



丁来强 wjo1212@163.com

# Python高效大数据工作流与任务调度

丁来强 (Lai Qiang Ding)



wjo1212



wjo1212@163.com



LaiQiangDing

- About Me

- Farther of a 4 years' boy



- About Me

- Worked for 10+ years.
- @Splunk



- Agenda

- Background
  - Definition
  - Role in Data Infra
- Requirement
  - Problem
  - Challenges
  - Requirement
- Solutions
  - Overview
  - Luigi
  - Airflow
- Demo



- You will learn:

- Role of workflow scheduler for data engineering in ecosystem.
- Challenges and key requirements.
- Solutions and general differences.
- Architecture, design and practices of using Airflow and Luigi **in Python**
- Pitfalls and common patterns in design to use a workflow scheduler

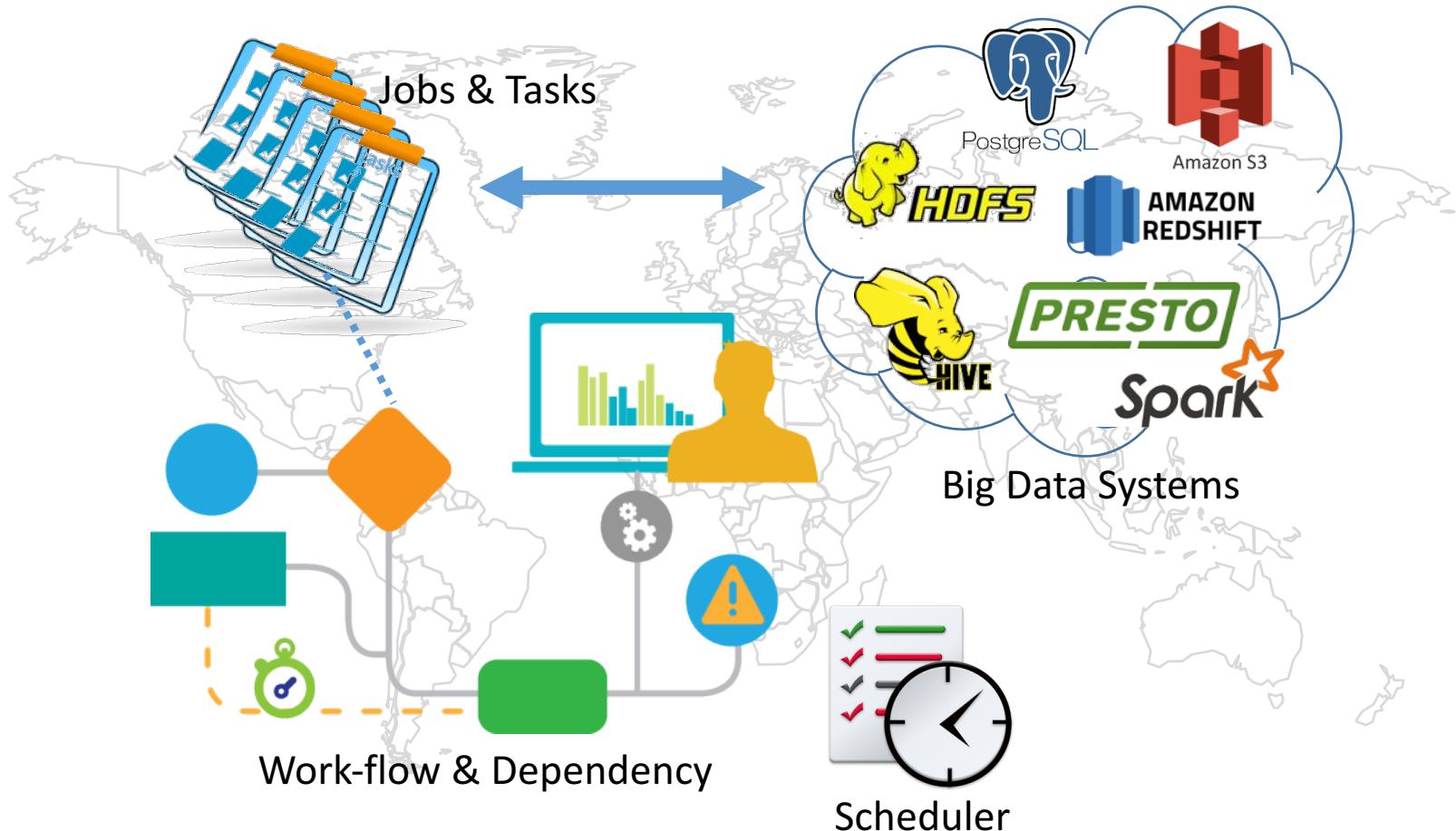


- Definition

## Big Data Workflow Scheduler

Schedule and manage dependencies of workflow of jobs in data infrastructure, mainly used in offline and near-line system.

# • Big Data Work-flow Scheduler



- Different with below categories:

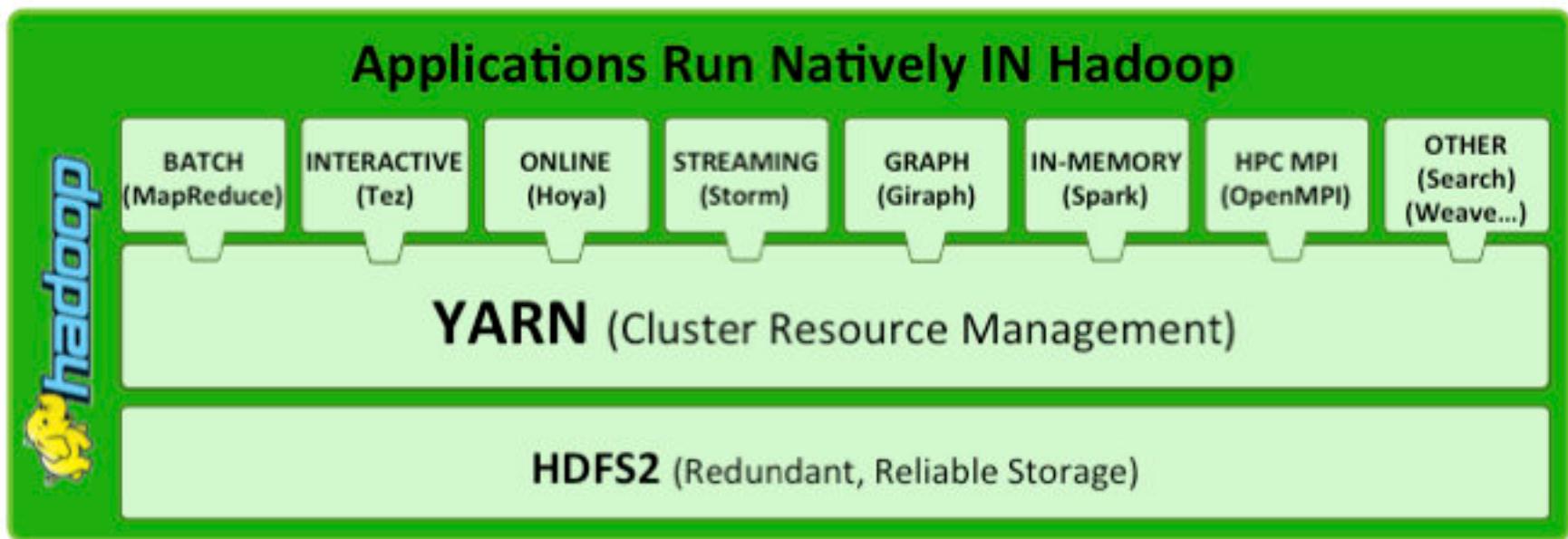
- BPM
  - Like Activiti
- Middleware workflow & SOA
  - Like AWS Simple Workflow
- Pure Data Driven Pipeline/API for Development
  - Like Apache Crunch, Apache Cascading, AWS Data Pipeline, Azure Data Factory
- Pure Streaming Process
  - Like Storm, Spark Streaming



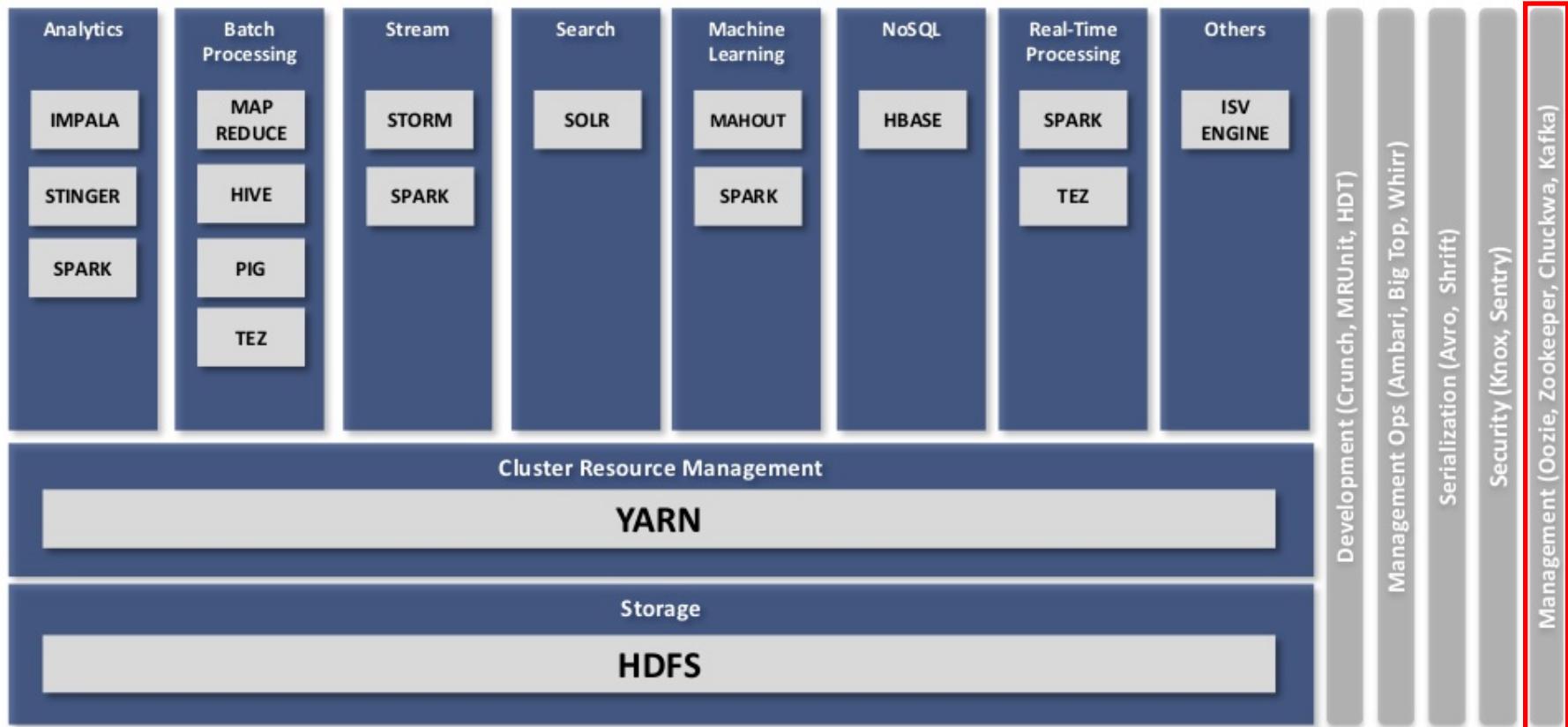
# Role in Data Infra

丁来强 wjo1212@163.com

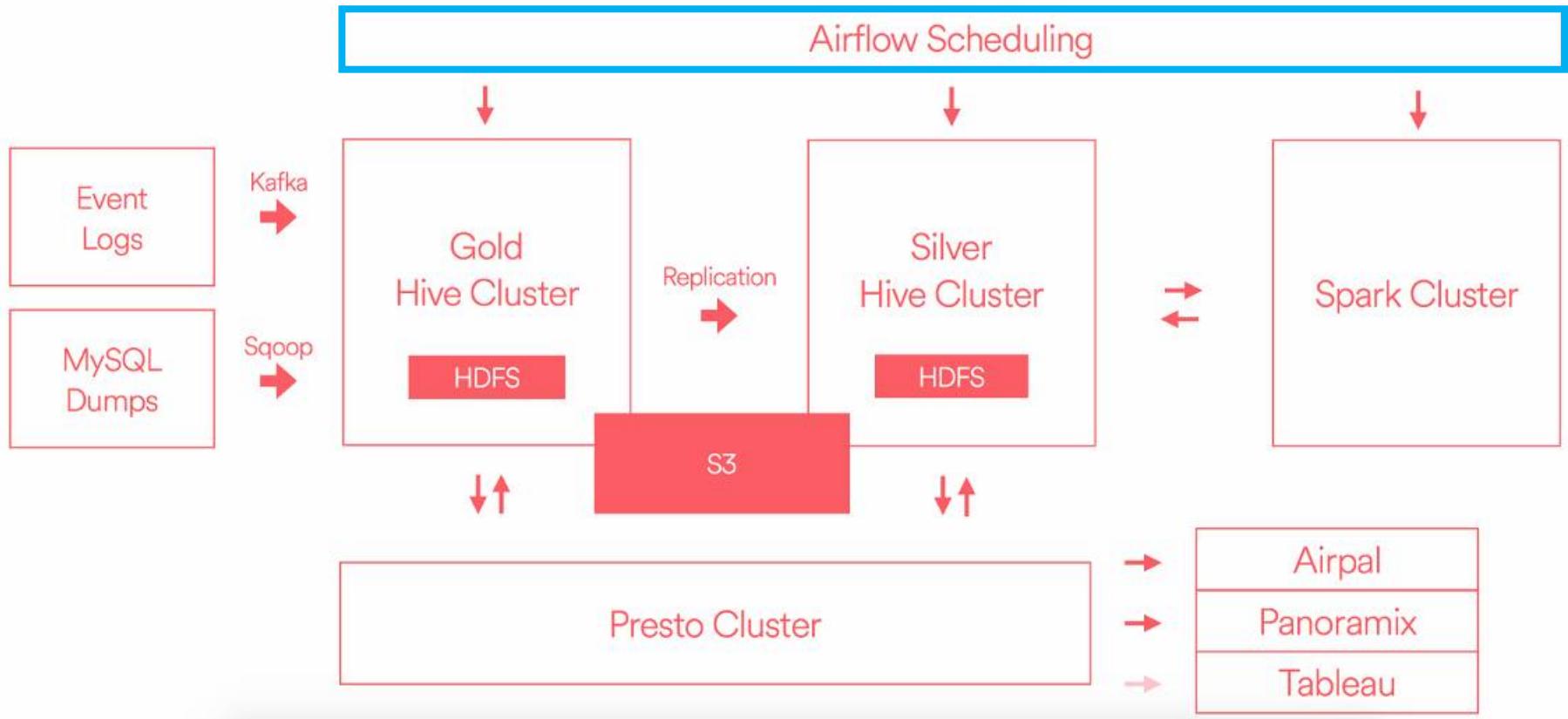
- Hadoop 2.0



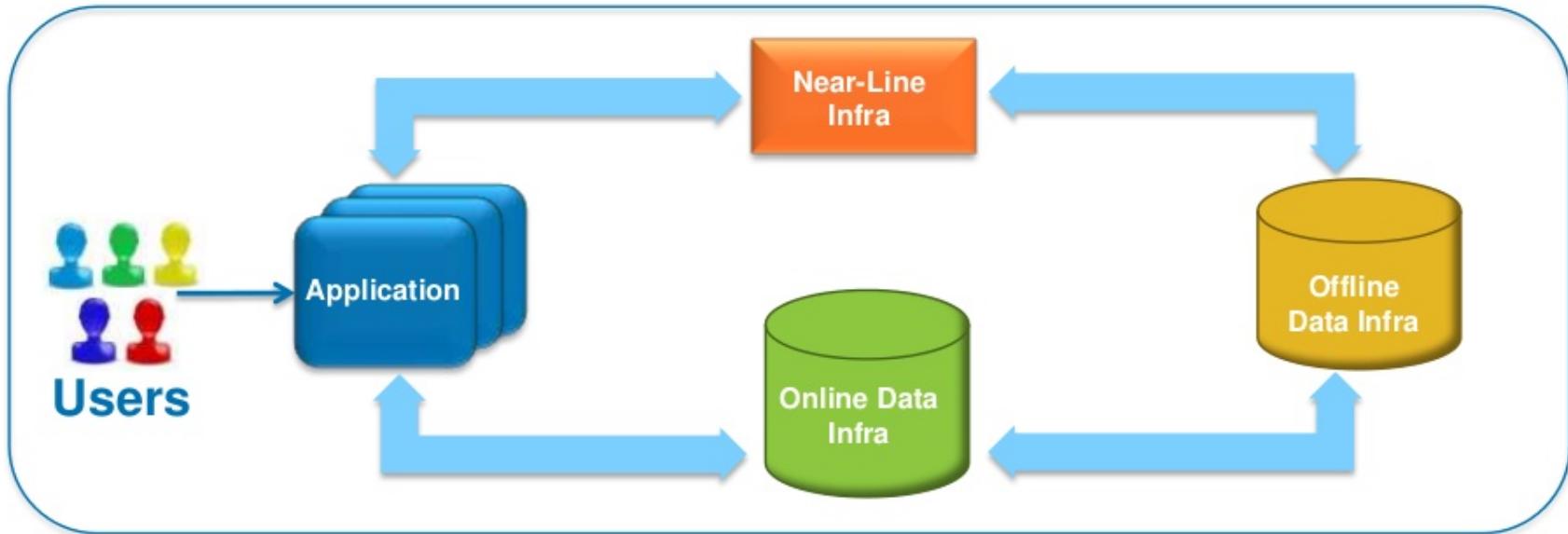
# ● Hadoop 2.0



# • Airbnb Data Infra

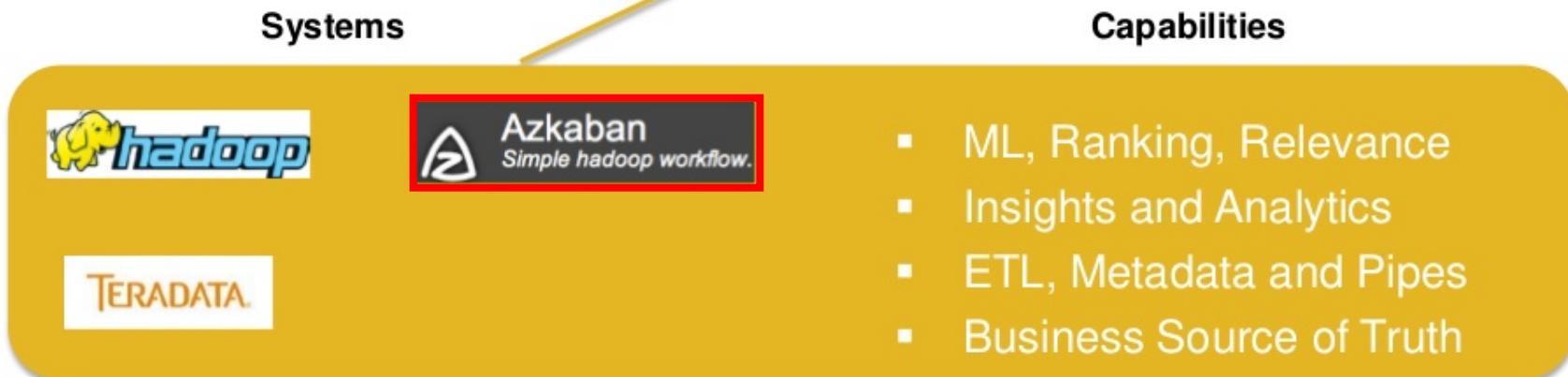
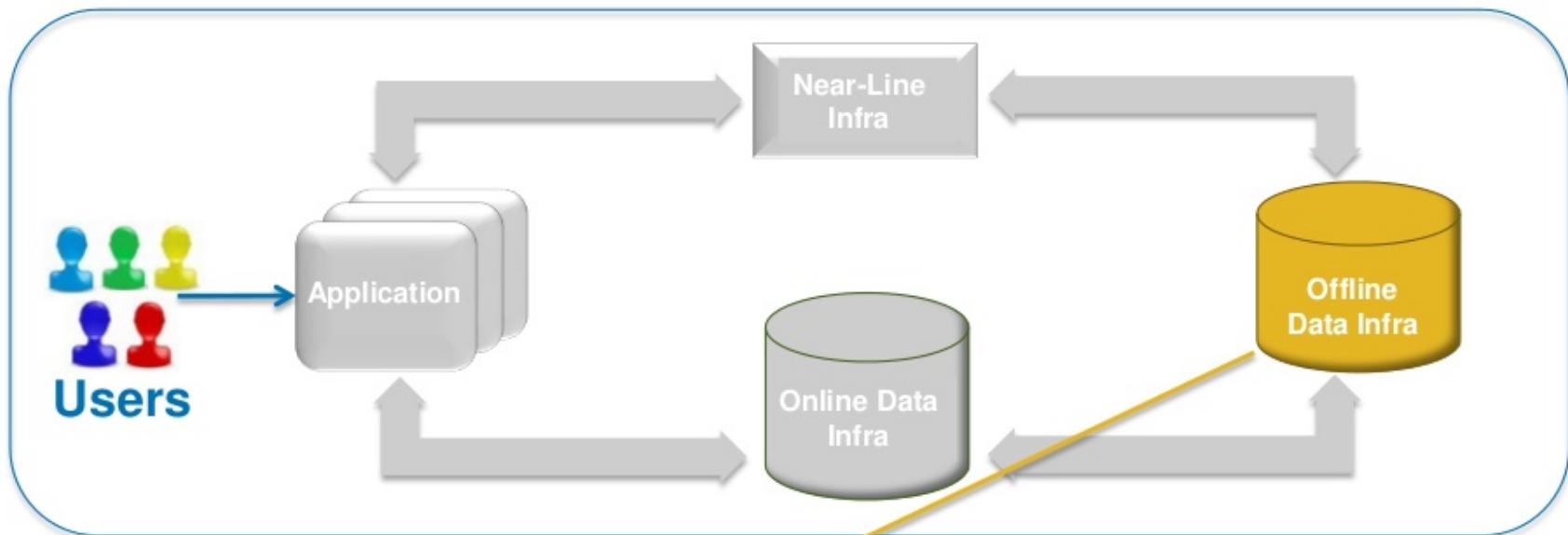


# LinkedIn Data Infra



Infrastructure	Latency & Freshness Requirements	Products
Online	Activity that should be reflected <b>immediately</b>	<ul style="list-style-type: none"><li>• Member Profiles</li><li>• Company Profiles</li><li>• Connections</li></ul> <ul style="list-style-type: none"><li>• Messages</li><li>• Endorsements</li><li>• Skills</li></ul>
Near-Line	Activity that should be reflected <b>soon</b>	<ul style="list-style-type: none"><li>• Activity Streams</li><li>• Profile Standardization</li><li>• News</li></ul> <ul style="list-style-type: none"><li>• Recommendations</li><li>• Search</li><li>• Messages</li></ul>
Offline	Activity that can be reflected <b>later</b>	<ul style="list-style-type: none"><li>• People You May Know</li><li>• Connection Strength</li><li>• News</li></ul> <ul style="list-style-type: none"><li>• Recommendations</li><li>• Next best idea...</li></ul>

# • LinkedIn Data Infra



- Data of workflow scheduler in Big Data



- 14 boxes dedicated for work-flow system
- 8,000 tasks daily



- Maintain 3 instances of work-flow system
- 2,500 flows, 30,000 jobs daily



- 2000+ tasks, 10,000+ Hadoop jobs daily





# What's the most important for a Big data workflow scheduler ?

# Dead Simple:

- Easy to use and configure





# Problems with big data job scheduling

丁来强 wjo1212@163.com

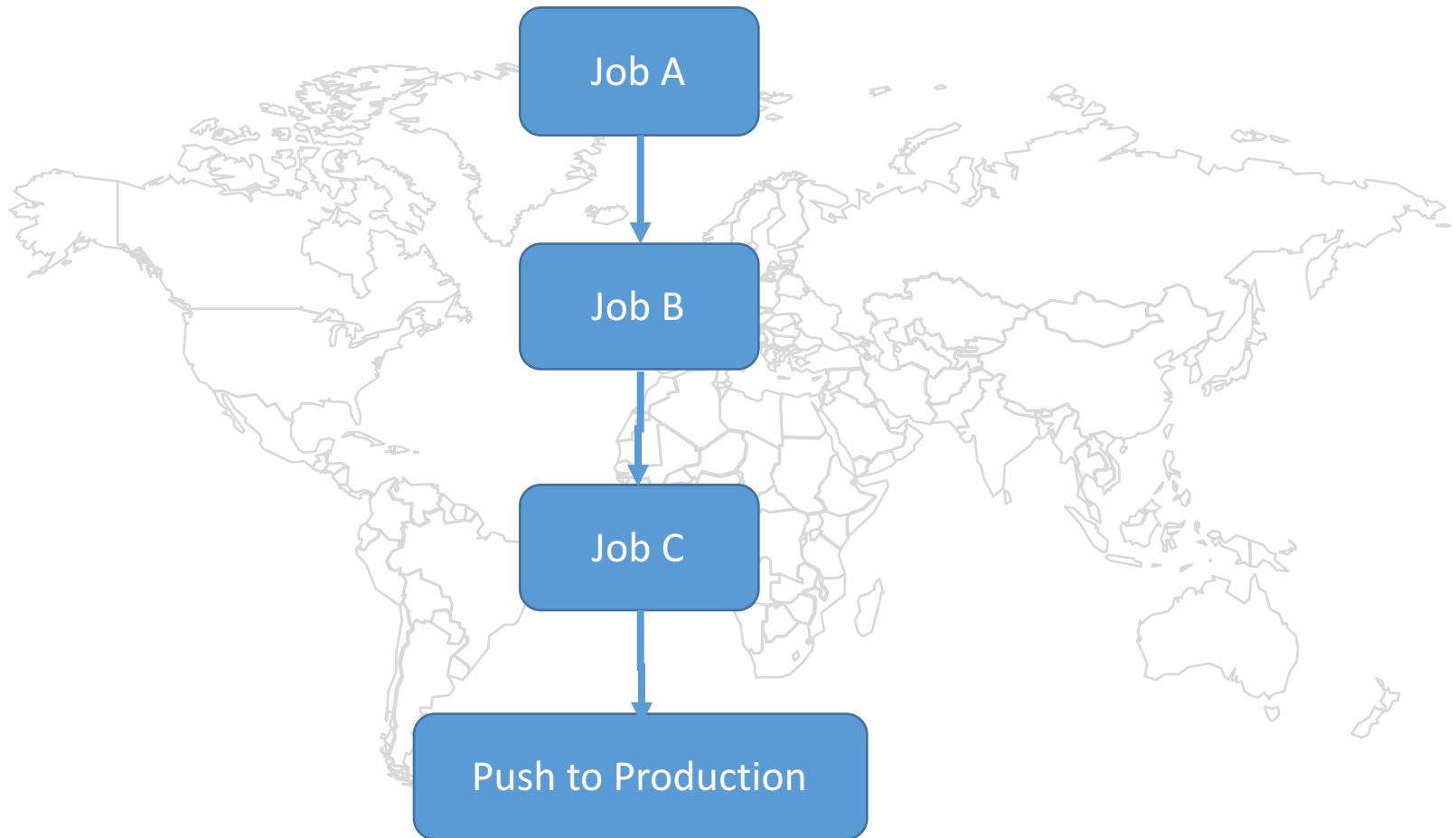
- Typical Challenge

- 数据工作流程复杂度越来越高
- 数据分析与批处理数据非常重要
- 大量时间花费在编写任务、检测与排错上

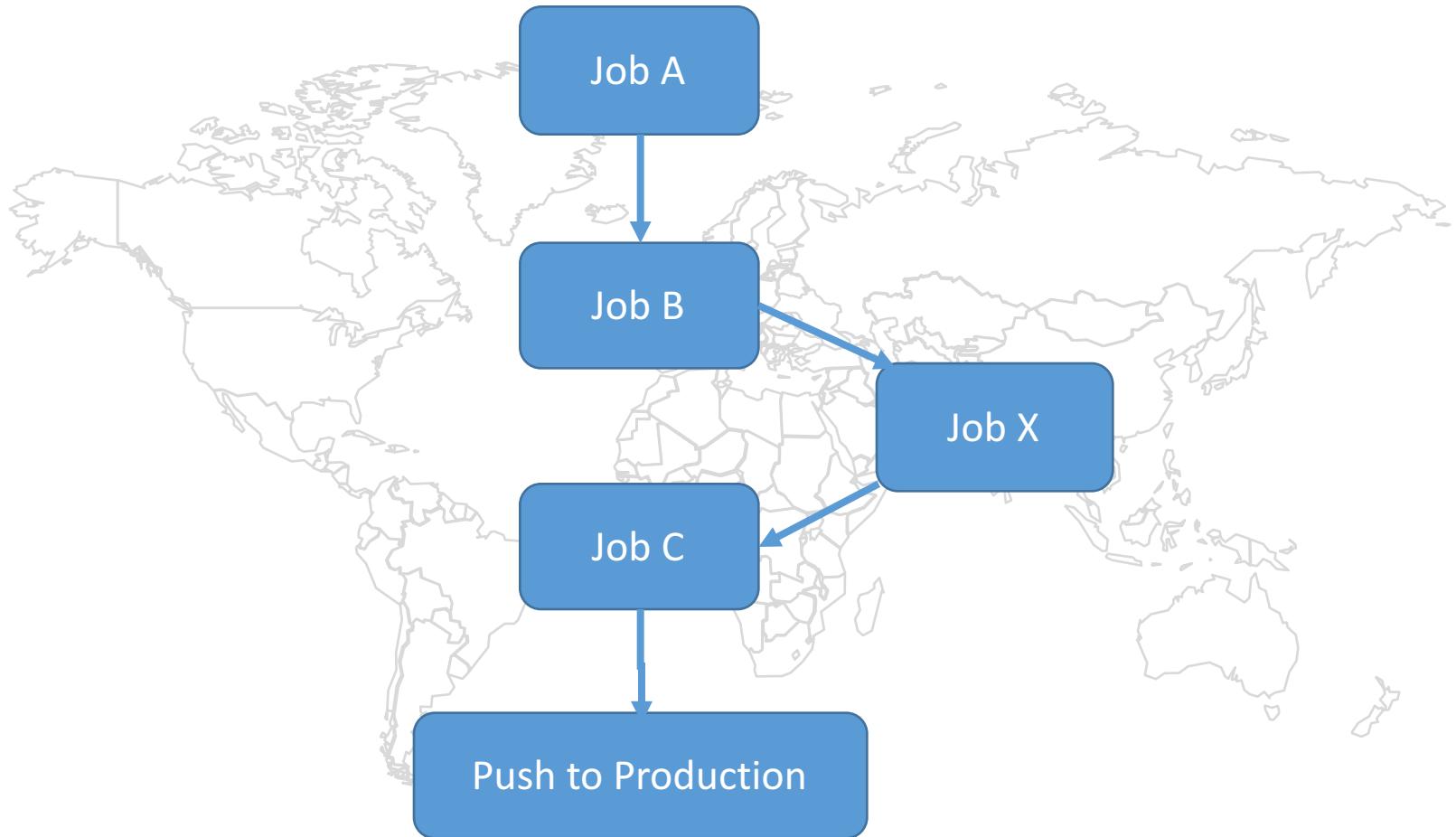


# Fragile process

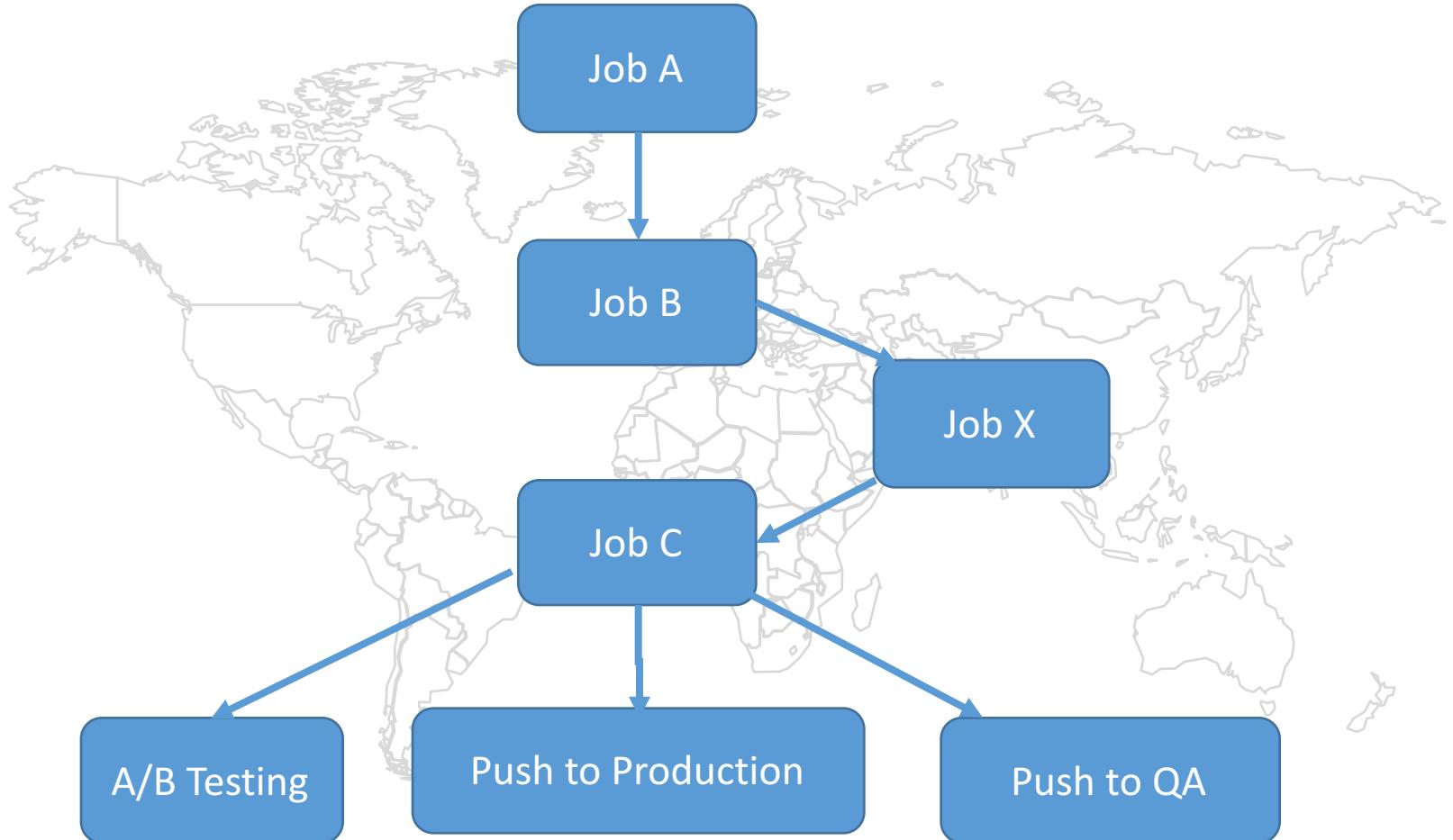
- Fragile process



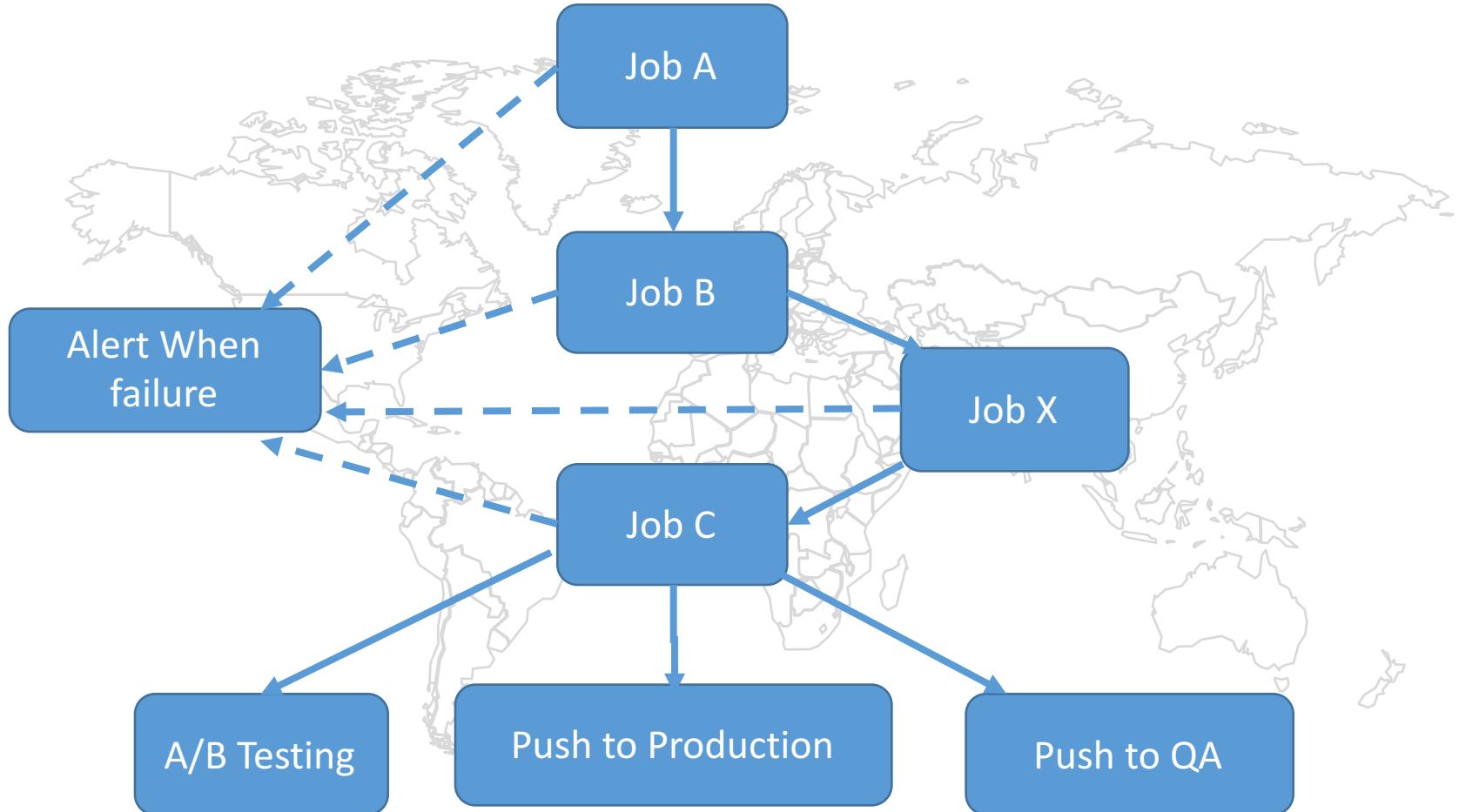
- Fragile process



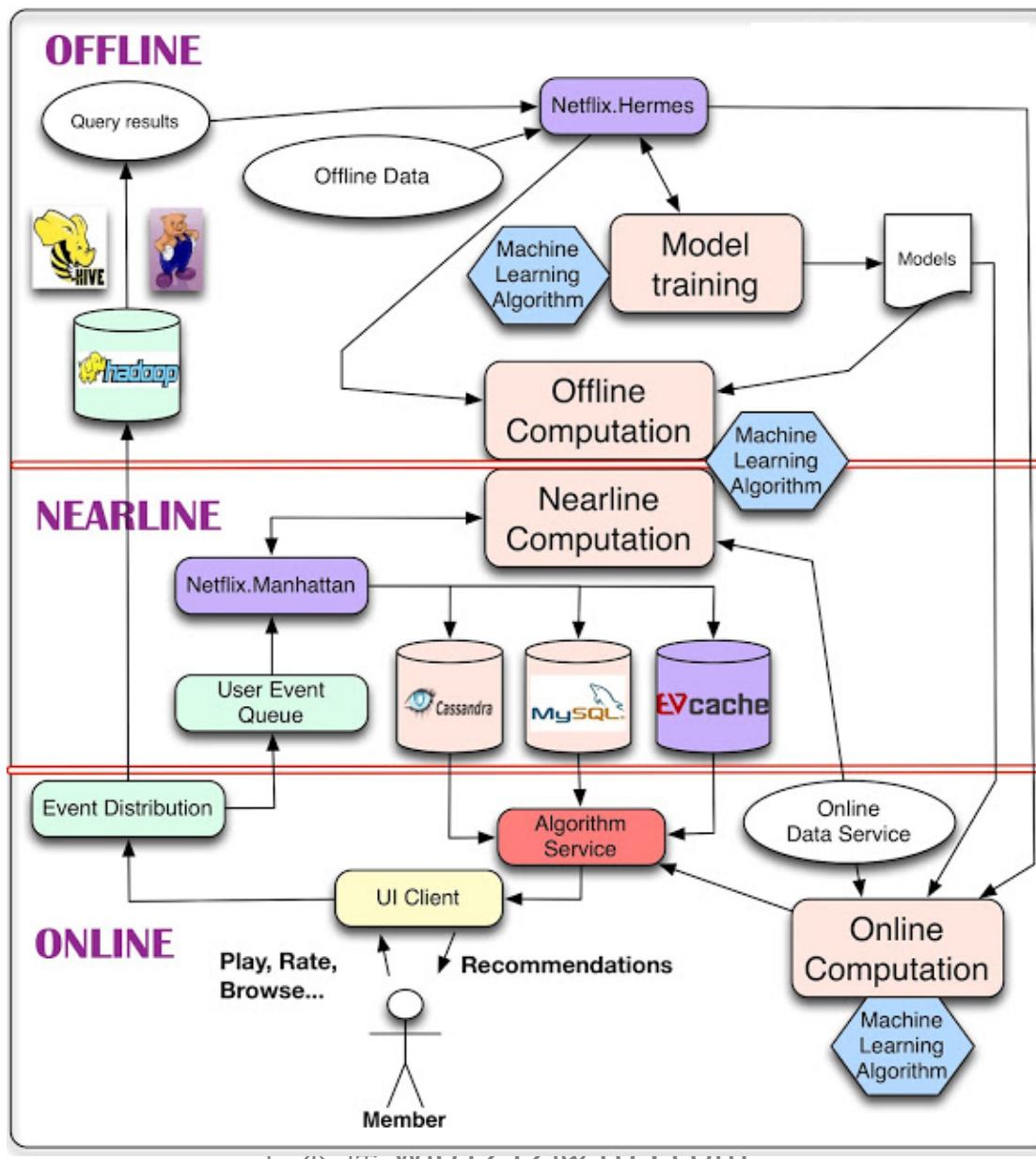
- Fragile process



- Fragile process



# ● Example: Netflix Recommendation System





# Fragile Failure Handling

丁来强 wjo1212@163.com

3. Some errors or bugs may exist in some jobs' logic



2. Job fails due to system or network  
may not be temporarily not available



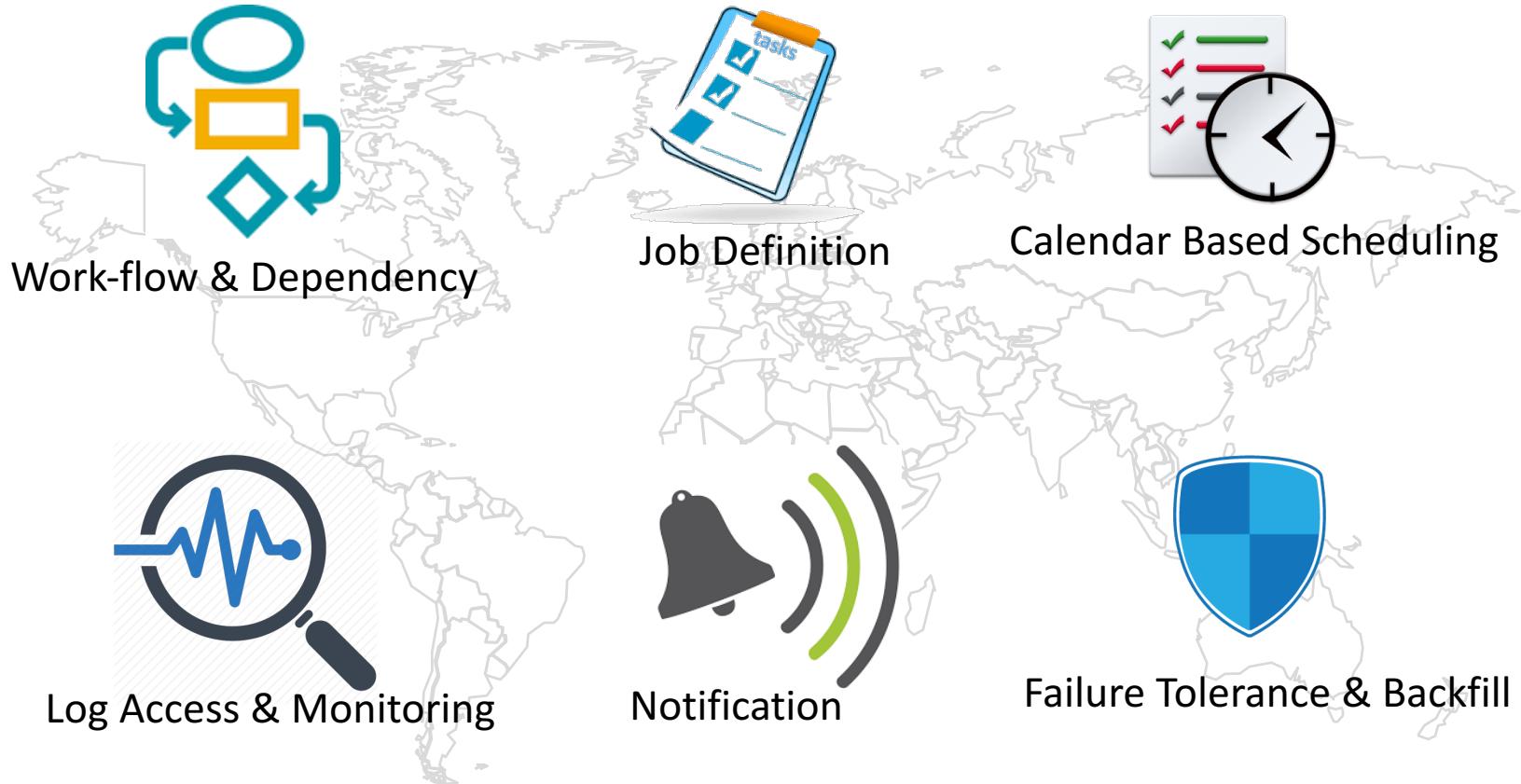
1. Scheduled triggers are skipped due to unavailability of sub-system

4. Performance is slow especially for some critical steps



丁来强 wjo1212@163.com

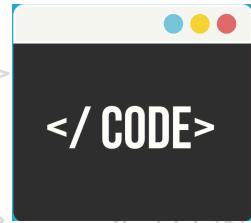
- Basic Needs



- Advanced Needs



Complex Rule



Programmatic



Operator OOB



Scalability

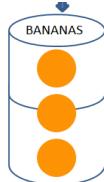


High Availability

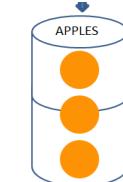
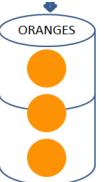


SLA Monitor & Alert

- Advanced Needs (cont')



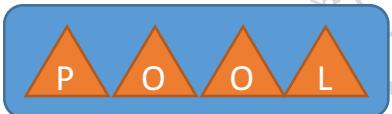
Queue (Affinity)



Data Profiling



Versioning



Pool (Limit concurrency + priority)



Event Driven Scheduler



Plugins



丁来强 wjo1212@163.com

- Options



Airflow

Luigi

Chronos



Azkaban

Oozie

cron

# ● Solution Overview

Basic Info	Luigi	Airflow	Azkaban	Oozie
Language	Python	Python	Java	Java
Github Stars	5,274	3,422	780	354
Contributors	256	178	37	18
Latest Version	2.3.1	1.7.1	3.1	4.2
History	4 years	1+ years	6+ years	6+ years
Invented by	Spotify	Airbnb	LinkedIn	Yahoo
Owned by	Spotify	Apache Incubator	Apache	Apache

- Azkaban



- Pros:
  - Born for Hadoop
    - Support all Hadoop, hive, pig versions
  - Easy to use Web UI:
    - Good Job visualization and monitoring
    - Flexible Module structure/Plugins
- Cons:
  - **Properties files based configuration**
  - Web UI only, No CLI and REST interfaces (need 3<sup>rd</sup> party AzkabanCLI)
  - Limited execution path control

# Azkaban GUI

Azkaban Local  
My Local Azkaban

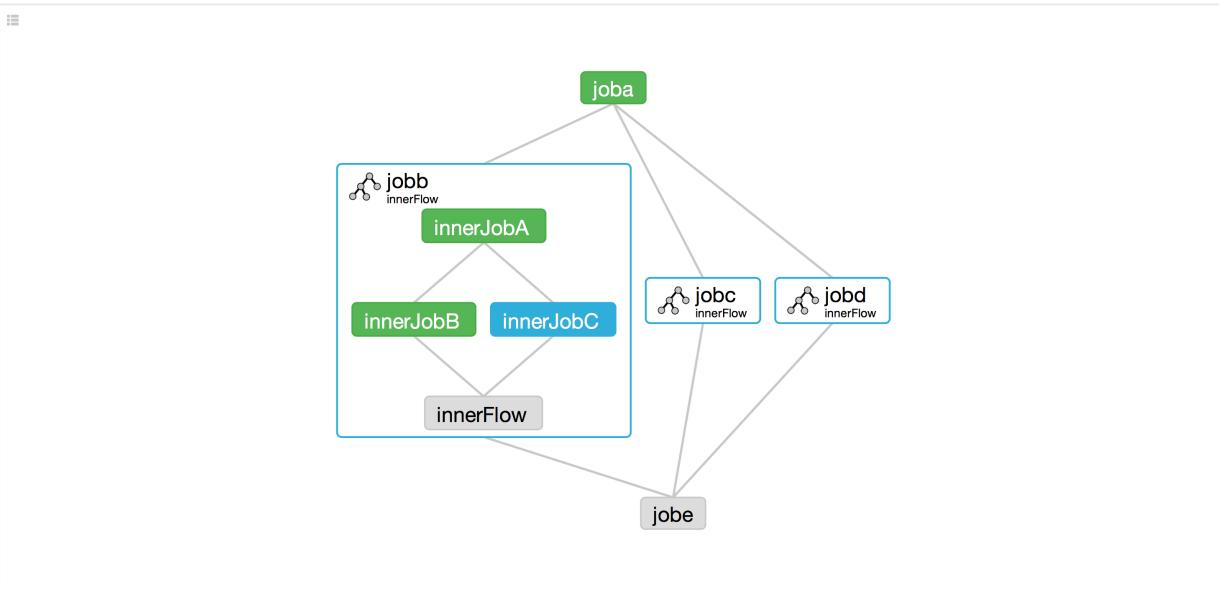
Projects Scheduling Executing History HDFS azkaban -

Flow Execution 27 RUNNING

Project embedded / Flow jobe / Execution 27

Submit User azkaban Duration 16 sec Start Time 2014-02-19 17:03:51s End Time -

Graph Job List Flow Log Stats Kill Pause



```

graph TD
    joba[joba] --> jobb[jobb]
    joba --> jobc[jobc]
    joba --> jobd[jobd]
    jobb --> innerJobA[innerJobA]
    jobb --> innerFlow[innerFlow]
    innerJobA --> innerJobB[innerJobB]
    innerJobA --> innerJobC[innerJobC]
    innerFlow --> jobe[jobe]
    innerJobB --> jobe
    innerJobC --> jobe
    jobc --> jobe
    jobd --> jobe
  
```

Test-embedded flows project.  
Last modified on 2014-01-19 12:18:49 by azkaban.

embeddedFlow2  
Test-embedded flows project.  
Last modified on 2014-01-19 12:17:00 by azkaban.

ID	Author	Created	Updated
27	azkaban	2014-01-21 11:58:47+0	2014-01-21 11:58:49+0
28	azkaban	2014-01-21 11:58:13+0	2014-01-21 11:58:15+0
29	azkaban	2014-01-21 11:58:05+0	2014-01-21 11:58:05+0

Elapsed	Status	Action
24 sec	<span style="background-color: green;">Success</span>	
1 sec	<span style="background-color: green;">Success</span>	
2 sec	<span style="background-color: red;">Failed</span>	
3 sec	<span style="background-color: red;">Failed</span>	

- Oozie 

- Pros:
  - Born for Hadoop
  - CLI, HTTP, JAVA API interfaces
  - Support extended Alert integration
- Cons:
  - **Higher learning curve**
  - PDL style XML based configuration
  - Limited Web UI (need Cloudera Hue)
  - No resource control



丁来强 wjo1212@163.com

- Overview

- Pros:
  - Programmatic by Python
  - Modeling is simple, Code is mature (~20K LOC)
  - Good support Hadoop (MR, logs, dist)
  - Test friendly, support local scheduler
- Cons:
  - Web UI is very limited
  - No built-in trigger (need cron)
  - Not design for large scaling (> 100K tasks)
  - No support distribution of execution

# ● Task Definition

```
class X(luigi.Task):  
    def requires(self): ...  
    def run(self): ...  
    def output(self): ...
```

Setup Dependencies:  
Return one or more Tasks

Logic

Output of the Task:  
Return one or more Targets

# ● Task Example

```
import luigi

class MyTask(luigi.Task):
    param = luigi.Parameter(default=42)

    def requires(self):
        return SomeOtherTask(self.param)

    def run(self):
        f = self.output().open('w')
        print >>f, "hello, world"
        f.close()

    def output(self):
        return luigi.LocalTarget('/tmp/foo/bar-%s.txt' % self.param)

if __name__ == '__main__':
    luigi.run()
```

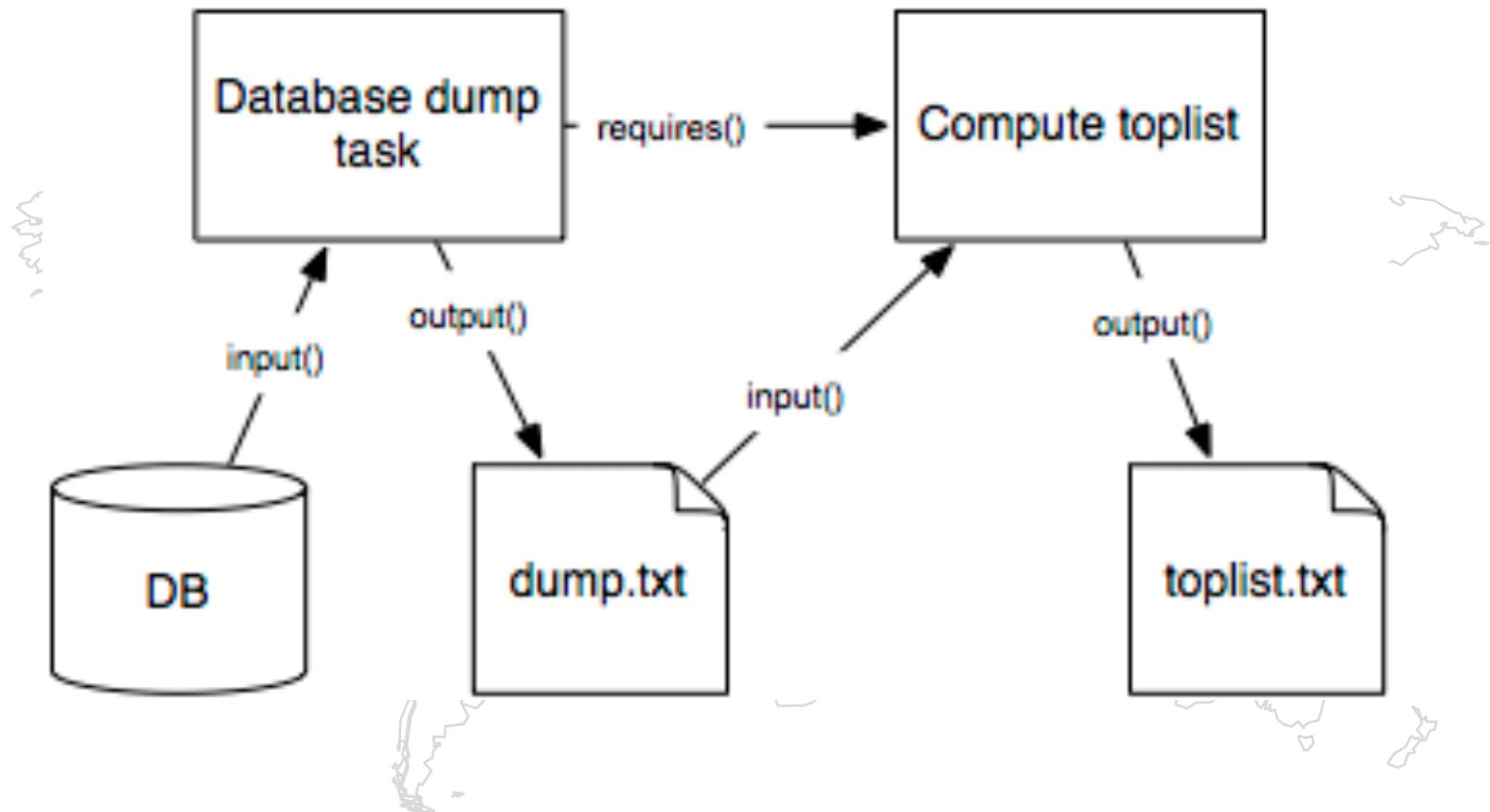
The business logic of the task

Where it writes output

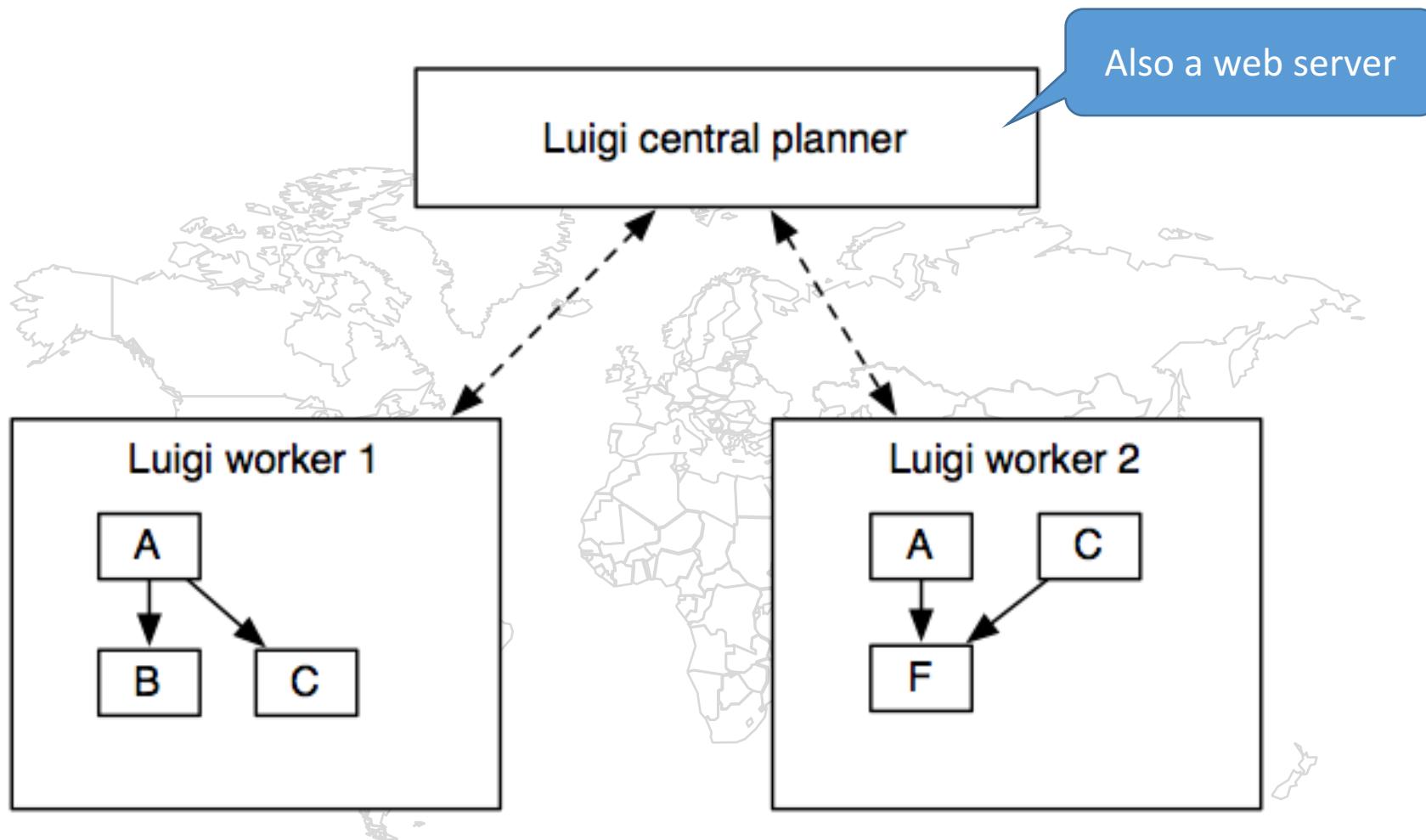
What other tasks it depends on

Parameters for this task

# • Task Execution



# • Architecture



- Architecture Notes

- Mainly manage the dependency and de-dup the task running.
- Mainly focus on data pipeline ETL.
- Limitations
  - No calendar trigger
  - Web UI is very limited
  - Too couple between worker and scheduler (not support > 100K tasks)
  - Execution is bundled on specific worker

# ● Web UI – execution status

Luigi Task Visualiser × localhost:8082/static/visualiser/index.html#

Luigi Task Status    Task List    Dependency Graph    Workers

5 CreateTrialEligibility	PENDING TASKS 6679	RUNNING TASKS 0	DONE TASKS 6946	FAILED TASKS 88
10 CreateUserCollectionSummary				
10 CreateCollectionTracksWithM				
24 BannerShownCleanedHour	UPSTREAM FAILED 351	DISABLED TASKS 0	UPSTREAM DISABLED 0	
10 AdUserMetricsImport				
10 AggregatedPlaysHive				
6 CreateReportingUsage				
10 LastLastActivityHive				
11 CreateUserPlaylistSummary				
3 FinancialUser				
30 AppAnnieDownloadsImport				
59 CreateArtistUsageBase				
9 Cass2HdfsCompact				
4 GCSToHDFSMain				
11 CreateUserCollectionSummary				

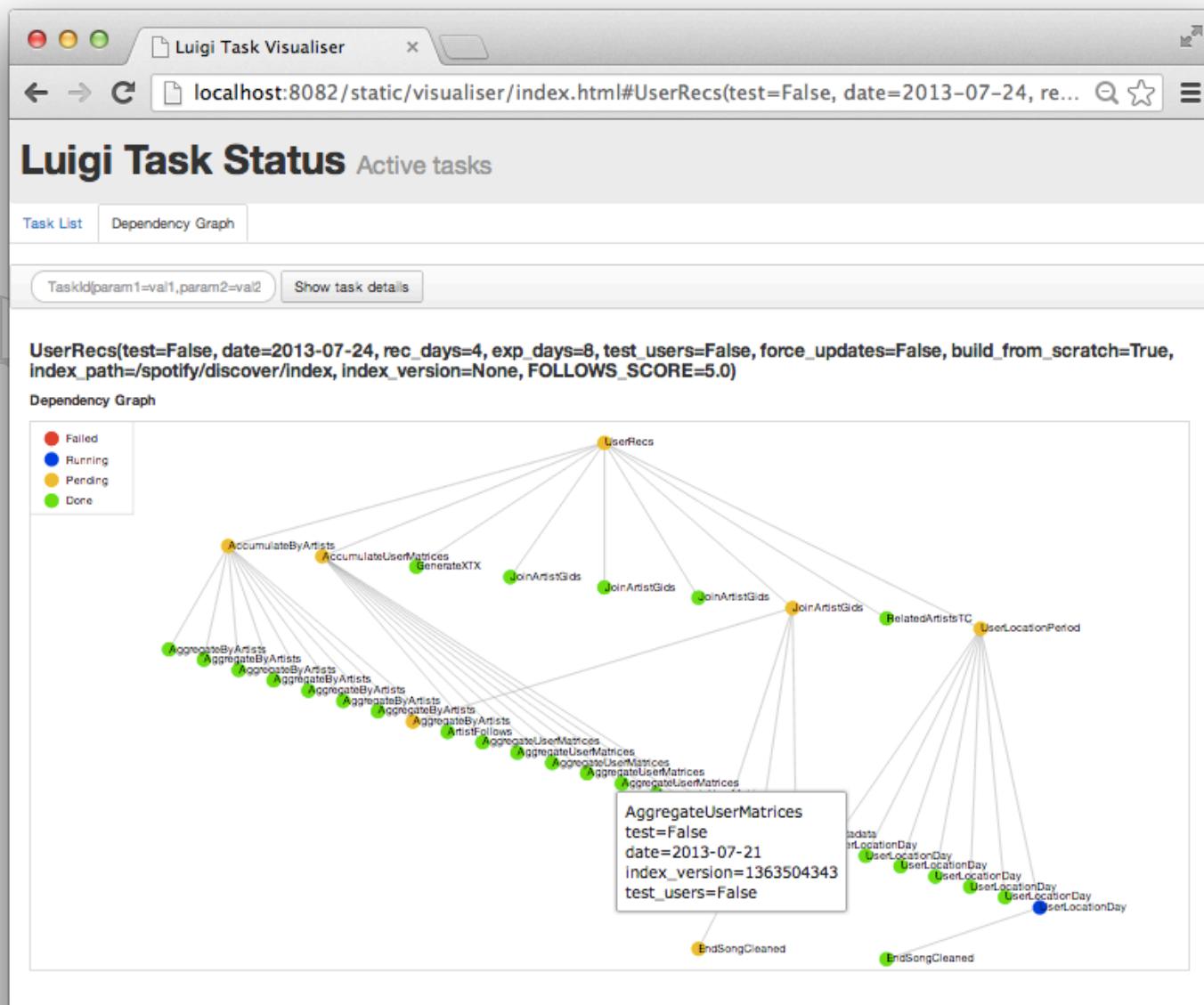
Displaying tasks of family **CreateReportingUsage**.

Show 10 entries    Filter table:  Filter on Server

Name	Details
DONE CreateReportingUsage	(test=False, date=2015-06-10, parallel=True)
DONE CreateReportingUsage	(test=False, date=2015-06-11, parallel=True)
PENDING CreateReportingUsage	(test=False, date=2015-06-13, parallel=True)
PENDING CreateReportingUsage	(test=False, date=2015-06-12, parallel=True)
UPSTREAM_FAILED CreateReportingUsage	(test=False, date=2015-06-14, parallel=True)
UPSTREAM_FAILED CreateReportingUsage	(test=False, date=2015-06-15, parallel=True)

Showing 1 to 6 of 6 entries (filtered from 14,064 total entries)    Previous  Next

- Web UI – DAG visualization



- Task and Targets Library

- Google Bigquery
- Hadoop jobs
- Hive queries
- Pig queries
- Scalding jobs
- Spark jobs
- Postgresql, Redshift, Mysql tables
- and more...



丁来强 wjo1212@163.com

- Overview

- Pros (we will see):
  - More General Flexible Architecture
  - Very compelling Web UI
  - Lots of cool features OOB, Rich Operator library
  - Fast growing adoption (30+ companies)
  - Test friendly (test mode and SequentialScheduler)
- Cons:
  - Coding quality is not so mature (UT coverage is not high)
  - No event driven scheduler (same to all others solutions)

- Airflow Tech Stack

- Python Code ( < 20K LOC )
- DB: SQLAlchemy
- Celery for distributed execution
- Web Server: Flask / gunicorn
- UI: d3.js / Highcharts / Pandas
- Templating: Jinjia2

# Airflow Web



Work-flow & Dependency

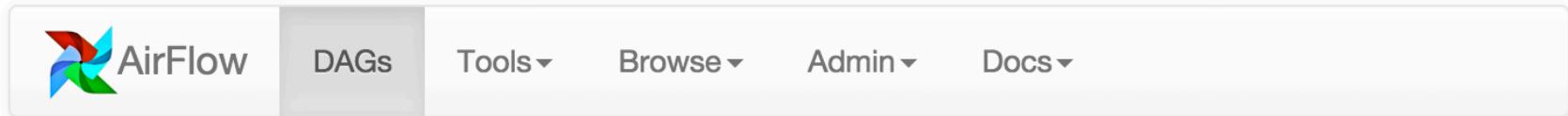


Log Access & Monitoring



Data Profiling

# ● Web UI – Overall status

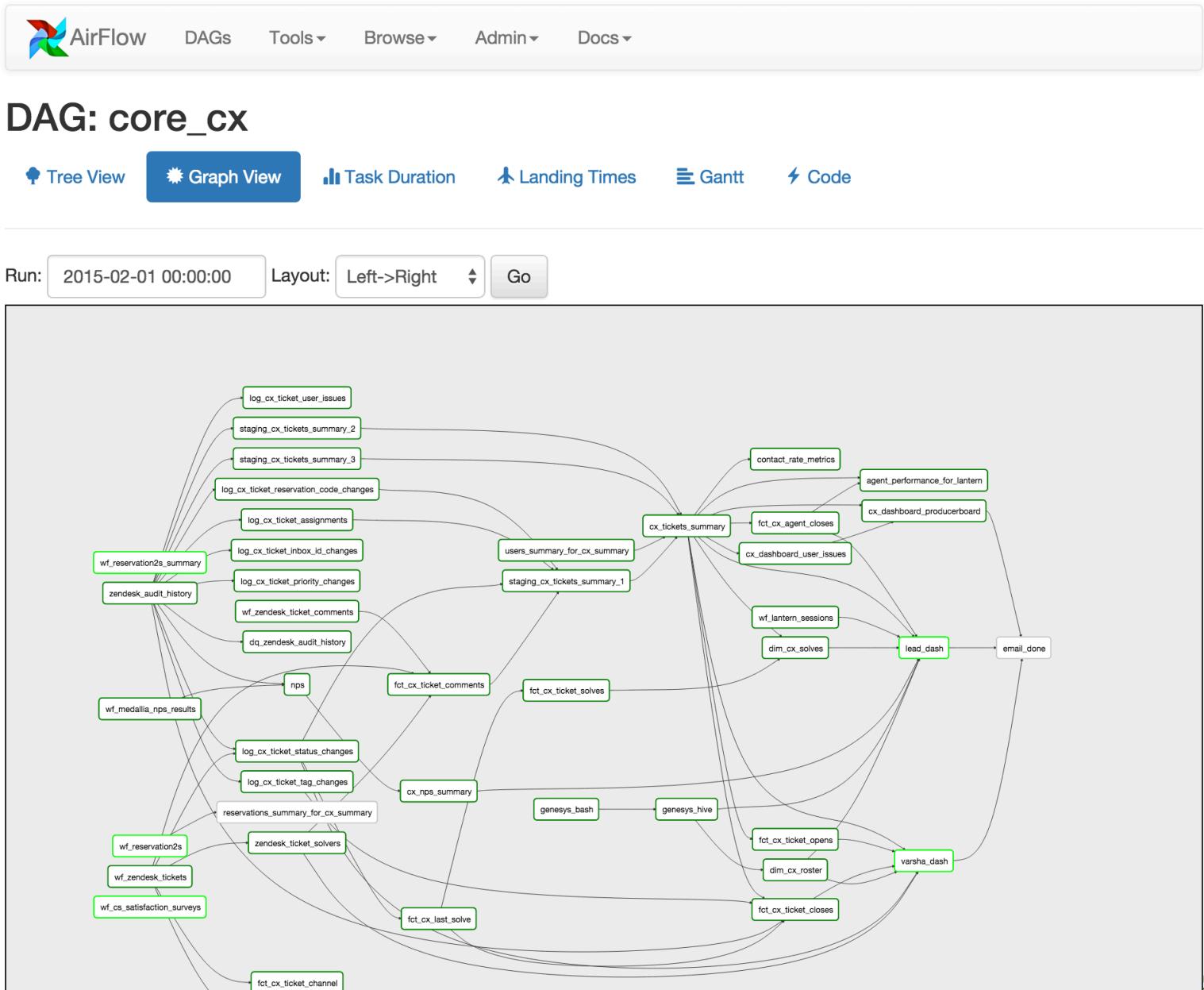


## DAGs

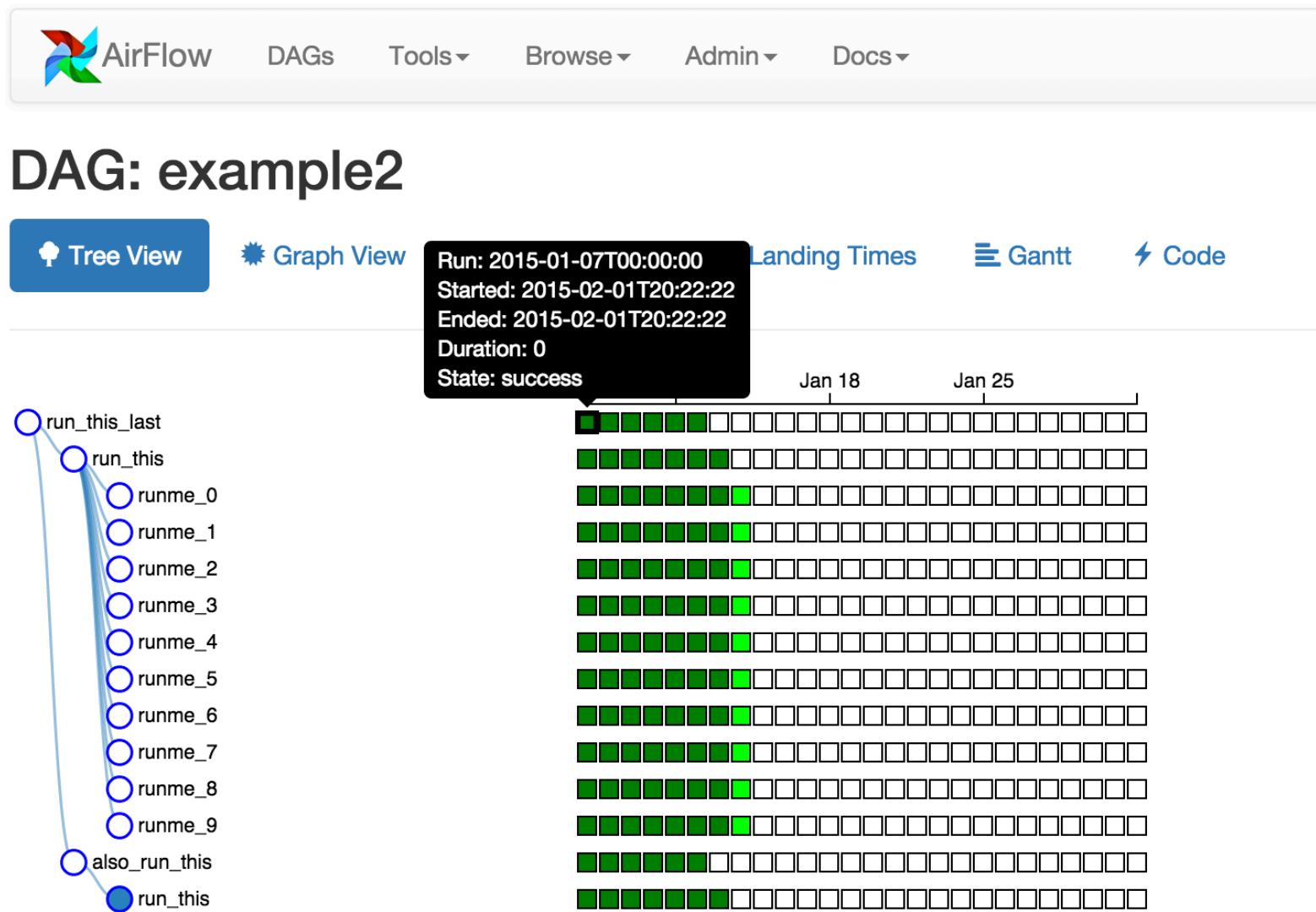
DAG	Filepath	Owner	Task by State	Links
example1	example_dags/example1.py	airflow	<span>80</span> <span>1</span> <span>0</span>	
example2	example_dags/example2.py	airflow	<span>128</span> <span>10</span> <span>0</span>	
example3	example_dags/example3.py	airflow	<span>138</span> <span>5</span> <span>0</span>	



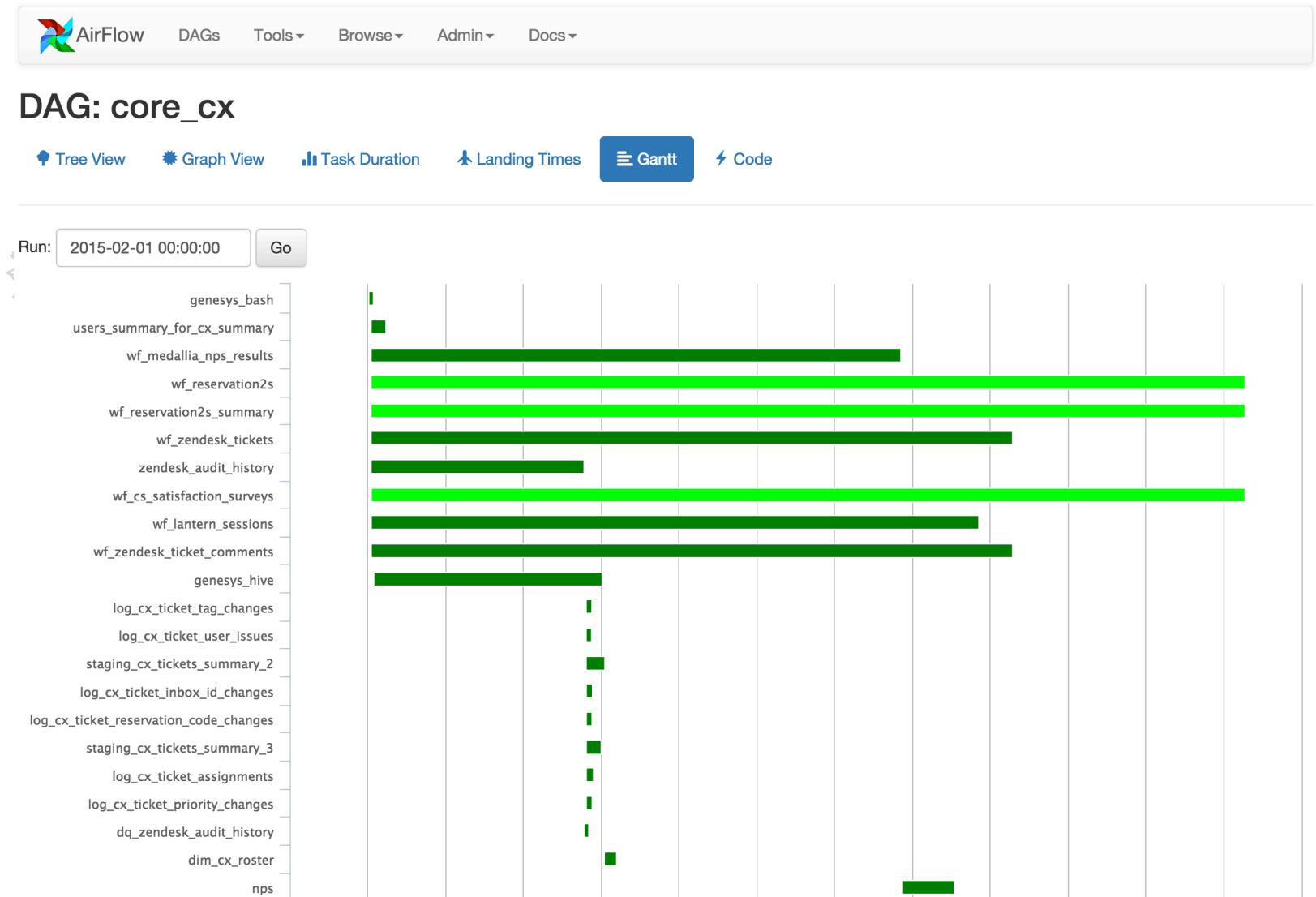
# ● Web UI – workflow visualization



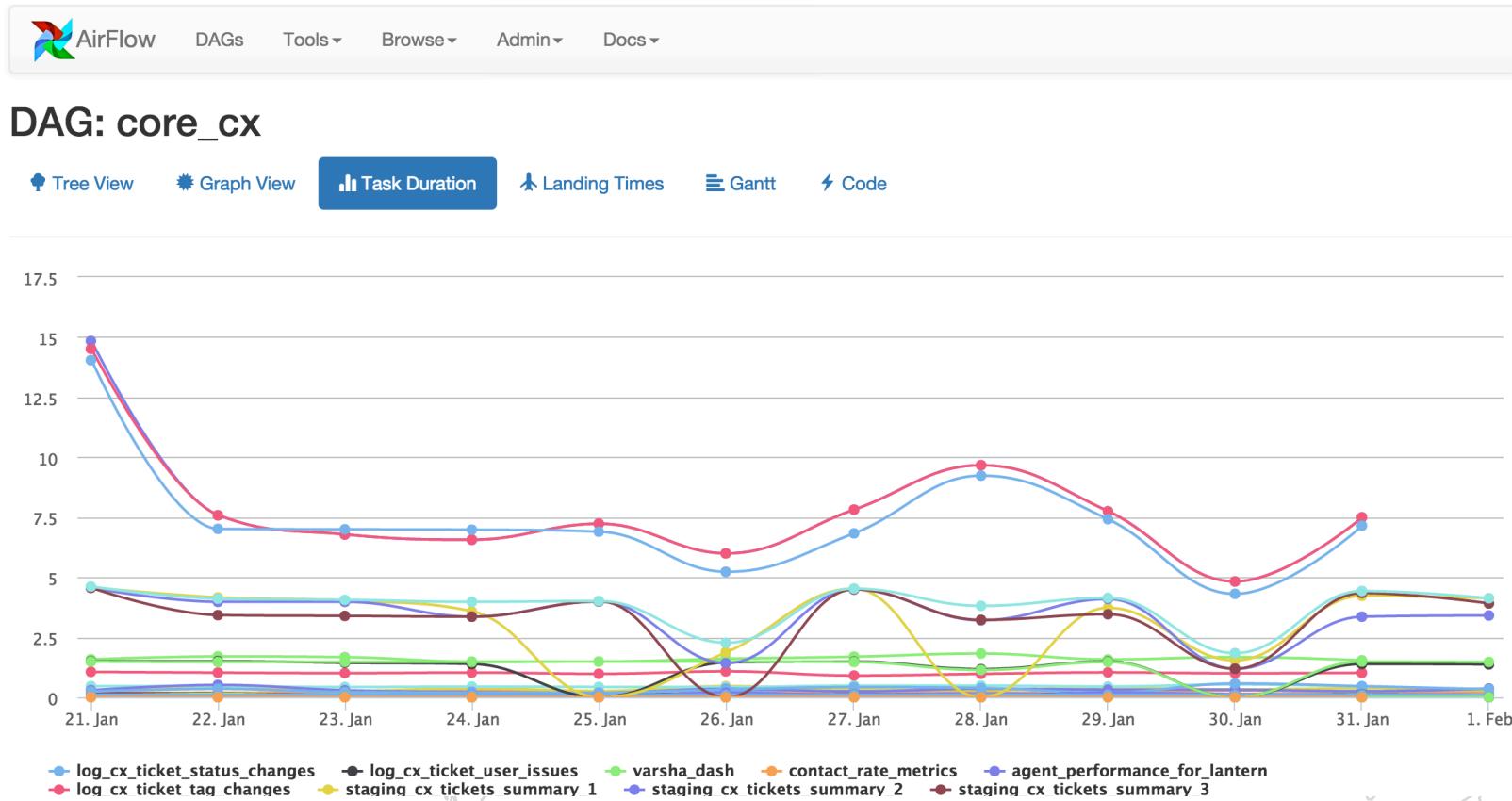
# • Web UI – execution history



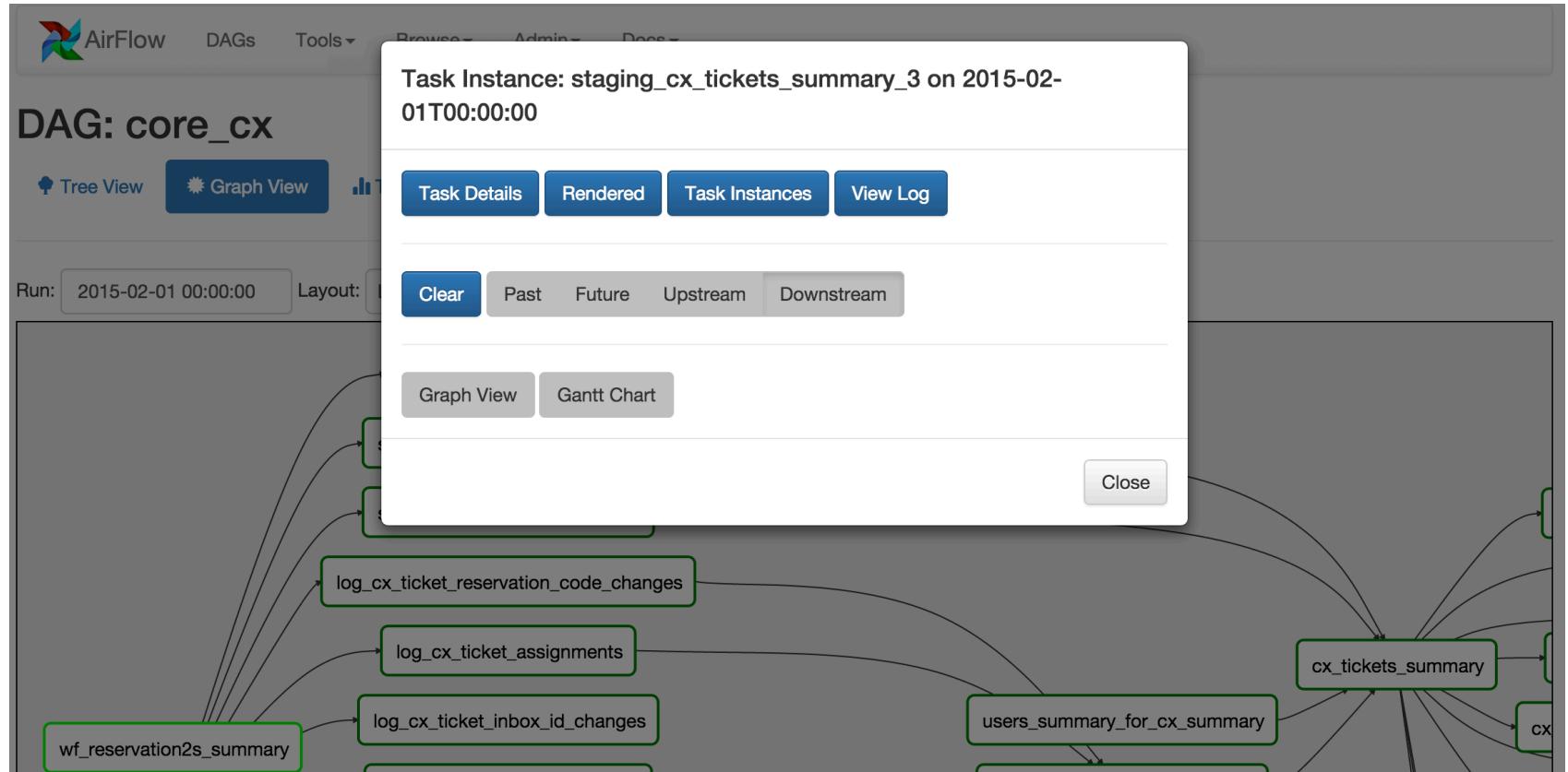
# ● Web UI – performance profile



# ● Web UI – Performance stats over time



# • Web UI – Deep dive for task execution

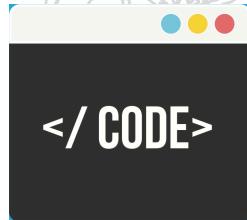




# Airflow Concepts



Work-flow & Dependency



Programmatic



Operator OOB

# ● DAG (Directed Acyclic Graph)

```
from airflow.operators import PythonOperator
from airflow.models import DAG
```

```
dag = DAG(
    dag_id='mydag',
    default_args={ 'owner': 'airflow' },
    schedule_interval='0 0 * * *',
    dagrun_timeout=timedelta(minutes=60))
```

DAG: a collection of tasks  
w/ scheduling settings

```
task1 = BashOperator(
    task_id='mydag_task1',
    bash_command='echo "{{ run_id }}" && sleep 1',
    dag=dag)
```

Task: an instance of BashOperator  
Support templating

```
def python_logic(param1): pass
```

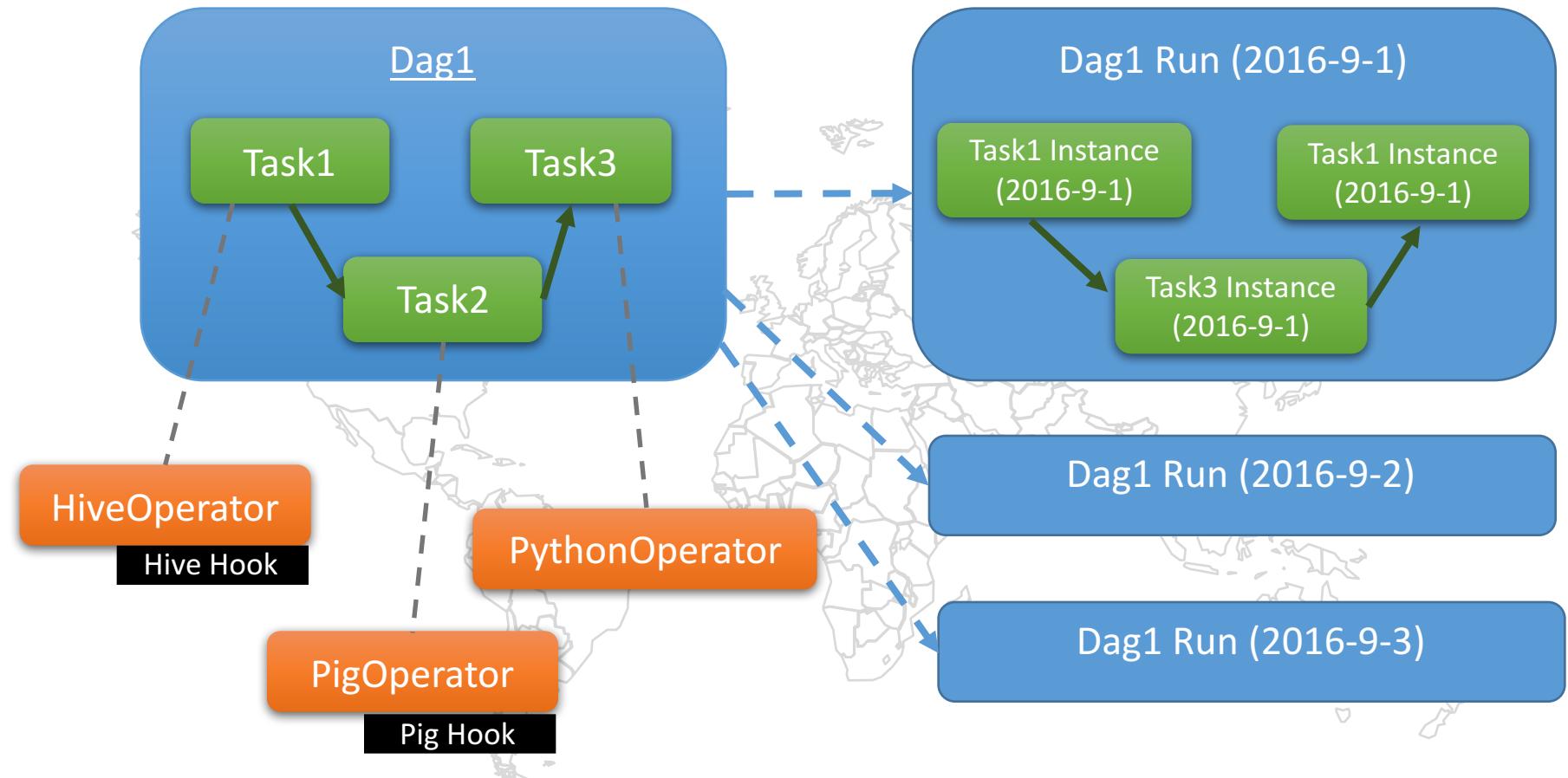
```
task2 = PythonOperator(
    task_id='mydag_task2',
    python_callable=python_logic,
    op_kwargs={'param1': 10},
    dag=dag)
```

An task of another kind of  
PythonOperator

Setup the dependencies

```
task1.set_upstream(task2)
```

# DAG execution



- Concepts – DAG, DAG Run

- DAG
  - A collection of Tasks
  - Setting of Calendar Scheduling

- Dag Run
  - A run instance of DAG with a scheduled date (ID: dag, start time and interval)

- Concepts – Operator, Task and TI

- Operator
  - Task templates
- Task
  - Instance of a Operator
- Task Instance (TI)
  - Belong to Dag Run
  - A run instance of a Task with a scheduled date (id: dag, task, start time and interval)

- Concepts - Operator

- Operator

- Task templates, general categories:
  - Sensor
  - Branching
  - Transformer
- Settings of Trigger Rules, retry etc.
- Use Hook for real operation w/ external systems

# ● Operator Library

- Google Bigquery, Could Storage
- AWS S3, EMR
- Spark SQL
- Docker
- Presto
- Sqoop
- Hive jobs
- Vertica
- Qubole
- SSH
- Hipchat, Slack, Email
- Postgresql, Redshift, Mysql, Oracle etc.
- and more...

# ● Parameterized Tasks

- Variables
  - Global parameters
- Connections
  - External system's connection string, confidential, extra parameters etc. Normally used by Hook.
- DAG Parameters/Macros
- Templating
  - Using Jinjia for batch or any places that fit
- Xcom
  - Share data between Tasks



# Architecture Insight



Scalability

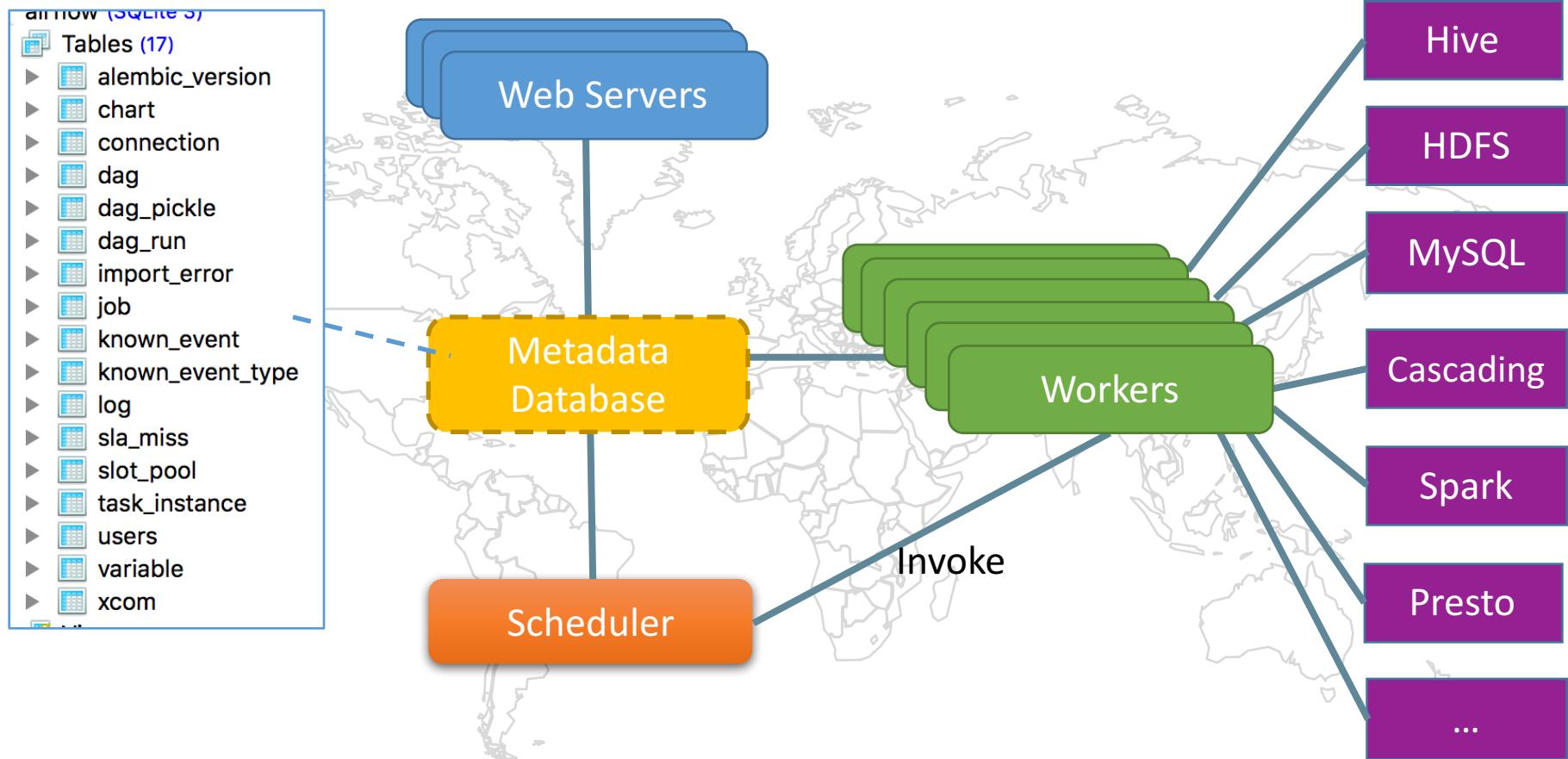


High Availability

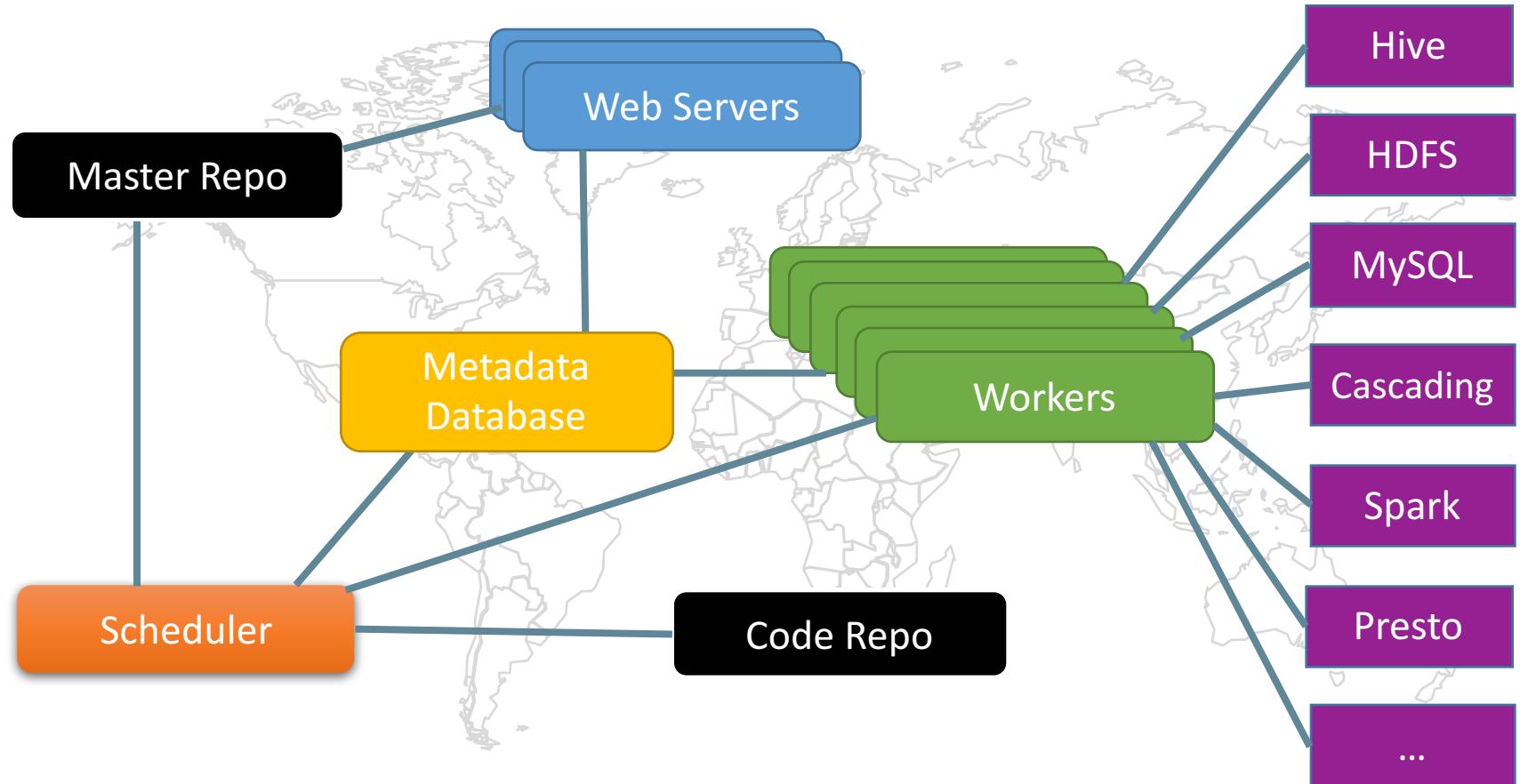


Versioning

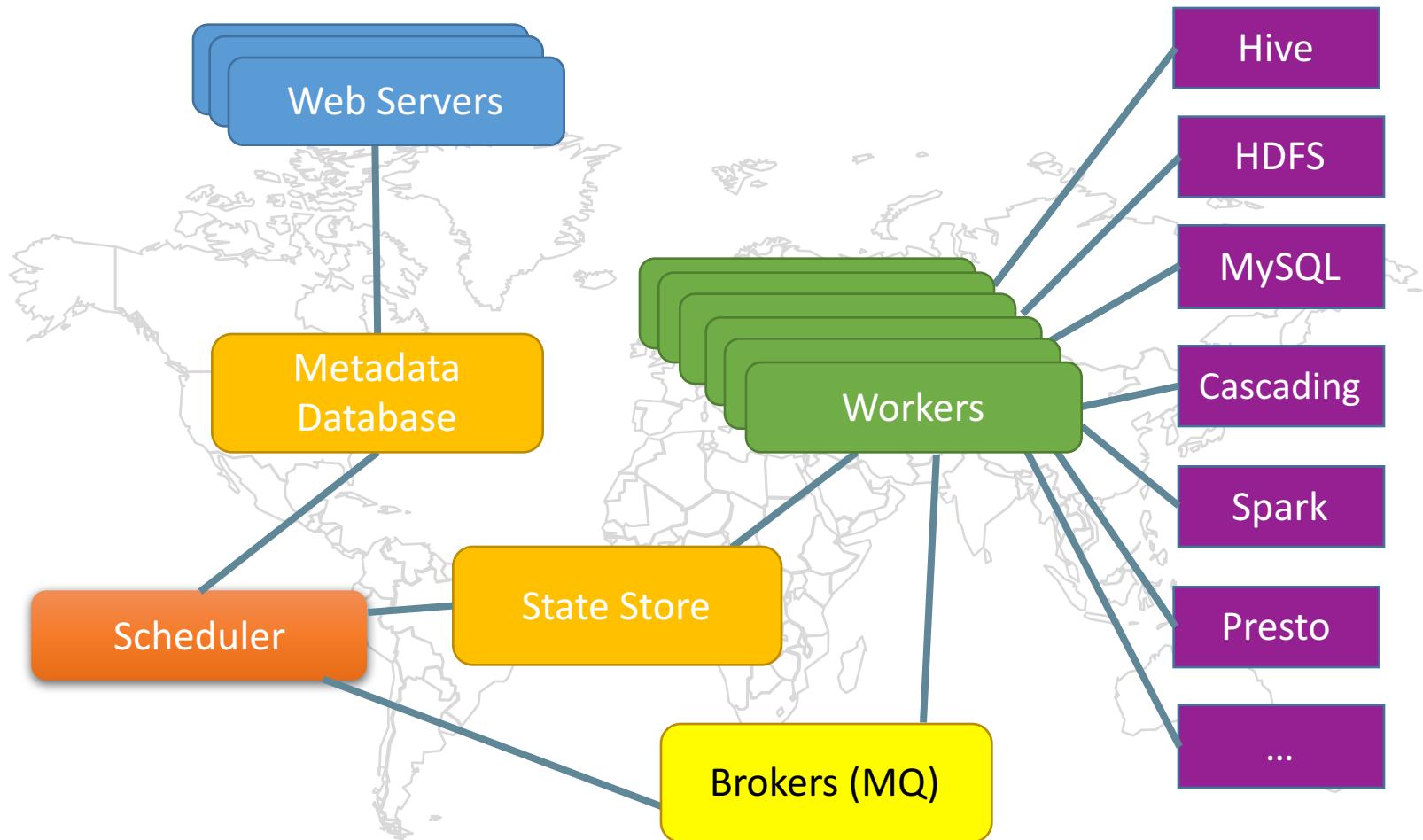
# ● Airflow Architecture (Local Scheduler)



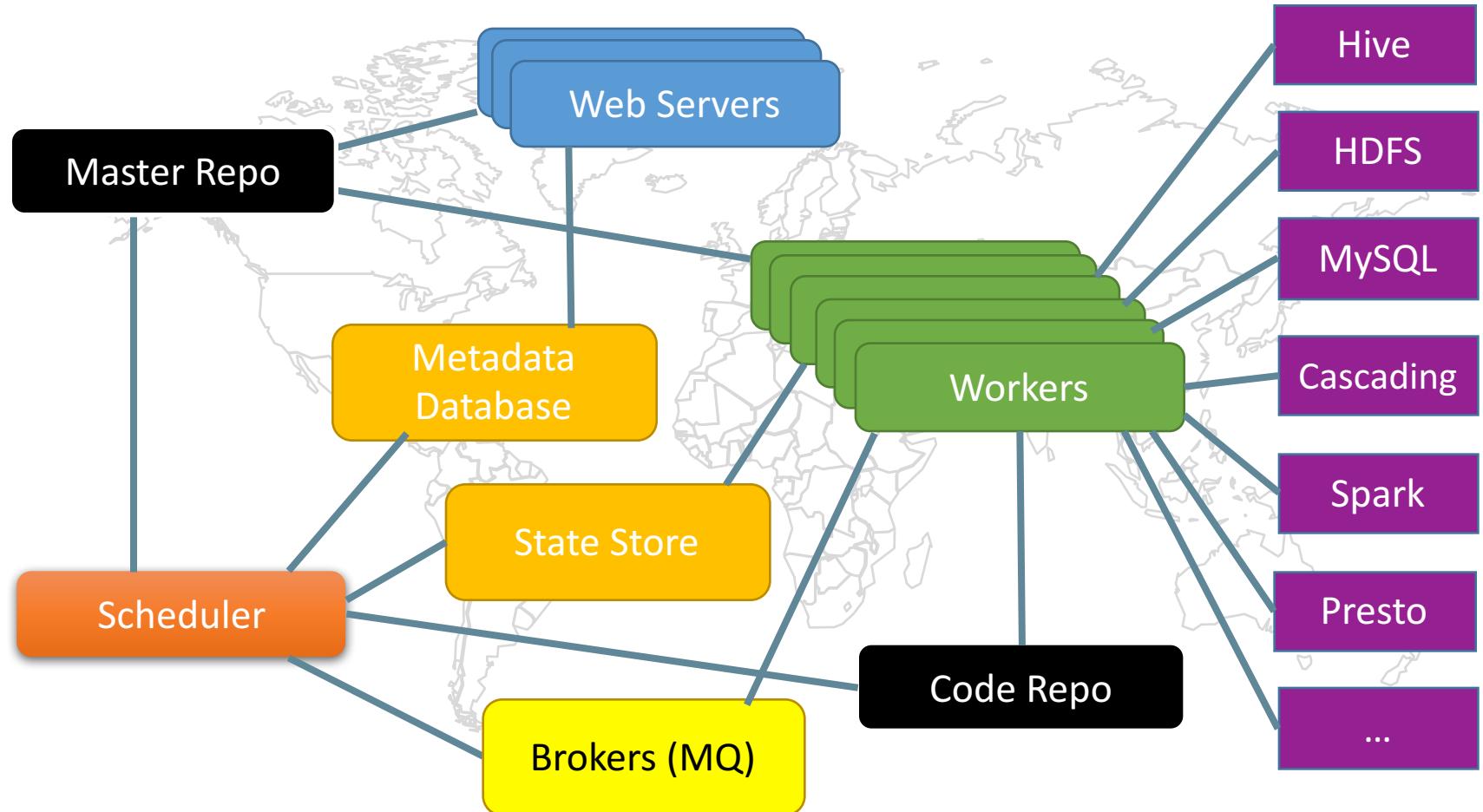
# • Local Scheduler – w/ version control



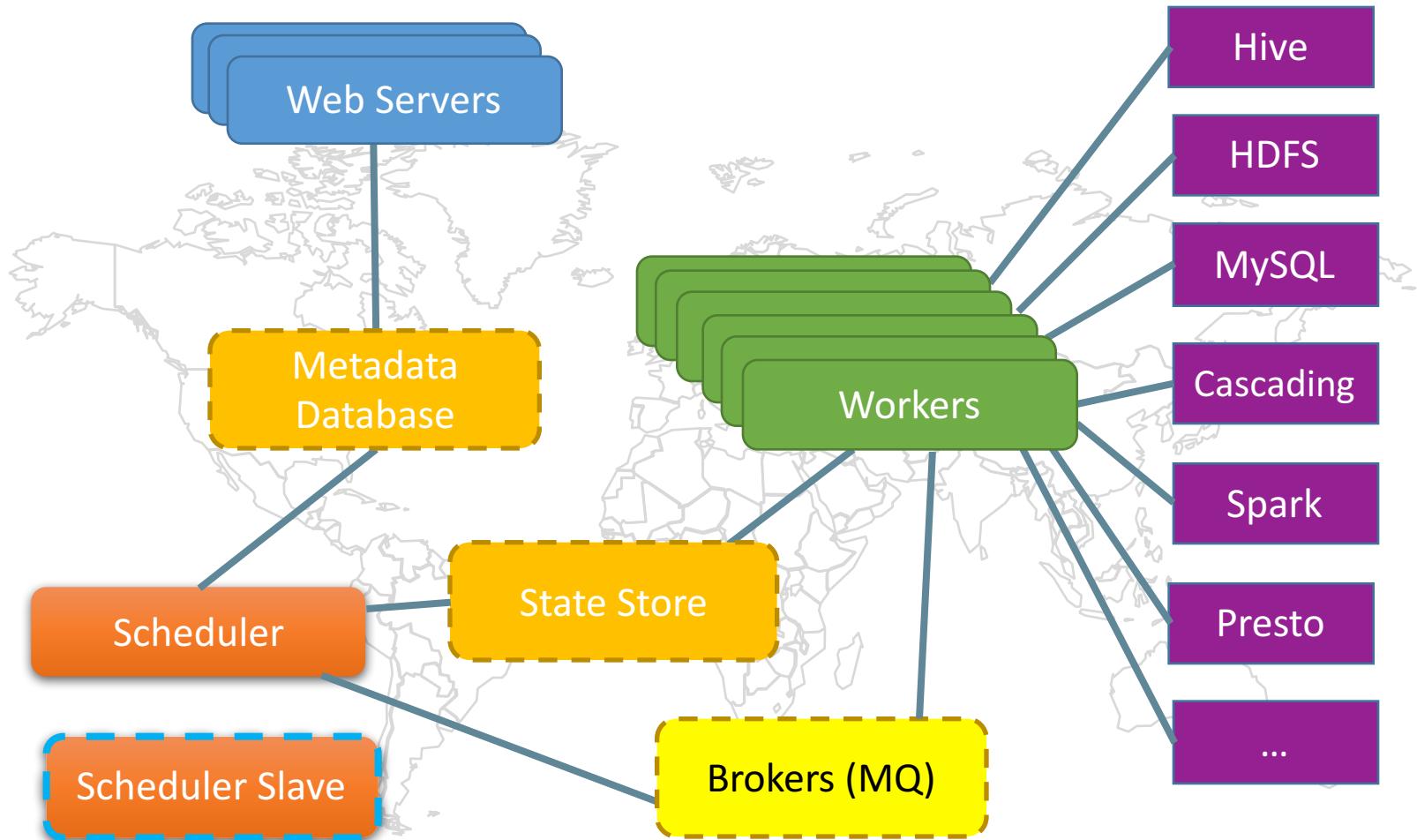
# • Airflow Architecture (Celery Scheduler)



# • Celery Scheduler – w/ version control



# • Airflow Architecture - HA



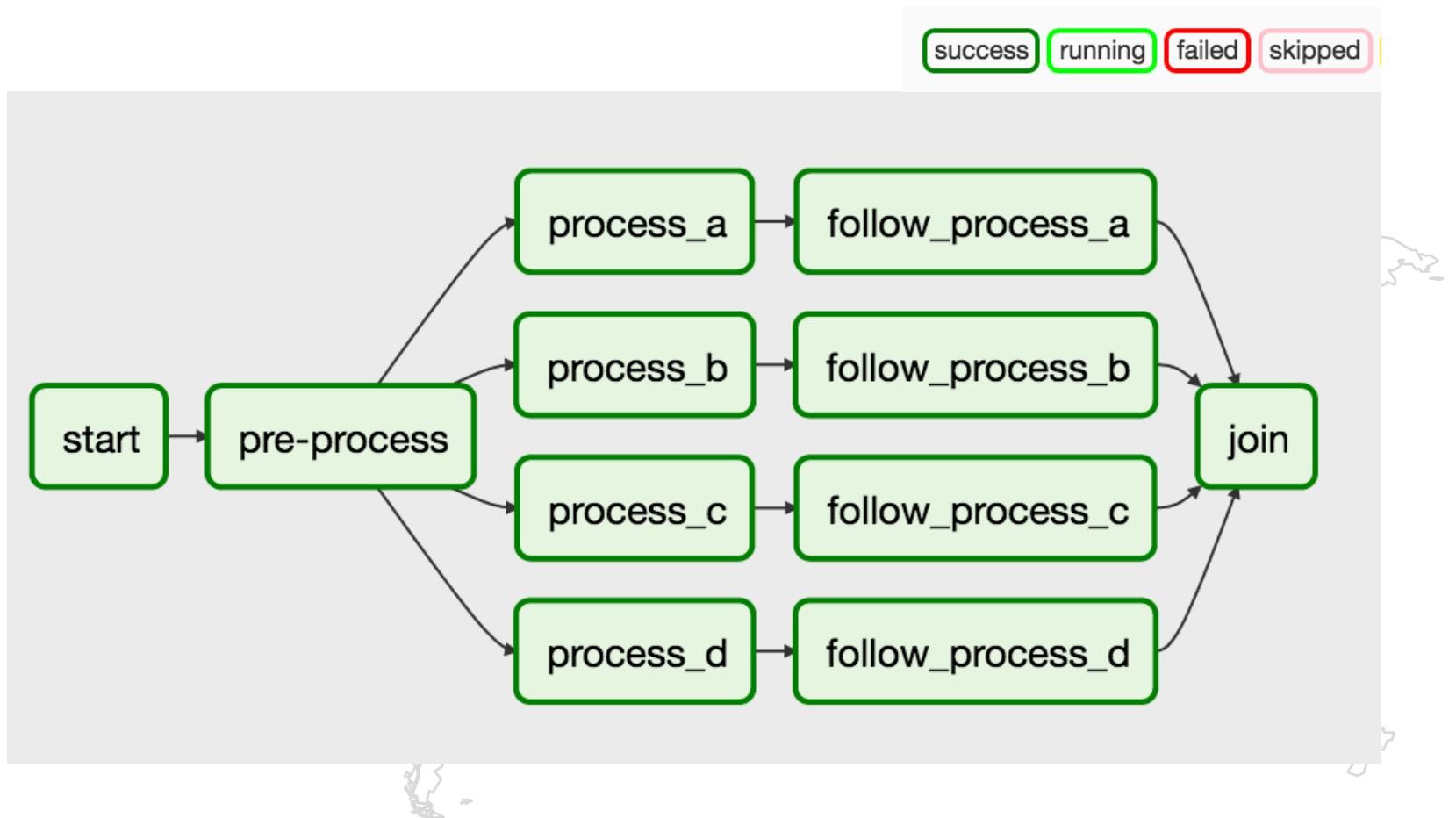


# Workflow Patterns

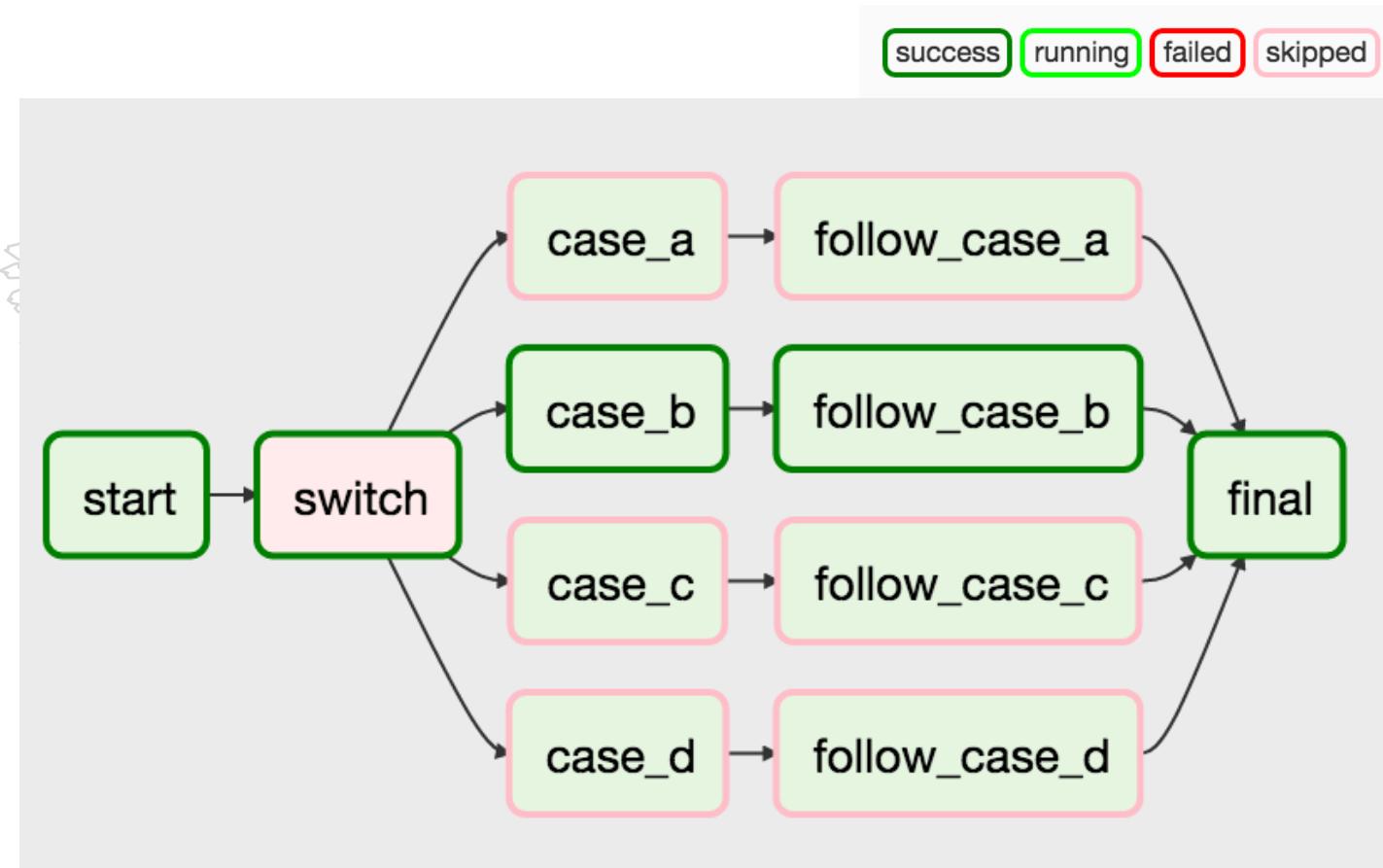


Complex Rule

- Process in parallel

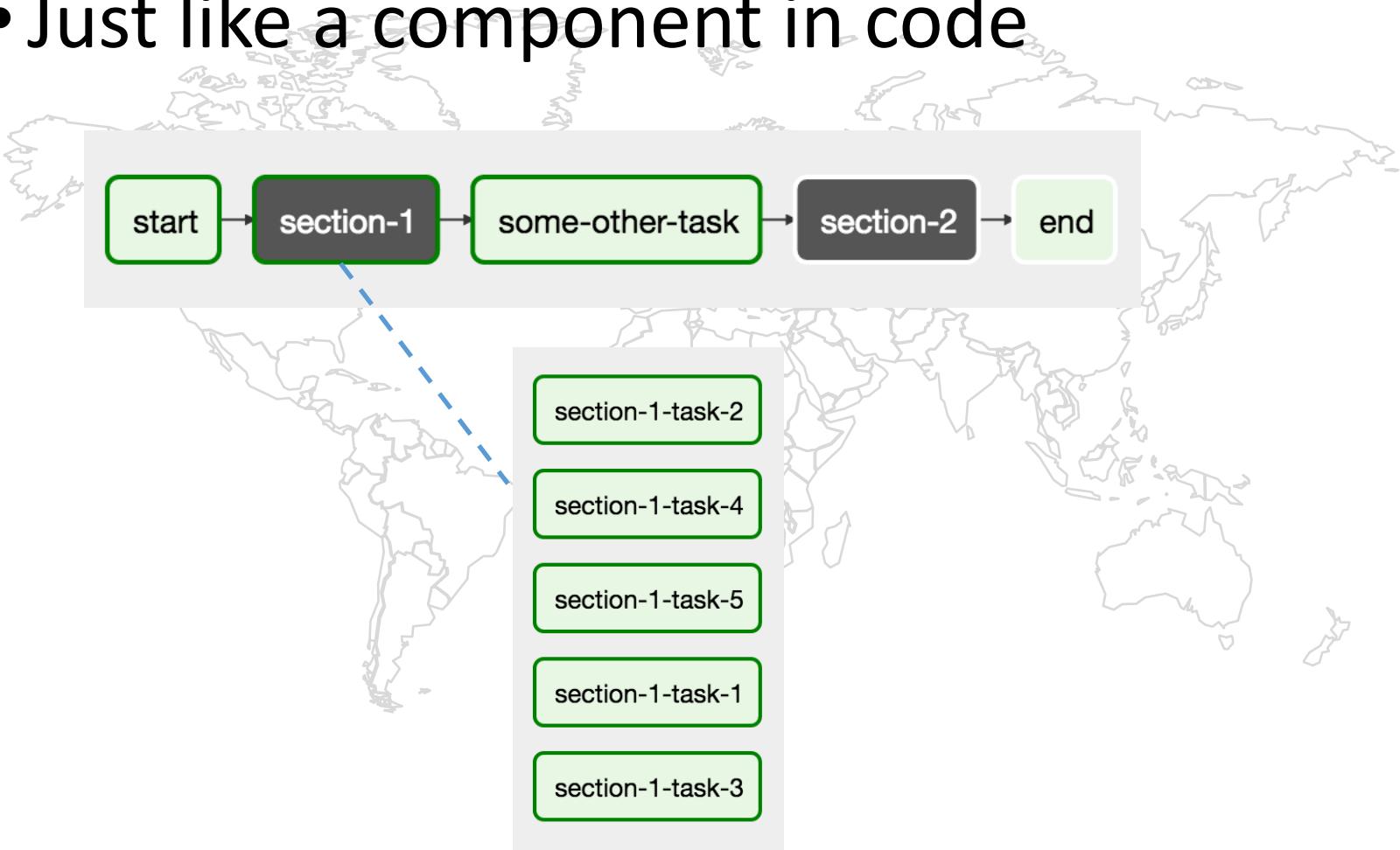


# ● Switch

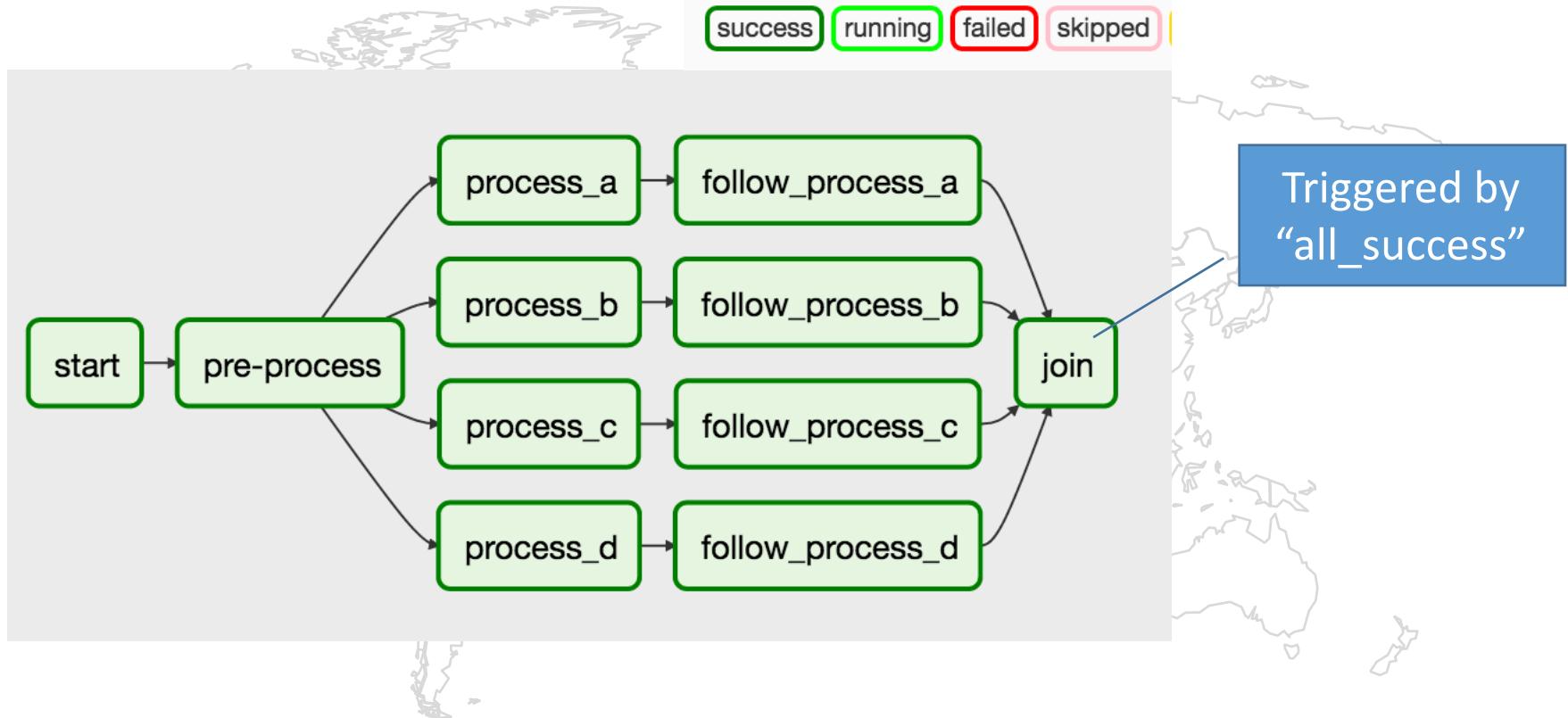


- Sub Dag

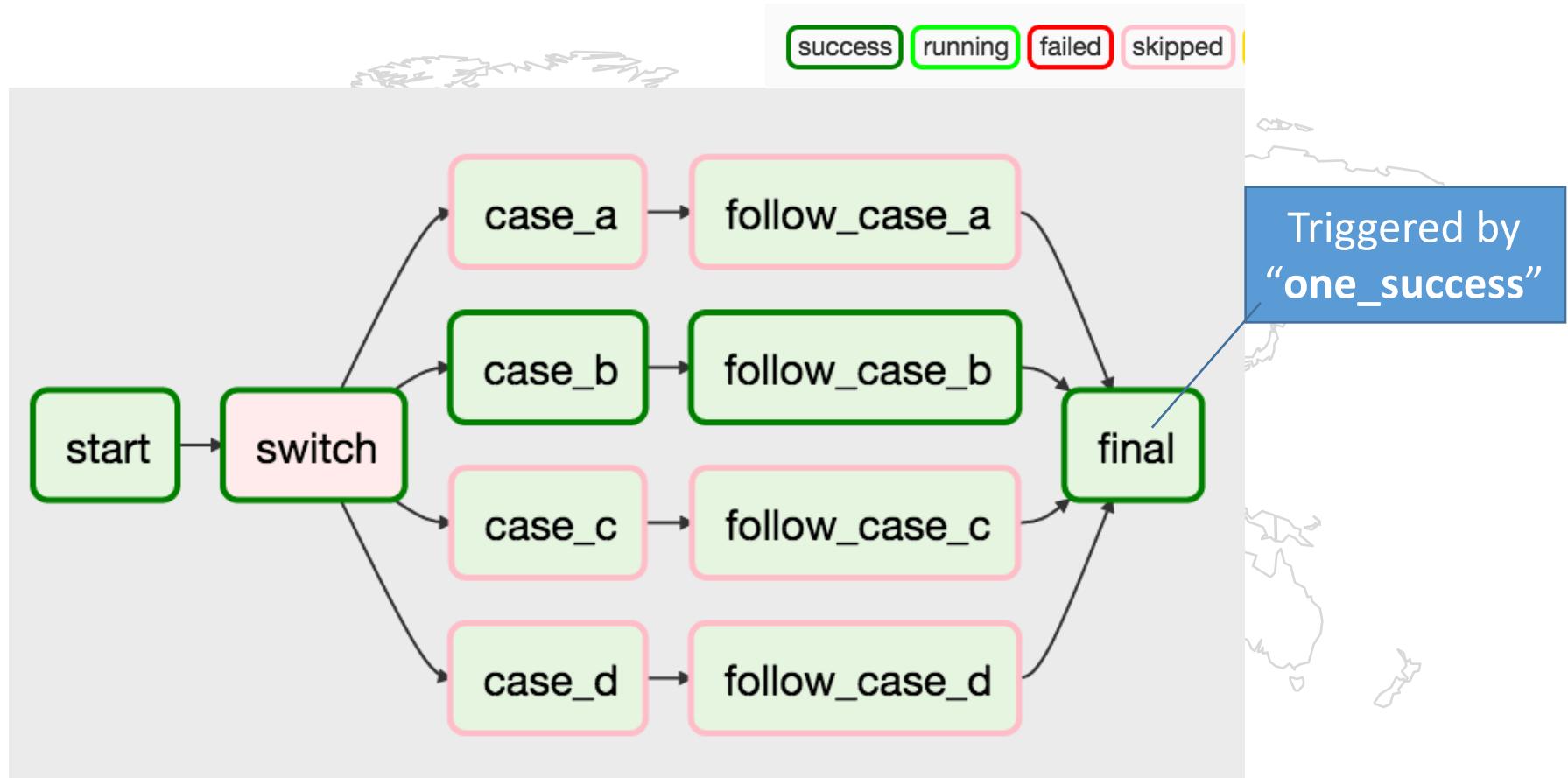
- Easier to control, re-use and test
- Just like a component in code



# • Trigger Rule – all success

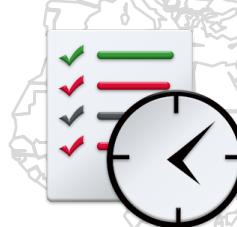


- Trigger Rule – one success



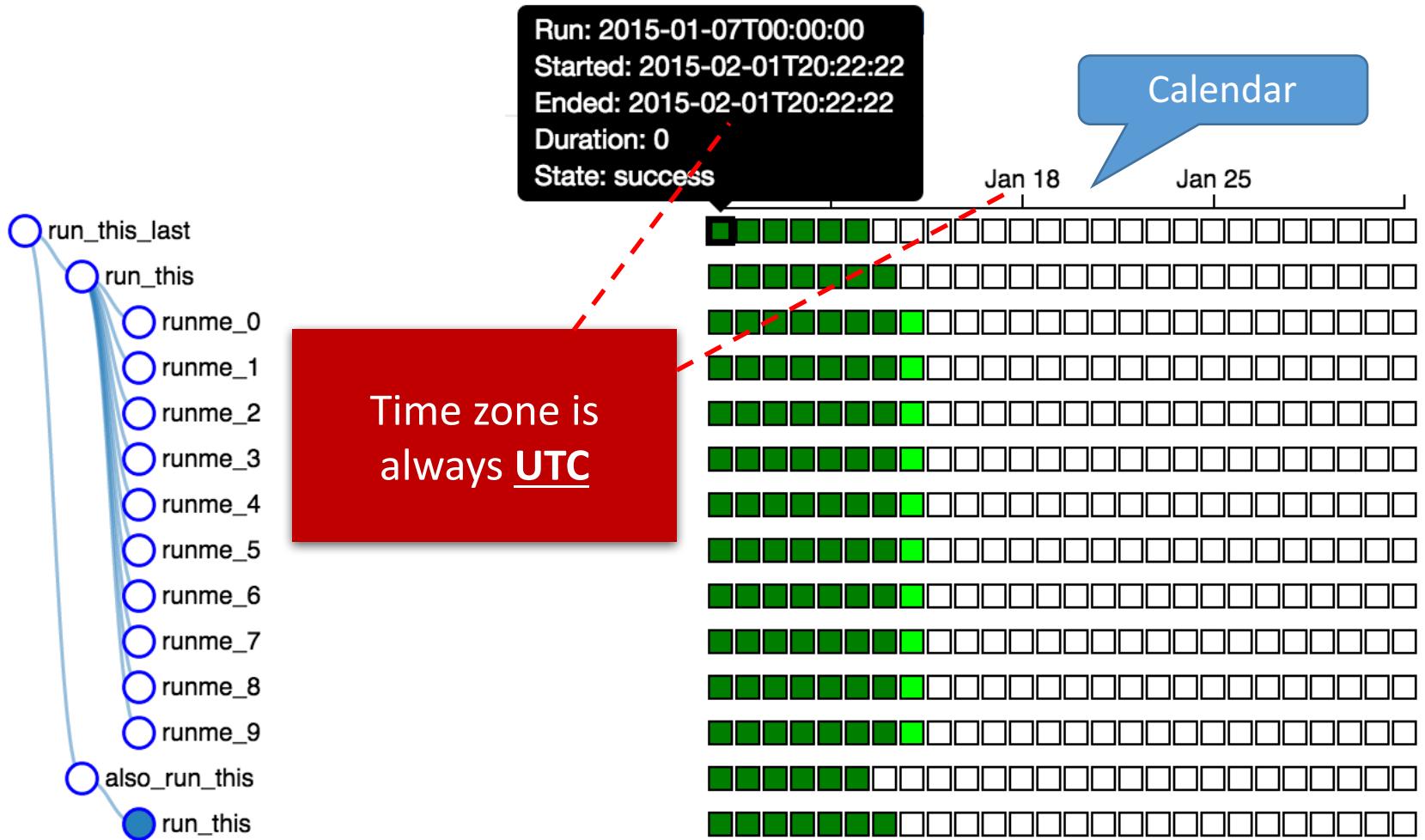


# Scheduling Practice



Calendar Based Scheduling

# • Calendar based scheduling (UTC)



# ● Scheduler – interval in workflow

```
dag = DAG(dag_id="mydag",
           default_args= {
               'start_date': datetime(2016,9,8),
           },
           schedule_interval='0 */4 * * *')
```

Run:

scheduled\_2016-09-09T16:00:00

- ✓ scheduled\_2016-09-09T16:00:00
- scheduled\_2016-09-09T12:00:00
- scheduled\_2016-09-09T08:00:00
- scheduled\_2016-09-09T04:00:00
- scheduled\_2016-09-09T00:00:00
- scheduled\_2016-09-08T20:00:00
- scheduled\_2016-09-08T16:00:00
- scheduled\_2016-09-08T12:00:00
- scheduled\_2016-09-08T08:00:00
- scheduled\_2016-09-08T04:00:00
- scheduled\_2016-09-08T00:00:00**

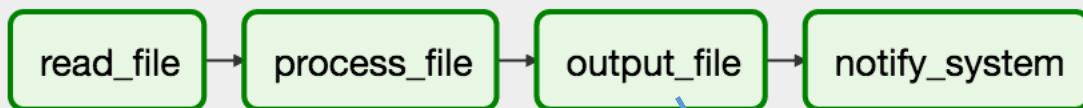
## Note

Every Dag Run will  
only start when next  
Dag Run's execution  
time meets

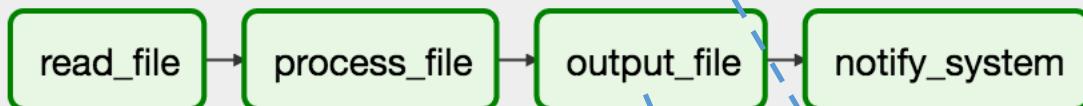
Run at 04:00:00

- Scheduler – recursive running

Run1 (09-01 00:00)



Run2 (09-01 04:00)



What if output\_file in Run1 and Run2 impact each others?

Idea

Try best to avoid this kind of design

# • Principle when defining Task

## Granularity

Task granularity should be proper

- Choose “Right size” for one task
- Task should execute simultaneously

## Atomic

Each Task should be atomic

- isolation from concurrent processing
- Either succeed or failure, no grey state
- Failure will not impact the system

# ● Idempotent Task

```
def do_job(*args, **kwargs):
    make_file1()
    make_file2()
    make_file3()
    do_final_process()

def clean_up(context):
    clean_file_if_exist(
        'file1_path',
        'file2_path',
        'file3_path')

task = PythonOperator(
    task_id='file_operation',
    provide_context=True,
    python_callable=do_job,
    on_retry_callback=clean_up,
    dag=dag)
```

Cleanup env  
when failure

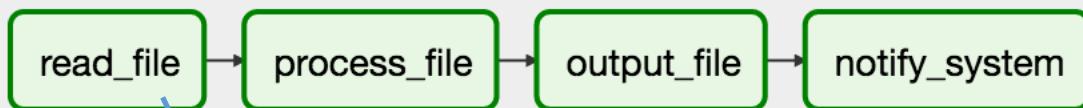
It's ideal case, in real cases...



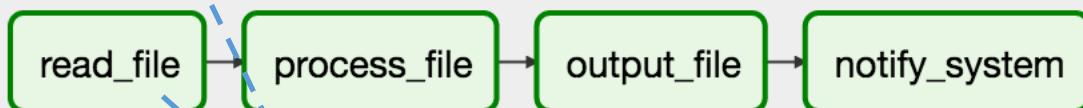
图样图森破

# • Scheduler – recursive dependency

Run1 (09-01 00:00)



Run2 (09-01 04:00)



What if `read_file` in `Run1` and `Run2` cannot run in parallel due to external system's limitation

Option 1

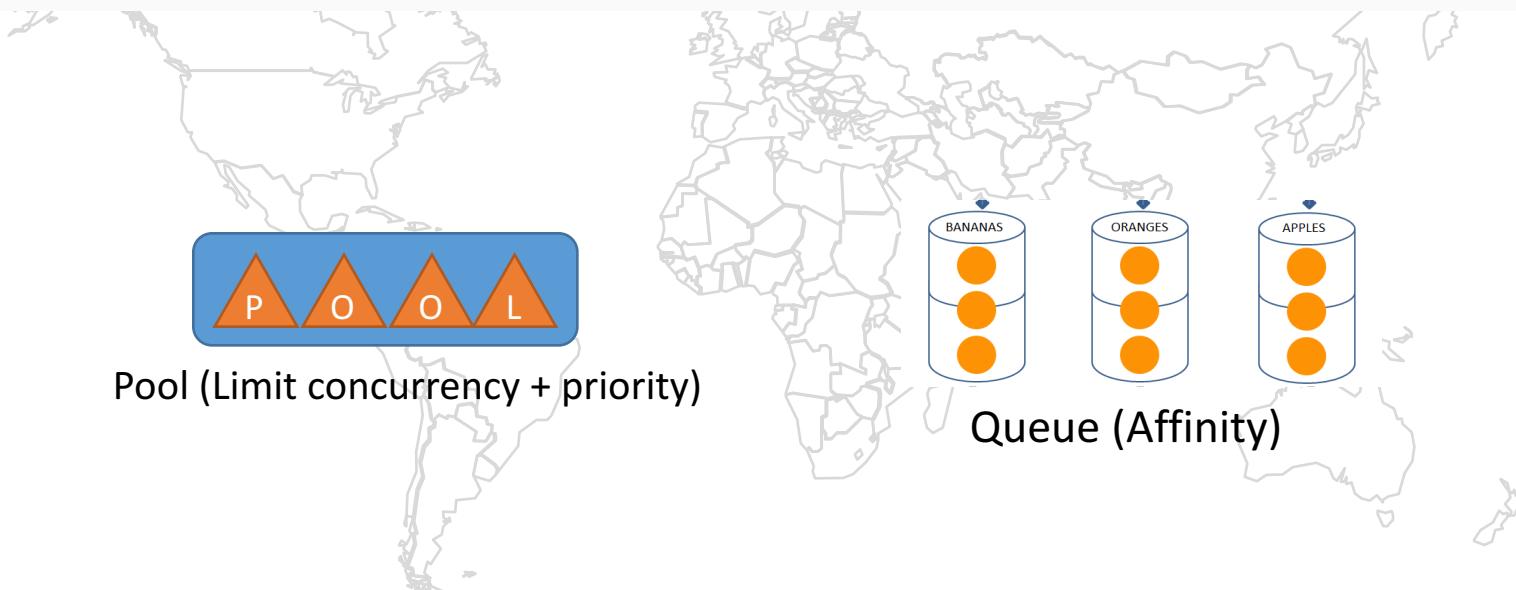
Assign a **pool** with **1 Slots** to for task `read_file`

Option 2

Turn on option “`depends_on_past`” for task `read_file`

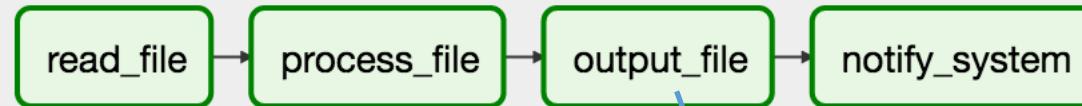
# ● Resource Control

	List (2)	Create	With selected▼			
				Slots	Used Slots	Queued Slots
<input type="checkbox"/>		Pool				
<input type="checkbox"/>	 	db connection		10	0	0
<input type="checkbox"/>	 	AWS CloudTrial Connections		10	0	0

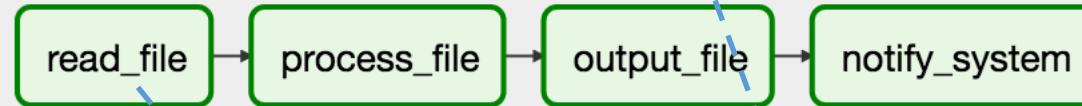


# • Scheduler – more recursive dependency

Run1 (09-01 00:00)



Run2 (09-01 04:00)



What if *read\_file* in *Run2* replies on *output\_file* in *Run1* due to restriction or necessary stateful design?

Option

Turn on option “`wait_for_downstream`” for task `read_file`  
(This will force to turn on “`depends_on_past`”)

# ● Scheduler – recursive dependency pitfall

```
dag = DAG(dag_id="mydag",
           default_args= {
               'start_date': datetime(2016,9,8),
           },
           schedule_interval='0 */4 * * *')-----
```

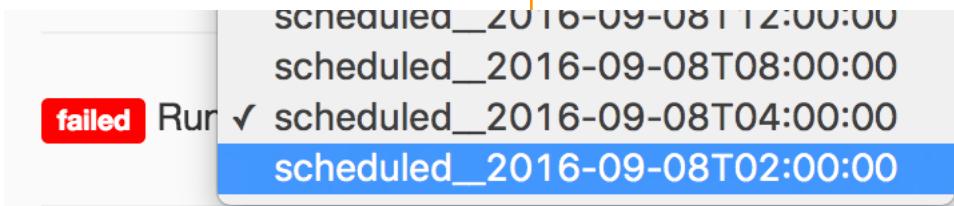
```
'@once'    just one time
 '@hourly': '0 * * * *',
 '@daily':  '0 0 * * *',
 '@weekly': '0 0 * * 0',
 '@monthly': '0 0 1 * *',
 '@yearly':  '0 0 1 1 *'
```

## Note

**start\_date** and **schedule\_interval** should be aligned

2016-09-08 00:00:00 is aligned  
 2016-09-08 02:00:00 is NOT aligned

This will make the DAG failure if the option  
**“depends\_on\_past”** is turned on



- Some other notes

- Update the dag id when changing the logic inside
- Using SLA alert for critical tasks
- Feature in plan:
  - Event Driven Scheduler
  - Mesos Scheduler
  - More operators
  - More syntax sugar



丁来强 wjo1212@163.com

- Now you've learned:

- Definition and ecosystem.
- Challenges and key requirements.
- Solutions and general comparisons.
- Most important part of Airflow and Luigi
  - Architecture, design, patterns, pitfalls and practices etc.

谢谢观看

 wjo1212

 wjo1212@163.com

 LaiQiangDing

