



Apache Airflow



نرم افزار مدیریت
جریان های کار

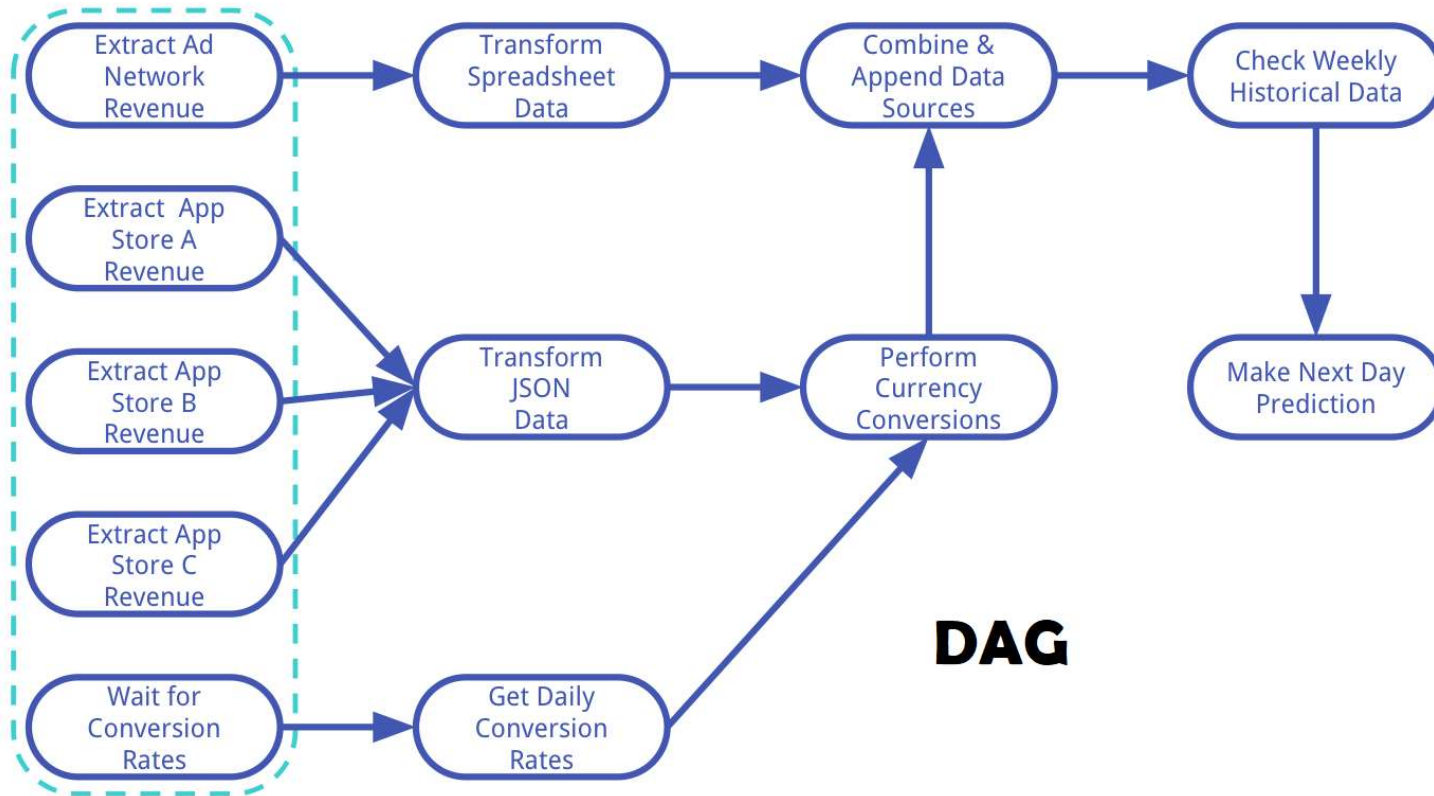
آنچه خواهیم دید

۱. آشنایی با مفاهیم پایه جریان‌های کار
۲. روش سنتی زمان‌بندی کارها (کران جابز)
۳. مفاهیم پایه ایرفلو و معماری آن
۴. بررسی محیط گرافیکی ایرفلو
۵. کارگاه عملی / کار با کران جابز
۶. کارگاه عملی / کار با ایرفلو





مدیریت جریان کار





ویژگیهای مورد نیاز WMS

- محیط گرافیکی (در کنار خط فرمان قدرتمند)
- انعطاف در تعریف خطوط پردازش داده
- SLA
- Backfill/Catch Up
- Metrics/Logs
- Alerting/ Tunable Retries
- قابل توسعه
- مقیاس پذیر / اجرای موازی تسکها
- مجموعه غنی از عملگرها و توابع آماده
- جامعه کاربری فعال
- زمان بند



بازنگرانی اصلی

Source :
Data Pipelines With Apache Airflow

Name	Originated at	Workflows defined in	Written in	Scheduling	Backfilling	User interface[2]
Airflow	Airbnb	Python	Python	Yes	Yes	Yes
Argo	Applatix	YAML	Go	3rd party[3]		Yes
Azkaban	LinkedIn	YAML	Java	Yes	No	Yes
Conductor	Netflix	JSON	Java	No		Yes
Luigi	Spotify	Python	Python	No	Yes	Yes
Make		Custom DSL	C	No	No	No
Metaflow	Netflix	Python	Python	No		No
Oozie		XML	Java	Yes	Yes	Yes



Apache Oozie

192.168.56.41:8000/oozie/

Dashboard Workflows Coordinators Bundles

Running

Submission	Status	Name	Progress	Submitter	Created	Last Modified	Run	Id	Action
Fri, 06 Jan 2017 21:20:20	RUNNING	map-reduce-wf	50%	yarn	Fri, 06 Jan 2017 21:20:15	Fri, 06 Jan 2017 21:20:20	0	0000019-170105175308386-oozie-oozi-W	Kill Suspend
Fri, 06 Jan 2017 17:14:52	SUSPENDED	pig-wf	50%	yarn	Fri, 06 Jan 2017 17:13:59	Fri, 06 Jan 2017 17:14:52	0	0000014-170105175308386-oozie-oozi-W	Kill Resume
Fri, 06 Jan 2017 17:12:25	SUSPENDED	pig-wf	50%	yarn	Fri, 06 Jan 2017 17:11:25	Fri, 06 Jan 2017 17:12:25	0	0000013-170105175308386-oozie-oozi-W	Kill Resume
Fri, 06 Jan 2017 00:07:06	SUSPENDED	map-reduce-wf	50%	yarn	Fri, 06 Jan 2017 00:06:36	Fri, 06 Jan 2017 00:07:06	0	0000004-170105175308386-oozie-oozi-W	Kill Resume
Fri, 06 Jan 2017 00:04:50	SUSPENDED	map-reduce-wf	50%	yarn	Fri, 06 Jan 2017 00:04:50	Fri, 06 Jan 2017 00:04:50	0	0000003-170105175308386-oozie-oozi-W	Kill Resume
Thu, 05 Jan 2017 23:52:23	SUSPENDED	map-reduce-wf	50%	hdfs	Thu, 05 Jan 2017 23:51:42	Thu, 05 Jan 2017 23:52:23	0	0000002-170105175308386-oozie-oozi-W	Kill Resume


192.168.56.41:8000/oozie/list_oozie_workflow/0000019-170105175308386-oozie-oozi-W//

Search the web and Windows

20:32 25-02-2017



Linkedin Azkaban

 **Azkaban Local**
\${version} My Local Azkaban

Projects | Scheduling | Executing | History | Flow Trigger Schedule

azkaban ▾

Flow flow_trigger

Schedule / Execute Flow

Project ET / Flow flow_trigger

Graph | Executions | **Flow Triggers** | Summary

Flow Trigger Instance Id	Submitting user	Start Time	End Time	Elapsed	Status	Action
acd4e5b1-21cd-41d1-b4e9-ddd1e4ddd9c7	azkaban	2018-07-31 17:22 00s	-	29 sec	RUNNING	
340392ba-263f-4cba-b578-32130bbcdd4d	azkaban	2018-07-31 17:21 25s	-	1m 3s	RUNNING	
39b6c9fe-5528-491e-8b16-bfc3681d596f	azkaban	2018-07-30 17:28 00s	-	23h 54m 29s	RUNNING	
296e76a1-6b6c-4660-9d89-8e2acd101a79	azkaban	2018-07-30 17:26 47s	2018-07-30 17:28 10s	1m 23s	SUCCEEDED	
6b43f3e3-4722-4d1f-b692-a634c6a8fa41	azkaban	2018-07-30 11:56 00s	2018-07-30 17:28 10s	5h 32m 10s	SUCCEEDED	
0aa7d799-7db1-4bd0-895d-6c4d26d4a286	azkaban	2018-07-30 11:54 00s	2018-07-30 17:28 10s	5h 34m 10s	SUCCEEDED	



Luigi

Luigi Task Status

Task ListDependency GraphWorkersResources

Running

TASK FAMILIES

1 AggregateArtists

28 Streams

PENDING TASKS
0

RUNNING TASKS
0

BATCH RUNNING TA...
0

DONE TASKS
29

FAILED TASKS
0

UPSTREAM FAILURE
0

DISABLED TASKS
0

UPSTREAM DISABLED
0

Displaying **DONE**, tasks of family **Streams**.

Show 10 entries

Filter table: Filter on Server ☐

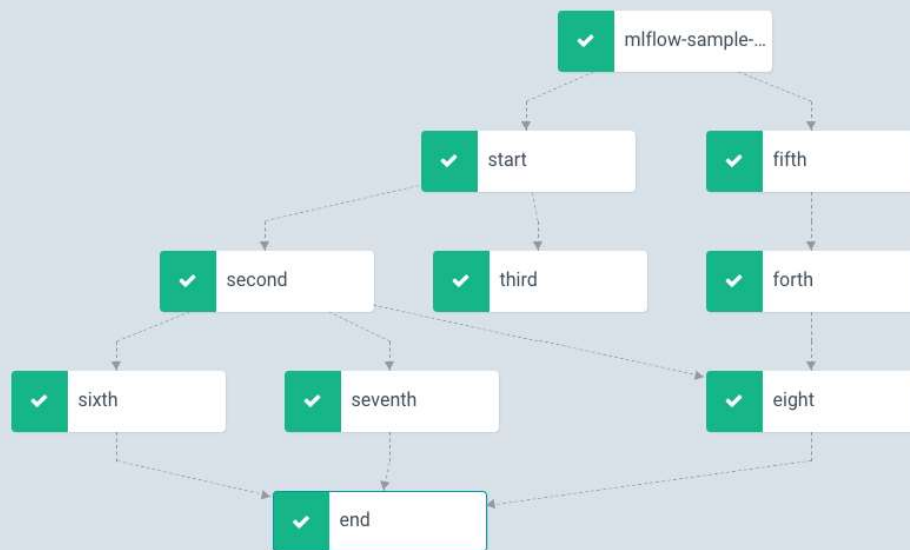
	Name	Details	Priority	Time	Actions
✓ DONE	Streams	date=2018-02-01	0	18/1/2018 10:50:01	
✓ DONE	Streams	date=2018-02-08	0	18/1/2018 10:50:01	
✓ DONE	Streams	date=2018-02-18	0	18/1/2018 10:50:01	
✓ DONE	Streams	date=2018-02-19	0	18/1/2018 10:50:01	
✓ DONE	Streams	date=2018-02-02	0	18/1/2018 10:50:01	



Argo

WORKFLOW DETAILS

Workflows / mlflow-sample-qzxcg2



SUMMARY

CONTAINERS

ARTIFACTS

NAME mlflow-sample-qzxcg2.end

TYPE Pod

PHASE ✓ Succeeded

START TIME 2019-03-03T01:22:26Z

END TIME 2019-03-03T01:22:30Z

DURATION 00:04 min

YAML

LOGS

Parameters

alpha 0.7

ratio 0.4

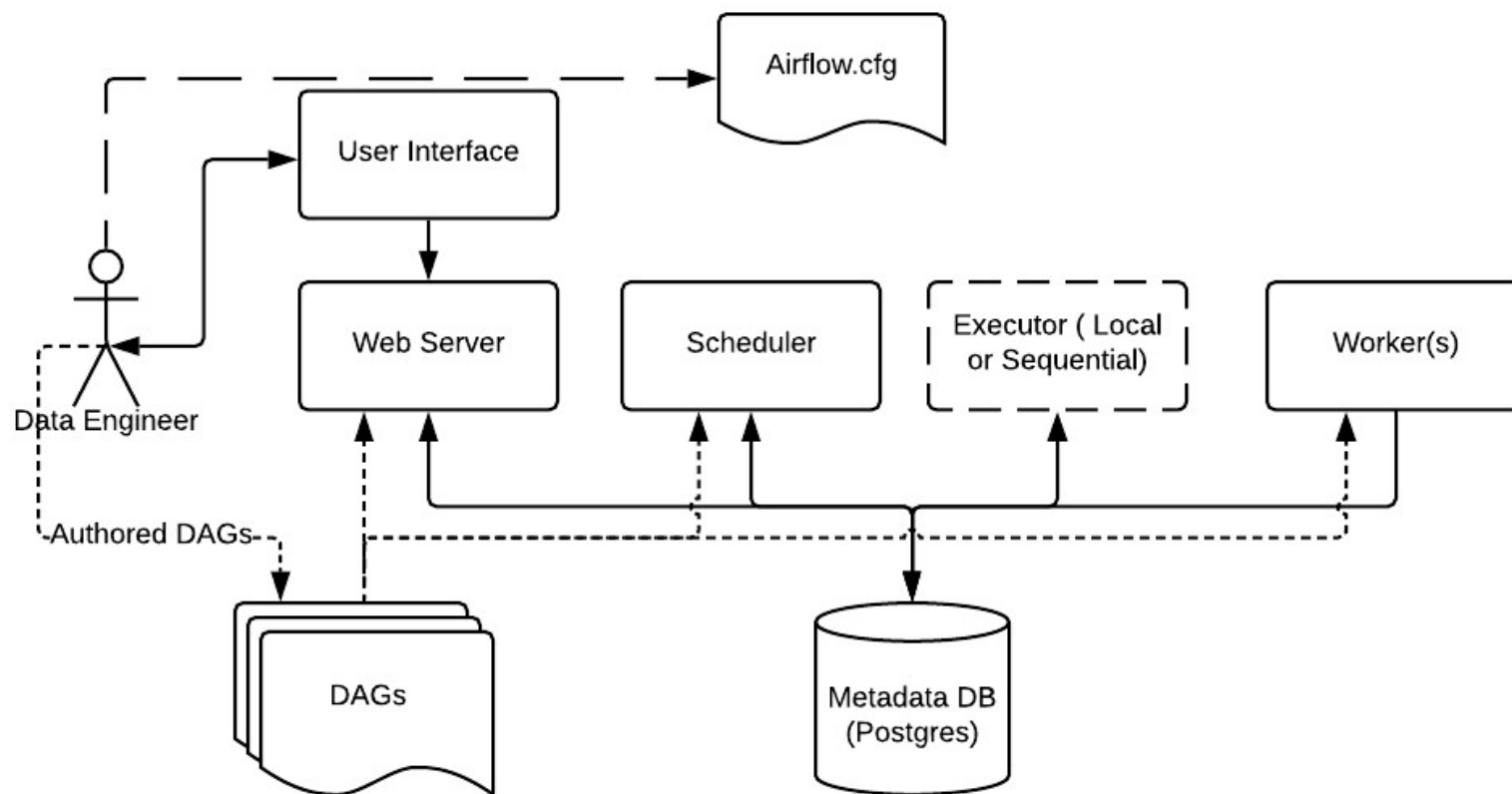


مفاهیم پایه ایرفلو

- › **DAG**: graph of operator usages (=tasks)
- › **Operator**: "Transformation" step
 - › **Sensor**: Operator which polls with frequency / timeout (e.g. LocalFileSensor)
 - › **Executor**: Trigger operation (e.g. HiveOperator, BashOperator, PigOperator, ...)
- › **Task**: Usage of Operator in DAG
 - › **Task Instance**: run of a Task at a point in time
- › **Hook**: Interface to external System (JDBCHook, HTTPHook, ...)



معماری ایرفلو





انواع زیرساخت‌های اجرای تسک‌ها



Apache
Airflow



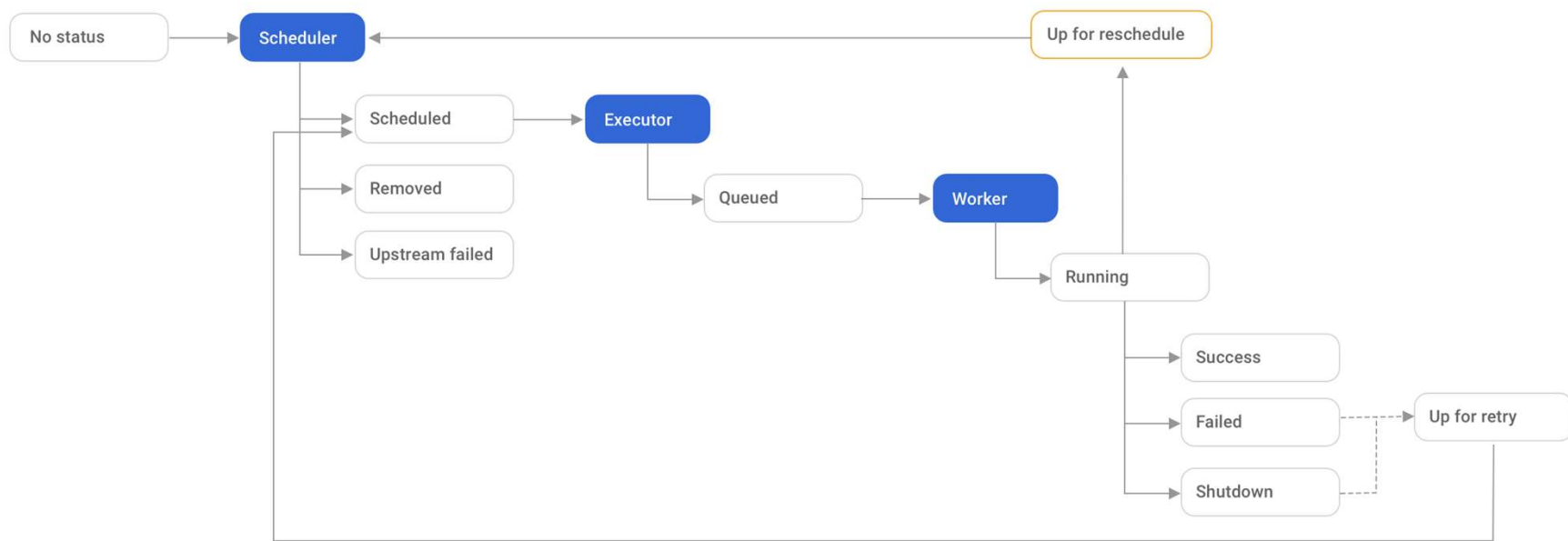
LocalExecutor
SequentialExecutor
DebugExecutor



Apache
MESOS™



چرخه حیات یک تسک

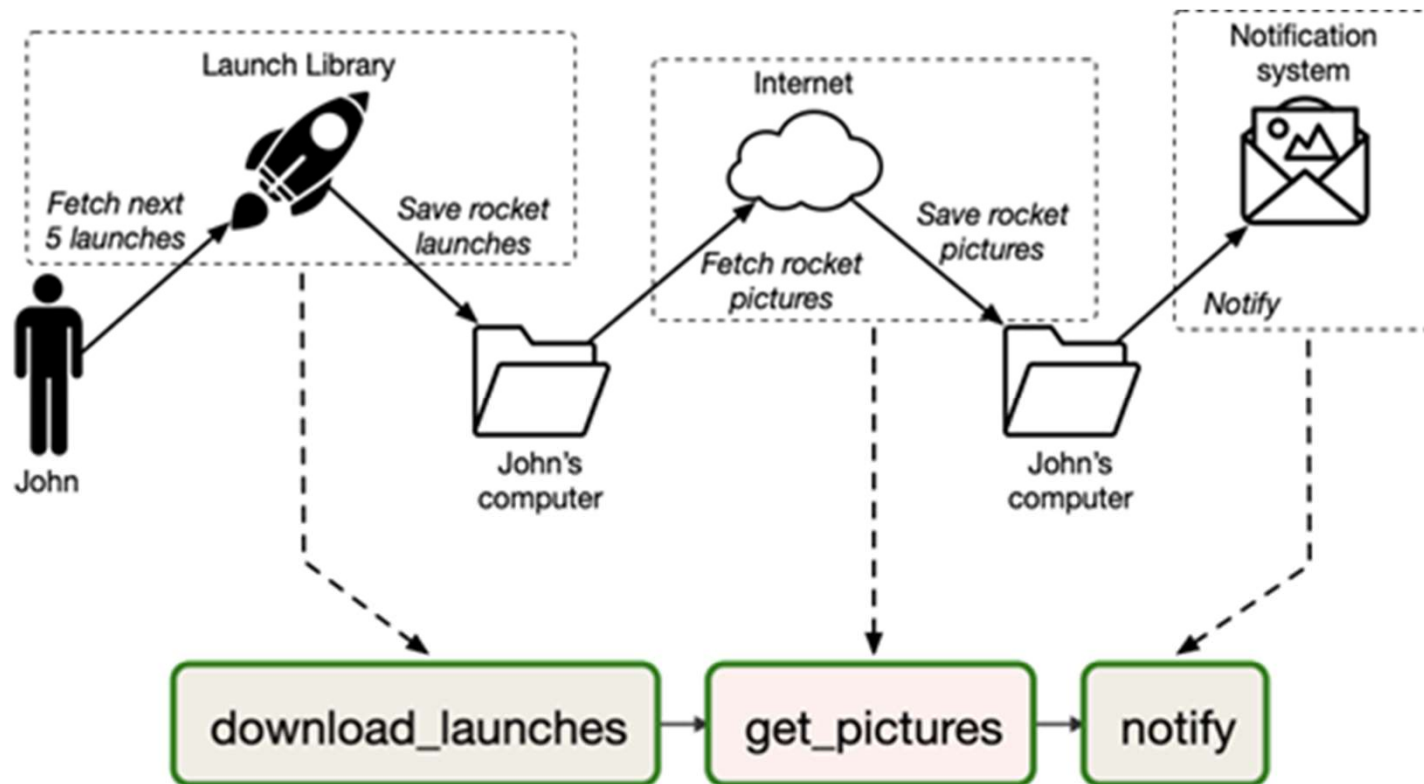


■ Component □ Task stage □ Task stage only for sensor — Stage transition --- Alternative stage transition

■ success ■ running ■ failed ■ skipped ■ upstream_failed ■ up_for_reschedule ■ up_for_retry ■ queued □ no_status



بررسی یک مثال





تعریف جریان کار

Listing 2.2 DAG for downloading and processing rocket launch data

```
1 import json
2 import pathlib
3
4 import airflow
5 import requests
6 from airflow import DAG
7 from airflow.operators.bash_operator import BashOperator
8 from airflow.operators.python_operator import PythonOperator
9
10 dag = DAG(
11     dag_id="download_rocket_launches",
12     start_date=airflow.utils.dates.days_ago(14),
13     schedule_interval=None,
14 )
15
16 download_launches = BashOperator(
17     task_id="download_launches",
18     bash_command="curl -o /tmp/launches.json 'https://launchlibrary.net/1.4/launch",
19     dag=dag,
20 )
```

A

B

C

D

E

F



تعریف جریان کار


```
40 get_pictures = PythonOperator(  
41     task_id="get_pictures",  
42     python_callable=_get_pictures,  
43     dag=dag,  
44 )  
45  
46 notify = BashOperator(  
47     task_id="notify",  
48     bash_command='echo "There are now $(ls /tmp/images/ | wc -l) images."',  
49     dag=dag,  
50 )  
51  
52 download_launches >> get_pictures >> notify
```

H

I



بررسی محیط ایرفلو

 Airflow

DAGs

Data Profiling ▾

Browse ▾

Admin ▾

Docs ▾

About ▾

2019-06-16 19:47:54 UTC

DAGs

Search:

Showing 1 to 1 of 1 entries

«

<

1

>

»

[Hide Paused DAGs](#)



Graph View

Off DAG: download_rocket_launches schedule: None

Graph View Tree View Task Duration Task Tries Landing Times Gantt Details Code Trigger DAG Refresh Delete

None Base date: 2019-06-17 19:04:15 Number of runs: 25 Run: Layout: Left->Right Go Search for...

BashOperator PythonOperator success running failed skipped up_for_reschedule up_for_retry queued no_status

download_launches → get_pictures → notify

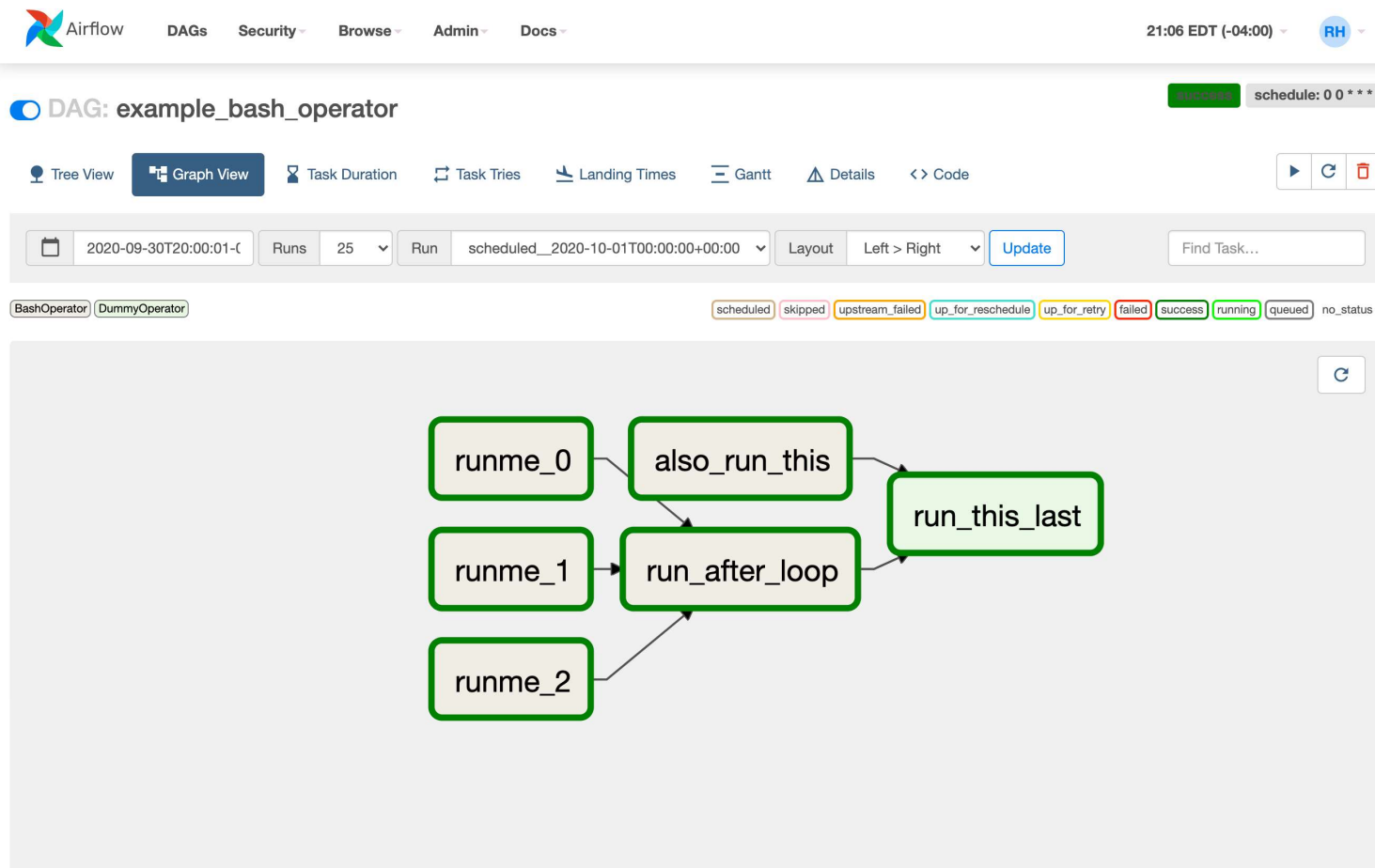
Operator types in DAG

DAG structure

State legend



Graph View





Task Details

notify on 2019-06-14T00:00:00+00:00

Task Instance Details

Rendered

Task Instances

View Log

Download Log (by attempts):

1

Run

Ignore All Deps

Ignore Task State

Ignore Task Deps

Clear

Past

Future

Upstream

Downstream

Recursive

Mark Failed

Past

Future

Upstream

Downstream

Mark Success

Past

Future

Upstream

Downstream



Task Details

```
*** Reading local file: /root/airflow/logs/download_rocket_launches/notify/2019-06-18T19:06:28.102026+00:00/1.log
[2019-06-18 19:06:58,698] {__init__.py:1139} INFO - Dependencies all met for <TaskInstance: download_rocket_launches.notify 2019-06-18T19:06:28.102026+00:00>
[2019-06-18 19:06:58,705] {__init__.py:1139} INFO - Dependencies all met for <TaskInstance: download_rocket_launches.notify 2019-06-18T19:06:28.102026+00:00>
[2019-06-18 19:06:58,705] {__init__.py:1353} INFO -

-----

[2019-06-18 19:06:58,705] {__init__.py:1354} INFO - Starting attempt 1 of 1
[2019-06-18 19:06:58,705] {__init__.py:1355} INFO -

-----

[2019-06-18 19:06:58,715] {__init__.py:1374} INFO - Executing <Task(BashOperator): notify> on 2019-06-18T19:06:28.102026+00:00
[2019-06-18 19:06:58,716] {base_task_runner.py:119} INFO - Running: ['airflow', 'run', 'download_rocket_launches', 'notify', '2019-06-18T19:06:28.102026+00:00']
[2019-06-18 19:06:59,871] {base_task_runner.py:101} INFO - Job 85: Subtask notify [2019-06-18 19:06:59,871] {__init__.py:51} INFO - Using executor LocalExecutor
[2019-06-18 19:07:00,126] {base_task_runner.py:101} INFO - Job 85: Subtask notify [2019-06-18 19:07:00,126] {__init__.py:305} INFO - Fetching task instance
[2019-06-18 19:07:00,153] {base_task_runner.py:101} INFO - Job 85: Subtask notify [2019-06-18 19:07:00,152] {cli.py:517} INFO - Running task
[2019-06-18 19:07:00,165] {bash_operator.py:81} INFO - Tmp dir root location:
/tmp
[2019-06-18 19:07:00,165] {bash_operator.py:90} INFO - Exporting the following env vars:
AIRFLOW_CTX_DAG_ID=download_rocket_launches
AIRFLOW_CTX_TASK_ID=notify
AIRFLOW_CTX_EXECUTION_DATE=2019-06-18T19:06:28.102026+00:00
AIRFLOW_CTX_DAG_RUN_ID=manual__2019-06-18T19:06:28.102026+00:00
[2019-06-18 19:07:00,165] {bash_operator.py:104} INFO - Temporary script location: /tmp/airflowtmpdhnvdwi/notify3obkuldp
[2019-06-18 19:07:00,165] {bash_operator.py:114} INFO - Running command: P484#yIS1 "There are now $(ls /tmp/images/ | wc -l) images."
[2019-06-18 19:07:00,173] {bash_operator.py:123} INFO - Output:
[2019-06-18 19:07:00,177] {bash_operator.py:127} INFO - There are now 5 images.
[2019-06-18 19:07:00,177] {bash_operator.py:131} INFO - Command exited with return code 0
[2019-06-18 19:07:03,692] {logging_mixin.py:95} INFO - [2019-06-18 19:07:03,692] {jobs.py:2562} INFO - Task exited with return code 0
```




زمان بندی

preset	meaning	cron
None	Don't schedule, use for exclusively "externally triggered" DAGs	
@once	Schedule once and only once	
@hourly	Run once an hour at the beginning of the hour	0 * * * *
@daily	Run once a day at midnight	0 0 * * *
@weekly	Run once a week at midnight on Sunday morning	0 0 * * 0
@monthly	Run once a month at midnight of the first day of the month	0 0 1 * *
@yearly	Run once a year at midnight of January 1	0 0 1 1 *



زمان‌بندی – کران جابز


```
1 # _____ minute (0 - 59)
2 # |_____ hour (0 - 23)
3 # | |_____ day of the month (1 - 31)
4 # | | |_____ month (1 - 12)
5 # | | | |_____ day of the week (0 - 6) (Sunday to Saturday;
6 # | | | | 7 is also Sunday on some systems)
7 # * * * * *
```

copy

```
0 2 * * * /bin/sh backup.sh
* * * * * /scripts/script.sh
*/10 * * * * /scripts/monitor.sh
0 */4 * * * /scripts/script.sh
0 4,17 * * sun,mon /scripts/script.sh
* * * * * sleep 30; /scripts/script.sh
```



متغیرها

Airflow

DAGs

Data Profiling ▾

Browse ▾

Admin ▾

Docs ▾

Variables



















List (9)

Create

Add Filter ▾

With selected ▾

Search

<input type="checkbox"/>		Key	Val
<input type="checkbox"/>	 	secret_password	*****
<input type="checkbox"/>	 	not_so_hidden	test value
<input type="checkbox"/>	 	secret	*****
<input type="checkbox"/>	 	password	*****
<input type="checkbox"/>	 	passwd	*****
<input type="checkbox"/>	 	api_key	*****
<input type="checkbox"/>	 	apikey	*****
<input type="checkbox"/>	 	authorization	*****
<input type="checkbox"/>	 	access_token	*****



XCom

On **DAG: my_dag**

[Graph View](#) [Tree View](#) [Task Duration](#) [Task Tries](#) [Landing Times](#) [Gantt](#) [Details](#) [Code](#) [Refresh](#) [Delete](#)

Task Instance: my_task


[Task Instance Details](#) [Rendered Template](#) [Log](#) [XCom](#)

XCom

Key	Value
start	2019-06-16T17:59:35.260928
end	2019-06-16T17:59:47.587656
queries	['foo', 'bar']
counts	[1]



Connections

 AirFlow

DAGs

Tools ▾

Browse ▾

Admin ▾

Docs ▾

List (4)

Create

With selected ▾

Configuration

Connections

Users

Reload DAGs



Pools





File Modifica Visualizza Cronologia Segnalibri Strumenti Aiuto

Visualize × Admin - Pools - Airflow × +

localhost:8080/admin/pool/

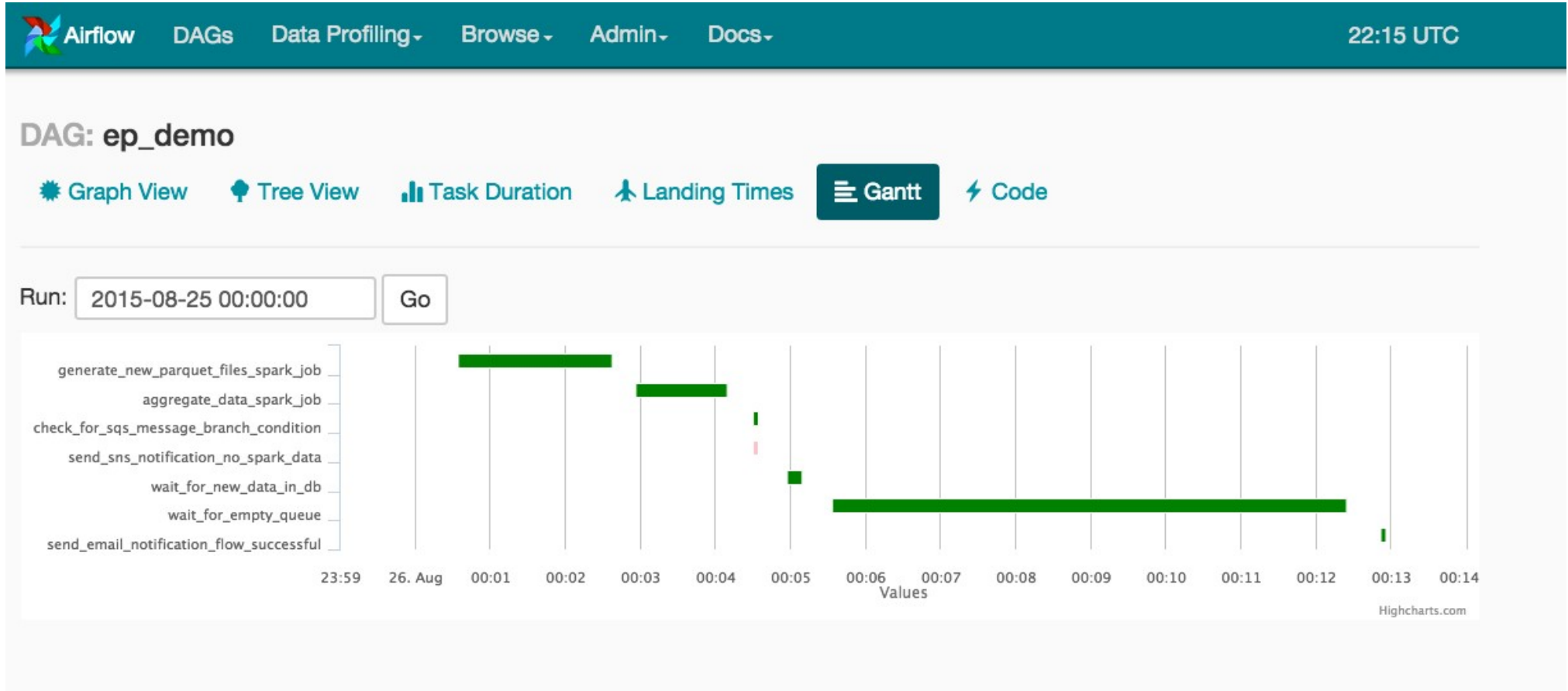
Airflow

List (3) Create With selected

<input type="checkbox"/>		Pool	Slots	Used Slots	Queued Slots
<input type="checkbox"/>	 	updates	50	0	77
<input type="checkbox"/>	 	crawler	1	6	18

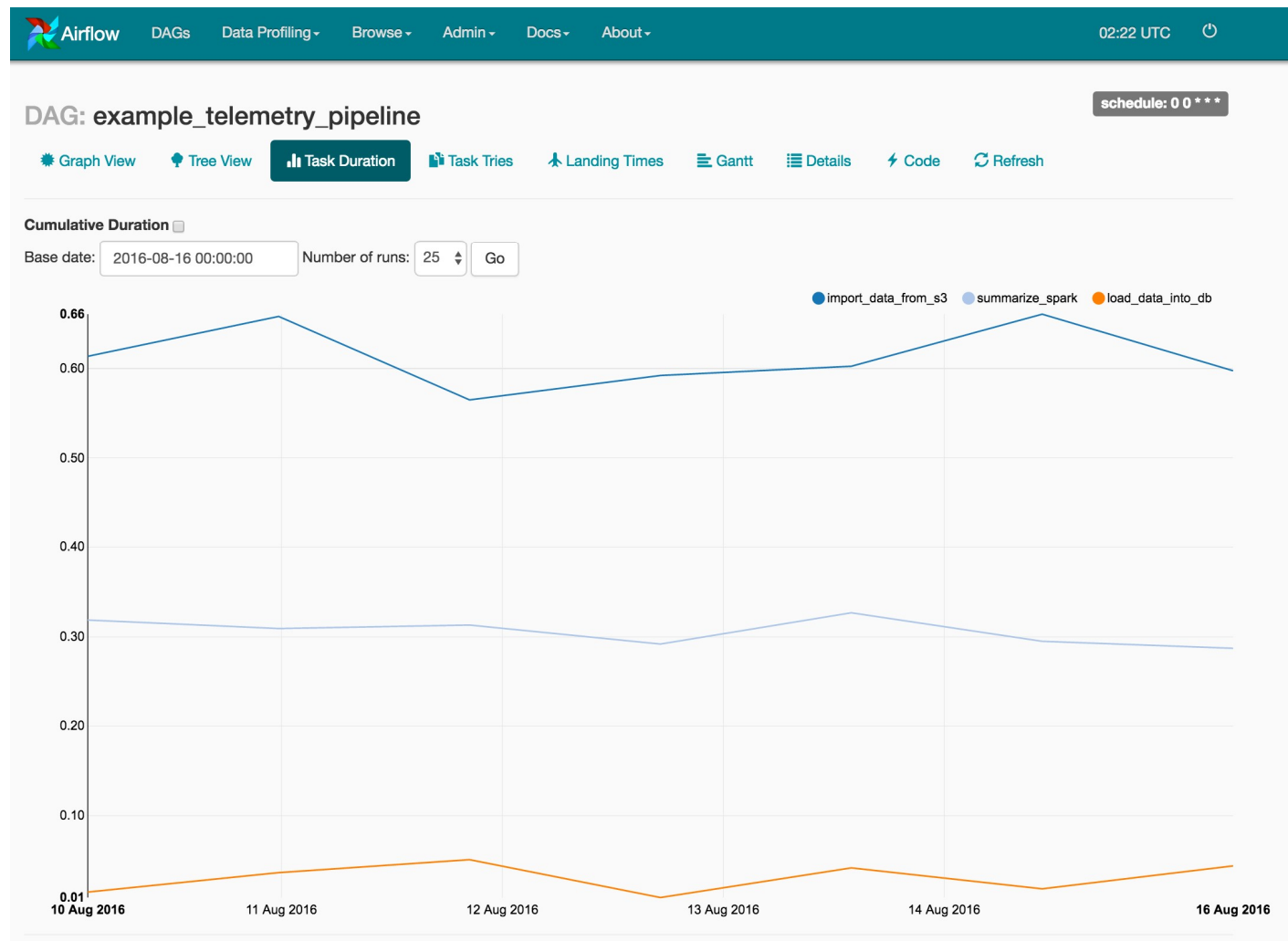


Gantt Chart





Task Duration Diagram





مثال عملی

• جریان کار

- دانلود فایل اکسل داد و ستد روزانه بورس ایران
- تبدیل فایل به CSV
- حذف فایل اکسل
- تبدیل فایل CSV به پارکت
- حذف فایل CSV

- کران جابز
- ایرفلو

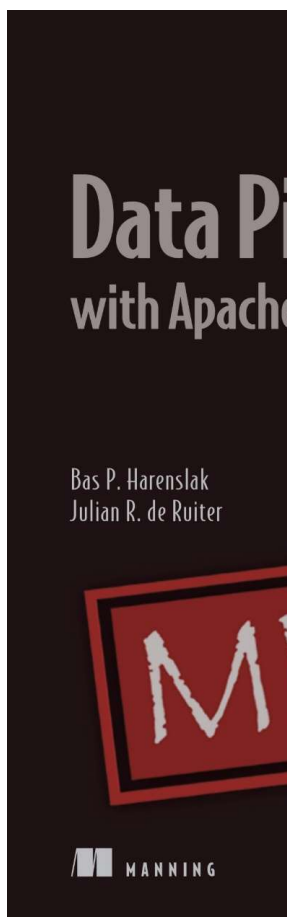


نصب ایرفلو

```
# airflow needs a home, ~/airflow is the default,  
# but you can lay foundation somewhere else if you prefer  
# (optional)  
export AIRFLOW_HOME=~/airflow  
  
# install from pypi using pip  
pip install apache-airflow  
  
# initialize the database  
airflow initdb  
  
# start the web server, default port is 8080  
airflow webserver -p 8080  
  
# start the scheduler  
airflow scheduler  
  
# visit localhost:8080 in the browser and enable the example dag in the home page
```



منابع



- <https://airflow.apache.org/docs/>
- <https://github.com/jghoman/awesome-e-apache-airflow>

