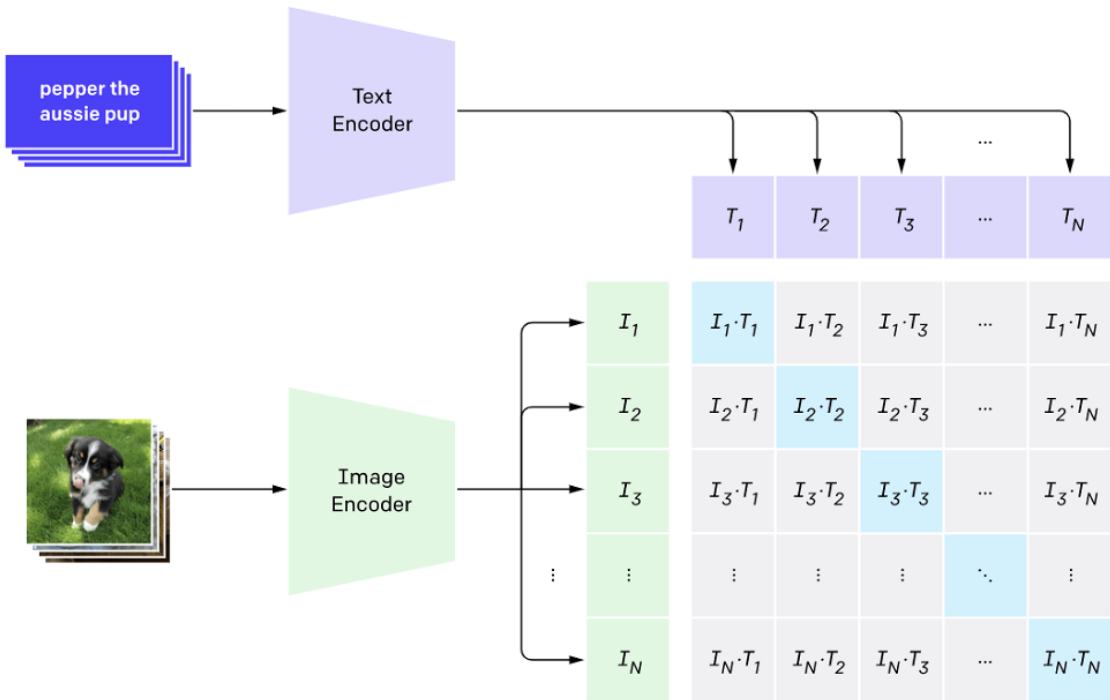
# How does the CLIP a

We can answer this question by understanding these three approaches: (1) contrastive pre-training, (2) dataset classifier creation from labeled text, and finally, (3) application of the zero-shot technique for classification.

Let's explain each of these three concepts.

## 1. Contrastive pre-training



Contrastive pre-training
Source: OpenAI

# 1. Contrastive pre-training

During this phase, a batch of 32,768 pairs of image and text is passed through the text and image encoders simultaneously to generate the vector representations of the text and the associated image, respectively.

The training is done by searching for each image, the closest text representation across the entire batch, which corresponds to maximizing cosine similarity between the actual N pairs that are maximally close.
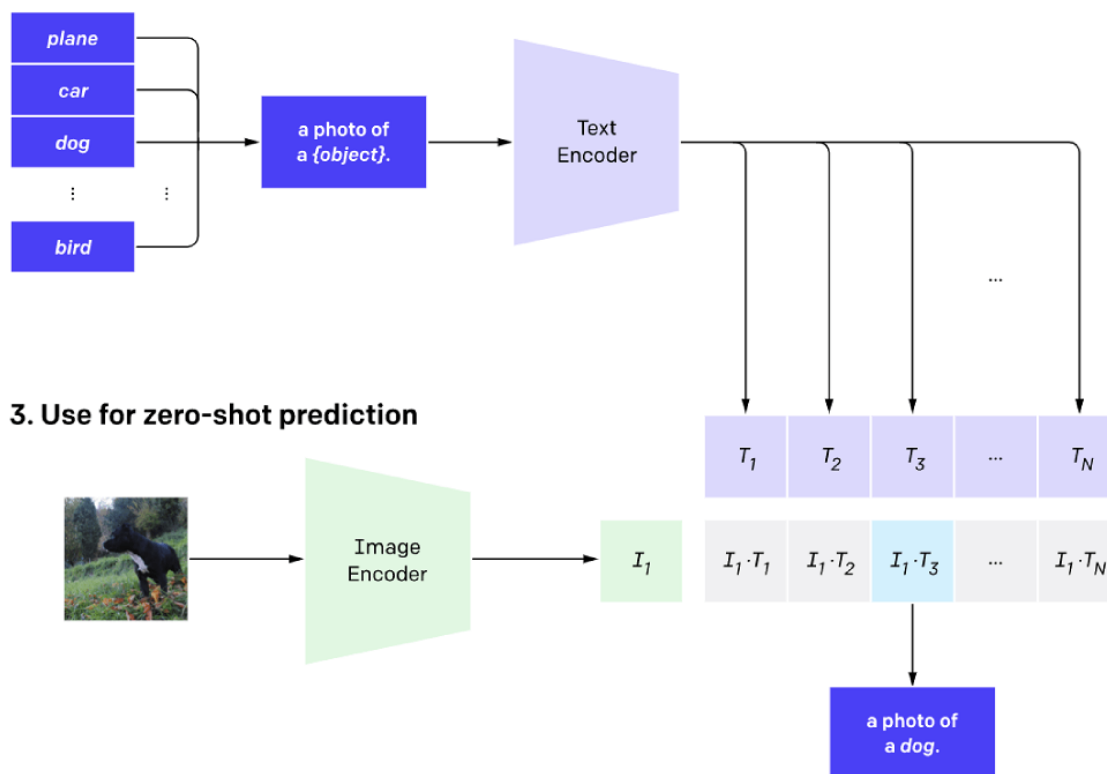
Also, it makes the actual images far awa    similarity.

Finally, a symmetric cross-entropy loss is optimized over the previously computed similarity scores.



Classification dataset creation and zero-shot prediction
Source: OpenAI

## 2. Create dataset classifier from label text

This second step section encodes all the labels/objects in the following context format: "**a photo of a {object}**. The vector representation of each context is generated from the text encoder.

If we have *dog, car*, and *plane* as the classes of the dataset, we will output the following context representations:

- a photo of a dog

- a photo of a car

- a photo of a plane

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.
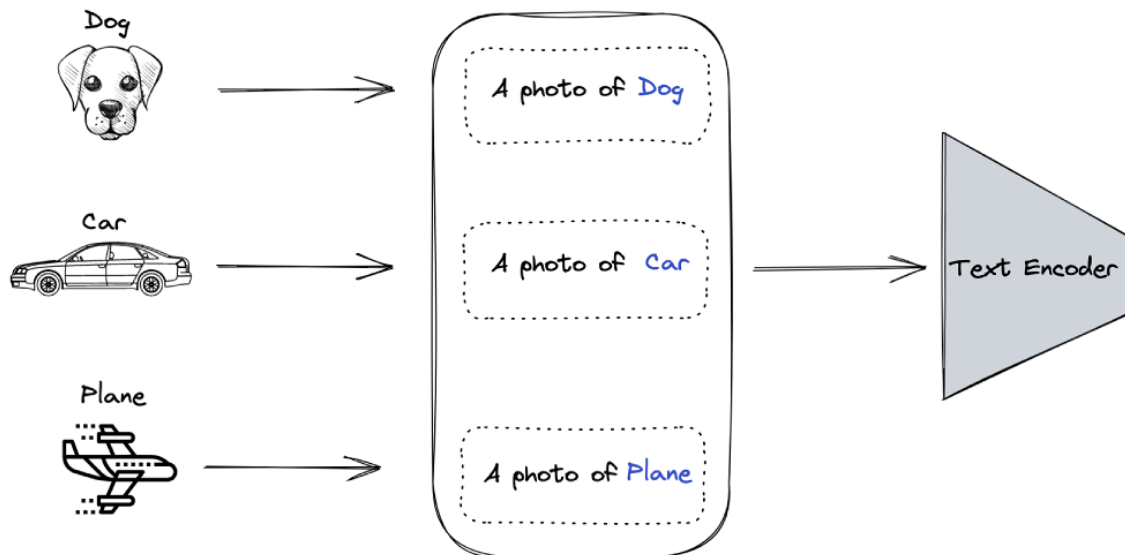
Image illustration of the context representations

## 3. Use of zero-shot prediction

We use the output of section 2 to predict which image vector corresponds to which context vector. The benefit of applying the zero-shot prediction approach is to make CLIP models generalize better on unseen data.

# Implementation of CLIP With Python

Now that we know the architecture of CLIP and how it works, this section will walk you through all the steps to successfully implement two real-world scenarios. First, you will understand how to perform an image search in natural language. Also, you will be able to perform an image-to-image search using.

At the end of the process, you will understand the benefits of using a vector database for such a use case.

**Don't miss the next one...**
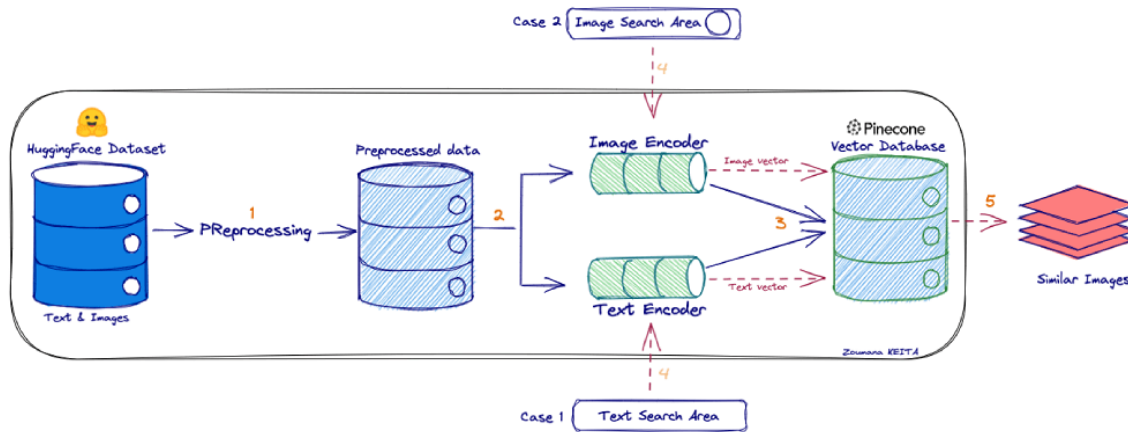
Get an email the next time we publish an article about machine learning and similarity search.

# General workflow of

*(Follow along with the Colab notebook!)*

The end-to-end process is explained through the workflow below. We start by collecting data from the Hugging Face dataset, which is then processed to further generate vector index vectors through the Image and Text Encoders. Finally, the Pinecone client is used to insert them to a vector index.

The user will then be able to search images based on either text or another image.



General workflow for image search

# Prerequisites

The following libraries are required to create the implementation.

### Install the libraries

```
1  %%bash
2  # Uncomment this if using it for the first time. -qqq for ZERO-OUT
3  pip3 -qqq install transformers torch datasets
4
5  # The following two libraries avoid the UnidentifiedImageError
6  pip3 -qqq install gdcm
7  pip3 -qqq install pydicom
8  pip -qqq install faiss-gpu
9  pip -qqq install pinecone
```

### Import the libraries

```
1    import os
2    import faiss
3    import torch
4    import skimage
5    import requests
6    import pinecone
7    import numpy as np
8    import pandas as pd
9    from PIL import Image
10   from io import BytesIO
11   import IPython.display
12   import matplotlib.pyplot as plt
13   from datasets import load_dataset
14   from collections import OrderedDict
15   from transformers import CLIPProcessor, CLIPModel, CLIPTokenizer
```

# Data acquisition and exploration

The conceptual captions dataset consists of around 3.3M images with two main columns: the image URL and its caption. You can find more details from the corresponding huggingface link.

```
1    # Get the dataset
2    image_data = load_dataset("conceptual_captions", split="train")
```

## Data preprocessing

Not all URLs in the dataset are valid. We fix that by testing and removing all erroneous URL entries.

```
1    def check_valid_URLs(imag
2      try:
3        response = requests
4        Image.open(BytesIO(
5        return True
6      except:
7        return False
```

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

```
 8    def get_image(image_URL):
 9        response = requests.get(image_URL)
10        image = Image.open(BytesIO(response.content)).convert("RGB")
11        return image
```

The following expression creates a new dataframe with a new column "is_valid" which is True
when the URL is valid or False otherwise.

```
1    # Transform dataframe
2    image_data_df["is_valid"] = image_data_df["image_url"].apply(check_va
3    # Get valid URLs
4    image_data_df = image_data_df[image_data_df["is_valid"]==True]
5    # Get image from URL
6    image_data_df["image"] = image_data_df["image_url"].apply(get_image)
```

The second step is to download the images from the URLs. This helps us avoid constant web
requests.

## Image and text embeddings implementation

The prerequisites to successfully implement the encoders are the model, the processor, and the
tokenizer.

The following function fulfills those requirements from the model ID and the device used for the
computation, either CPU or GPU.

```
 1    def get_model_info(model_ID, device):
 2    # Save the model to device
 3        model = CLIPModel.from pretrained(model ID).to(device)
 4        # Get the processor
 5        processor = CLIPProc
 6    # Get the tokenizer
 7        tokenizer = CLIPToken
 8            # Return model, p.
 9        return model, proces
10    # Set the device
```

**Don't miss the next one...**

Get an email the next time we publish an article about
machine learning and similarity search.

```
11   device = "cuda" if torch.cuda.is_available() else "cpu"
12   # Define the model ID
13   model_ID = "openai/clip-vit-base-patch32"
14   # Get model, processor & tokenizer
15   model, processor, tokenizer = get_model_info(model_ID, device)
```

## Text embeddings

We start by generating the embedding of a single text before applying the same function across the entire dataset.

```
1    def get_single_text_embedding(text):
2    inputs = tokenizer(text, return_tensors = "pt")
3        text_embeddings = model.get_text_features(**inputs)
4        # convert the embeddings to numpy array
5        embedding_as_np = text_embeddings.cpu().detach().numpy()
6    return embedding_as_np
7    def get_all_text_embeddings(df, text_col):
8    df["text_embeddings"] = df[str(text_col)].apply(get_single_text_embed
9    return df
10   # Apply the functions to the dataset
11   image_data_df = get_all_text_embeddings(image_data_df, "caption")
```

The first five rows look like this:

| | image_url | caption | is_valid | image | text_embeddings |
|---|---|---|---|---|---|
| 0 | http://lh6.ggpht.com/-lvRtNLNcG8o/TpFyrudaT6I/... | a very typical bus station | True | <PIL.Image.Image image mode=RGB size=800x534 a... | [[0.25922304, -0.08825898, 0.020317025, -0.127... |
| 1 | http://78.media.tumblr.com/3b133294bdc7c7784b7... | sierra looked stunning in this top and this sk... | True | <PIL.Image.Image image mode=RGB size=500x441 a... | [[0.0041467994, 0.18943565, -0.123970225, 0.30... |
| 2 | https://media.gettyimages.com/photos/young-con... | young confused girl standing in front of a war... | True | <PIL.Image.Image image mode=RGB size=490x612 a... | [[-0.28737983, -0.34814143, -0.04288538, 0.401... |
| 3 | https://thumb1.shutterstock.com/display_pic_wi... | interior design of modern living room with fir... | True | <PIL.Image.Image image mode=RGB size=450x470 a... | [[0.56064534, -0.15138063, -0.43740302, -0.339... |
| 4 | https://thumb1.shutterstock.com/display_pic_wi... | cybernetic scene isolated on white background . | True | <PIL.Image.Image image mode=RGB size=450x470 a... | [[0.035292536, 0.24262792, -0.12724756, -0.210... |

Format of the vector ind

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

## Image embeddings

The same process is used for image embeddings but with different functions.

```python
def get_single_image_embedding(my_image):
image = processor(
        text = None,
        images = my_image,
        return_tensors="pt"
        )["pixel_values"].to(device)
embedding = model.get_image_features(image)
# convert the embeddings to numpy array
    embedding_as_np = embedding.cpu().detach().numpy()
    return embedding_as_np
 def get_all_images_embedding(df, img_column):
    df["img_embeddings"] = df[str(img_column)].apply(get_single_image
    return df
image_data_df = get_all_images_embedding(image_data_df, "image")
```

The final format of the text and image vector index looks like this:

| | image_url | caption | is_valid | image | text_embeddings | img_embeddings |
|---|---|---|---|---|---|---|
| 0 | http://lh6.ggpht.com/-IvRtNLNcG8o/TpFyrudaT6I/... | a very typical bus station | True | <PIL.Image.Image image mode=RGB size=800x534 a... | [[0.25922304, -0.08825898, 0.020317025, -0.127... | [[-0.0034022853, -0.053583913, 0.35247508, 0.3... |
| 1 | http://78.media.tumblr.com/3b133294bdc7c7784b7... | sierra looked stunning in this top and this sk... | True | <PIL.Image.Image image mode=RGB size=500x441 a... | [[0.0041467994, 0.18943565, -0.123970225, 0.30... | [[-0.25019708, -0.1325763, 0.09706805, 0.97886... |
| 2 | https://media.gettyimages.com/photos/young-con... | young confused girl standing in front of a war... | True | <PIL.Image.Image image mode=RGB size=490x612 a... | [[-0.28737983, -0.34814143, -0.04288538, 0.401... | [[-0.36655784, 0.3118331, -0.13266361, 0.34909... |
| 3 | https://thumb1.shutterstock.com/display_pic_wi... | interior design of modern living room with fir... | True | <PIL.Image.Image image mode=RGB size=450x470 a... | [[0.56064534, -0.15138063, -0.43740302, -0.339... | [[-0.17221001, -0.29784596, -0.10141284, -0.06... |
| 4 | https://thumb1.shutterstock.com/display_pic_wi... | cybernetic scene isolated on white background . | True | <PIL.Image.Image image mode=RGB size=450x470 a... | [[0.035292536, 0.24262792, -0.12724756, -0.210... | [[0.18897031, -0.0012195408, -0.6513251, -0.12... |

Vector index with image and captions embeddings (Image by Author)

# Vector storage approach—Local vector index Vs. A cloud-based vector index

In this section, we will explore two differe̶n̶t̶... for performing the searches: The first is ... Pinecone. Both approaches use the cos...

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

# Using local dataframe as vector index

The helper function ***get_top_N_images*** generates similar images for the two scenarios illustrated in the workflow above: text-to-image search or image-to-image search.

```python
from sklearn.metrics.pairwise import cosine_similarity
def get_top_N_images(query, data, top_K=4, search_criterion="text"):
    # Text to image Search
    if(search_criterion.lower() == "text"):
      query_vect = get_single_text_embedding(query)
    # Image to image Search
    else:
      query_vect = get_single_image_embedding(query)
    # Relevant columns
    revevant_cols = ["caption", "image", "cos_sim"]
    # Run similarity Search
    data["cos_sim"] = data["img_embeddings"].apply(lambda x: cosine_si
    data["cos_sim"] = data["cos_sim"].apply(lambda x: x[0][0])
    """
    Retrieve top_K (4 is default value) articles similar to the query
    """
    most_similar_articles = data.sort_values(by='cos_sim', ascending=
    return most_similar_articles[revevant_cols].reset_index()
```

Let's understand how we perform the recommendation.

→ The user provides either a text or an image as a search criterion, but the model performs a text-to-image search by default.

→ In line 17, a cosine similarity is performed between each image vector and the user's input vector.

→ Finally, in line 24, sort the result based on the similarity score in descending order, and we return the most similar images by exclud
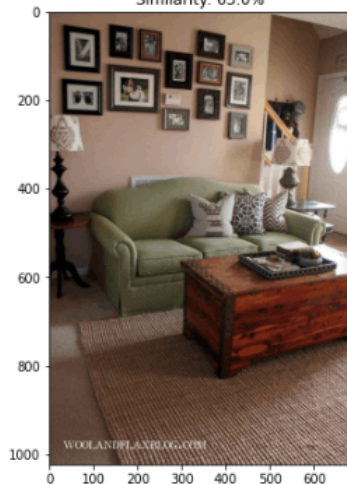
**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

# Example of searches

This helper function makes it easy to ha
images. Each image will have the corres

```python
1   def plot_images_by_side(top_images):
2     index_values = list(top_images.index.values)
3     list_images = [top_images.iloc[idx].image for idx in index_values]
4     list_captions = [top_images.iloc[idx].caption for idx in index_value
5     similarity_score = [top_images.iloc[idx].cos_sim for idx in index_va
6     n_row = n_col = 2
7     _, axs = plt.subplots(n_row, n_col, figsize=(12, 12))
8     axs = axs.flatten()
9     for img, ax, caption, sim_score in zip(list_images, axs, list_captio
10        ax.imshow(img)
11        sim_score = 100*float("{:.2f}".format(sim_score))
12        ax.title.set_text(f"Caption: {caption}\nSimilarity: {sim_score}
13    plt.show()
```

## Text-to-image

→ First, the user provides the text that is used for the search.

→ Second, we run a similarity search.

→ Third, we plot the images recommended by the algorithm.

```python
1   query_caption = image_data_df.iloc[10].caption
2   # Print the original query text
3   print("Query: {}".format(query_caption))
4   # Run the similarity search
5   top_images = get_top_N_images(query_caption, image_data_df)
6   # Plot the recommended images
7   plot_images_by_side(top_images)
```

Line 3 generates the following text:

*Query: actor arrives for the premiere of*

Line 9 produces the plot below.

Don't miss the next one...

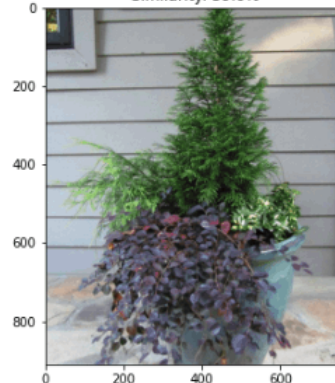Get an email the next time we publish an article about machine learning and similarity search.

Caption: actor and daughters uk premiere held
Similarity: 27.0%

Caption: pop artist attends the 3rd annual at guest house
Similarity: 27.0%

Caption: actor arrives to the premiere
Similarity: 25.0%

Caption: actor attends the premiere of season
Similarity: 25.0%

Images corresponding to the text: "actor arrives for the premiere of the film"

## Image-to-image

The same process applies. The only difference this time is that the user provides an image instead of a caption.

```
1    # Get the query image an
2    query_image = image_data_
3    query_image
```
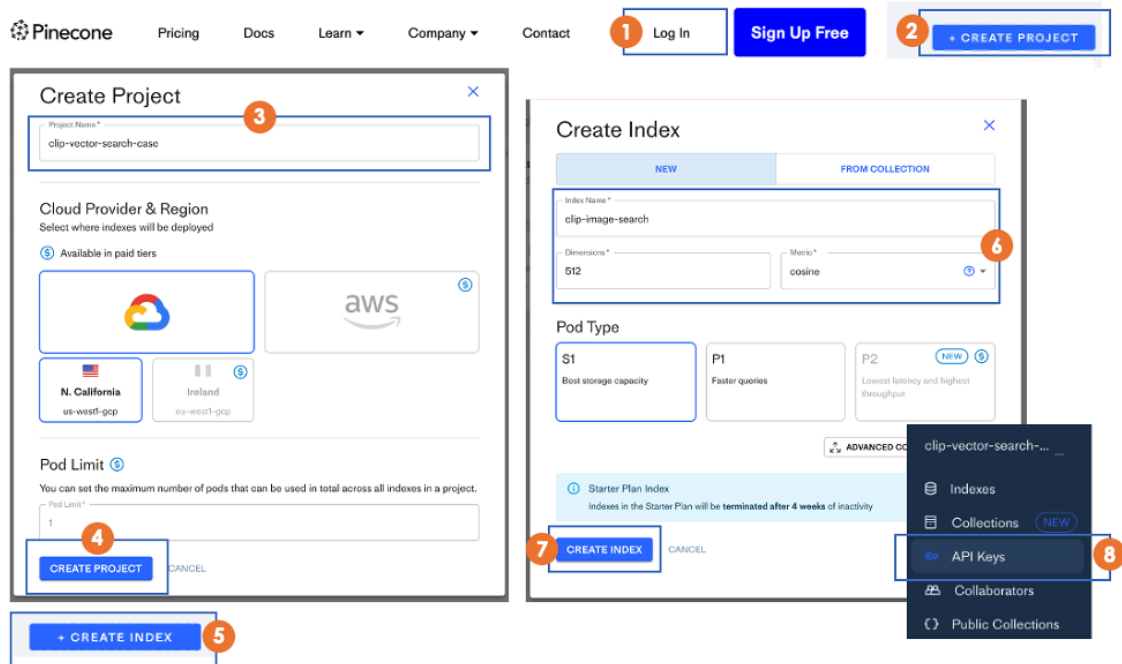
**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

Original image of search (image at the index)

```
1   # Run the similarity search and plot the result
2   top_images = get_top_N_i
3   # Plot the result
4   plot_images_by_side(top_:
```

**Don't miss the next one...**

Get an email the next time we publish an article about
machine learning and similarity search.

We run the search by specifying the search_criterion which is "image" in line 2.

The final result is shown below.



Images corresponding to the image-to-image search (Image by Author)

We can observe that some of the images are less similar which introduces noise in the recommendation. We can reduce that noise by specifying a threshold level of similarity. For instance, consider all the images with at least 60% similarity.

# Leveraging the power of a managed vector index using Pinecone

Pinecone provides a fully-managed, eas[y]
high-performance vector search applica[]

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

This section will walk you through the steps from acquiring your API credentials to implementing the search engine.

## Acquire your Pinecone API

Below are the eight steps to acquire your API credentials, starting from the Pinecone website.



Eight main steps to acquire your Pinecone Client API

## Configure the vector index

From the API, we can create the index that allows us to perform all the create, update, delete, and insert actions.

```
1   pinecone.init(
2       api_key = "YOUR_API_KEY",
3       environment="YOUR_ENV"
4   )
5   my_index_name = "clip-ima
6   vector_dim = image_data_
7
8   if my_index_name not in
9     # Create the vectors di
10    pinecone.create_index(na
```

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

```
11                              dimension=vector_dim,
12                              metric="cosine", shards=1,
13                              pod_type='s1.x1')
14   # Connect to the index
15   my_index = pinecone.Index(index_name = my_index_name)
```

- ***pinecone.init*** section initializes the pinecone workspace to allow future interactions.

- from lines 8 to 9 we specify the name we want for the vector index, and also the dimension of the vectors, which is 512 in our scenario.

- from lines 11 to 16 we create the index if it does not already exist.

The result of the following instruction shows that we have no data in the index.

```
1   my_index.describe_index_stats()
```

The only information we have is the dimension, which is 512.

```
1   {'dimension': 512,
2           'index_fullness': 0.0,
3             'namespaces': {},
4   'total_vector_count': 0}
```

## Populate the database

Now that we have configured the Pinecone database, the next step is to populate it with the following code.

```
image_data_df["vector_id
image_data_df["vector_id
# Get all the metadata
final_metadata = []
```

```
 5    for index in range(len(image_data_df)):
 6      final_metadata.append({
 7          'ID':  index,
 8          'caption': image_data_df.iloc[index].caption,
 9          'image': image_data_df.iloc[index].image_url
10      })
11    image_IDs = image_data_df.vector_id.tolist()
12    image_embeddings = [arr.tolist() for arr in image_data_df.img_embedd:
13    # Create the single list of dictionary format to insert
14    data_to_upsert = list(zip(image_IDs, image_embeddings, final_metadata
15    # Upload the final data
16    my_index.upsert(vectors = data_to_upsert)
17    # Check index size for each namespace
18    my_index.describe_index_stats()
```

*Let's understand what is going on here.*

~~The data to upsert requires three components: the unique identifiers (IDs) of each observation,~~ the list of embeddings being stored, and the metadata containing additional information about the data to store.

→ From lines 5 to 12, the metadata is created by storing the "ID", "caption" and "URL" of each observation.

→ On lines 14 and 15, we generate a list of IDs, and convert the embeddings into a list of lists.

→ Then, we create a list of dictionaries mapping the IDs, embeddings, and metadata.

→ The final data is upserted to the index with the *.upsert()* function.

Similarly to the previous scenario, we can check that all vectors have been upserted via `my_index.describe_index_stats()` .

## Start the query

All that remains is to query our index usi... Both will use the following syntax:

```
1   my_index.query(my_query_embedding, top_k=N, include_metadata=True)
```

→ *my_query_embedding* is the embedding (as a list) of the query (caption or image) provided by the user.

→ *N* corresponds to the top number of results to return.

→ *include_metadata=True* means that we want the query result to include metadata.

## Text to image

```
1   # Get the query text
2   text_query = image_data_df.iloc[10].caption
3
4   # Get the caption embedding
5   query_embedding = get_single_text_embedding(text_query).tolist()
6
7   # Run the query
8   my_index.query(query_embedding, top_k=4, include_metadata=True)
```

Below is the JSON response returned from the query



text-to-imag

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

From the "matches" attribute, we can observe the top four most similar images returned by the query.

## Image-to-image

The same approach applies to image-to-image search.

```
1   image_query = image_data_df.iloc[43].image
```

This is the image provided by the user as the search criteria.



Query image

```
# Get the text embedding
query_embedding = get_si...
```

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

```python
# Run the query
my_index.query(query_embedding, top_k=4, include_metadata=True)
```



image-to-image query result (Image by Author)

Once you've finished don't forget to delete your index to free up your resources with:

```python
1   pinecone.delete_index(my_index)
```

# What are the advantages to using a Pinecone over a local pandas dataframe?

This approach using Pinecone has seve

→ **Simplicity**: the querying approach is
the full responsibility of managing the ve

→ **Speed**: Pinecone approach is faster, v

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.

→ **Scalability**: vector index hosted on Pinecone is scalable with little-to-no user effort from us. The first approach would become increasingly complex and slow as we scale.

→ **Lower chance of information loss**: the vector index based on Pinecone is hosted in the cloud with backups and high information security. The first approach is too high risk for production use-cases.

→ **Web-service friendly**: the result provided by the query is in JSON format and can be consumed by other applications, making it a better fit for web-based applications.

# Conclusion

Congratulations, you have just learned how to fully implement an image search application using both image and natural language. I hope the benefits highlighted are valid enough to take your project to the next level using vector databases.

Multiple resources are available at our Learning Center to further your learning.

The source code for the article is available here.

# References

Code Notebook

[1] A. Radford, J. W. Kim, et al., Learning Transferable Visual Models From Natural Language Supervision (2021)

[2] A. Vaswani, et al., Attention Is All You Need (2017), NeurIPS

Share via: 𝕏 in Y

**Zoumana Keita**

Technical Writer

PRODUCT

Overview

Documentation

Trust and Security

SOLUTIONS

Search

Generative AI

Customers

RESOURCES

Learning Center

Community

Pinecone Blog

Support Center

System Status

COMPANY

About

Partners

Careers

Newsroom

Contact

LEGAL

Terms

Privacy

Cookies

**Don't miss the next one...**

Get an email the next time we publish an article about machine learning and similarity search.