

Департамент образования и науки города Москвы
Государственное автономное образовательное учреждение высшего
образования города Москвы

«Московский городской педагогический университет»

Институт цифрового образования

Департамент информатики, управления и технологий

Макарова Виктория Сергеевна

Бизнес кейс «Rocket»

Проектный практикум по разработке ETL-решений

Направление подготовки

38.03.05 Бизнес-информатика

Профиль подготовки

Аналитика данных и эффективное управление

Курс обучения: 4

Форма обучения: очная

Преподаватель: кандидат технических наук,

доцент Босенко Тимур Муртазович

Москва

2025

Задачи

Общее задание. Создать исполняемый файл с расширением .sh, который автоматизирует выгрузку данных из контейнера в основную ОС данных, полученные в результате работы DAG в Apache Airflow.

1.1 Спроектировать верхнеуровневую архитектуру аналитического решения задания Бизнес-кейса «Rocket» в draw.io. Необходимо использовать:

Source Layer - слой источников данных.

Storage Layer - слой хранения данных.

Business Layer - слой для доступа к данным пользователей.

1.2 Спроектировать архитектуру DAG Бизнес-кейса «Rocket» в draw.io. Необходимо использовать:

Source Layer - слой источников данных.

Storage Layer - слой хранения данных.

Business Layer - слой для доступа к данным пользователей.

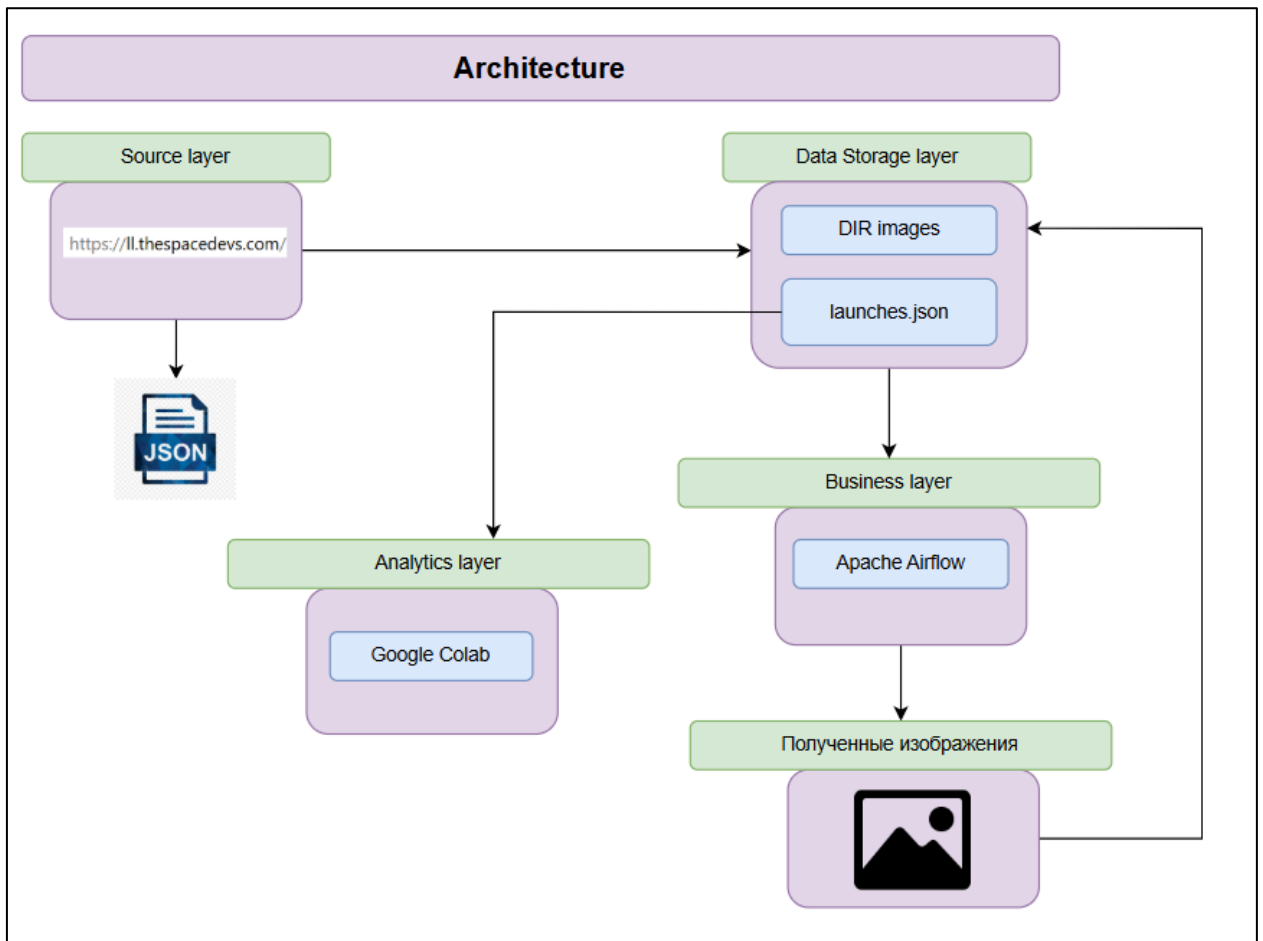
1.3 Построить диаграмму Ганта работы DAG в Apache Airflow.

Индивидуальное задание Вариант 7

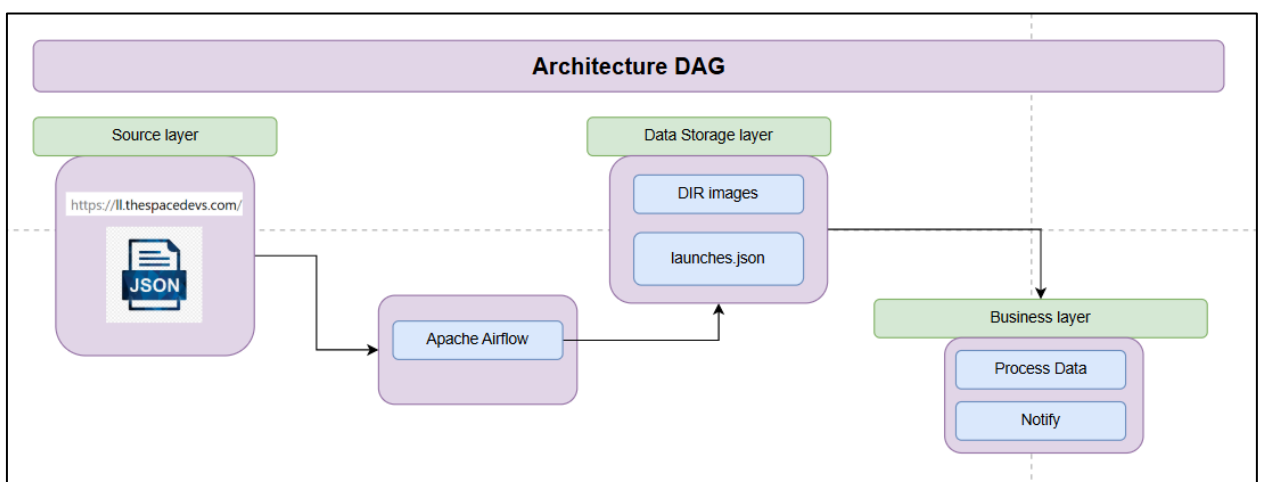
1. Создать отчет по скачиванию изображений на основе данных JSON.
2. Оценить необходимость использования BashOperator и PythonOperator для таких задач.
3. Изучить альтернативы для получения данных о стартах ракет с использованием API.

Ход работы

Представлена архитектура аналитического решения задания Бизнес-кейса «Rocket» в draw.io.



Представлена архитектура DAG Бизнес-кейса «Rocket» в draw.io

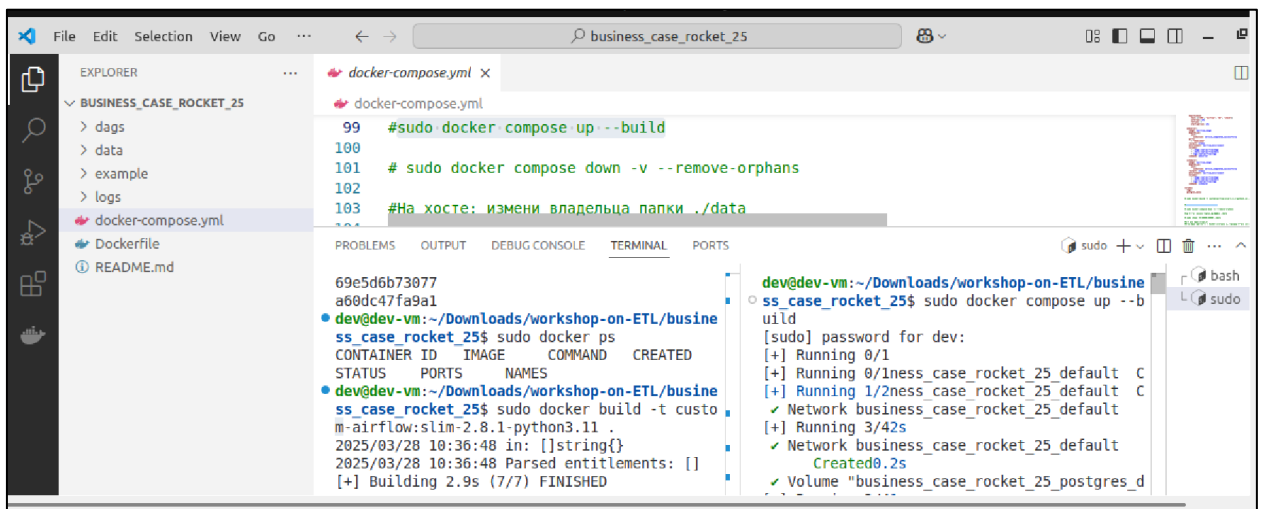


```
dev@dev-vm:~/Downloads/workshop-on-ETL/business_case_rocket_25$ sudo chown -R 50000:50000 ./data
[sudo] password for dev:
dev@dev-vm:~/Downloads/workshop-on-ETL/business_case_rocket_25$
```

Проверка активных контейнеров и остановка их работы.

```
dev@dev-vm: ~/Downloads/workshop-on-ETL/business_case_rocket_25$ docker ps
CONTAINER ID   IMAGE                                COMMAND                  CREATED        STATUS        PORTS                               NAMES
69e5d6b73077   custom-airflow:slim-2.8.1-python3.11  "/usr/bin/dumb-init ..." 6 days ago    Up About a minute    0.0.0.0:8080->8080/tcp, [::]:8080->8080/tcp    business_case_umbrella_25-webserver-1
a60dc47fa9a1   custom-airflow:slim-2.8.1-python3.11  "/usr/bin/dumb-init ..." 6 days ago    Up About a minute    8080/tcp                                     business_case_umbrella_25-scheduler-1
dev@dev-vm: ~/Downloads/workshop-on-ETL/business_case_rocket_25$ sudo docker stop $(sudo docker ps -q)
69e5d6b73077
a60dc47fa9a1
dev@dev-vm: ~/Downloads/workshop-on-ETL/business_case_rocket_25$ sudo docker ps
CONTAINER ID   IMAGE                                COMMAND                  CREATED        STATUS        PORTS                               NAMES
dev@dev-vm: ~/Downloads/workshop-on-ETL/business_case_rocket_25$
```

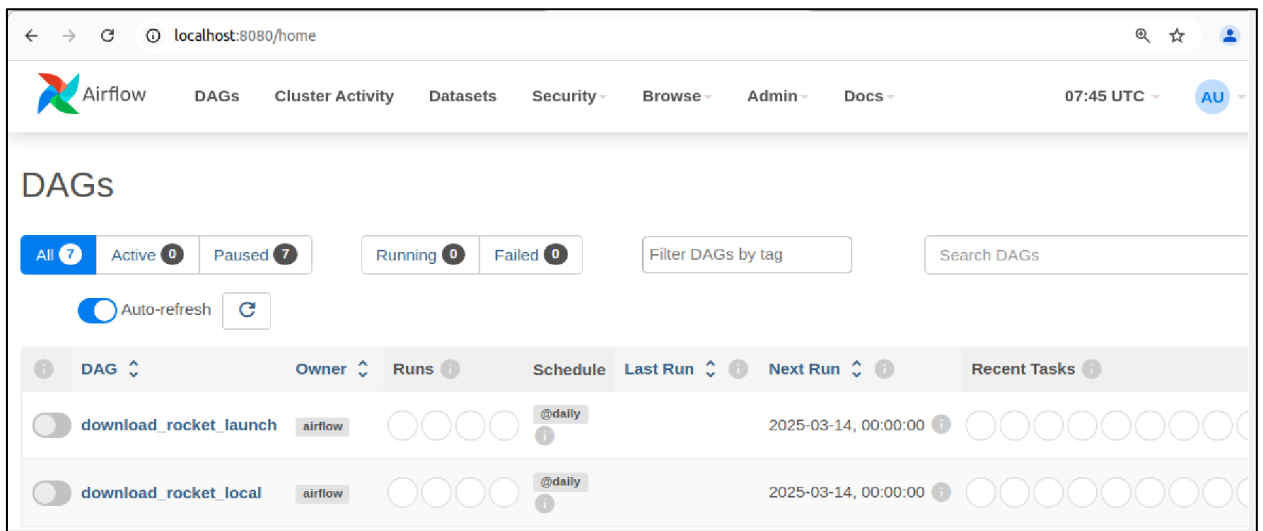
Билдим и запускаем контейнеры.



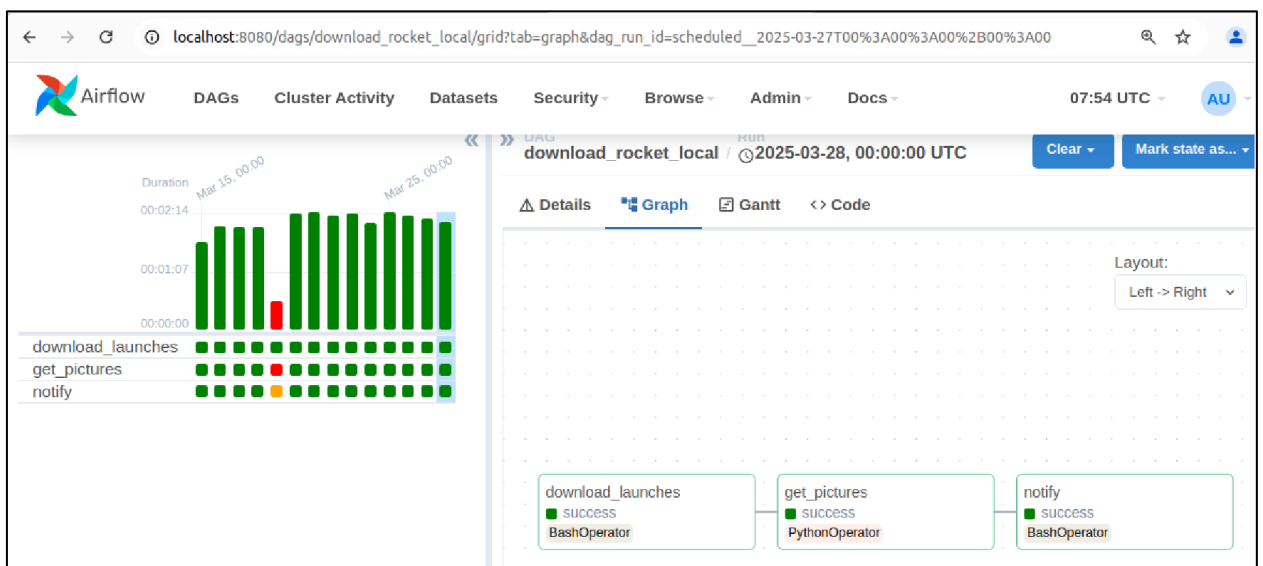
Код используемого дага

```
$ copy_files.sh download_rocket_local.py x
dags > download_rocket_local.py
1 import json
2 import pathlib
3
4 import airflow.utils.dates
5 import requests
6 import requests.exceptions as requests_exceptions
7 from airflow import DAG
8 from airflow.operators.bash import BashOperator
9 from airflow.operators.python import PythonOperator
10
11 dag = DAG(
12     dag_id="download_rocket_local",
13     description="Download rocket pictures of recently launched rockets.",
14     start_date=airflow.utils.dates.days_ago(14),
15     schedule_interval="@daily"
16 )
17
18 # Изменение пути для скачивания JSON-файла в папку data
19 download_launches = BashOperator(
20     task_id="download_launches",
21     bash_command="curl -o /opt/airflow/data/launches.json -L 'https://ll.thespacedevs.com/2.0.0/launch/upcoming', # noqa: E501
22     dag=dag,
23 )
24
25 def _get_pictures():
26     # Обеспечиваем существование директории для изображений в папке data
27     images_dir = "/opt/airflow/data/images"
28     pathlib.Path(images_dir).mkdir(parents=True, exist_ok=True)
29
30     # Загружаем все картинки из launches.json
31     with open("/opt/airflow/data/launches.json") as f:
32         launches = json.load(f)
33         image_urls = [launch["image"] for launch in launches["results"]]
34         for image_url in image_urls:
35             try:
36                 response = requests.get(image_url)
37                 image_filename = image_url.split("/")[-1]
38                 target_file = f"{images_dir}/{image_filename}"
39                 with open(target_file, "wb") as f:
40                     f.write(response.content)
41                 print(f"Downloaded {image_url} to {target_file}")
42             except requests.exceptions.MissingSchema:
```

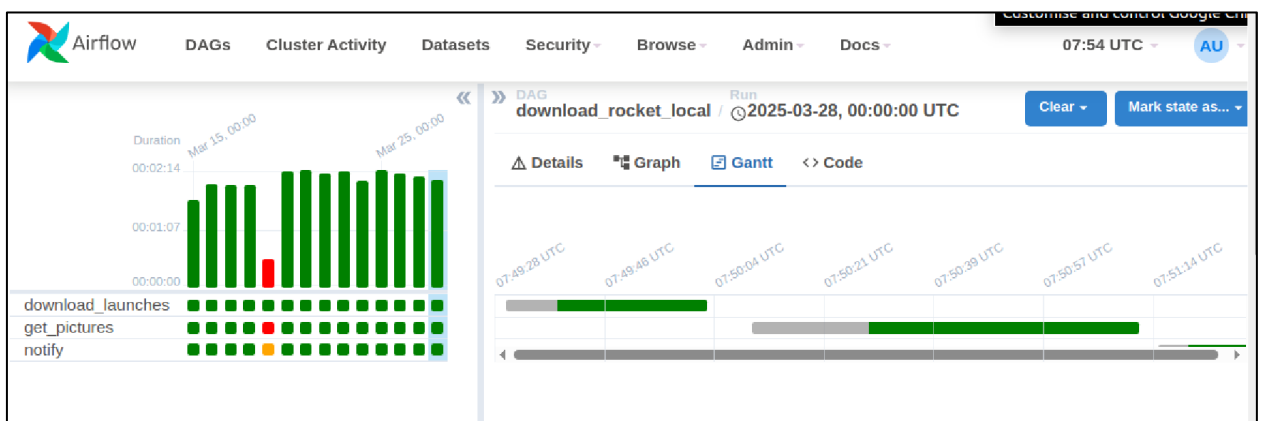
Проверяем доступность AirFlow.



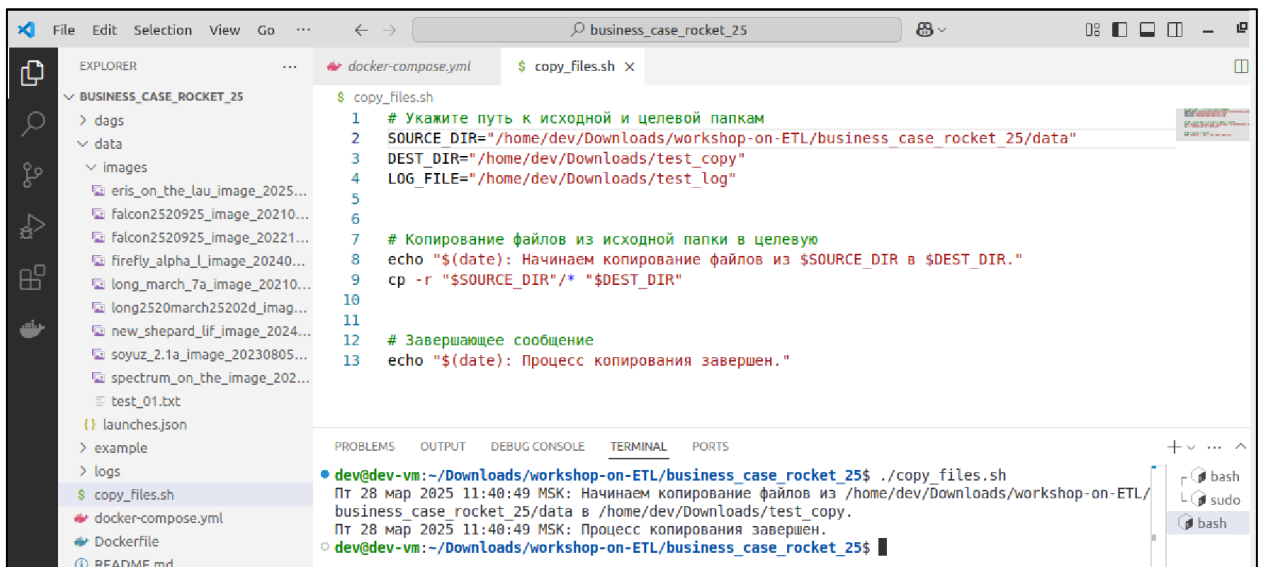
Запускаем даг `download_rocket_local`. Даг отработал успешно.



На рисунке представлена диаграмма Ганта выполнения дага.



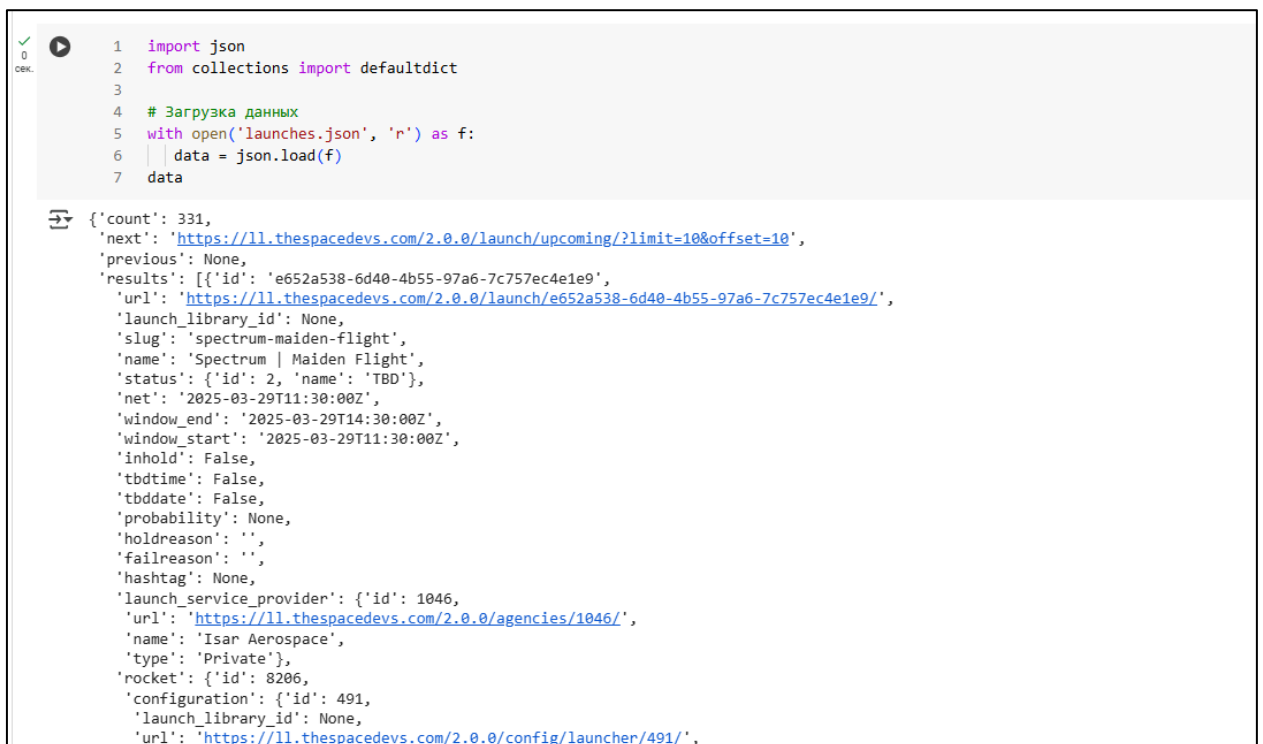
Далее создаем исполняемый файл `copy_files.sh`, который автоматизирует выгрузку данных из контейнера в основную ОС данных, полученные в результате работы DAG в Apache Airflow.



Задание 1

Переходим к аналитике и созданию отчета в GoogleColab

Импортируем файл launches.json и выведем данные этого файла



Создадим отчет о предстоящих запусках

```
1 import json
2 from datetime import datetime
3
4
5 # Анализ данных
6 total_launches = data['count']
7 upcoming_launches = len(data['results'])
8
9 # Обработка каждого запуска
10 launch_details = []
11 for launch in data['results']:
12     # Парсинг даты
13     net_time = datetime.strptime(launch['net'], '%Y-%m-%dT%H:%M:%SZ')
14
15     details = {
16         'name': launch['name'],
17         'status': launch['status']['name'],
18         'date': net_time.strftime('%Y-%m-%d'),
19         'time': net_time.strftime('%H:%M:%S'),
20         'provider': launch['launch_service_provider']['name'],
21         'rocket': launch['rocket']['configuration']['full_name'],
22         'window': f"{launch['window_start']} to {launch['window_end']}"
23     }
24     launch_details.append(details)
25
26 # Генерация отчета
27 report = f"""
28 === ОТЧЕТ О ПРЕДСТОЯЩИХ ЗАПУСКАХ ===
29 Всего предстоящих запусков в системе: {total_launches}
30 Анализируемых в этом отчете: {upcoming_launches}
31
32 Детали ближайших запусков:
33 """
34
35 for idx, launch in enumerate(launch_details, 1):
36     report += f"""
37 {idx}. {launch['name']}
```

```
=== ОТЧЕТ О ПРЕДСТОЯЩИХ ЗАПУСКАХ ===
Всего предстоящих запусков в системе: 331
Анализируемых в этом отчете: 10

Детали ближайших запусков:

1. Spectrum | Maiden Flight
  - Ракета: Spectrum
  - Провайдер: Isar Aerospace
  - Дата: 2025-03-29
  - Время: 11:30:00 UTC
  - Статус: TBD
  - Окно запуска: 2025-03-29T11:30:00Z to 2025-03-29T14:30:00Z

2. Long March 7A | Unknown Payload
  - Ракета: Long March 7A
  - Провайдер: China Aerospace Science and Technology Corporation
  - Дата: 2025-03-29
  - Время: 16:05:00 UTC
  - Статус: Go
  - Окно запуска: 2025-03-29T15:57:00Z to 2025-03-29T17:15:00Z

3. Firefly Alpha | FLTA006 (Message in a Booster)
  - Ракета: Firefly Alpha
  - Провайдер: Firefly Aerospace
  - Дата: 2025-03-30
  - Время: 13:37:00 UTC
  - Статус: TBD
  - Окно запуска: 2025-03-30T13:37:00Z to 2025-03-30T15:16:00Z

4. Falcon 9 Block 5 | Starlink Group 6-80
  - Ракета: Falcon 9 Block 5
  - Провайдер: SpaceX
  - Дата: 2025-03-30
  - Время: 19:20:20 UTC
  - Статус: TBD
  - Окно запуска: 2025-03-30T19:16:00Z to 2025-03-30T23:45:00Z
```

Выведем статистику по используемым провайдерам.

Статистика по провайдерам

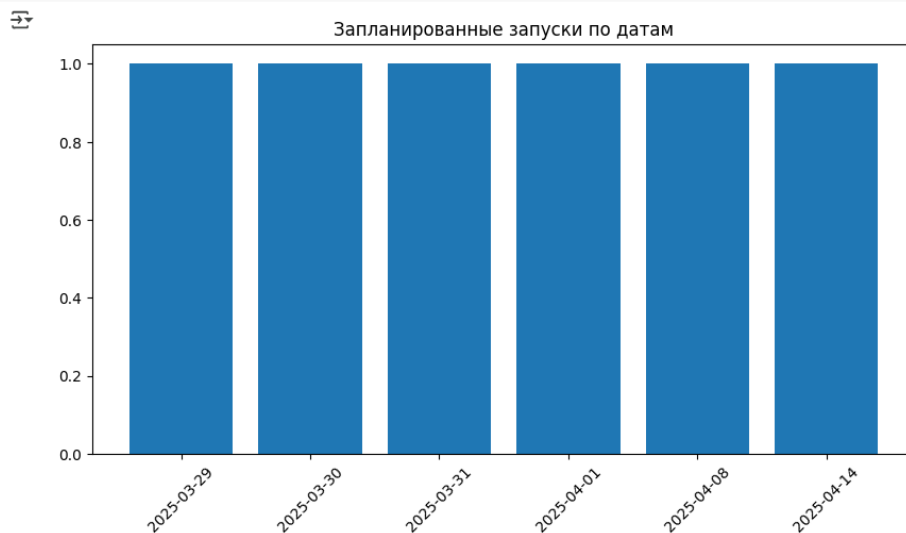
```
1 from collections import Counter
2 providers = Counter([l['provider'] for l in launch_details])
3 providers
```

```
Counter({'Isar Aerospace': 1,
        'China Aerospace Science and Technology Corporation': 2,
        'Firefly Aerospace': 1,
        'SpaceX': 3,
        'Gilmour Space Technologies': 1,
        'Russian Federal Space Agency (ROSCOSMOS)': 1,
        'Blue Origin': 1})
```

Таким образом больше всего провайдеров SpaceX : 3

Построим график по ближайшим запланированным запускам ракет.

```
1 import matplotlib.pyplot as plt
2
3 dates = [l['date'] for l in launch_details]
4 plt.figure(figsize=(10,5))
5 plt.bar(dates, [1]*len(dates))
6 plt.xticks(rotation=45)
7 plt.title('Запланированные запуски по датам')
8 plt.show()
```



Выведем расширенный отчет о предстоящих запусках


```
1 import json
2 from datetime import datetime
3 from collections import Counter
4 import matplotlib.pyplot as plt
5 import calendar
6
7
8 # Основная обработка данных
9 launches = data['results']
10 total_launches = data['count']
11
12 # 1. Распределение по статусам запусков
13 status_counts = Counter([launch['status']['name'] for launch in launches])
14
15 # 2. Частота запусков по дням недели
16 weekday_counts = Counter()
17 for launch in launches:
18     date = datetime.strptime(launch['net'], '%Y-%m-%dT%H:%M:%SZ')
19     weekday_counts[calendar.day_name[date.weekday()]] += 1
20
21 # 3. Топ самых активных провайдеров
22 provider_counts = Counter([launch['launch_service_provider']['name'] for launch in launches])
23 top_providers = provider_counts.most_common(5)
24
25 # 4. Средняя продолжительность окон запуска
26 window_durations = []
27 for launch in launches:
28     start = datetime.strptime(launch['window_start'], '%Y-%m-%dT%H:%M:%SZ')
29     end = datetime.strptime(launch['window_end'], '%Y-%m-%dT%H:%M:%SZ')
30     duration = (end - start).total_seconds() / 3600 # в часах
31     window_durations.append(duration)
32 avg_window = sum(window_durations) / len(window_durations)
33
34 # Генерация текстового отчета
35 report = f"""
36 === РАСШИРЕННЫЙ ОТЧЕТ О ПРЕДСТОЯЩИХ ЗАПУСКАХ ===
37 Всего предстоящих запусков: {total_launches}
38 Анализируется записей: {len(launches)}
39
```

```
=== РАСШИРЕННЫЙ ОТЧЕТ О ПРЕДСТОЯЩИХ ЗАПУСКАХ ===
Всего предстоящих запусков: 331
Анализируется записей: 10

1. РАСПРЕДЕЛЕНИЕ ПО СТАТУСАМ:
-----
TBD: 5 (50.0%)
Go: 5 (50.0%)

2. ЧАСТОТА ЗАПУСКОВ ПО ДНЯМ НЕДЕЛИ:
-----
Tuesday: 4 запусков
Saturday: 2 запусков
Sunday: 2 запусков
Monday: 2 запусков

3. ТОП-5 АКТИВНЫХ ПРОВАЙДЕРОВ:
-----
SpaceX: 3 запусков
China Aerospace Science and Technology Corporation: 2 запусков
Isar Aerospace: 1 запусков
Firefly Aerospace: 1 запусков
Gilmour Space Technologies: 1 запусков

4. СРЕДНЯЯ ПРОДОЛЖИТЕЛЬНОСТЬ ОКОН ЗАПУСКА:
-----
Среднее: 2.02 часов
Минимальное: 0.00 часов
Максимальное: 4.66 часов
```

Построение визуализаций на основе отчета.

```

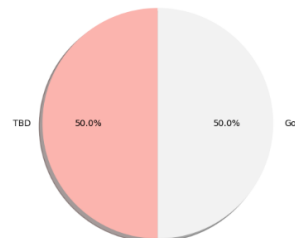
70
71 # Визуализация данных в виде pie-чартов
72 plt.figure(figsize=(18, 12))
73
74 # Цветовые палитры
75 status_colors = plt.cm.Pastell(np.linspace(0, 1, len(status_counts)))
76 weekday_colors = plt.cm.Set3(np.linspace(0, 1, len(weekday_counts)))
77 provider_colors = plt.cm.Accent(np.linspace(0, 1, len(top_providers)))
78
79 # График 1: Распределение по статусам
80 plt.subplot(2, 2, 1)
81 plt.pie(status_counts.values(), labels=status_counts.keys(), autopct='%1.1f%%',
82         colors=status_colors, startangle=90, shadow=True)
83 plt.title('Распределение по статусам', fontsize=12, pad=20)
84
85 # График 2: Дни недели
86 plt.subplot(2, 2, 2)
87 plt.pie(weekday_counts.values(), labels=weekday_counts.keys(), autopct='%1.1f%%',
88         colors=weekday_colors, startangle=90, wedgeprops={'edgecolor': 'white'})
89 plt.title('Запуски по дням недели', fontsize=12, pad=20)
90
91 # График 3: Топ провайдеров
92 plt.subplot(2, 2, 3)
93 providers, counts = zip(*top_providers)
94 plt.pie(counts, labels=providers, autopct='%1.1f%%',
95         colors=provider_colors, startangle=45, explode=[0.05]*len(top_providers))
96 plt.title('Топ-5 провайдеров', fontsize=12, pad=20)
97
98 # График 4: Продолжительность окон (оставляем гистограмму, так как это количественные данные)
99 plt.subplot(2, 2, 4)
100 n, bins, patches = plt.hist(window_durations, bins=10, color='skyblue', edgecolor='navy')
101 plt.title('Распределение продолжительности окон', fontsize=12)
102 plt.xlabel('Часы')
103 plt.ylabel('Количество запусков')
104
105 # Добавляем процентные значения к гистограмме
106 total = len(window_durations)
107 for i in range(len(patches)):
108     plt.text(bins[i], n[i]+0.5, f'{n[i]/total:.1%}',
109             ha='center', va='bottom')
110
111 plt.tight_layout(pad=3.0)
112 plt.savefig('launches_pie_analysis.png', dpi=300, bbox_inches='tight')
113 plt.show()

```

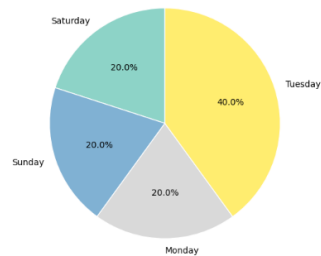
SpaceX: 3 запусков
 China Aerospace Science and Technology Corporation: 2 запусков
 Isar Aerospace: 1 запуск
 Firefly Aerospace: 1 запуск
 Gilmour Space Technologies: 1 запуск
 4. средняя продолжительность окон запуска:

 Среднее: 2.82 часов
 Минимальное: 0.88 часов
 Максимальное: 4.66 часов

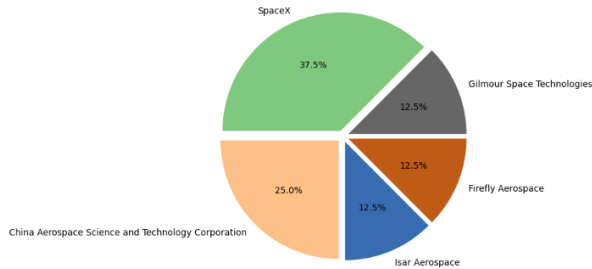
Распределение по статусам



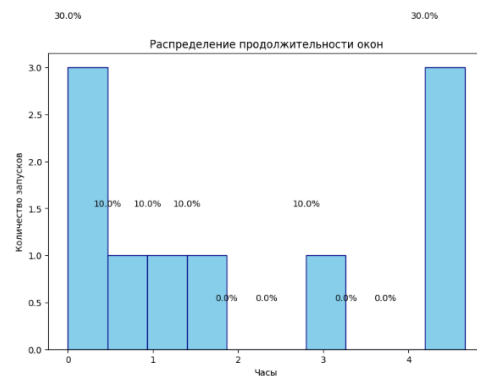
Запуски по дням недели



Топ-5 провайдеров



Распределение продолжительности окон



Таким образом:

- Успешность и неуспешность запуска ракет распределяется в разрезе 50/50.
- Самый частый день недели по запускам Четверг (приходится на него 4 запуска).
- Самый активный провайдер SpaceX (3 запуска).
- Средняя продолжительность окон запуска около 2 часов.

Задание 2

Оценить необходимость использования BashOperator и PythonOperator для таких задач.

DAG использует оба оператора:

- BashOperator для загрузки JSON через curl
- PythonOperator для обработки JSON и загрузки изображений

Это рациональное разделение:

- BashOperator идеален для простых HTTP-запросов (curl/wget)
- PythonOperator необходим для сложной логики (парсинг JSON, обработка изображений)

BashOperator подходит для простого скачивания файлов, работой с файловой системой, запуска CLI-утилит.

PythonOperator подходит там, где используется сложная логика: парсинг JSON, анализ данных, работа с Python-библиотеками, использование Xcom для передачи данных между задачами.

Задание 3

Изучить альтернативы для получения данных о стартах ракет с использованием API.

NASA Launch Schedule API

Источник: [NASA API Portal](#)

Особенности:

- ✓ Официальные данные NASA о предстоящих запусках
- ✓ Бесплатный доступ (до 1000 запросов/час)
- ✓ Включает миссии SpaceX, ULA, Rocket Lab (партнёрские запуски)

SpaceX API

Источник: [GitHub / r-spacex/SpaceX-API](https://github.com/r-spacex/SpaceX-API)

Особенности:

- ✓ Полные данные о запусках SpaceX (Falcon 9, Starship)
- ✓ История всех миссий с 2006 года
- ✓ Бесплатно и без ограничений

Launch Library 2 (альтернатива The Space Devs)

Источник: [Launch Library 2](https://github.com/spaceup/launch_library)

Особенности:

- ✓ Более стабильная версия API (v2)
- ✓ Подробные данные о ракетах, миссиях, площадках
- ✓ Бесплатно (но требуется API-ключ)

ВЫВОД:

Использование данного решения на основе AirFlow позволит компании операционную эффективность вследствие автоматизации процесса (загрузки, обработки, логгирования), также предоставляет возможность масштабировать данный процесс, подключая новые источники данных. Процесс стандартизован, что позволяет гарантированно выполнять по расписанию.

Компания будет получать улучшенное качество данных за счет логгирования, сохранения данных по шаблону (папки, наименования файлов). Вследствие чего эти данные будут использоваться для аналитики и предсказания будущих запусков.