

БУЛДАКОВ ДМИТРИЙ

6233-010402D

ХОД РАБОТЫ

Пайплайн для инференса данных

Перед началом работы требовалось создать собственный образ (или взять существующий) контейнера с необходимыми библиотеками для дальнейшей работы с первым и вторым заданием. Получился следующий докер-файл.

```
FROM tensorflow/tensorflow:latest

RUN pip install fpdf numpy pandas scikit-learn

LABEL Maintainer="vapaov.tensorflow_sklearn"

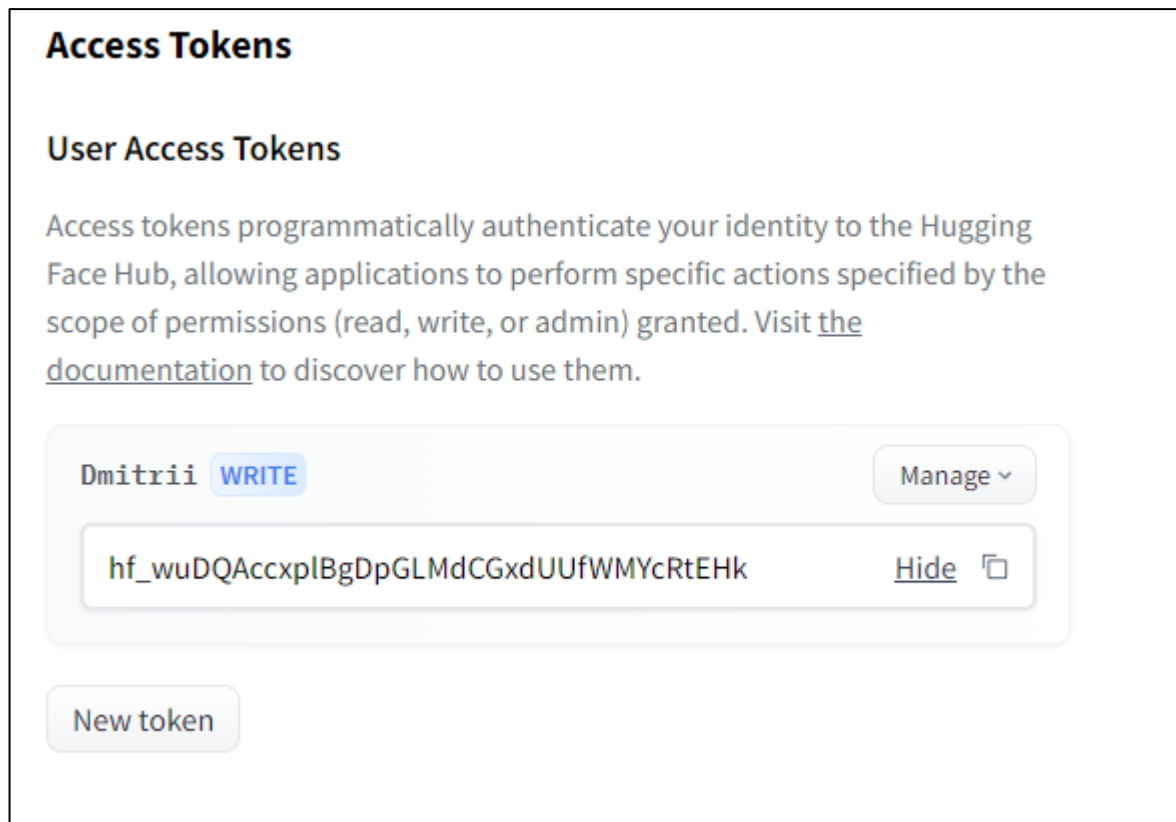
WORKDIR /usr/app/src
```

Затем образ собрали, присвоили тэг и отправили в докер-хаб.

```
C:\Users\Odmityr>cd prerequisites
C:\Users\Odmityr\prerequisites>docker build . -t custom_container
[+] Building 28.3s (4/7)                                docker:default
[+] Building 14.7s (8/8) FINISHED                        docker:default
-> [internal] load build definition from Dockerfile
-> [internal] load .dockerignore
-> [internal] load metadata for docker.io/tensorflow/tensorflow:latest
-> resolve docker.io/tensorflow/tensorflow:latest@sha256:4689c724a7d65a7d289cc2ae536fa3cdeb636b2df3e23da3a44f
-> sha256:6a8c4ad35b6db8ef0bea19569618f584cb4a8adf5bed962e3db9aa4c0010d1 4.54kB / 4.54kB
-> sha256:5cc77b1d99a1e092e1bf48f5ead9028d5ac80e76729a10f02f74cb72a82d870 769B / 769B
-> sha256:345716e6b3b788c2ad108607e9745804ebf3242a7000c7b55cbb99c1a971d 861B / 861B
-> sha256:4689c724a7d65a7d289cc2ae536fa3cdeb636b2df3e23da3a44fdd66ac3af46f 2.83kB / 2.83kB
-> sha256:aec8e9d3d3972efa43bf4ee3cda659c0f787f8f59c82fb3e48c87cb22a12e 29.54MB / 29.54MB
-> sha256:5cd1e9a0284bf5a49731ecf8d7436d199fef8adc18745d487843cae92d93e9b1 163B / 163B
-> sha256:f190ee1fa1d7baf03b64c1a00e441777c2ba26c3478a1fd46049b8415f919a 37.63MB / 37.63MB
-> sha256:909fb368ac91957ba7f421488cc96a1a9d732795d6c754f9d8cdbc4792370 133.67MB / 133.67MB
-> extracting sha256:aec8e9d3d3972efa43bf4ee3cda659c0f787f8f59c82fb3e48c87cb22a12e 2.15
-> sha256:bcb6c16aad4fac3c21d24d39c503ae697731ba183e0bef0ac576ff529fc1baec 970B / 970B
-> sha256:77cb63a520e5d19dd7e7e1a3b29b35aca75db372592dd4b14d144ed9e6ff5273 128B / 128B
-> sha256:3268aa633e510fe5595e2faade355f5ca532d32987ed8087a6e1ffc92b98f9b 37.27MB / 37.27MB
-> extracting sha256:345716e6b3b788c2ad108607e9745804ebf3242a7000c7b55cbb99c1a971d 0.0B
-> extracting sha256:5cc77b1d99a1e092e1bf48f5ead9028d5ac80e76729a10f02f74cb72a82d870 0.0B
-> extracting sha256:5cd1e9a0284bf5a49731ecf8d7436d199fef8adc18745d487843cae92d93e9b1 0.0B
-> sha256:c6cbc212b10dd3c38925320760e6bbe7fb3047ddae8ce92fb893dde185e95e1 297.40MB / 297.40MB
-> sha256:f2e46a0a82df72a17e81b6c585aae77a18510b5e96186888c75b54ef277bf82 1.13kB / 1.13kB
-> extracting sha256:3268aa633e510fe5595e2faade355f5ca532d32987ed8087a6e1ffc92b98f9b 1.0B
-> sha256:6c6bc212b10dd3c38925320760e6bbe7fb3047ddae8ce92fb893dde185e95e1 1.13kB / 1.13kB
-> extracting sha256:5cd1e9a0284bf5a49731ecf8d7436d199fef8adc18745d487843cae92d93e9b1 9.3B
-> extracting sha256:bcb6c16aad4fac3c21d24d39c503ae697731ba183e0bef0ac576ff529fc1baec 0.0B
-> extracting sha256:77cb63a520e5d19dd7e7e1a3b29b35aca75db372592dd4b14d144ed9e6ff5273 0.0B
-> extracting sha256:3268aa633e510fe5595e2faade355f5ca532d32987ed8087a6e1ffc92b98f9b 3.0B
-> extracting sha256:6c6bc212b10dd3c38925320760e6bbe7fb3047ddae8ce92fb893dde185e95e1 18.55
-> extracting sha256:f2e46a0a82df72a17e81b6c585aae77a18510b5e96186888c75b54ef277bf82 0.0B
-> extracting sha256:831122755b02c47aa9ad5a9ad0da7717c9cbfe7f8a5230d24f6747d35e9925 0.0B
-> [2/3] RUN pip install fpdf numpy pandas scikit-learn
-> [3/3] WORKDIR /usr/app/src
-> exporting to image
-> exporting layers
-> writing image sha256:cbacc2c0ae692af9055f32ea793f4aad8c4ffde03a99b530036af4763be7ede
-> naming to docker.io/library/custom_container

What's Next?
View summary of image vulnerabilities and recommendations → docker scout quickview
C:\Users\Odmityr\prerequisites>
```


В рамках первого задания будем использовать API Hugging Face. Для доступа к API требуется регистрация и получение специального токена.



DAG-файл получился длинным и адекватно в отчёт его не вставить. Можно посмотреть в репозитории: video_to_summary_pdf.py. Для пайплайна были написаны специальные функции.

1_audio_to_text.py – переводит звуковую дорожку файла в текст при помощи сервиса openai/whisper-small

```
1 import requests
2 url = "https://api-inference.huggingface.co/models/openai/whisper-small"
3 headers = {"Authorization": "Bearer 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"}
4
5 with open('/data/audio.aac', "rb") as file:
6     data = file.read()
7     response = requests.post(url, headers=headers, data=data)
8     result = response.json()
9     text_from_audio = result['text']
10    text_file = open("/data/text.txt", "w+")
11    text_file.write(text_from_audio)
12    text_file.close()
```

2_text_to_summary.py – кратко пересказывает написанное в текстовом файле. Используется сервис slauw87/bart_summarisation.

```
import requests
url = "https://api-inference.huggingface.co/models/slauw87/bart_summarisation"
headers = {"Authorization": "Bearer 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"}

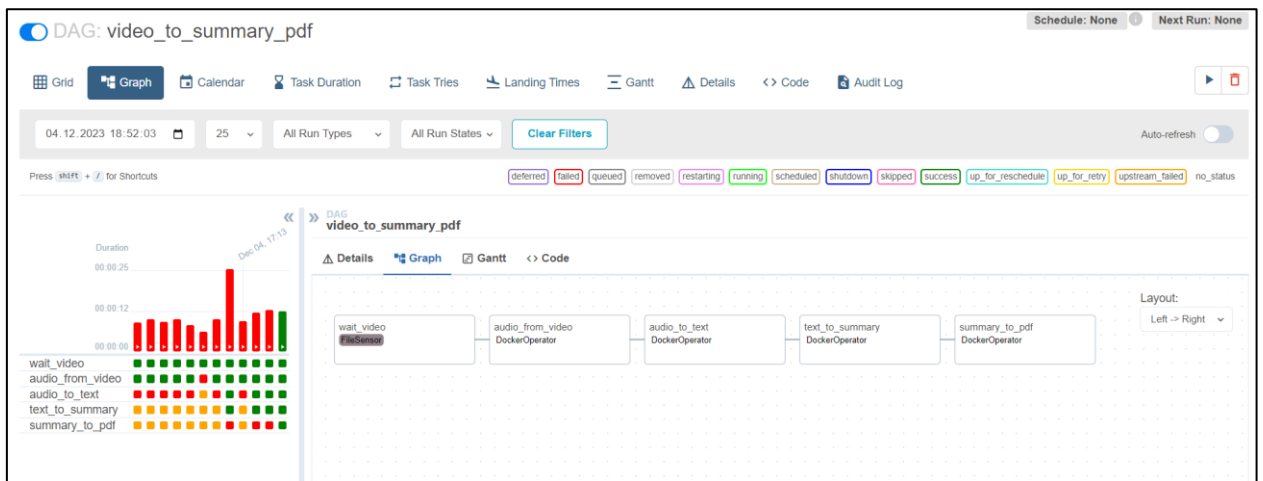
with open('/data/text.txt', "rb") as file:
    data = file.read()
    response = requests.post(url, headers=headers, json={'inputs': f"{data}"})
    result = response.json()
    summary = result[0]['summary_text']
    text_file = open("/data/summary.txt", "w+")
    text_file.write(summary)
    text_file.close()
```

3_summary_to_pdf.py – с помощью библиотеки FPDF формирует файл формата PDF.

```
from fpdf import FPDF

file = open("/data/summary.txt", "r")
pdf = FPDF()
pdf.add_page()
for text in file:
    pdf.set_font('Helvetica', '', size=14)
    pdf.multi_cell(w=100, h=10, border = 1, txt = text, align = 'J')
pdf.output("/data/pdf_result.pdf")
```

После всех настроек и нескольких попыток, пайплайн успешно запустился.



В результате получился PDF файл с следующим содержанием:

Mark will buy a new shirt for Mark. Mark and Mark will have a coffee and then go shopping. Mark has a meeting at 12.30, so they have a lot of time to go shopping before the next meeting. Mark will be back in 10 minutes.

Пайплайн для обучения модели

Для обучения модели написал небольшой DAG.

```
from datetime import datetime
from airflow import DAG
from docker.types import Mount
from airflow.providers.docker.operators.docker import DockerOperator

default_args = {
    'owner': 'airflow',
    'start_date': datetime(2023, 12, 3),
}

dag = DAG(
    'train_mnist_model',
    default_args=default_args,
    schedule_interval=None,
)

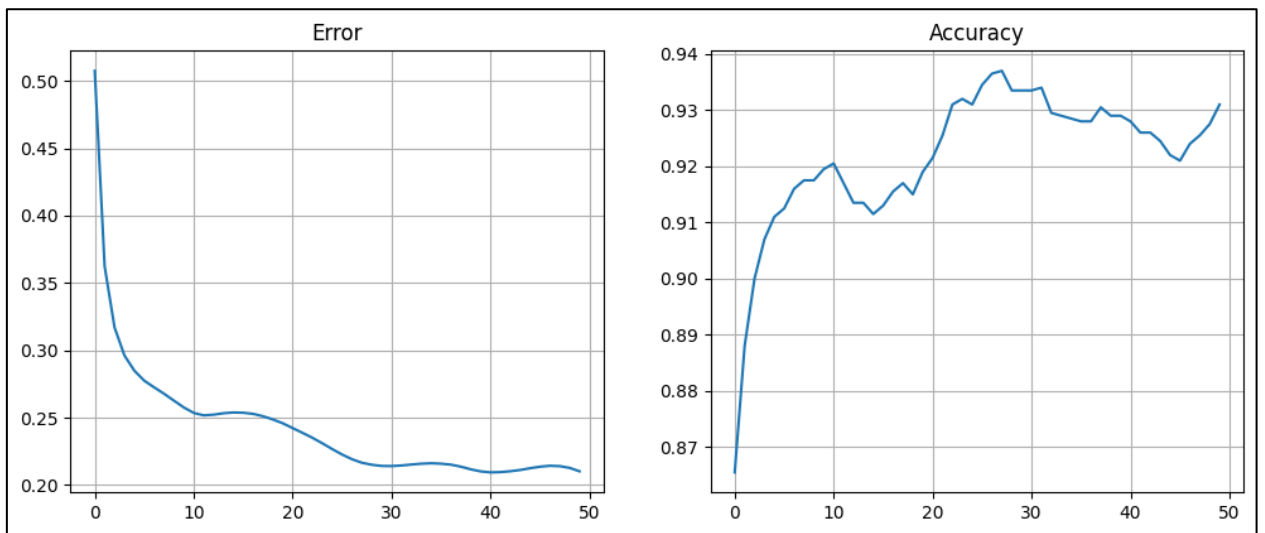
train_model = DockerOperator(
    task_id='train_model',
    image='vapaov/custom_container:1.0',
    command='python /data/nn_mnist.py',
    mounts=[Mount(source='/data', target='/data', type='bind')],
    docker_url="tcp://docker-proxy:2375",
    dag=dag,
)

train_model
```

Я взял нейросеть, написанную мною в рамках курса «Нейронные сети и машинное обучение». Она имеет довольно простой алгоритм и невысокую точность, но, тем не менее, пайплайн я захотел протестировать именно на ней. Обучение происходит на MNIST-датасете. К сожалению гитхаб не позволяет загружать файлы больше 25 МБ.

Ссылка на датасет: <https://www.kaggle.com/datasets/oddrational/mnist-in-csv>

На скриншоте ниже можно наблюдать графики функции потерь и точности.



В `mnist_nn.py` с помощью библиотеки `logging` реализовал логирование. После успешного запуска пайплайна логи можно посмотреть в файле `log.log`

