

**Руководство Администратора ETL
Системы бизнес-анализа BI (СБА)**

10-01-2023

Сокращение	Описание
БД	База данных
ЕХД	Единое хранилище данных
Dag, даг	Directed Acyclic Graph это основной элемент Airflow, собирающий задачи вместе, организованный с помощью зависимостей и взаимосвязей, чтобы указать, как они должны выполняться.
ETL	трёхэтапный процесс управления данными, который означает «извлечение, преобразование, загрузка»
ПО	Программное обеспечение

Перечень используемых терминов	Описание
Администратор	Работник, должностные обязанности которого заключаются в обеспечении штатной работы Системы, проведение аварийных работ.

Основной функционал ЕХД.

Целевая аудитория.

Настоящее руководство администратора предназначено для пользователей программного обеспечения – программы ЭВМ – «Система бизнес-анализа ВІ» (далее – СБА).

Назначение документа.

Документ содержит описание последовательности действий в основных операциях, которые должны производить администраторы СБА.

Необходимая подготовка специалиста по администрированию СБА.

Основные компетенции Администратора СБА:

- навыки работы с ОС Linux на уровне продвинутого пользователя;
- навыки работы с Docker;
- навыки работы с БД Postgres на уровне продвинутого пользователя;
- навыки работы с Apache Airflow на уровне продвинутого пользователя.

Авторизация СБА

1) Авторизация в БД

Для входа в ЕХД СБА необходимо создать подключение к БД через SQL-клиент.

Параметры подключения к БД:

Настройки соединения

Свойства соединения с PostgreSQL



Настройки соединения PostgreSQL

Server

Connect by: Host URL

URL: jdbc:postgresql://172.20.200.51:5432/RA_DW_DEV

Хост: 172.20.200.51 Порт: 5432

База данных: RA_DW_DEV

Аутентификация

Аутентификация: Database Native

Пользователь:

Пароль: Сохранять пароль

Advanced

Роль сессии:

Локальный клиент: PostgreSQL

Учетная запись предоставляется по требованию.

2) Авторизация в интерфейсе управления процессами загрузки и обработки данных Airflow.

Для входа в интерфейс управления процессами необходимо:

- 1) В браузере перейти по адресу: <http://172.20.200.41:8080/login/>
- 2) Учетные данные Администратору предоставляются дополнительно.

Описание установки и настройки инструментов.

Реализация процедур ETL реализована с использованием ПО Apache Airflow v.2.4.

Фреймворк Apache Airflow развернут на виртуальной машине под управлением операционной системы Ubuntu Server. Установка ПО сделана с использованием Docker контейнеризации.

Apache Airflow развернут на машине *172.20.200.41* вход на сервер осуществляется через TTY консоль, SSH протокол. Пользователь *airflow*.

Для работы с контейнерами Airflow необходимо перейти в папку установка ПО

```
cd /opt/airflow
```

```
airflow@ecs-dwh-vm1:/opt/airflow$ ls -l
total 100
-rw-r--r--  1 root root  49634 Oct  3 21:56 airflow.cfg
drwxrwxr-x  7 root root   4096 Jan 22 12:21 dags
-rw-rw-r--  1 root root  10739 Jan 18 23:41 docker-compose.yaml
-rw-r--r--  1 root root  10732 Jan 11 10:07 docker-compose.yaml.bak
-rw-r--r--  1 root root    150 Jan 18 23:39 Dockerfile
drwxr-xr-x  2 root root   4096 Oct 27 22:49 libs
drwxrwxr-x 20 root root   4096 Jan 22 12:38 logs
drwxrwxr-x  2 root root   4096 Sep 25 23:07 plugins
drwxr-xr-x  2 root root   4096 Jan 30 03:00 tmp_files
airflow@ecs-dwh-vm1:/opt/airflow$
```

Проверка статуса сервисов.

```
docker container ls
```

```
airflow@ecs-dwh-vm1:/opt/airflow$ docker container ls
CONTAINER ID   IMAGE                                COMMAND                  CREATED        STATUS              PORTS                               NAMES
e672e2e29299  apache/ra_airflow:1.0.1            "/usr/bin/dumb-init ..." 6 days ago    Up 6 days (unhealthy) 8080/tcp                            airflow-airf
low-triggerer-1
0b7c1f072ef8  apache/ra_airflow:1.0.1            "/usr/bin/dumb-init ..." 6 days ago    Up 6 days (healthy)   0.0.0.0:8080->8080/tcp, :::8080->8080/tcp airflow-airf
low-webserver-1
89b662567f06  apache/ra_airflow:1.0.1            "/usr/bin/dumb-init ..." 6 days ago    Up 6 days (unhealthy) 8080/tcp                            airflow-airf
low-worker-1
e8fb68ad3f6a  apache/ra_airflow:1.0.1            "/usr/bin/dumb-init ..." 6 days ago    Up 6 days (unhealthy) 8080/tcp                            airflow-airf
low-scheduler-1
116c742bf287  postgres:13                         "docker-entrypoint.s..." 6 days ago    Up 6 days (healthy)   5432/tcp                            airflow-post
gres-1
1e73124342f5  redis:latest                         "docker-entrypoint.s..." 6 days ago    Up 6 days (healthy)   6379/tcp                            airflow-redi
s-1
```

Запуск Apache Airflow

```
Docker compose up
```

Остановка Apache Airflow

```
Docker compose down
```

WEB интерфейс Apache Airflow доступен через браузер, по адресу <http://172.20.200.41:8080/home>

В Apache Airflow преднастроены соединения к используемым БД.

List Connection			
Search ▾			
<input type="button" value="+"/> <input type="button" value="Actions ▾"/> <input type="button" value="←"/>			
<input type="checkbox"/>	Conn Id ↓	Conn Type ↓	Description ↓
<input type="checkbox"/>	BITRIX_24_MY_SQL	mysql	
<input type="checkbox"/>	CONN_PREPROD_CRM	postgres	Preprod CRM
<input type="checkbox"/>	CONN_RA_PG_DEV	postgres	рублево-Архангельское PasgreSQL Dev

Так же установлен ряд служебных переменных (Variables).

List Variable			
Search ▾			
<input type="button" value="+"/> <input type="button" value="Actions ▾"/> <input type="button" value="←"/>			
<input type="checkbox"/>	Key ↓	Val ↓	Description ↓
<input type="checkbox"/>	B24	{"CONNECTION_ID": "BITRIX_2..."	Bitrix24
<input type="checkbox"/>	CRM	{"CONNECTION_ID": "CONN_P..."	Параметры источника CRM
<input type="checkbox"/>	DWH	{"CONNECTION_ID": "CONN_R..."	Параметры DWH
<input type="checkbox"/>	DWH_CONNECTION_ID	CONN_RA_PG_DEV	DWH connection
<input type="checkbox"/>	PROFIT_URL	https://pb12354.profitbase.ru/expo...	Ссылка на отчет "Profit"

Для загрузки данных написан ряд Airflow DAGs, код DAGs расположен в /opt/airflow/dags.

Переиспользуемый в DAGs код вынесен в отдельные python модули /opt/airflow/dags/lib.

Конфигурационные параметры для загрузки данных расположены в yaml файлах.

/opt/airflow/dags/config

Конфигурация содержит данные об источнике, целевой таблице, названия целевой схемы, [перечень целевых столбцов], флаг необходимости очистки

таблицы перед вставкой, словарь для подстановок source-target (преобразования типов, применение функций и т.д.).

Пример конфигурации для ETL

```
b24_deal:
  source: DWH
  src_table: stage.b24_deal
  dest_schema: facts
  dest_fields:
  dest_truncate_yn: n
  substs:
    ddu_payment_method: ddu_payment_method::int
```

Процессы загрузки запускаются автоматически по расписанию начиная с 0:00 (UTC)

DAGs

All 57 Active 13 Paused 44

DAG	Owner	Runs	Schedule	Last Run
<input type="checkbox"/> FCT_INCR_CRM_THREAD_1	airflow	81 / 31	@daily	2023-01-30, 21:29:15
<input type="checkbox"/> FCT_INIT_B24_DICTS	airflow	46 / 4	@daily	2023-01-29, 03:00:00
<input type="checkbox"/> FCT_INIT_B24_SMART	airflow	45 / 4	0 2 * * *	2023-01-29, 05:00:00
<input type="checkbox"/> FCT_INIT_B24_SMART_128	airflow	9 / 2	0 2 * * *	2023-01-29, 05:00:00
<input type="checkbox"/> FCT_INIT_CRM_COMMON	airflow	39 / 51	@daily	2023-01-30, 21:34:03
<input type="checkbox"/> FCT_INIT_CRM_DICTS	airflow	45 / 45	@daily	2023-01-29, 03:00:00
<input type="checkbox"/> FCT_INIT_PROFIT	airflow	18 / 35	@daily	2023-01-29, 03:00:00
<input type="checkbox"/> STG_INCR_B24_THREAD_1	airflow	34 / 54	@daily	2023-01-30, 21:26:05
<input type="checkbox"/> STG_INCR_CRM_THREAD_1	airflow	31 / 80	@daily	2023-01-29, 03:00:00
<input type="checkbox"/> STG_INIT_B24_DICTS	airflow	46 / 9	@daily	2023-01-29, 03:00:00
<input type="checkbox"/> STG_INIT_B24_SMART	airflow	41 / 11	@daily	2023-01-29, 03:00:00
<input type="checkbox"/> STG_INIT_B24_SMART_128	airflow	9 / 0	@daily	2023-01-29, 03:00:00
<input type="checkbox"/> STG_INIT_CRM_DICTS	airflow	78 / 34	@daily	2023-01-29, 03:00:00

Результаты работы DAGs доступны через веб-интерфейс. Доступны как общая информация о статусе загрузки, так и детальная информация (logs)

Более подробную информацию о работе с Apache Airflow можно получить в официальной документации к проекту <https://airflow.apache.org/docs/>

Администрирование пользователей СБА.

Создание пользователя:

Для создания нового пользователя в Системе требуется войти в БД под учетной записью Администратора и выполнить скрипт создания пользователя указав имя и пароль. (пример маски имени пользователя: Фамилия_ИО/Belov_EA)

```
CREATE USER username WITH PASSWORD '*****';
```

Редактировании пользователя:

Для редактирования пользователя требуется войти в БД и выполнить скрипт указав имя пользователя и прописав дополнительные модификаторы в поле «OPTION».

```
ALTER USER username WITH [OPTION];
```

Блокировка пользователя:

Для блокировки пользователя требуется войти в БД и выполнить скрипт, указав имя пользователя.

```
ALTER USER username NOLOGON;
```

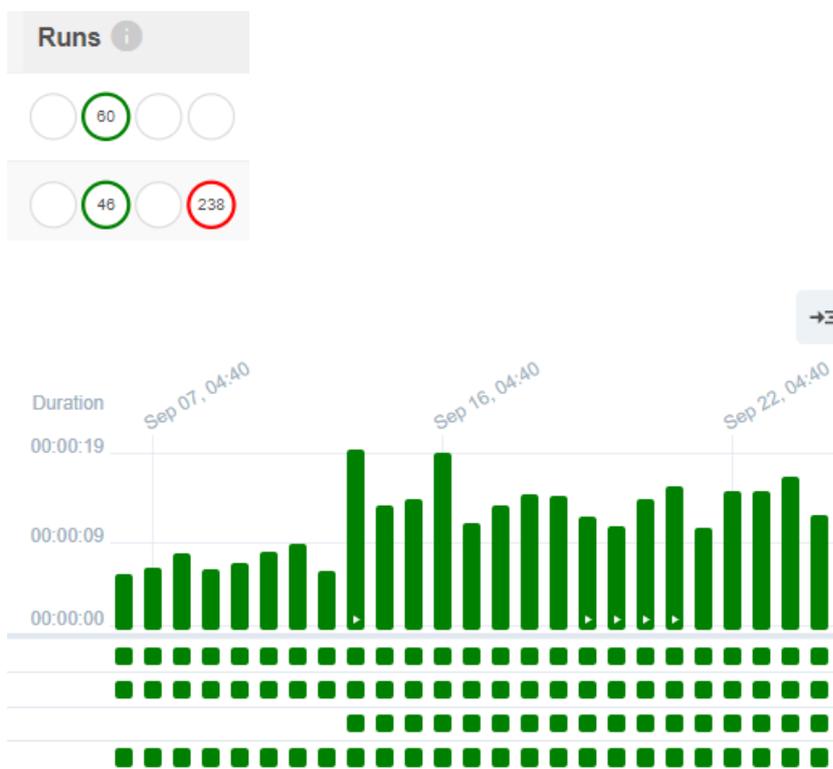
Администрирование процессов загрузки и обработки данных Airflow

Включение/выключение процесса.

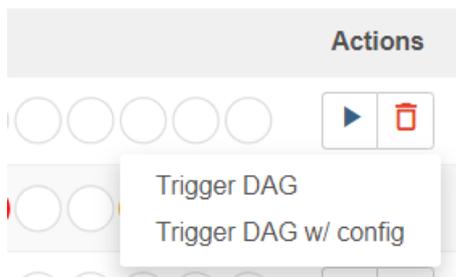
Для работы в Системе необходимо авторизоваться в интерфейсе. В главном меню нажать кнопку DAG. Нажать на кнопку:



Контроль успешности запусков DAG осуществляется через поле RUNS



Для ручного старта отработки DAG, необходимо нажать на стрелочку и выбрать Trigger DAG



Для понимания структуры DAG, можно перейти на следующее окно.

