# Лабораторная работ 2. Airflow and docker.

## Грибанов Данил, 6233

## Ход работы

Все файлы (dockerfile, docker-compose и прочее) находятся в папке airflow (так же dags и скрипты для запуска соотв. частей).

Перед запуском самой системы (докера airflow) следует подготовить следующее:
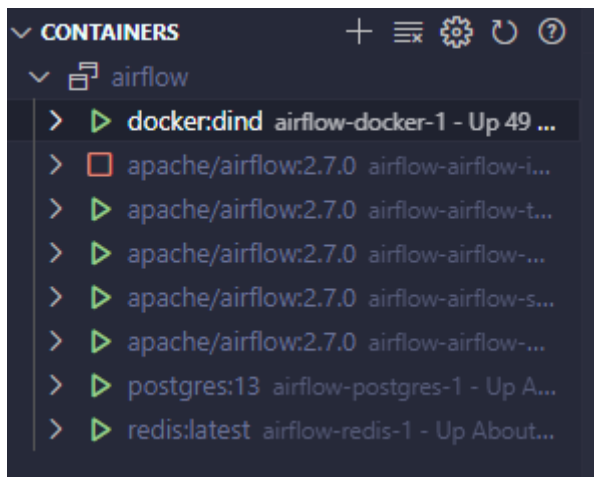
Создадим хранилище с данными для запуска и обучения всех систем:

```
313
314    volumes:
315      postgres-db-volume:
316      airflow-data-volume:
317        driver: local
318        driver_opts:
319          type: none
320          o: bind
321          device: "${AIRFLOW_PROJ_DIR:-.}/data"
```
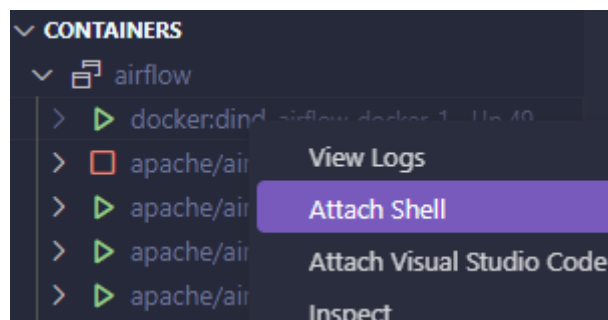
И добавим данное хранилище с папкой докер-файлов нужных в dind:

```
services:

  docker:
    image: docker:dind          You, 48 minutes ago • Uncommitted changes
    privileged: true
    # Build additional dockers in dind TODO: Is there some way to do init prebuild for dind?
#     command: 'cd /dockerfiles/huggy-face && docker build . -t huggy_face_image'
    environment:
      DOCKER_TLS_CERTDIR: ""
    volumes:
      - airflow-data-volume:/data
      - ${AIRFLOW_PROJ_DIR:-.}/dockerfiles:/dockerfiles

  postgres:
    image: postgres:13
```

Запустим docker-compose.yaml для сборки airflow и нужного:

Дождемся инициализации (пару минут), и зайдем в docker:dind:



Теперь подготовим здесь докер для использования нейронных сетей:

cd /dockerfiles/huggy-face && docker build . -t huggy_face_image

Дождемся сборки (до 10 минут). Теперь мы готовы к запуску основных DAG разработанных для лабораторных работ. Код DAG к соотв. пунктам задания имеет комментарии и пояснения.

*<u>Видео в краткий пересказ</u>*

Airflow/dags/airflow_lab2.py

```python
1.  import os
2.  from datetime import datetime
3.  from airflow import DAG
4.  from airflow.providers.docker.operators.docker import DockerOperator
5.  from airflow.sensors.filesystem import FileSensor
6.
7.  from docker.types import Mount
8.
9.  default_args = {
10.     'owner': 'airflow',
11.     'start_date': datetime(2023, 1, 1),
12.     'retries': 1,
13. }
14.
15.
16. dag = DAG(
17.     'audio_to_text_to_summary_to_pdf',
18.     default_args=default_args,
19.     description='DAG for extracting audio, transforming to text, summarizing, and saving
    as PDF',
```

```python
20.      schedule_interval=None,
21. )
22.
23. # TODO: Connection could be done via PythonAPI - but I didnt found HOW - so, do this in
    Web instead...
24. #file_connection = Connection(
25. #     conn_id="file_connection",
26. #     conn_type="fs",
27. #     description="Connection to file-path",
28. #)
29.
30. wait_for_new_file = FileSensor(
31.      task_id='wait_for_new_file',
32.      poke_interval=10,  # Interval to check for new files (in seconds)
33.      filepath='/opt/airflow/data/lab2',  # Target folder to monitor
34.      fs_conn_id='file_connection',
35.      dag=dag,
36. )
37.
38. extract_audio = DockerOperator(
39.      task_id='extract_audio',
40.      image='jrottenberg/ffmpeg',
41.      docker_url="tcp://docker:2375", # For Dind usage case
42.      mount_tmp_dir=False,
43.      network_mode='bridge',
44.      entrypoint='bash',
45.      command=['-c', 'cd /data/lab2 && for single_video in ./*.mp4; do ffmpeg -y -i
    "${single_video}" -ss 1 -to 5 -vn "./../lab2_output/${single_video}.wav"; done'],
46.      mounts=[
47.          Mount(source='/data', target='/data', type='bind'),
48.      ],
49.      dag=dag,
50. )
51.
52. audo2text = DockerOperator(
53.      task_id='audio2text',
54.      image='huggy_face_image',
55.      docker_url="tcp://docker:2375", # For Dind usage case
56.      mount_tmp_dir=False,
57.      network_mode='bridge',
58.      entrypoint='bash',
59.      command=['-c', "python /data/audio2text.py"],
60.      mounts=[
61.          Mount(source='/data', target='/data', type='bind'),
62.      ],
63.      dag=dag,
64. )
65.
66. text2summary = DockerOperator(
67.      task_id='text2summary',
68.      image='huggy_face_image',
69.      docker_url="tcp://docker:2375", # For Dind usage case
70.      mount_tmp_dir=False,
71.      network_mode='bridge',
72.      entrypoint='bash',
73.      command=['-c', "python /data/text2summary.py"],
74.      mounts=[
75.          Mount(source='/data', target='/data', type='bind'),
76.      ],
77.      dag=dag,
78. )
79.
80. wait_for_new_file >> extract_audio >> audo2text >> text2summary
```
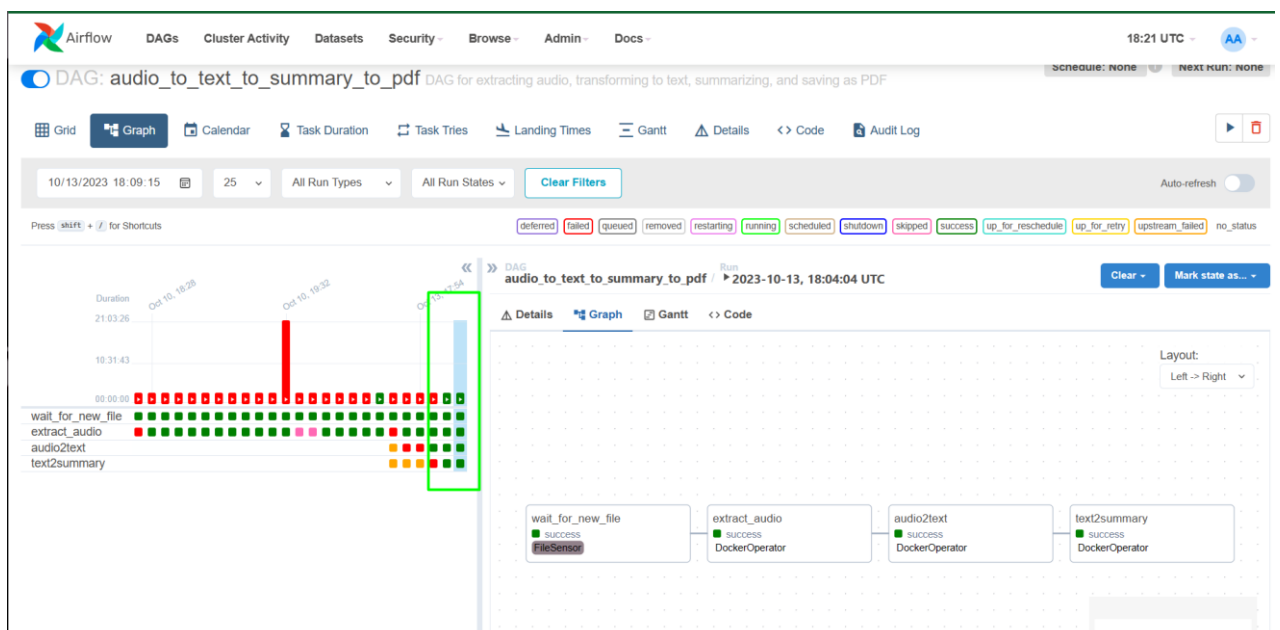
*Код audio2text:*

```python
1.  from transformers import pipeline
```

```
2.  import glob
3.  import os
4.
5.  for audio_file_path in glob.glob('/data/lab2_output/*.wav'):
6.      _, filename = os.path.split(audio_file_path)
7.      output_txt_file_path = f'/data/lab2_output/text/{filename}.txt'
8.      if os.path.isfile(output_txt_file_path):
9.          continue # skip, already exist
10.
11.     pipe = pipeline("automatic-speech-recognition", "openai/whisper-tiny")
12.     res = pipe(audio_file_path)
13.
14.     with open(output_txt_file_path, "w") as text_file:
15.         text_file.write(res['text'])
```

*Код text2summary:*

```
1.  from transformers import pipeline
2.  import glob
3.  import os
4.
5.  for txt_file_path in glob.glob('/data/lab2_output/text/*.txt'):
6.      _, filename = os.path.split(txt_file_path)
7.      output_summary_file_path = f'/data/lab2_output/summary/{filename}'
8.      if os.path.isfile(output_summary_file_path):
9.          continue # skip, already exist
10.
11.     with open(txt_file_path, 'r') as fr:
12.         text = fr.read()
13.
14.     summarizer = pipeline("summarization", max_length=9) # Since our input small
15.     text_summ = summarizer(text)
16.
17.     with open(output_summary_file_path, "w") as text_file:
18.         text_file.write(text_summ[0]['summary_text'])
```

## Обучение сети при обнаружении нового файла на примере MNIST

Для тестирования такой ситуации, я разбил обучающую выборку MNIST на 10 случайных частей. В папке airflow/data/lab2_nn_train/full_data – полные файлы данных, сама загрузка и слежка идет за папкой airflow/data/lab2_nn_train/data, куда можно закидывать или удалять файла для тестирования.

В качестве обучения взят набор данных MNIST и пару легких сверточных сетей.

Airflow/dags/airflow_train_lab2.py

```python
1.  import os
2.  from datetime import datetime
3.  from airflow import DAG
4.  from airflow.providers.docker.operators.docker import DockerOperator
5.  from airflow.sensors.filesystem import FileSensor
6.
7.  from docker.types import Mount
8.
9.  default_args = {
10.     'owner': 'airflow',
11.     'start_date': datetime(2023, 1, 1),
12.     'retries': 1,
13. }
14.
15.
16. dag = DAG(
17.     'train_nn',
18.     default_args=default_args,
19.     description='DAG train NN',
20.     schedule_interval=None,
21. )
22.
23. wait_for_new_file = FileSensor(
24.     task_id='wait_for_new_train_file',
25.     poke_interval=10,  # Interval to check for new files (in seconds)
26.     filepath='/opt/airflow/data/lab2_nn_train/data',  # Target folder to monitor
27.     fs_conn_id='file_train_connection',
28.     dag=dag,
29. )
30.
31. train_nn = DockerOperator(
32.     task_id='train_nn_on_updated_data',
33.     image='huggy_face_image',
34.     docker_url="tcp://docker:2375", # For Dind usage case
35.     mount_tmp_dir=False,
36.     network_mode='bridge',
37.     entrypoint='bash',
38.     command=['-c', "python /data/lab2_nn_train/train_nn_mnist.py"],
39.     mounts=[
40.         Mount(source='/data', target='/data', type='bind'),
41.     ],
42.     dag=dag,
43. )
44.
```

```
45.
46. wait_for_new_file >> train_nn
```

## Код обучения:

```python
1.  import argparse
2.  import torch
3.  import torch.nn as nn
4.  import torch.nn.functional as F
5.  import torch.optim as optim
6.  from torch.utils.data import Dataset
7.  from torchvision import transforms
8.  from torch.optim.lr_scheduler import StepLR
9.
10. import numpy as np
11. import glob
12. import os
13. import sys
14. from datetime import datetime
15.
16. SAVE_MODEL_PATH = '/data/lab2_nn_train/models'
17. DATA_PATH = '/data/lab2_nn_train/data'
18. TEST_FILENAME = 'test_data.npz'
19. TRAIN_PREFIX = 'train_'
20.
21.
22. BATCH_SIZE = 32
23. EPOCHS = 15
24. LR = 1e-4
25. LR_STEP = 0.7
26. CUDA = False
27. SEED = 2023
28. LOG_INTERVAL = 10
29. SAVE_MODEL = True
30.
31.
32. class MnistDataset(Dataset):
33.
34.     def __init__(self, x, y, transform=None):
35.         assert len(x) == len(y)
36.         self.x = x
37.         self.y = y
38.         self.transform = transform
39.
40.     def __len__(self):
41.         return len(self.x)
42.
43.     def __getitem__(self, idx):
44.         x_sample = self.x[idx]
45.
46.         if self.transform:
47.             x_sample = self.transform(x_sample)
48.
49.         return x_sample, self.y[idx]
50.
51.
52. class Net(nn.Module):
53.     def __init__(self):
54.         super(Net, self).__init__()
55.         self.conv1 = nn.Conv2d(1, 32, 3, 1)
56.         self.conv2 = nn.Conv2d(32, 64, 3, 1)
57.         self.dropout1 = nn.Dropout(0.25)
58.         self.dropout2 = nn.Dropout(0.5)
59.         self.fc1 = nn.Linear(9216, 128)
60.         self.fc2 = nn.Linear(128, 10)
61.
62.     def forward(self, x):
```

```python
63.         x = self.conv1(x)
64.         x = F.relu(x)
65.         x = self.conv2(x)
66.         x = F.relu(x)
67.         x = F.max_pool2d(x, 2)
68.         x = self.dropout1(x)
69.         x = torch.flatten(x, 1)
70.         x = self.fc1(x)
71.         x = F.relu(x)
72.         x = self.dropout2(x)
73.         x = self.fc2(x)
74.         output = F.log_softmax(x, dim=1)
75.         return output
76.
77.
78. def train(model, device, train_loader, optimizer, epoch):
79.     model.train()
80.     for batch_idx, (data, target) in enumerate(train_loader):
81.         data, target = data.to(device), target.to(device)
82.         optimizer.zero_grad()
83.         output = model(data)
84.         loss = F.nll_loss(output, target)
85.         loss.backward()
86.         optimizer.step()
87.         if batch_idx % LOG_INTERVAL == 0:
88.             print('Train Epoch: {} [{}/{} ({:.0f}%)]\tLoss: {:.6f}'.format(
89.                 epoch, batch_idx * len(data), len(train_loader.dataset),
90.                 100. * batch_idx / len(train_loader), loss.item()))
91.
92.
93. def test(model, device, test_loader):
94.     model.eval()
95.     test_loss = 0
96.     correct = 0
97.     with torch.no_grad():
98.         for data, target in test_loader:
99.             data, target = data.to(device), target.to(device)
100.             output = model(data)
101.             test_loss += F.nll_loss(output, target, reduction='sum').item()  #
    sum up batch loss
102.             pred = output.argmax(dim=1, keepdim=True)  # get the index of the max
    log-probability
103.             correct += pred.eq(target.view_as(pred)).sum().item()
104.
105.     test_loss /= len(test_loader.dataset)
106.
107.     print('\nTest set: Average loss: {:.4f}, Accuracy: {}/{} ({:.0f}%)\n'.format(
108.         test_loss, correct, len(test_loader.dataset),
109.         100. * correct / len(test_loader.dataset)))
110.
111.
112.     def main():
113.         use_cuda = CUDA and torch.cuda.is_available()
114.
115.         torch.manual_seed(SEED)
116.         exp_folder_path = os.path.join('/data/lab2_nn_train/res',
    str(datetime.now()))
117.         os.makedirs(exp_folder_path, exist_ok=True)
118.         sys.stdout = open(os.path.join(exp_folder_path, "print.log"), 'w')
119.
120.         if use_cuda:
121.             device = torch.device("cuda")
122.         else:
123.             device = torch.device("cpu")
124.
125.         train_kwargs = {'batch_size': BATCH_SIZE}
126.         test_kwargs = {'batch_size': BATCH_SIZE}
127.         if use_cuda:
128.             cuda_kwargs = {'num_workers': 1,
```

```python
                                'pin_memory': True,
                                'shuffle': True}
                train_kwargs.update(cuda_kwargs)
                test_kwargs.update(cuda_kwargs)

            transform=transforms.Compose([
                transforms.ToTensor(),
                transforms.Normalize((0.1307,), (0.3081,))
                ])
            x_train_list = []
            y_train_list = []

            for single_train_file in glob.glob(os.path.join(DATA_PATH,
        f'{TRAIN_PREFIX}*.npz')):
                single_train_loaded = np.load(single_train_file)
                x_train_list.append(single_train_loaded['x_train'])
                y_train_list.append(single_train_loaded['y_train'])

            x_train_np = np.concatenate(x_train_list, axis=0)
            y_train_np = np.concatenate(y_train_list, axis=0)
            dataset_train = MnistDataset(
                x=x_train_np, y=y_train_np,
                transform=transform
            )

            test_loaded = np.load(os.path.join(DATA_PATH, TEST_FILENAME))
            dataset_test = MnistDataset(
                x=test_loaded['x_test'], y=test_loaded['y_test'],
                transform=transform
            )
            train_loader = torch.utils.data.DataLoader(dataset_train,**train_kwargs)
            test_loader = torch.utils.data.DataLoader(dataset_test, **test_kwargs)

            model = Net().to(device)
            optimizer = optim.Adadelta(model.parameters(), lr=LR)

            scheduler = StepLR(optimizer, step_size=1, gamma=LR_STEP)
            for epoch in range(1, EPOCHS + 1):
                train(model, device, train_loader, optimizer, epoch)
                test(model, device, test_loader)
                scheduler.step()

            if SAVE_MODEL:
                torch.save(model.state_dict(), os.path.join(exp_folder_path,
        'mnist_cnn.pt'))


        if __name__ == '__main__':
            main()
```