

> Kohchekt > 1 ypok > Airflow

- > Планировщики задач
- > Cron
- > Airflow
- > Граф и DAG

> Планировщики задач

Так или иначе, нам часто приходится выполнять наши задачи по некоторому расписанию, например:

- 1. Отправка менеджеру "свежего" отчёта по основным метрикам
- 2. Обновление базы данных
- 3. Даже банальные git pull и git commit можно делать по расписанию

Очевидно, что каждая из этих задач может звучать тривиально, но на деле это может быть большое количество команд, которые необходимо выполнять последовательно раз в день, а может и раз в 5 минут.

Да, вполне можно написать скрипты, исполнение которых приведет к желаемому результату, но даже простое исполнение скрипта один раз в определенный интервал может стать проблемой, например, в меру нашей забывчивости.

Планировщики задач способны избавить нас от выполнения "рутинных" и однообразных задач по расписанию.

Сам по себе, **планировщик задач** - это некоторая утилита, чаще всего с довольно тривиальным синтаксисом, способная раз в определенное время (или, например, по наступлению определенного состояния) выполнять установленную команду. Наиболее популярный сейчас планировщик - **Cron**.

> Cron

Cron (Command Run ON) - это системная утилита в операционных системах семейства Unix, которая используется для выполнения определенных действий по расписанию.

В файл **crontab** помещаются инструкции, написанные на специальном синтаксисе, понятном **cron**.

Для описания инструкции Вам достаточно указать расписание, по которому следует выполнять команду, а затем саму команду.

Расписание указывается в виде следующего выражения:

```
минута час день(месяца) месяц день(недели)
```

Каждое значение может быть описано двумя цифрами или заданным шаблоном. День недели также можно описывать английскими сокращениями (MON, TUE, WED, THU, FRI, SAT, SUN).

Шаблоны, которые могут быть использованы в выражении:

- любое значение
- , несколько значений
- диапазон значений
- оператор который позволяет указать значения, которые будут повторяться в течение определенного интервала между ними. К примеру, если в поле минута указано */3, задача будет исполняться каждую третью минуту

Синтаксис указания расписания, по началу, может казаться довольно громоздким и непонятным, но на деле - это один из самых удобных способов. Быстро и удобно составить необходимое расписание может помочь <u>сайт</u>.

Далее, указывается исполняемая команда, тут уже всё ограничивается исключительно Вашим полетом фантазии, это может быть почти любая команда.

Пример описываемой инструкции:

/3 * * 5 * echo 'Three minutes passed in May' - каждые три минуты в 5 месяце (мае) будет выводится сообщение о том, что прошло три минуты.

Это также может быть выполнение python-скрипта.

С помощью cron мы можем выполнять даже последовательность инструкций. Например, каждый день, в 12 часов, мы хотим подготавливать отчет, с помощью скрипта **prepare.py** через 5 минут (12:05), мы хотим отправлять наш отчёт и для этого у нас есть специальный скрипт **send.py**, для этого можем подготовить инструкцию:

00 12 * * * python3 prepare.py - каждый день в 12 часов запускается скрипт, который подготавливает наш отчет

о<u>5 12 * * * python</u>3 send.py - каждый день в 12:05 отчет будет отправляться менеджеру

Основными командами, позволяющими Вам работать с crontab, будут:

```
crontab -1 - просмотр всех записанных инструкций crontab -e - редакция инструкций
```

Со всеми доступными командами Вы можете ознакомиться, выполнив man crontab Несмотря на все достоинства Cron, также остается существенное количество **недостатков**:

- Отсутствие какой-либо визуализации. Когда задач становится сильно много, очень сложно понимать что, когда и как выполняется, исключительно по crontab -1 Вам довольно неудобно фильтровать Ваши задачи, строить последовательности их выполнения
- Отсутствие какой-либо интерактивности. Чтобы посмотреть исполняемые файлы, Вам придется довольно долго бегать от crontab -1 до выполняемых

скриптов(хорошо, если они хотя бы в одной директории или структурированны), в попытке создать картину происходящего

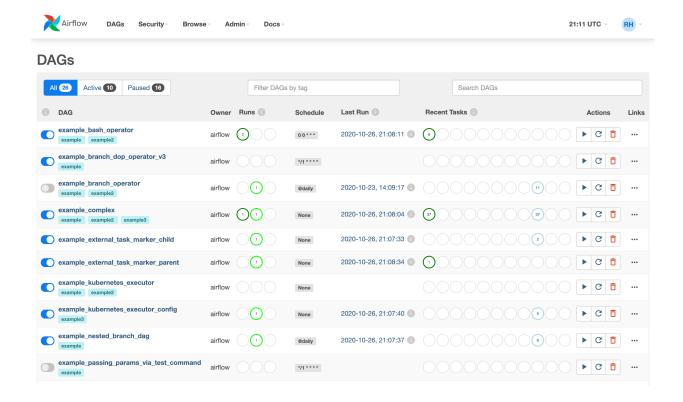
- **Работа в терминале.** Редко, конечно, для аналитика работа в терминале становится серьёзной проблемой, но тем не менее, она требует гораздо больше ресурсов и усилий
- Отсутствие сообщений об ошибках. В случаях, когда что-либо не может выполнится и падает с ошибкой, cron не делает никакого уведомления об этом. Поэтому на выслеживание ошибки может уйти несколько рабочих часов (вспоминаем первые два пункта)

Решить большую часть проблем Cron может Airflow.

> Airflow

Airflow - это библиотека(даже комплекс библиотек), позволяющая очень легко и удобно работать с расписанием и мониторингом выполняемых задач.

Интерфейс Airflow выглядит следующим образом



Здесь мы видим большую подпись **DAGs**:

DAG - это основной юнит работы с Airflow, мы обсудим его подробнее позже, для начала можем считать, что это некоторая глобальная задача, решаемая путем последовательного выполнения более мелких, редуцированных задач.

Чуть детальнее рассмотрим интерфейс:

На главной страничке у нас перечислены все доступные ДАГи, вкладочки **AII**, **Active** и **Paused** позволяют фильтровать ДАГи в соответствии с состоянием их выполнения, у каждого ДАГа сначала стоит переключатель, отвечающий за то, активен ли DAG или нет, затем идет название, владелец, информация о запусках и их состояниях, расписание(в формате Cron), информация по последним более маленьких задачах, составляющих нашу большую и некоторые хот-кеи для работы с ДАГом: запуск мгновенно, перезагрузить и удалить.

Чаще всего, для того, чтобы добавить DAG необходимо просто загрузить его в специальную выделенную папку в git - репозитории. О подробном устройстве и создании самих ДАГов Вы узнаете в следующих лекциях.

> Граф и DAG

Прежде всего, вспомним, что такое граф.

Граф - это некоторая математическая абстракция, которая задается множеством вершин и множеством дуг между ними.

Мы же не будем сильно погружаться в математическую теорию графов, и для нас граф - это описание маленьких задач, которые мы будем называть тасками и описание правил перехода между ними.

(спойлер: DAG - это тоже граф, но немного особенный)

В общем случае, граф выглядит примерно следующим образом



Такой граф называется обыкновенным, здесь мы видим 4 вершины(таски) и две дуги(перехода) между тасками 0 и 1, и тасками 0 и 2. В случае обыкновенного графа наличие перехода между двумя тасками обозначает возможность перехода в обе стороны, то есть, если между 0 и 1 есть переход, это означает, что мы можем переходить из 0 в 1 из 1 в 0.

С точки зрения математики всё, конечно, здорово, но на практике мы сталкиваемся со следующей проблемой: наличие двух переходов между двумя вершинами в обе стороны сразу же приводит нас к циклу. Циклом называется последовательность переходов в результате которой мы можем вернутся в начало (т. е. в изначальную таску).

В нашем случае, цикл, например:

0 -> 1, 1 -> 0, мы вышли из таски 0 и в неё же смогли вернуться.

Почему это проблема?

У нас может не быть явного правила для завершения выполнения DAGa. Помним о том, что DAG - некоторая глобальная задача, а таски - всего лишь более мелкие, необходимы для выполнения глобальной.

На практике:

Представьте, что у нас есть задача: выгрузить данные, подготовить их и выслать менеджеру.

- Таска 1 выгрузка данных
- Таска 0 подготовка и обработка
- Таска 2 отправка данных

Начинаем выполнение нашей глобальной задачи из таски 1:

Мы выгрузили данные, перешли в таску 0, подготовили и обработали данные, но, вместо того, чтобы просто перейти в таску 2, мы можем вернуться в таск 1 и снова начать выгружать данные. И такое может повторяться сколько угодно раз. Аналогичных примеров можно придумать довольно много.

Поэтому Airflow накладывает некоторые ограничения на загружаемые глобальные задачи:

- 1. Граф глобальной задачи должен быть ацикличным (не должен содержать ни одного цикла)
- 2. Граф глобальной задачи должен быть направленным (между тасками путь в каждую сторону должен быть указан напрямую)

Объединяем эти правила и переводим на английский язык, получаем

DAG - Directed acyclic graph, направленный, ацикличный граф.

Приведем пример ДАГа, построенного в соответствии с нашей задачей:



Как видим, после выполнения таски 1 мы однозначно попадём в таску 0, а после таски 0 в таску 2, а после выполнения таски 2 наш DAG будет считаться выполненным.

Кстати, присутствие отдельной таски 3, которая никак не используется, не вносит никаких препятствий для выполнения нашего DAG. Таким образом, мы можем "отключать" недоработанные или на данный момент ненужные таски.

Отойдём от теоретической модели, что же на практике в Airflow?

Для каждой глобальной задачи составляется DAG, DAG состоит из подзадач, которые необходимо выполнить для завершения глобальной задачи.

В веб-интерфейсе мы можем наблюдать DAG целиком, а также состояние выполнения каждой конкретной таски, можем устанавливать количество попыток выполнения таски, а также, в случае чего, видеть, что и с какой ошибкой упало.

Пример довольно непростого ДАГа, который может встретиться при работе:

