
Deep Dive into Reverse Image Search Engine with Towhee

In the previous “Build a Reverse Image Search Engine in Minutes” codelab, we built and prototyped a proof-of-concept reverse image search engine. Now, let’s optimize our algorithm, feed it with large-scale image datasets, and deploy it as a micro-service with Towhee.

Preparation

duration: 1

Install Dependencies

First we need to install dependencies such as pymilvus, towhee, transformers, opencv-python and fastapi.

```
1 $ python -m pip -q install pymilvus towhee transformers fastapi opencv-python
```

Prepare the data

There is a subset of the ImageNet dataset (100 classes, 10 images for each class) is used in this demo, and the dataset is available via Github. The dataset is same as our previous codelab: “Build a Reverse Image Search Engine in Minutes”, and to make things easy, we’ll repeat the important code blocks below; if you have already downloaded data, please move on to next section.

The dataset is organized as follows:

- **train**: directory of candidate images;
- **test**: directory of the query images;
- **reverse_image_search.csv**: a csv file containing an **id***, **path***, and ***label*** for each image;

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

To use the dataset for image search, let’s first define some helper functions:

- **read_images(results)**: read images by image IDs;

- **ground_truth(path)**: ground-truth for each query image, which is used for calculating mHR(mean hit ratio) and mAP(mean average precision);

```

1 import cv2
2 import pandas as pd
3 from towhee._types.image import Image
4
5 df = pd.read_csv('reverse_image_search.csv')
6 df.head()
7
8 id_img = df.set_index('id')['path'].to_dict()
9 label_ids = {}
10 for label in set(df['label']):
11     label_ids[label] = list(df[df['label']==label].id)
12
13 def read_images(results):
14     imgs = []
15     for re in results:
16         path = id_img[re.id]
17         imgs.append(Image(cv2.imread(path), 'BGR'))
18     return imgs
19
20 def ground_truth(path):
21     label = path.split('/')[2]
22     return label_ids[label]

```

Create a Milvus Collection

Before getting started, please make sure you have installed milvus. Let's first create a `reverse_image_search` collection that uses the L2 distance metric and an IVF_FLAT index.

```

1 from pymilvus import connections, FieldSchema, CollectionSchema,
   DataType, Collection, utility
2
3 connections.connect(host='127.0.0.1', 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     ]
15     schema = CollectionSchema(fields=fields, description='reverse image
16                             search')

```

```
14     collection = Collection(name=collection_name, schema=schema)
15
16     index_params = {
17         'metric_type': 'L2',
18         'index_type': "IVF_FLAT",
19         'params': {"nlist": 2048}
20     }
21     collection.create_index(field_name="embedding", index_params=
        index_params)
22     return collection
```

Improve the Model

duration: 3

In “**Build a Reverse Image Search Engine in Minutes**” codelab, we evaluated the search engine with **mHR** and **mAP** and have already tried to improve the metrics by normalizing the embedding features and using a more complex model. The experiment shows a significant boost in the metrics.

Let's run some more experiments, with the end goal of improving the search results. First, we begin with a benchmark that compares **ResNet** and **VGG** models to the increasingly popular Vision Transformer (**ViT**), to see whether we can achieve better accuracy or get a reasonable trade-off between the accuracy and performance. We will then try to fix some bad cases by a preceding object detection model.

Model Benchmark: VGG vs ResNet vs EfficientNet

Three models will be included in the benchmark: **VGG16**, **resnet50**, and **efficientnet-b2**. We won't include too many models in this tutorial, but we encourage you to play around with different models. We can't include too much models in this notebook, for it might be too much time-consuming for the readers. But you can add your own interested model to the benchmark and try the notebook on your own machine, for example, **VGG19** (4096-d), **resnet101** (2048-d), or **efficient-b7** (2560-d). The following metrics will be included in the benchmark:

- **mHR (recall@K)**: This metric describes how many actual relevant results were returned out of all ground-truth relevant results by the search engine. For example, if we have put 100 pictures of cats into the search engine and then query the image search engine with another picture of cats. The total relevant result is 100, and the actual relevant results are the number of cat images in the top 100 results returned by the search engine. If there are 80 images about cats in the search result, the hit ratio is 80/100;

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

We'll use a helper class to compute runtime:

```
1 import time
2
3 class Timer:
4     def __init__(self, name):
5         self._name = name
6
7     def __enter__(self):
8         self._start = time.time()
9         return self
10
11    def __exit__(self, *args):
12        self._interval = time.time() - self._start
13        print('%s: %.2fs'%(self._name, self._interval))
```

Then to run some models to evaluate:

```
1 import towhee
2
3 model_dim = {
4     'vgg16': 4096,
5     'resnet50': 2048,
6     'tf_efficientnet_b2': 1408
7 }
8
9 for model in model_dim:
10     collection = create_milvus_collection(model, model_dim[model])
11
12     with Timer(f'{model} load'):
13         (
14             towhee.read_csv('reverse_image_search.csv')
15             .runas_op['id', 'id'](func=lambda x: int(x))
16             .image_decode['path', 'img']()
17             .image_embedding.timm['img', 'vec'](model_name=model)
18             .tensor_normalize['vec', 'vec']()
19             .to_milvus['id', 'vec'](collection=collection, batch
20                                 =100)
21         )
22     with Timer(f'{model} query'):
23         ( towhee.glob['path']('./test/*/*.JPEG')
24           .image_decode['path', 'img']()
25           .image_embedding.timm['img', 'vec'](model_name=model)
26           .tensor_normalize['vec', 'vec']()
27           .milvus_search['vec', 'result'](collection=collection,
28                                       limit=10)
29           .runas_op['path', 'ground_truth'](func=ground_truth)
30           .runas_op['result', 'result'](func=lambda res: [x.id
```

```

29         for x in res])
30         .with_metrics(['mean_hit_ratio', '
31                        mean_average_precision'])
32         .evaluate['ground_truth', 'result'](model)
33         .report()
34     )

```

vgg16 load: 34.15s

	mean_hit_ratio	mean_average_precision
vgg16	0.652	0.849296

vgg16 query: 6.10s

resnet50 load: 16.73s

	mean_hit_ratio	mean_average_precision
resnet50	0.781	0.917373

resnet50 query: 2.81s

tf_efficientnet_b2 load: 21.88s

	mean_hit_ratio	mean_average_precision
tf_efficientnet_b2	0.818	0.925592

tf_efficientnet_b2 query: 3.18s

Vision Transformer

Now, let's try using a transformer-based model. We'll first initialize the embedding model for use in a `vit_embedding` function which can be directly called to generate the embedding vector. This function can then be used via `runas_op` in `DataCollection`.

```

1  import torch
2  from transformers import ViTFeatureExtractor, ViTModel
3
4  feature_extractor = ViTFeatureExtractor.from_pretrained('google/vit-
5      large-patch32-384')
6  model = ViTModel.from_pretrained('google/vit-large-patch32-384')
7  device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
8  model.to(device)
9
10 def vit_embedding(img):
11     img = img.cv2_to_rgb()
12     inputs = feature_extractor(img, return_tensors="pt")
13     outputs = model(inputs['pixel_values'].to(device))
14     return outputs.pooler_output.detach().cpu().numpy().flatten()
15
16 collection = create_milvus_collection('huggingface_vit', 1024)

```

```

17
18 with Timer('ViT load'):
19     (
20         towhee.read_csv('reverse_image_search.csv')
21         .runas_op['id', 'id'](func=lambda x: int(x))
22         .image_decode['path', 'img']()
23         .runas_op['img', 'vec'](func=vit_embedding)
24         .tensor_normalize['vec', 'vec']()
25         .to_milvus['id', 'vec'](collection=collection, batch=100)
26     )
27
28 with Timer('ViT query'):
29     (
30         towhee.glob['path']('./test/*/*.JPEG')
31         .image_decode['path', 'img']()
32         .runas_op['img', 'vec'](func=vit_embedding)
33         .tensor_normalize['vec', 'vec']()
34         .milvus_search['vec', 'result'](collection=collection,
35                                         limit=10)
36         .runas_op['path', 'ground_truth'](func=ground_truth)
37         .runas_op['result', 'result'](func=lambda res: [x.id for x
38                                                         in res])
39         .with_metrics(['mean_hit_ratio', 'mean_average_precision'])
40         .evaluate['ground_truth', 'result']('huggingface_vit')p
41         .report()
42     )

```

ViT load: 63.23s

	mean_hit_ratio	mean_average_precision
huggingface_vit	0.905	0.96266

ViT query: 7.25s

We ran the experiment above on our own and found that **ViT-large** performed the best in this dataset. **EfficientNet-B2** also performs well (with less than half the runtime of **ViT-large**). We encourage you to try using your own datasets, model architectures, and training techniques for comparison.

Models	Mean Hit Ratio	Mean Average Precision	Interval (Query)
VGG16	0.652	0.849	6.15s
ResNet50	0.781	0.917	2.96s
EfficientNet-B2	0.818	0.925	3.35s
ViT-large	0.907	0.965	7.27s

Note: There are also predefined operators of ViT models on towhee hub. You can try the other ViT models with the image-embedding/timm.

Dimensionality Reduction

duration: 2

In a production system it is often practical to minimize the embedding dimension in order to minimize memory consumption. Random projection is a dimensionality reduction method for a set vectors in Euclidean space. Since this method is fast and requires no training, we'll try this technique on embeddings generated by the EfficientNet-B2 model:

```
1 import numpy as np
2
3 projection_matrix = np.random.normal(scale=1.0, size=(1408, 512))
4
5 def dim_reduce(vec):
6     return np.dot(vec, projection_matrix)
7
8 collection = create_milvus_collection('tf_efficientnet_b2_512', 512)
9
10 # load embeddings into milvus
11 (
12     towhee.read_csv('reverse_image_search.csv')
13     .runas_op['id', 'id'](func=lambda x: int(x))
14     .image_decode['path', 'img']()
15     .image_embedding.timm['img', 'vec'](model_name='
16         tf_efficientnet_b2')
17     .runas_op['vec', 'vec'](func=dim_reduce)
18     .tensor_normalize['vec', 'vec']()
19     .to_milvus['id', 'vec'](collection=collection, batch=100)
20 )
21 # query and evaluation
22 ( towhee.glob['path']('./test/*/*.JPEG')
23     .image_decode['path', 'img']()
```



```

24     .image_embedding.timm['img', 'vec'](model_name='
      tf_efficientnet_b2')
25     .runas_op['vec', 'vec'](func=dim_reduce)
26     .tensor_normalize['vec', 'vec']()
27     .milvus_search['vec', 'result'](collection=collection, limit
      =10)
28     .runas_op['path', 'ground_truth'](func=ground_truth)
29     .runas_op['result', 'result'](func=lambda res: [x.id for x in
      res])
30     .with_metrics(['mean_hit_ratio', 'mean_average_precision'])
31     .evaluate['ground_truth', 'result']('tf_efficientnet_b2_512')
32     .report()
33 )

```

	mean_hit_ratio	mean_average_precision
tf_efficientnet_b2_512	0.78	0.915922

The dimension of embedding vectors is reduced from 1408 to 512, thereby reducing memory usage by around 60%. Despite this, it maintains a reasonable performance (91.6% mAP for reduced vectors vs 92.5% for full vectors).

Object Detection with YOLO

duration: 2

Finally, we can try cropping candidate bounding boxes reverse image search, i.e. use YOLOv5 to get the object of the image before image feature vector extraction, and then use that object to represent the image data for insertion and search.

`get_object` function is used to get the image of the largest object detected by YOLOv5, or the image itself if there is no object, then insert the resulting image into Milvus, and finally do the search. Object detection is very common in product search.

```

1  yolo_collection = create_milvus_collection('yolo', 2048)
2  resnet_collection = create_milvus_collection('resnet', 2048)
3
4  def get_object(img, boxes):
5      if len(boxes) == 0:
6          return img
7      max_area = 0
8      for box in boxes:
9          x1, y1, x2, y2 = box
10         area = (x2-x1)*(y2-y1)
11         if area > max_area:
12             max_area = area

```

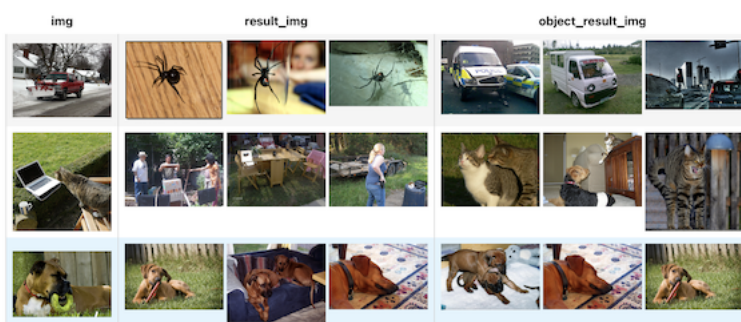
```

13         max_img = img[y1:y2,x1:x2,:]
14     return max_img
15
16     with Timer('resnet load'):
17         (towhee.read_csv('reverse_image_search.csv')
18          .runas_op['id', 'id'](func=lambda x: int(x))
19          .image_decode['path', 'img']()
20          .image_embedding.timm['img', 'vec'](model_name='resnet50')
21          .tensor_normalize['vec', 'vec']()
22          .to_milvus['id', 'vec'](collection=resnet_collection, batch
23                                =100)
24
25     with Timer('yolo+resnet load'):
26         (towhee.read_csv('reverse_image_search.csv')
27          .runas_op['id', 'id'](func=lambda x: int(x))
28          .image_decode['path', 'img']()
29          .object_detection.yolov5['img', ('boxes', 'class', 'score')]()
30          .runas_op[('img', 'boxes'), 'object'](func=get_object)
31          .image_embedding.timm['object', 'object_vec'](model_name='
32                resnet50')
33          .tensor_normalize['object_vec', 'object_vec']()
34          .to_milvus['id', 'object_vec'](collection=yolo_collection,
35                                         batch=100)

```

resnet load: 16.98s yolo+resnet load: 33.97s

Using a preceding object detection with yolo can fix the bad cases, but it also makes the search engine much slower.



Making Our Image Search Engine Production Ready

duration: 2

To put the image search engine into production, we need to feed it with a large-scale dataset and deploy a microservice to accept incoming queries.

Optimize for large-scale dataset

When the dataset becomes very large, as huge as tens of millions of images, it faces two significant problems:

1. embedding feature extractor and Milvus data loading needs to be fast so that we can finish the search index in time;
2. There are corrupted images or images with wrong formats in the dataset. It is impossible to clean up all such bad cases when the dataset is huge. So the data pipeline needs to be very robust to such exceptions.

Towhee supports parallel execution to improve performance for large-scale datasets, and also has `exception_safe` execution mode to ensure system stability.

Improve Performance with Parallel Execution

We are able to enable parallel execution by simply calling `set_parallel` within the pipeline. It tells towhee to process the data in parallel. Here is an example that enables parallel execution on a pipeline using ViT model. It can be seen that the execution speed below is nearly three times faster than before. And note that please clean up the GPU cache before runing with parallel.

```
1 collection = create_milvus_collection('test_resnet101', 2048)
2 with Timer('resnet101 load'):
3     (
4         towhee.read_csv('reverse_image_search.csv')
5         .runas_op['id', 'id'](func=lambda x: int(x))
6         .image_decode['path', 'img']()
7         .image_embedding.timm['img', 'vec'](model_name='resnet101')
8         .tensor_normalize['vec', 'vec']()
9         .to_milvus['id', 'vec'](collection=collection, batch=100)
10    )
11
12 collection_parallel = create_milvus_collection('test_resnet101_parallel', 2048)
13 with Timer('resnet101+parallel load'):
14     (
15         towhee.read_csv('reverse_image_search.csv')
16         .runas_op['id', 'id'](func=lambda x: int(x))
17         .set_parallel(3)
18         .image_decode['path', 'img']()
19         .image_embedding.timm['img', 'vec'](model_name='resnet101')
20         .tensor_normalize['vec', 'vec']()
21         .to_milvus['id', 'vec'](collection=collection_parallel,
22                                batch=100)
23    )
```

resnet101 load: 21.04s resnet101+parallel load: 14.42s

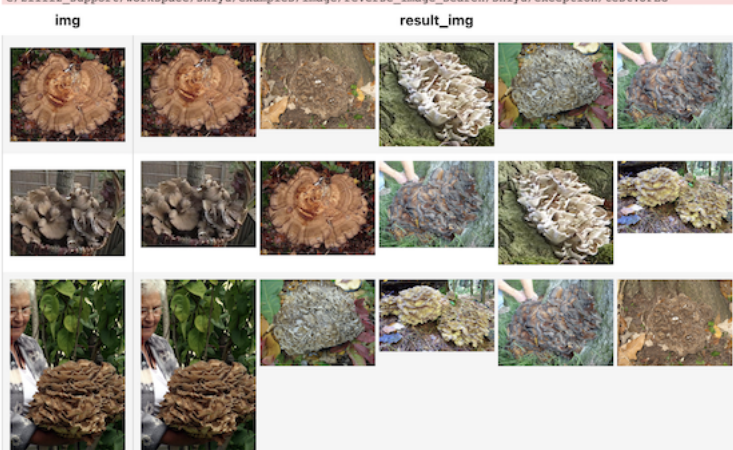
Exception Safe Execution

When we have large-scale image data, there may be bad data that will cause errors. Typically, the users don't want such errors to break the production system. Therefore, the data pipeline should continue to process the rest of the images and report the errors.

Towhee supports an exception-safe execution mode that allows the pipeline to continue on exceptions and represent the exceptions with `Empty` values. The user can choose how to deal with the `Empty` values at the end of the pipeline. During the query below, there are four images in total, one of them is broken, it just prints an error message instead of terminating because it has `exception_safe` and `drop_empty`, as you can see, `drop_empty` deletes `empty` data.

```
1 (
2     towhee.glob['path']('./exception/*.JPEG')
3     .exception_safe()
4     .image_decode['path', 'img']()
5     .image_embedding.timm['img', 'vec'](model_name='resnet50')
6     .tensor_normalize['vec', 'vec']()
7     .milvus_search['vec', 'result'](collection=resnet_collection,
8         limit=5)
9     .runas_op['result', 'result_img'](func=read_images)
10    .drop_empty()
11    .select['img', 'result_img']()
12    .show()
```

2022-05-17 17:44:14,341 - 140286030959168 - image_decode_cv2.py-image_decode_cv2:64 - ERROR: Read image ..
e/zilliz_support/workspace/shiyu/examples/image/reverse_image_search/shiyu/exception/test.JPEG



Deploy as a Microservice

duration: 2

The data pipeline used in our experiments can be converted to a function with `towhee.api` and `as_function()`, and we can also convert the data pipeline into a RESTful API with `serve()`, it generates FastAPI services from towhee pipelines.

Insert Image Data

```
1 import time
2 import towhee
3 from fastapi import FastAPI
4 from pymilvus import connections, Collection
5
6 app = FastAPI()
7 connections.connect(host='127.0.0.1', port='19530')
8 milvus_collection = Collection('resnet50')
9
10 @towhee.register(name='get_path_id')
11 def get_path_id(path):
12     timestamp = int(time.time()*10000)
13     id_img[timestamp] = path
14     return timestamp
15
16 @towhee.register(name='milvus_insert')
17 class MilvusInsert:
18     def __init__(self, collection):
19         self.collection = collection
20
21     def __call__(self, *args, **kwargs):
22         data = []
23         for iterable in args:
24             data.append([iterable])
25         mr = self.collection.insert(data)
26         self.collection.load()
27         return str(mr)
28
29 with towhee.api['file']() as api:
30     app_insert = (
31         api.image_load['file', 'img']()
32         .save_image['img', 'path'](dir='tmp/images')
33         .get_path_id['path', 'id']()
34         .image_embedding.timm['img', 'vec'](model_name='resnet50')
35         .tensor_normalize['vec', 'vec']()
36         .milvus_insert(['id', 'vec'], 'res')(collection=
37             milvus_collection)
38         .select['id', 'path']()
```

```
38         .serve('/insert', app)
39     )
```

Search Similar Image

```
1  with towhee.api['file']() as api:
2      app_search = (
3          api.image_load['file', 'img']()
4          .image_embedding.timm['img', 'vec'](model_name='resnet50')
5          .tensor_normalize['vec', 'vec']()
6          .milvus_search['vec', 'result'](collection=milvus_collection)
7          .runas_op['result', 'res_file'](func=lambda res: str([id_img[x.
8              id] for x in res]))
9          .select['res_file']()
10         .serve('/search', app)
```

Count Numbers

```
1  with towhee.api() as api:
2      app_count = (
3          api.map(lambda _: milvus_collection.num_entities)
4          .serve('/count', app)
5      )
```

Start Server

```
1  import uvicorn
2  import nest_asyncio
3
4  nest_asyncio.apply()
5  uvicorn.run(app=app, host='0.0.0.0', port=8000)
```

Finally to start FastAPI, there are three services `/insert`, `/search` and `/count`, you can run the following commands to test:

```
1  # upload an image and search
2  $ curl -X POST "http://0.0.0.0:8000/search" --data-binary @test/lion/
    n02129165_13728.JPEG -H 'Content-Type: image/jpeg'
3  # upload an image and insert
4  $ curl -X POST "http://0.0.0.0:8000/insert" --data-binary @test/banana
    /n07753592_323.JPEG -H 'Content-Type: image/jpeg'
5  # count the collection
6  $ curl -X POST "http://0.0.0.0:8000/count"
```