
Build a Reverse Video Search Engine in minutes

@Milvus.io

6/13/2022

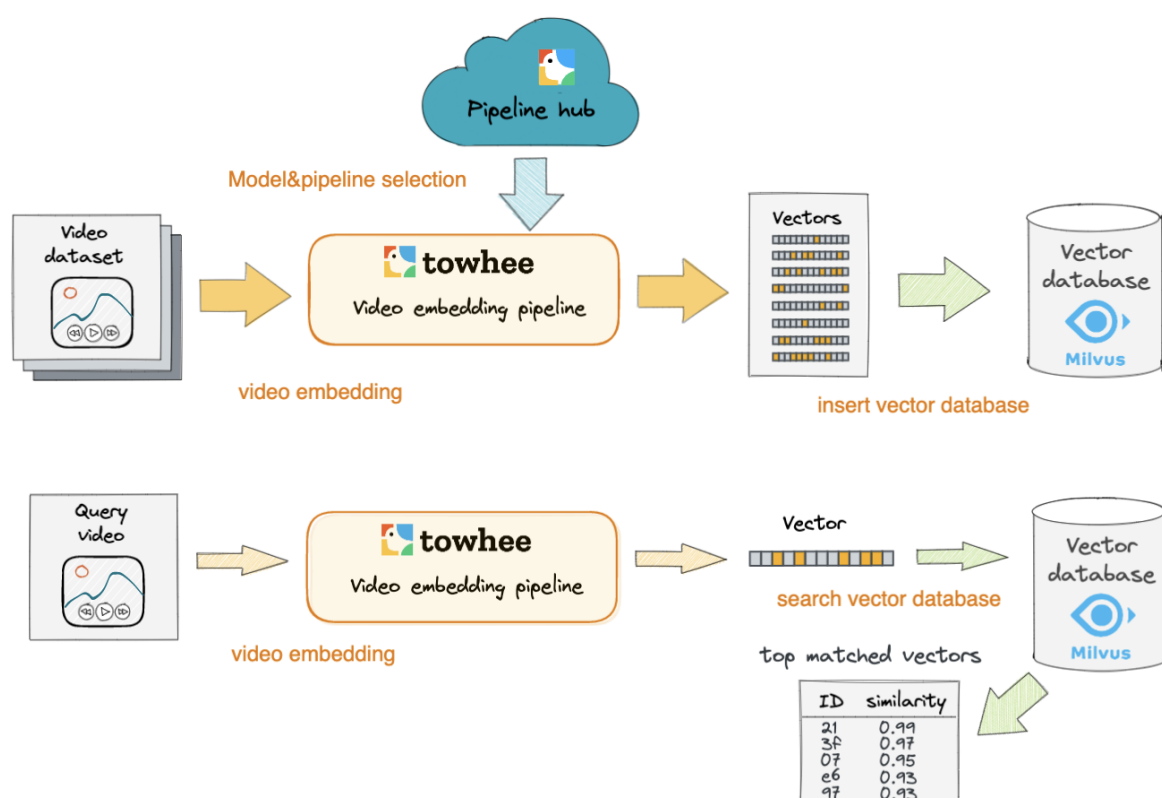
Build a Reverse Video Search Engine in minutes

Introduction

This codelab will show how to build a reverse-video-search engine from scratch using Milvus and Towhee. We will go through the procedure of building a reverse-video-search engine and evaluate its performance.

What is Reverse Video Search?

Reverse video search is similar like reverse image search. In simple words, it takes a video as input to search for similar videos. As we know that video-related tasks are harder to tackle, video models normally do not achieve as high scores as other types of models. However, there are increasing demands in AI applications in video. Reverse video search can effectively discover related videos and improve other applications.



What are Milvus & Towhee?

- Milvus is the most advanced open-source vector database built for AI applications and supports nearest neighbor embedding search across tens of millions of entries.
- Towhee is a framework that provides ETL for unstructured data using SoTA machine learning models.

Preparation

Install packages

Make sure you have installed required python packages(pymilvus, towhee, towhee.models, pillow, python, gradio):

```
1 $ python -m pip install -q pymilvus towhee towhee.models pillow ipython gradio
```

Prepare data

This tutorial will use a small data extracted from Kinetics400. You can download the subset from Github.

The data is organized as follows:

- **train:** candidate videos, 20 classes, 10 videos per class (200 in total)
- **test:** query videos, same 20 classes as train data, 1 video per class (20 in total)
- **reverse_video_search.csv:** a csv file containing an **id***, **path***, and **label*** for each video in train data

First to download the data and unzip it:

```
1 $ curl -L https://github.com/towhee-io/examples/releases/download/data/reverse_video_search.zip -O
2 $ unzip -q -o reverse_video_search.zip
```

Let's take a quick look:

```
1 import pandas as pd
2
3 df = pd.read_csv('./reverse_video_search.csv')
4 df.head(3)
```

	id	path	label
0	0	./train/country_line_dancing/bTbC3w_NlvM.mp4	country_line_dancing
1	1	./train/country_line_dancing/n2dWtEmNn5c.mp4	country_line_dancing
2	2	./train/country_line_dancing/zta-lv-xK7l.mp4	country_line_dancing

For later steps to easier get videos & measure results, we build some helpful functions in advance:

- **ground_truth:** get ground-truth video ids for the query video by its path

```

1 id_video = df.set_index('id')['path'].to_dict()
2 label_ids = {}
3 for label in set(df['label']):
4     label_ids[label] = list(df[df['label']==label].id)
5
6
7 def ground_truth(path):
8     label = path.split('/')[-2]
9     return label_ids[label]
```

Start Milvus

Before getting started with the engine, we also need to get ready with Milvus. Please make sure that you have started a Milvus service (Milvus Guide). Here we prepare a function to work with a Milvus collection with the following parameters:

- L2 distance metric
- IVF_FLAT index.

```

1 from pymilvus import connections, FieldSchema, CollectionSchema,
   DataType, Collection, utility
2
3 connections.connect(host='localhost', port='19530')
4
5 def create_milvus_collection(collection_name, dim):
6     if utility.has_collection(collection_name):
7         utility.drop_collection(collection_name)
8
9     fields = [
10     FieldSchema(name='id', dtype=DataType.INT64, description='ids',
11                 is_primary=True, auto_id=False),
12     FieldSchema(name='embedding', dtype=DataType.FLOAT_VECTOR,
13                 description='embedding vectors', dim=dim)
14     ]
```

```
13     schema = CollectionSchema(fields=fields, description='reverse video
14         search')
15     collection = Collection(name=collection_name, schema=schema)
16     # create IVF_FLAT index for collection.
17     index_params = {
18         'metric_type': 'L2',
19         'index_type': "IVF_FLAT",
20         'params': {"nlist": 400}
21     }
22     collection.create_index(field_name="embedding", index_params=
23         index_params)
24     return collection
```

Load Video Embeddings into Milvus

We first generate embeddings for videos with X3D model and then insert video embeddings into Milvus. Towhee provides a method-chaining style API so that users can assemble a data processing pipeline with operators.

```
1 import towhee
2
3 collection = create_milvus_collection('x3d_m', 2048)
4
5 dc = (
6     towhee.read_csv('reverse_video_search.csv')
7     .runas_op['id', 'id'](func=lambda x: int(x))
8     .video_decode.ffmpeg['path', 'frames'](sample_type='
9         uniform_temporal_subsample', args={'num_samples': 16})
10    .action_classification['frames', ('labels', 'scores', 'features')]
11    .pytorchvideo(
12        model_name='x3d_m', skip_preprocess=True)
13    .to_milvus['id', 'features'](collection=collection, batch=10)
14 )
```

Pipeline Explanation

Here are some details for each line of the assemble pipeline:

- `towhee.read_csv()`: read tabular data from csv file
- `.runas_op['id', 'id'](func=lambda x: int(x))`: for each row from the data, convert data type of the column id to int
- `.video_decode.ffmpeg['path', 'frames']()`: an embedded Towhee operator reading video as frames with specified sample method and number of samples. learn more
- `.action_classification['frames', ('labels', 'scores', 'features')].pytorchvideo()`: an embedded Towhee operator applying specified model to video frames,

which can be used to generate video embedding. [learn more](#)

- `.to_milvus['id', 'features']()`: insert video embedding into Milvus collection

Query Similar Videos from Milvus

Now all embeddings of candidate videos have been inserted into Milvus collection, we can query embeddings across the collection for nearest neighbors.

To get query embeddings, we should go through same pre-insert steps for each input video. Because Milvus returns video ids and vector distances, we use the `id_video` dictionary to get corresponding video paths based on ids.

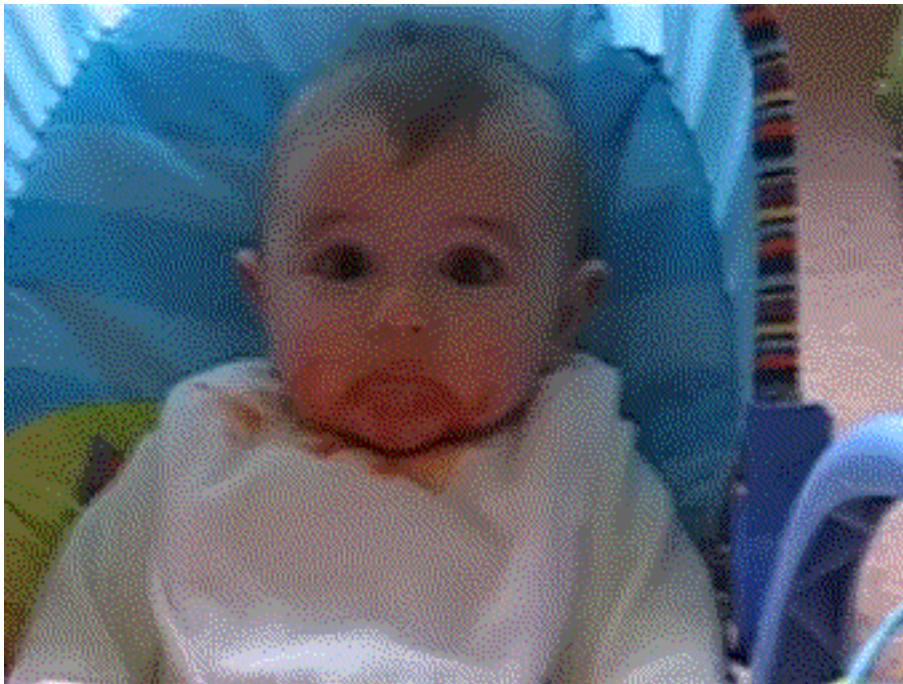
```
1 collection = Collection('x3d_m')
2
3 query_path = './test/eating_carrots/ty4UQlowp0c.mp4'
4
5 res_paths = (
6     towhee.glob['path'](query_path)
7     .video_decode.ffmpeg['path', 'frames'](sample_type='
8         uniform_temporal_subsample', args={'num_samples': 16})
9     .action_classification['frames', ('labels', 'scores', 'features
10         ')].pytorchvideo(
11         model_name='x3d_m', skip_preprocess=True)
12     .milvus_search['features', 'result'](collection=collection,
13         limit=10)
14     .runas_op['result', 'res_path'](func=lambda res: [id_video[x.id
15         ] for x in res])
16     .select['res_path']().to_list()[0].res_path
17 )
```

To display in the notebook, we convert videos to gifs. The code below first loads each video from its path and then gets full video frames with the embeded Towhee operator `.video_decode.ffmpeg()`. Finally converted gifs are saved under the directory `tmp_dir`. The section below is just help to show a search example.

```
1 import os
2 from IPython import display
3 from PIL import Image
4
5 tmp_dir = './tmp'
6 os.makedirs(tmp_dir, exist_ok=True)
7
8 def video_to_gif(video_path):
9     gif_path = os.path.join(tmp_dir, video_path.split('/')[-1][:4] + '
10         .gif')
11     frames = (
12         towhee.glob(video_path)
```

```
12         .video_decode.ffmpeg(sample_type='
13             uniform_temporal_subsample', args={'num_samples': 16})
14         .to_list()[0]
15     )
16     imgs = [Image.fromarray(frame) for frame in frames]
17     imgs[0].save(fp=gif_path, format='GIF', append_images=imgs[1:],
18                 save_all=True, loop=0)
19     return gif_path
20
21 html = 'Query video "{}": <br/>'.format(query_path.split('/')[-2])
22 query_gif = video_to_gif(query_path)
23 html_line = '<img src("{})"> <br/>'.format(query_gif)
24 html += html_line
25 html += 'Top 3 search results: <br/>'
26
27 for path in res_paths[:3]:
28     gif_path = video_to_gif(path)
29     html_line = '<img src("{}" style="display:inline;margin:1px"/>'.
30                 format(gif_path)
31     html += html_line
32 display.HTML(html)
```

Query video “eating_carrots”:



Top 3 search results:





Evaluation

We have just built a reverse video search engine. But how's its performance? We can evaluate the search engine against the ground truths.

In this section, we'll measure the performance with 2 metrics - mHR and mAP:

- **mHR (recall@K):**

- Mean Hit Ratio describes how many actual relevant results are returned out of all ground truths.
- Since Milvus return results with topK, we can also call this metric *recall@K*, where K is the count of searched results. When returned results are as many as ground truths, the hit ratio is equivalent to accuracy and we can take it as *accuracy@K* as well.
- For example, there are 100 archery videos in the collection. Then querying the engine with another archery video returns 70 archery videos out of 80 results. In this case, the number of ground truths is 100 and hit (correct) results are 70. So the hit ratio is 70/100.

- **mAP:**

- Average precision describes whether all of the relevant results are ranked higher than irrelevant results.

```
1 benchmark = (
```

```

2     towhee.glob['path']('./test/**/*.mp4')
3     .video_decode.ffmpeg['path', 'frames'](sample_type='
        uniform_temporal_subsample', args={'num_samples': 16})
4     .action_classification['frames', ('labels', 'scores', 'features
        ')].pytorchvideo(
5         model_name='x3d_m', skip_preprocess=True)
6     .milvus_search['features', 'result'](collection=collection,
        limit=10)
7     .runas_op['path', 'ground_truth'](func=ground_truth)
8     .runas_op['result', 'result'](func=lambda res: [x.id for x in
        res])
9     .with_metrics(['mean_hit_ratio', 'mean_average_precision'])
10    .evaluate['ground_truth', 'result']('x3d_m')
11    .report()
12 )

```

	mean_hit_ratio	mean_average_precision
x3d_m	0.57	0.751615

Optimization

We can see from above evaluation report, the current performance is not satisfactory. What can we do to improve the search engine? Of course we can fine-tune deep learning network with our own train data. Using more types of embeddings or filters by video tags/description/captions and audio can definitely enhance the search engine as well. But in this tutorial, I will just recommend some very simple options to make improvements.

Normalize embeddings

A quick optimization is normalizing all embeddings. Then the L2 distance will be equivalent to cosine similarity, which measures the similarity between two vectors using the angle between them, which ignores the magnitude of the vectors. We use the `.tensor_normalize['vec', 'vec']()` provided by Towhee to simply normalize all embeddings.

```

1 collection = create_milvus_collection('x3d_m', 2048)
2
3 dc = (
4     towhee.read_csv('reverse_video_search.csv')
5     .runas_op['id', 'id'](func=lambda x: int(x))
6     .video_decode.ffmpeg['path', 'frames'](sample_type='
        uniform_temporal_subsample', args={'num_samples': 16})
7     .action_classification['frames', ('labels', 'scores', 'features')
        ].pytorchvideo(

```

```

8         model_name='x3d_m', skip_preprocess=True)
9     .tensor_normalize['features', 'features']()
10    .to_milvus['id', 'features'](collection=collection, batch=10)
11 )
12
13 benchmark = (
14     towhee.glob['path']('./test/*/*.mp4')
15     .video_decode.ffmpeg['path', 'frames'](sample_type='
16         uniform_temporal_subsample', args={'num_samples': 16})
17     .action_classification['frames', ('labels', 'scores', 'features
18         ')].pytorchvideo(
19         model_name='x3d_m', skip_preprocess=True)
20     .tensor_normalize['features', 'features']()
21     .milvus_search['features', 'result'](collection=collection,
22         limit=10)
23     .runas_op['path', 'ground_truth'](func=ground_truth)
24     .runas_op['result', 'result'](func=lambda res: [x.id for x in
25         res])
26     .with_metrics(['mean_hit_ratio', 'mean_average_precision'])
27     .evaluate['ground_truth', 'result']('x3d_m_norm')
28     .report()
29 )

```

	mean_hit_ratio	mean_average_precision
x3d_m_norm	0.66	0.79646

With vector normalization, we have increased the mHR to 0.66 and mAP to about 0.8, which look better now.

Change model

There are more video models using different networks. Normally a more complicated or larger model will show better results while cost more. You can always try more models to tradeoff among accuracy, latency, and resource usage. Here I show the performance for the reverse video search engine using a SOTA model with multiscale vision transformer as backbone.

```

1 collection = create_milvus_collection('mvit_base', 768)
2
3 dc = (
4     towhee.read_csv('reverse_video_search.csv')
5     .runas_op['id', 'id'](func=lambda x: int(x))
6     .video_decode.ffmpeg['path', 'frames'](sample_type='
7         uniform_temporal_subsample', args={'num_samples': 32})
8     .action_classification['frames', ('labels', 'scores', 'features')
9         ].pytorchvideo(
10         model_name='mvit_base_32x3', skip_preprocess=True)

```

```

9     .tensor_normalize['features', 'features']()
10    .to_milvus['id', 'features'](collection=collection, batch=10)
11 )
12
13 benchmark = (
14     towhee.glob['path']('./test/*/*.mp4')
15     .video_decode ffmpeg['path', 'frames'](sample_type='
        uniform_temporal_subsample', args={'num_samples': 32})
16     .action_classification['frames', ('labels', 'scores', 'features
        ')].pytorchvideo(
17         model_name='mvit_base_32x3', skip_preprocess=True)
18     .tensor_normalize['features', 'features']()
19     .milvus_search['features', 'result'](collection=collection,
        limit=10)
20     .runas_op['path', 'ground_truth'](func=ground_truth)
21     .runas_op['result', 'result'](func=lambda res: [x.id for x in
        res])
22     .with_metrics(['mean_hit_ratio', 'mean_average_precision'])
23     .evaluate['ground_truth', 'result']('mvit_base')
24     .report()
25 )

```

	mean_hit_ratio	mean_average_precision
mvit_base	0.785	0.864911

Switching to MVIT model increases the mHR to 0.79 and mAP to 0.86, which are much better than X3D model. However, both insert and search time have increased. It's time for you to make trade-off between latency and accuracy. You're always encouraged to play around with this tutorial.

Release a Showcase

We've learnt how to build a reverse video search engine. Now it's time to add some interface and release a showcase. Towhee provides `towhee.api()` to wrap the data processing pipeline as a function with `as_function()`. So we can build a quick demo with this `milvus_search_function` with Gradio.

```

1 import gradio
2
3 with towhee.api() as api:
4     milvus_search_function = (
5         api.video_decode ffmpeg(
6             sample_type='uniform_temporal_subsample', args={'
                num_samples': 32})
7         .action_classification.pytorchvideo(model_name='
                mvit_base_32x3', skip_preprocess=True)
8         .runas_op(func=lambda x: x[-1])
9         .tensor_normalize()

```

```
10         .milvus_search(collection='mvid_base', limit=3)
11         .runas_op(func=lambda res: [id_video[x.id] for x in res])
12         .as_function()
13     )
14
15
16     interface = gradio.Interface(milvus_search_function,
17                                  inputs=gradio.Video(source='upload'),
18                                  outputs=[gradio.Video(format='mp4') for _
19                                           in range(5)]
19     )
20
21     interface.launch(inline=True, share=True)
```