



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

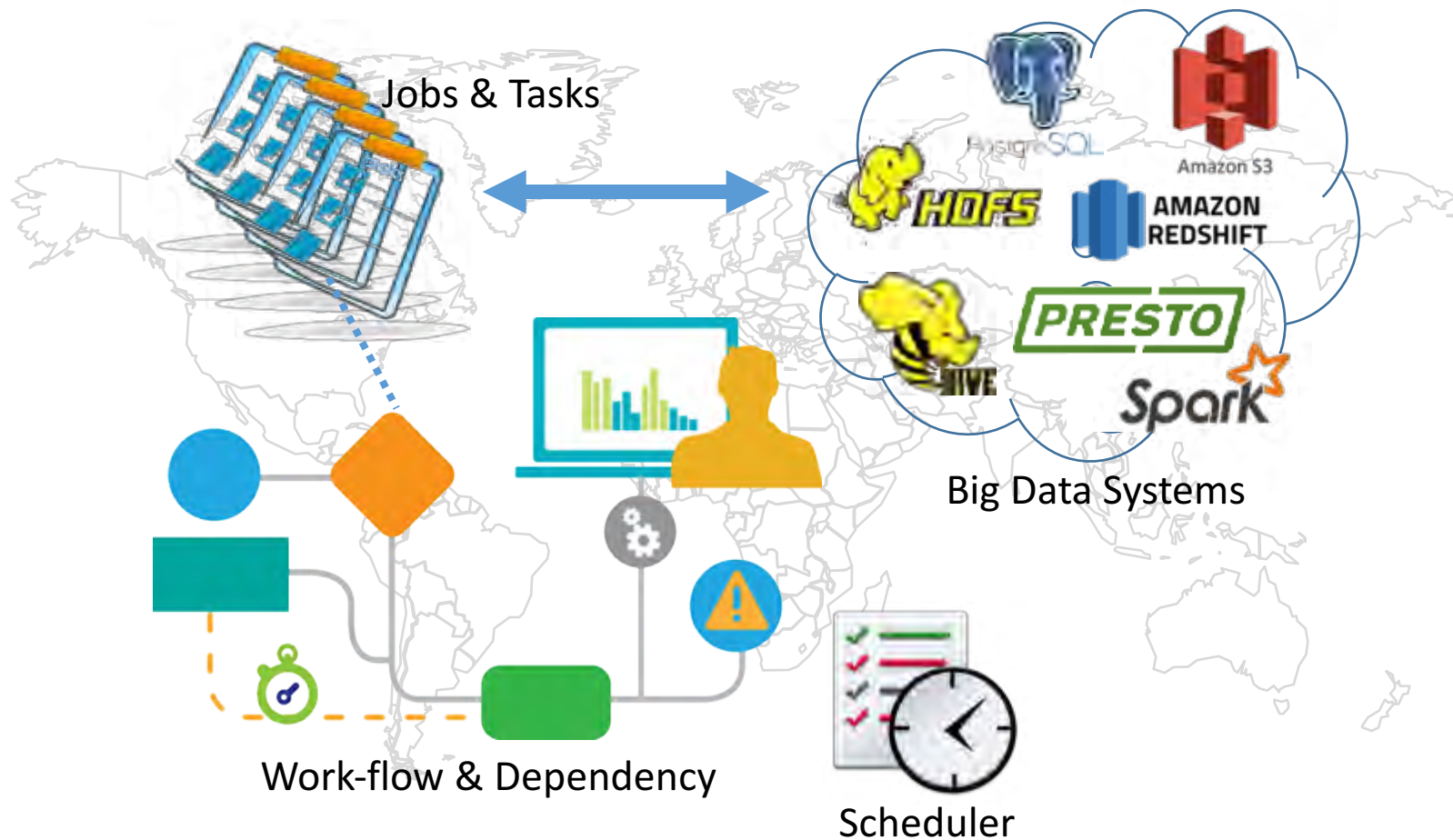


- Definition

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



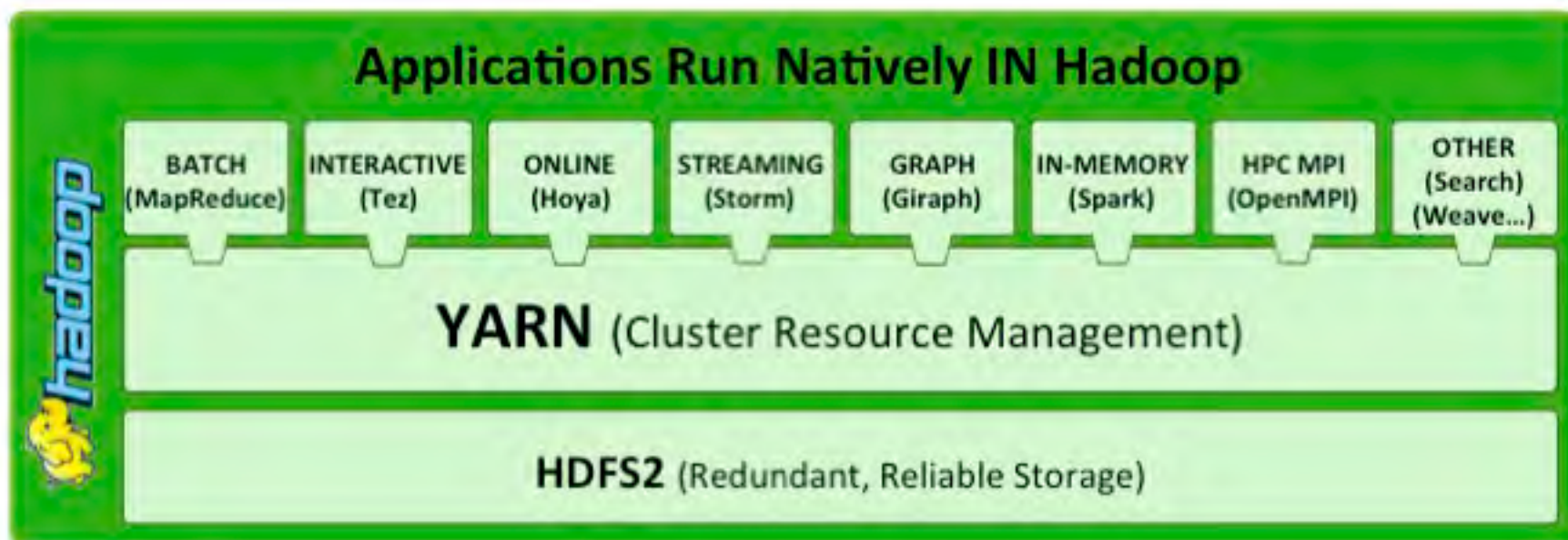
- Different with below categories:

- 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

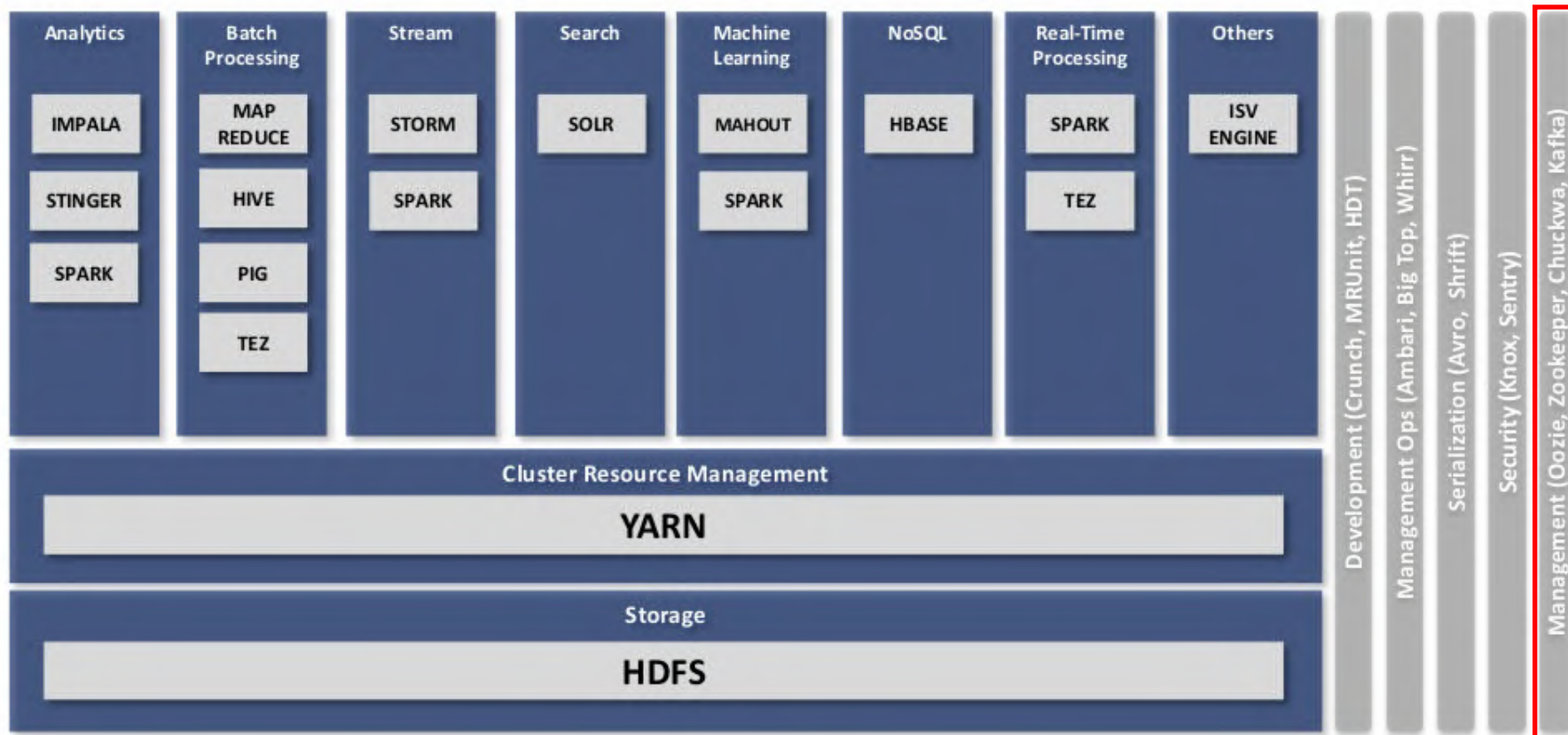


Role in Data Infra

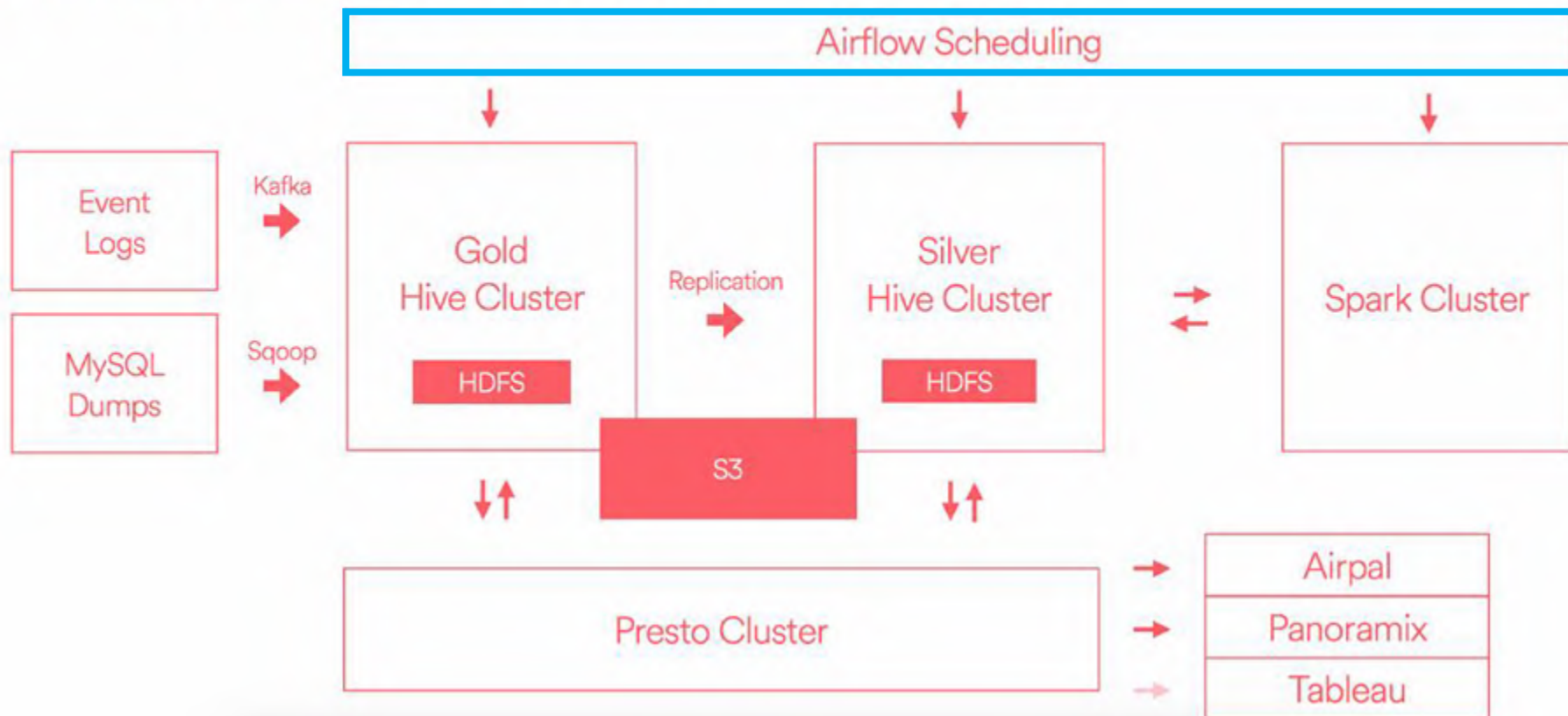
- Hadoop 2.0



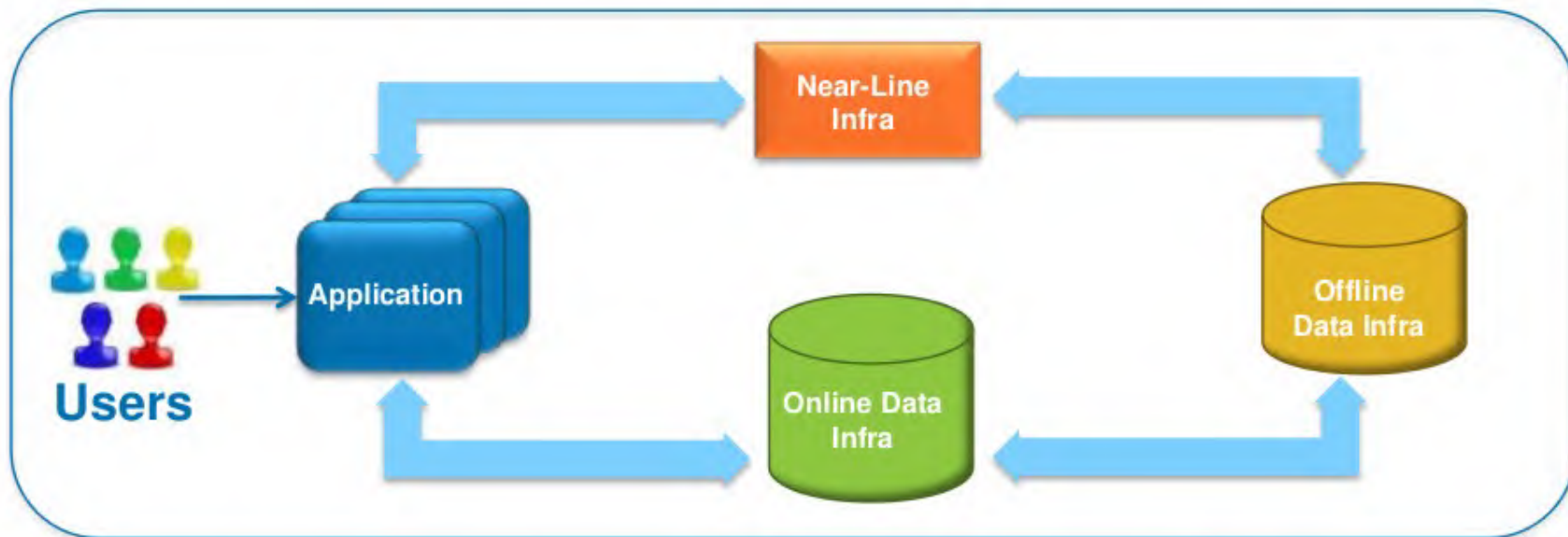
● Hadoop 2.0



● Airbnb Data Infra

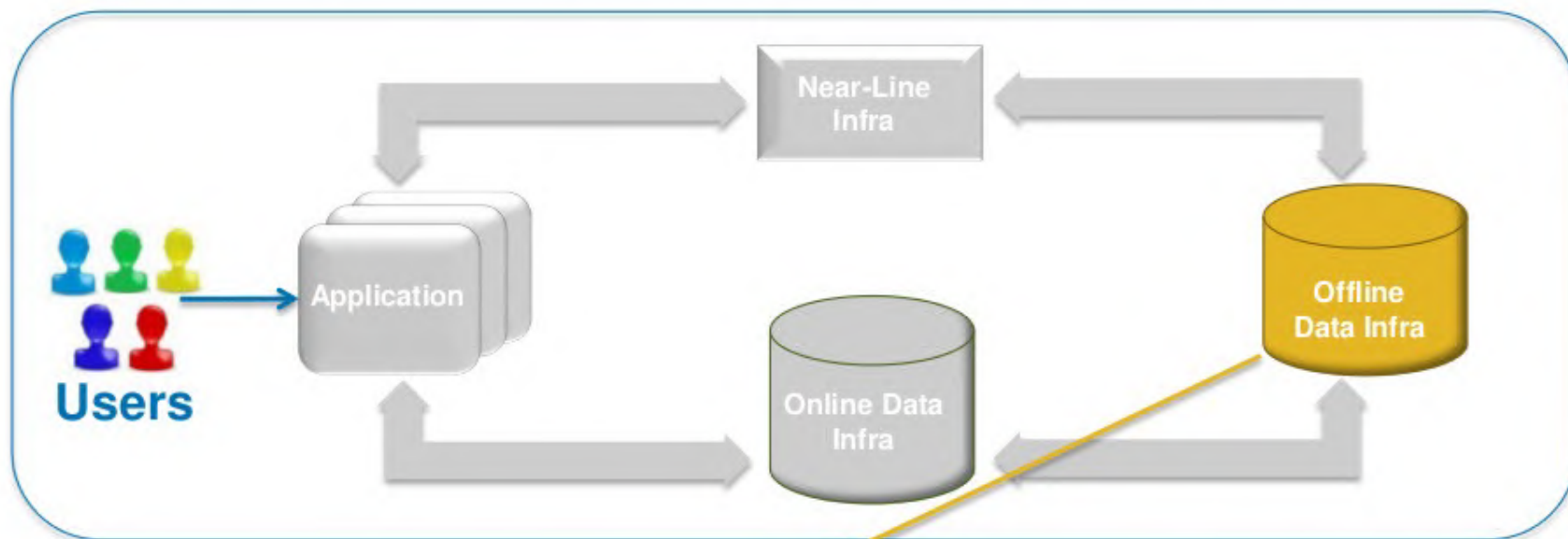


● LinkedIn Data Infra



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

● LinkedIn Data Infra



Systems



Capabilities

- ML, Ranking, Relevance
- Insights and Analytics
- ETL, Metadata and Pipes
- Business Source of Truth

- 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

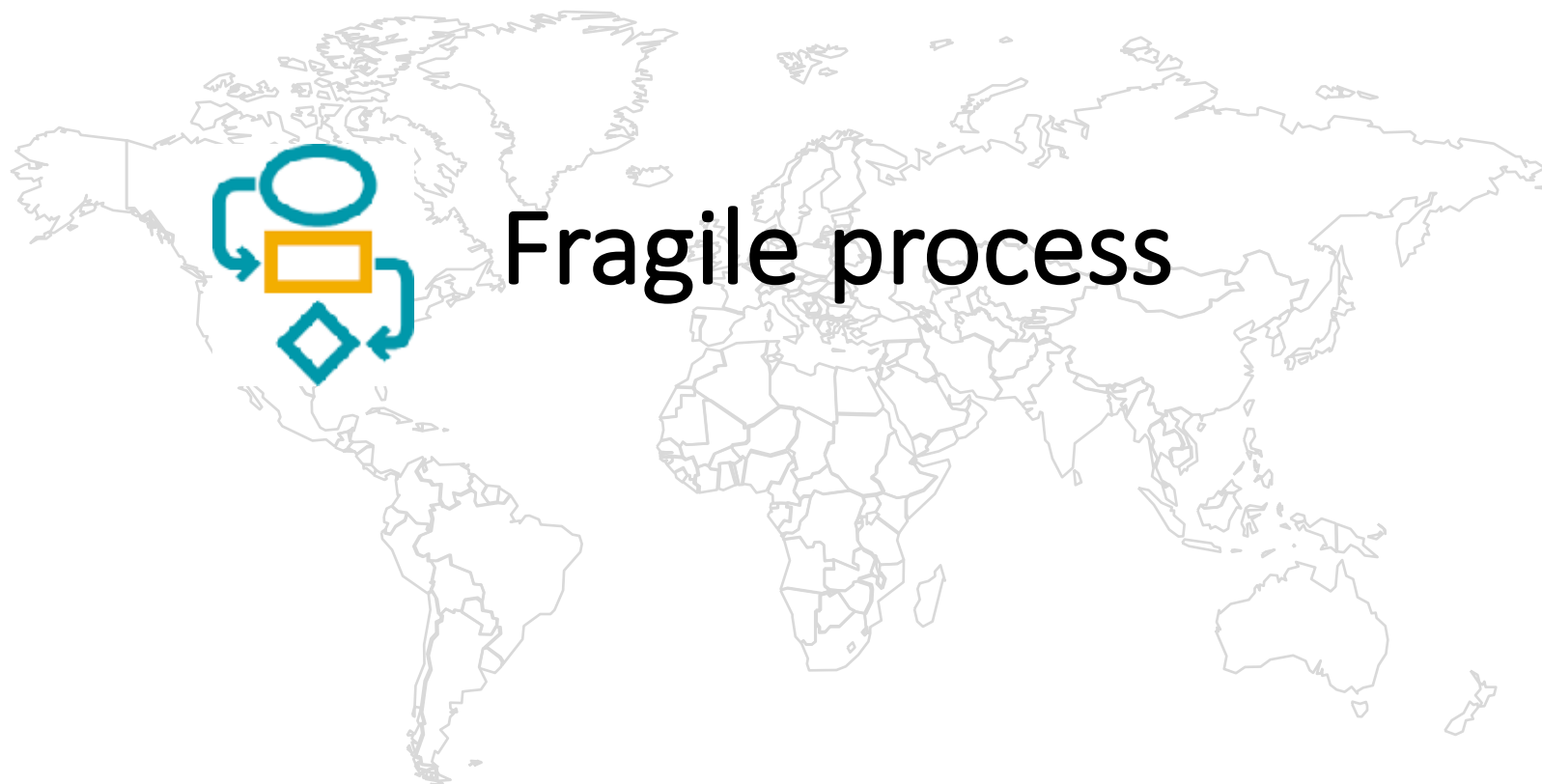




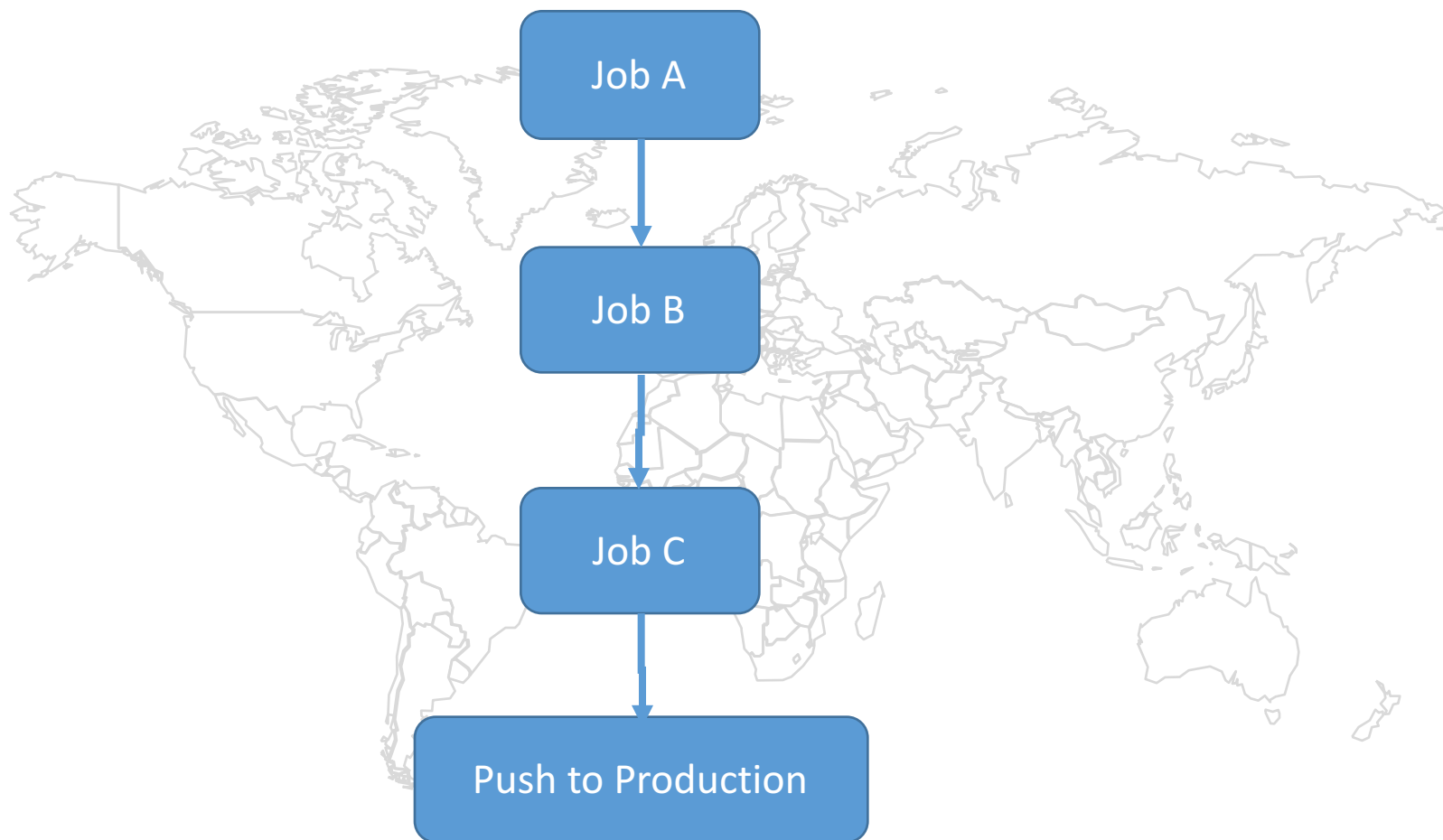
Problems with big data job scheduling

● Typical Challenge

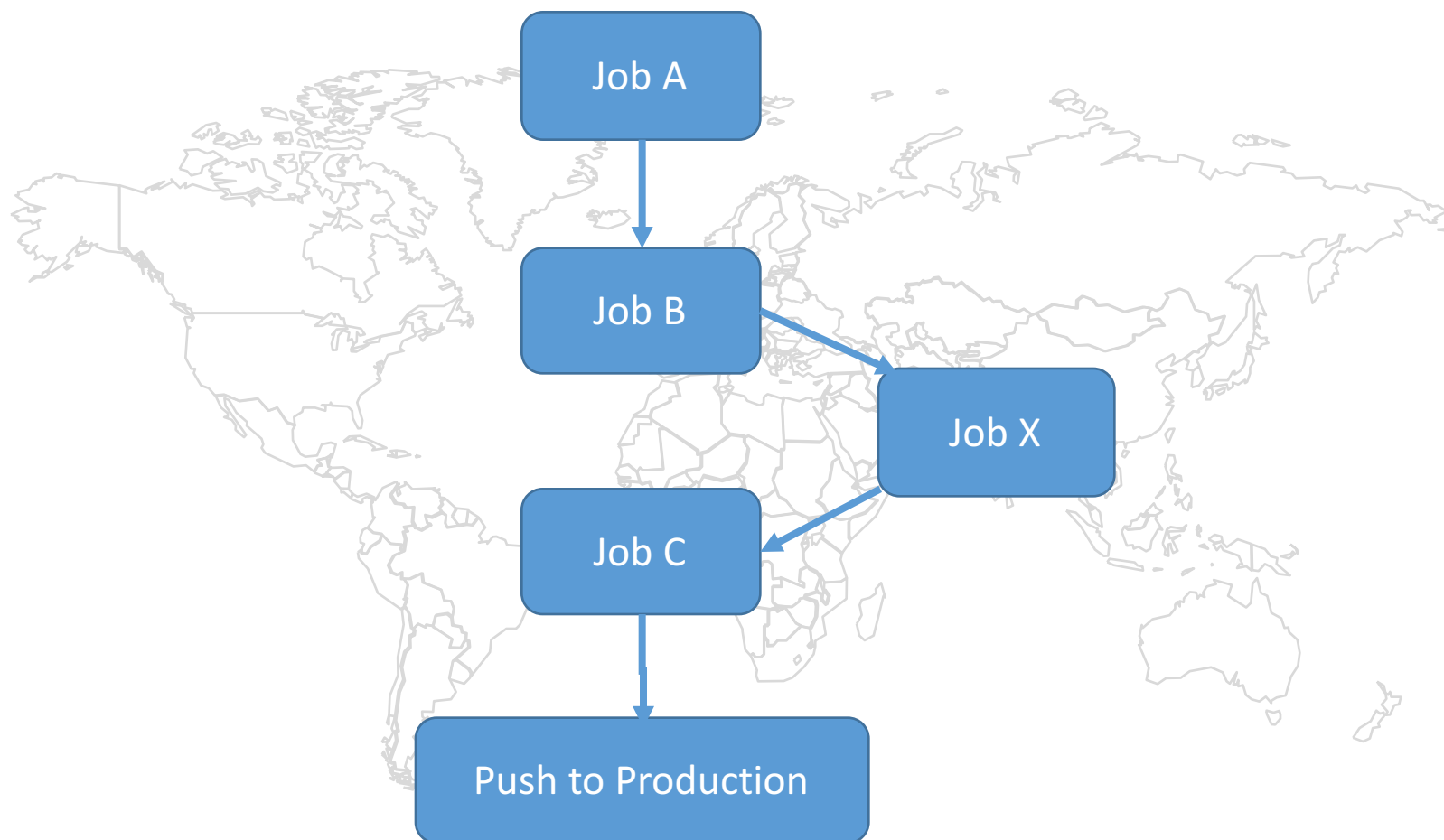
- 数据工作流程复杂度越来越高
- 数据分析与批处理数据非常重要
- 大量时间花费在编写任务、检测与排错上



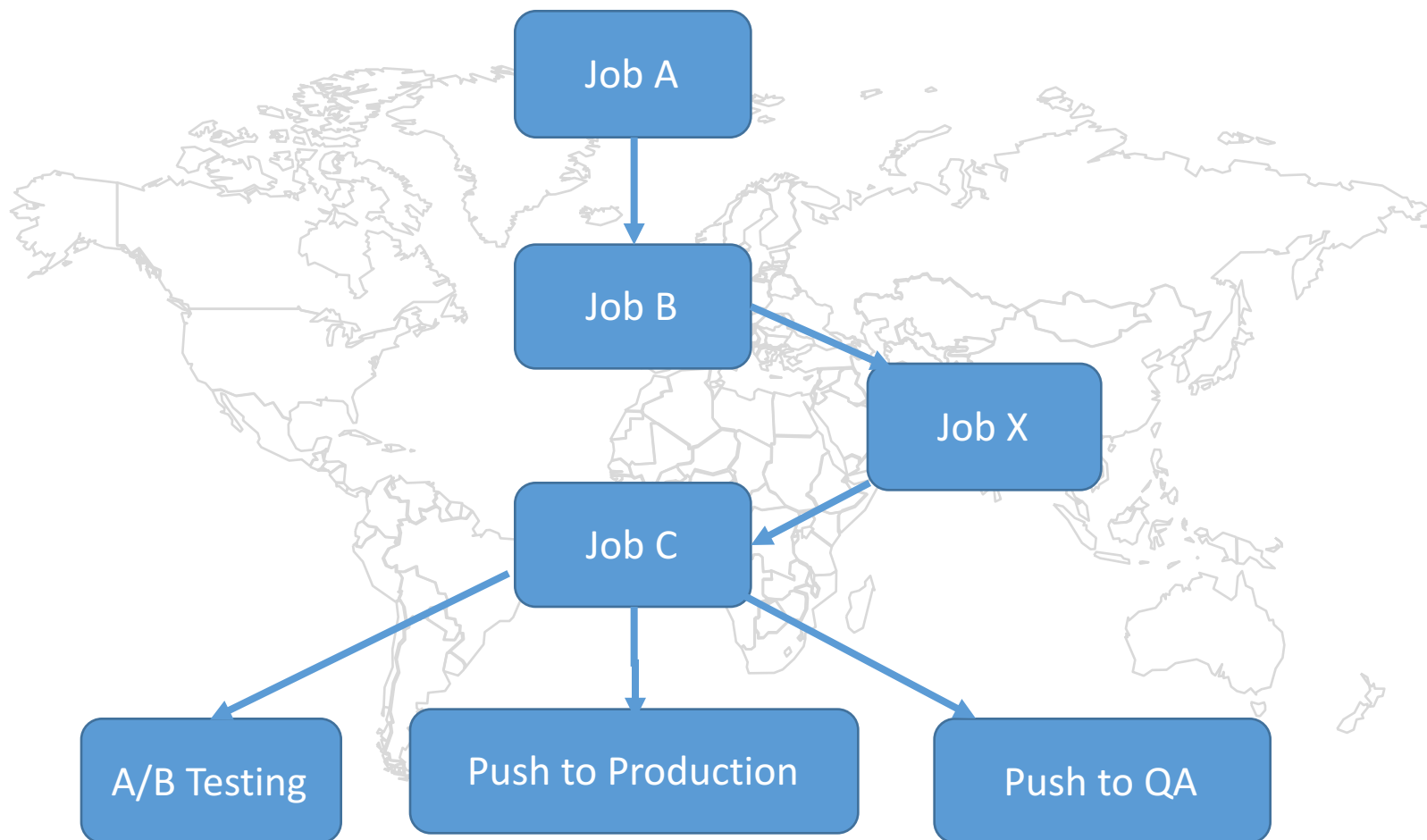
- Fragile process



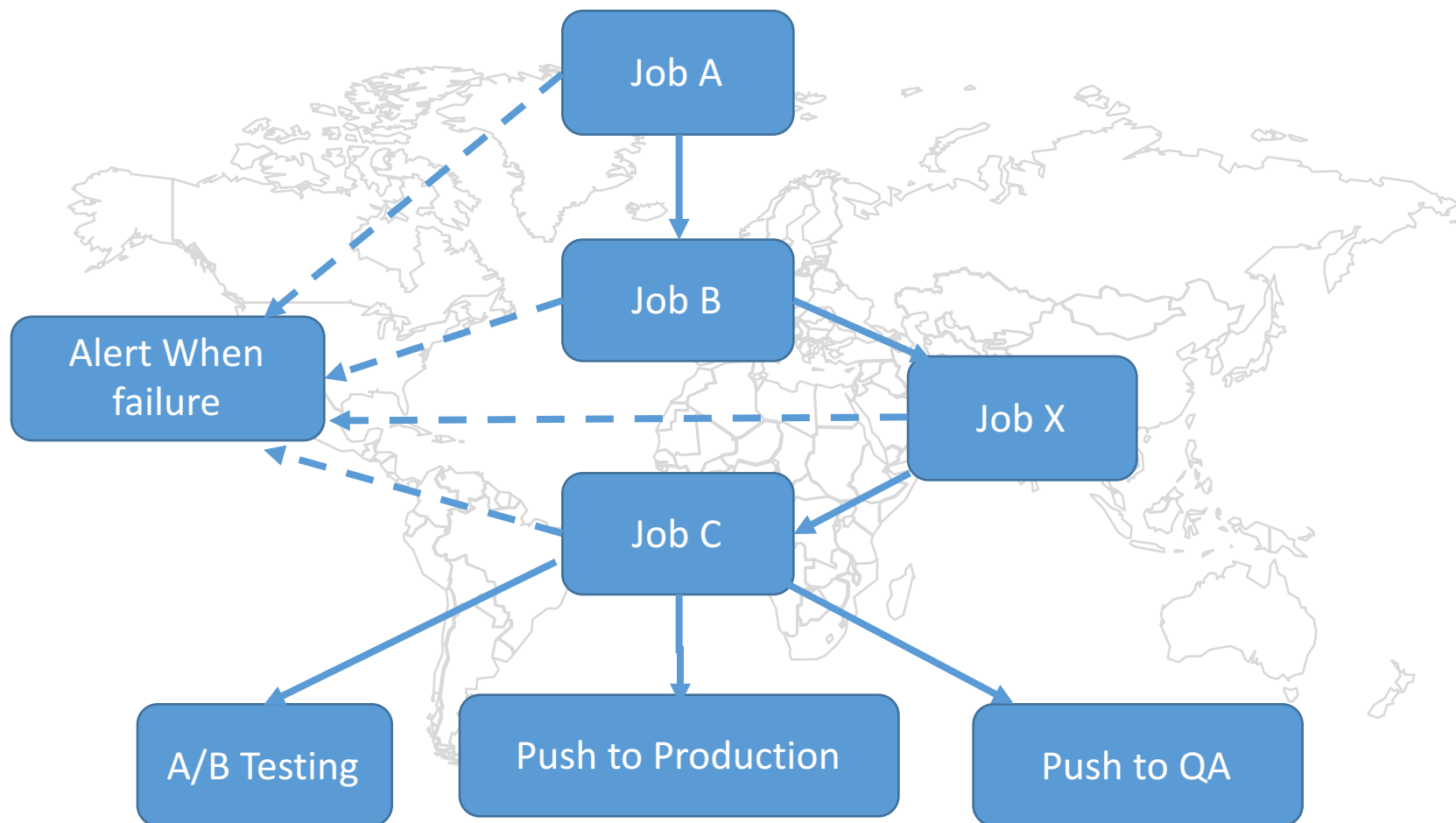
- Fragile process



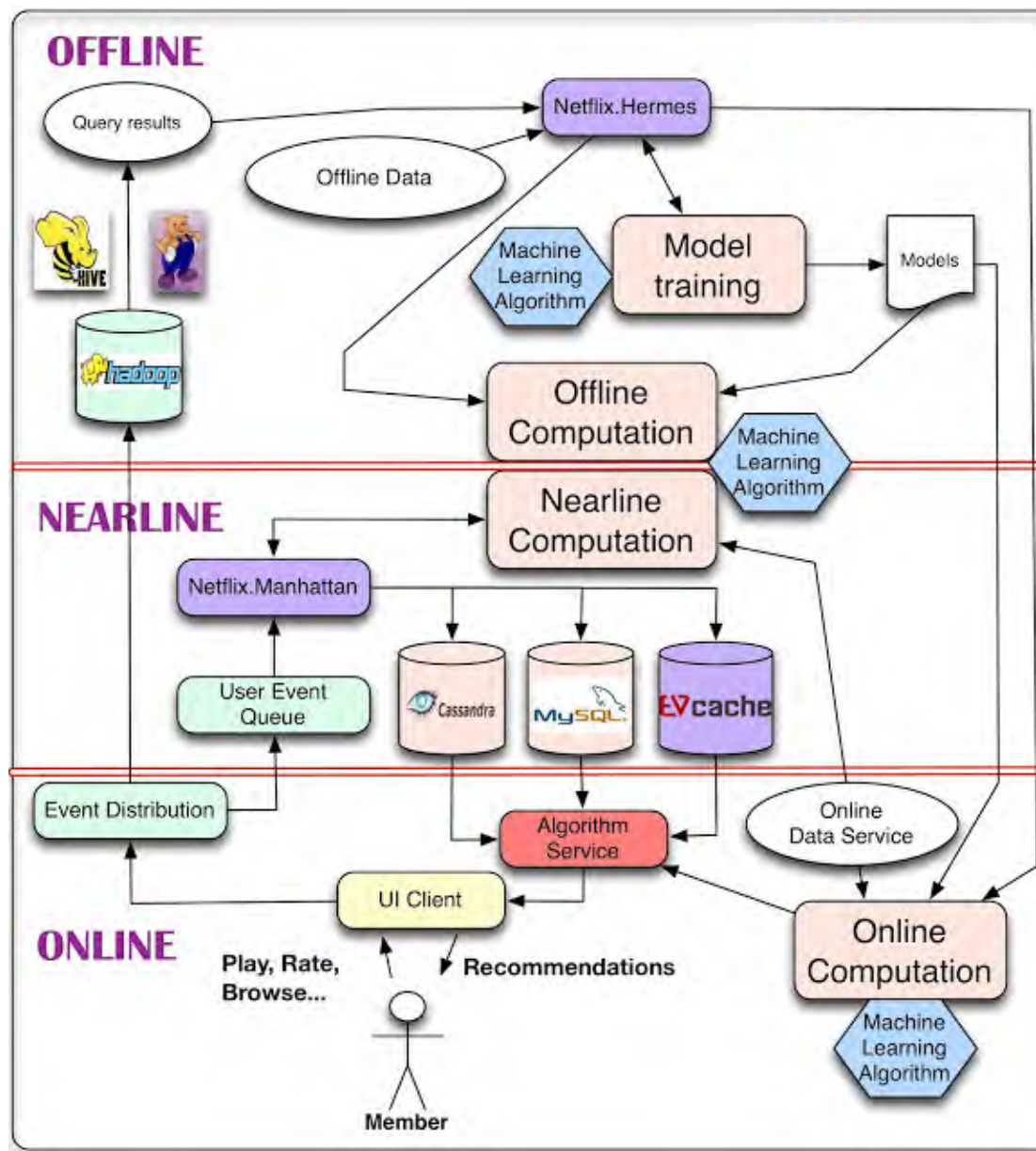
- Fragile process



- Fragile process



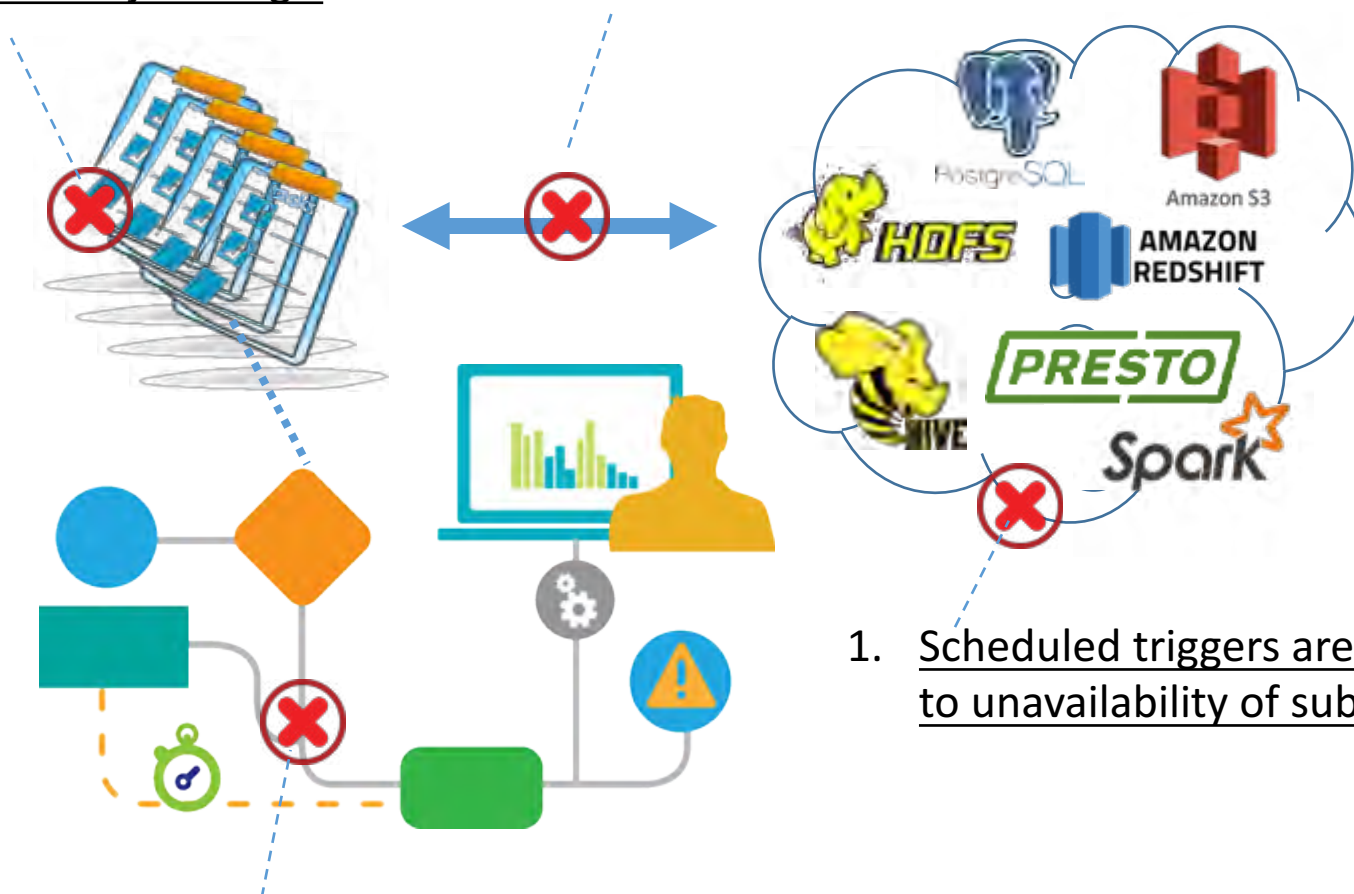
● Example: Netflix Recommendation System





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

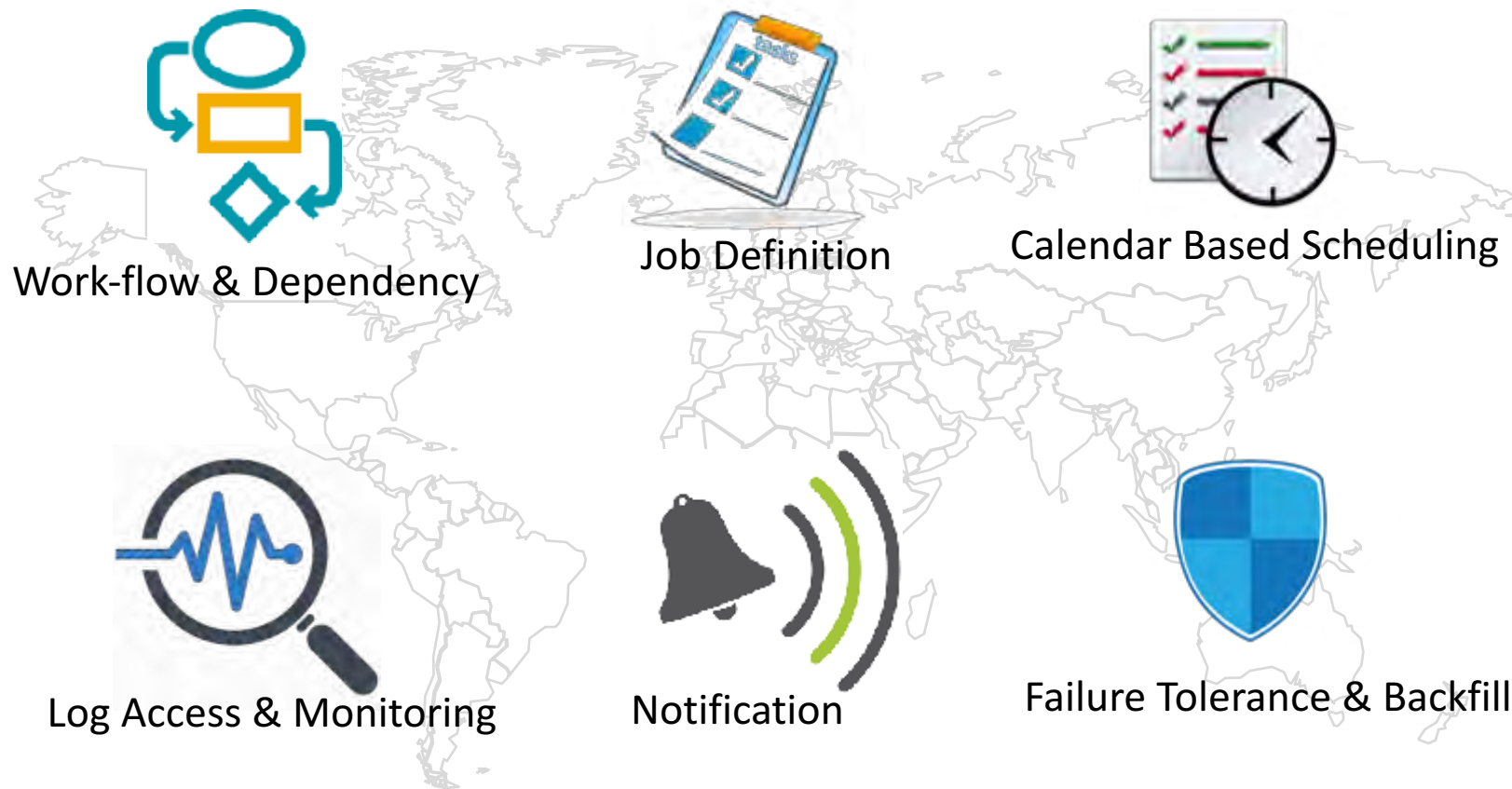


1. Scheduled triggers are skipped due to unavailability of sub-system

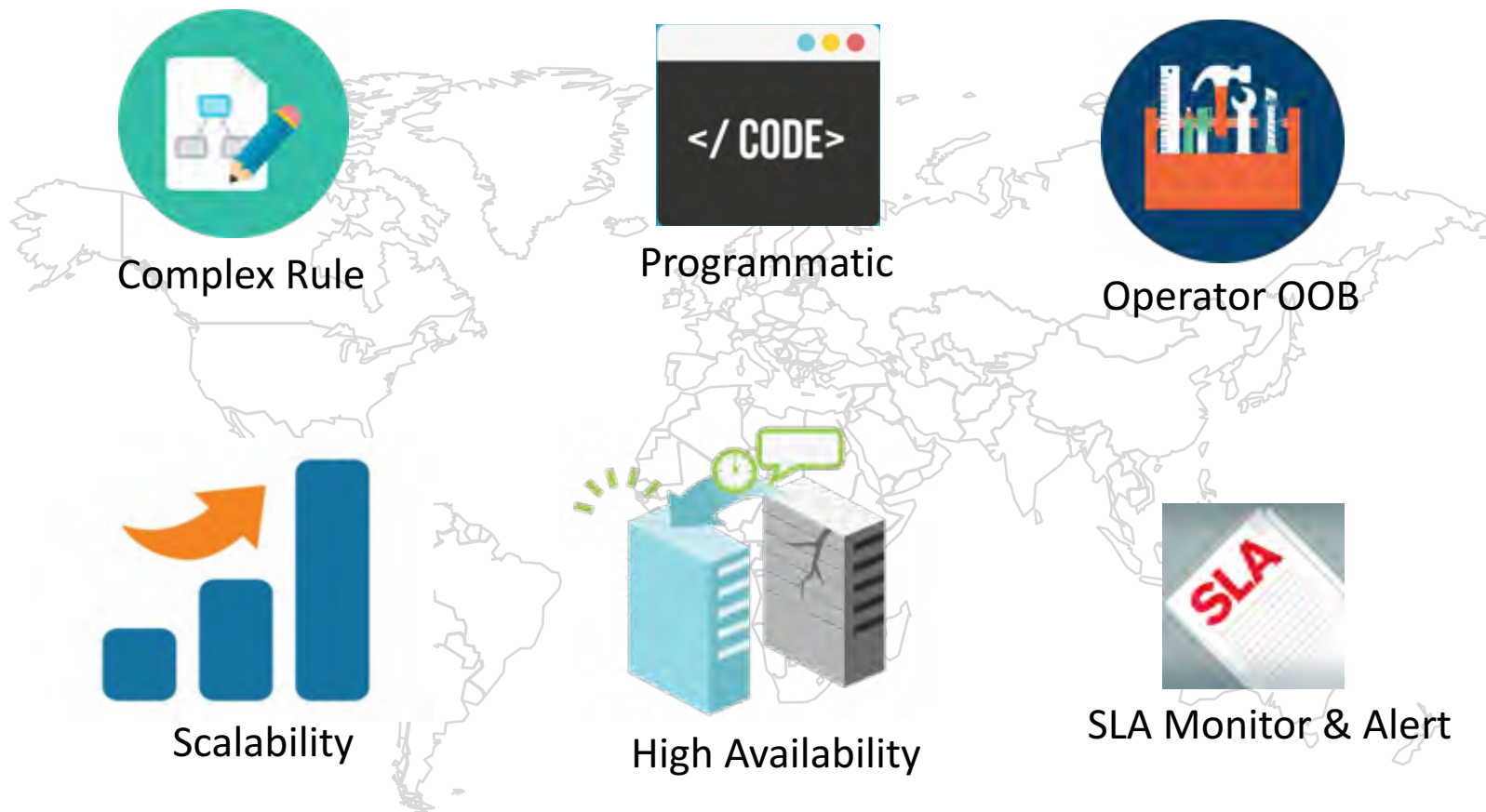
4. Performance is slow especially for some critical steps



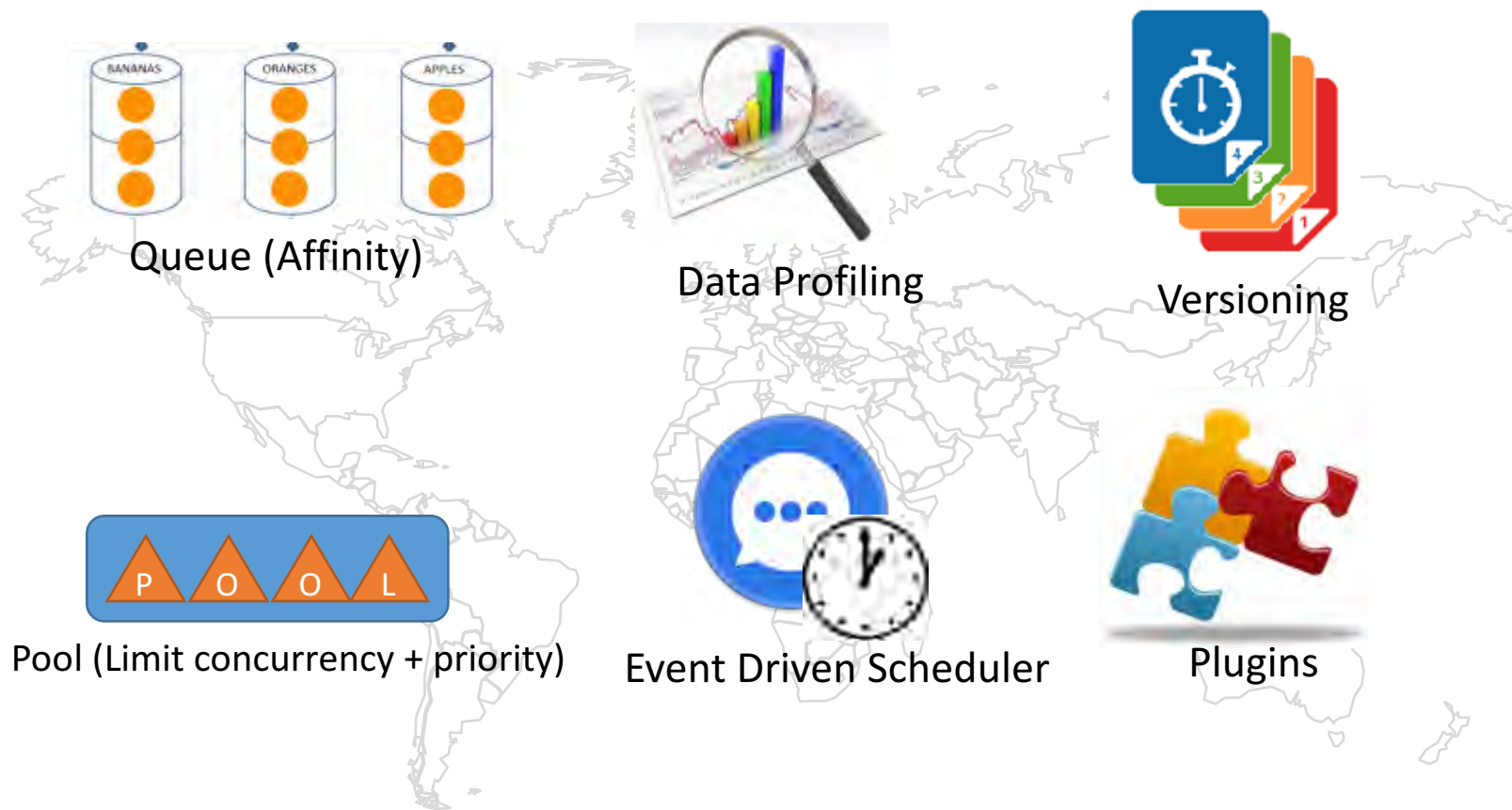
• Basic Needs



• Advanced Needs



Advanced Needs (cont')





Solution Overview

- Options



Airflow



Azkaban

Wt9t

Chronos



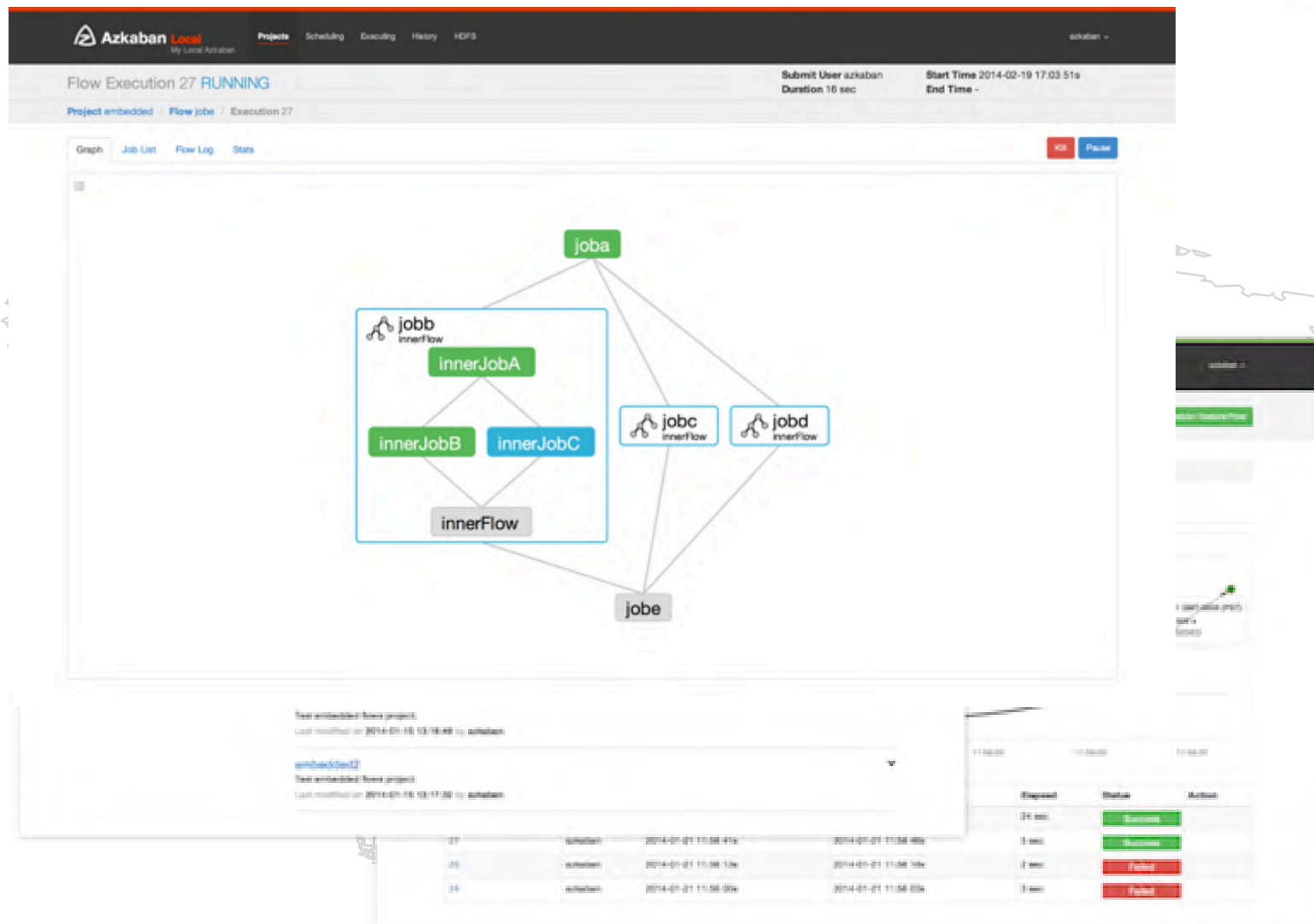
cron

● 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 |

- 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 3rd party AzkabanCLI)
 - Limited execution path control

● Azkaban GUI



● 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



● 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()
```

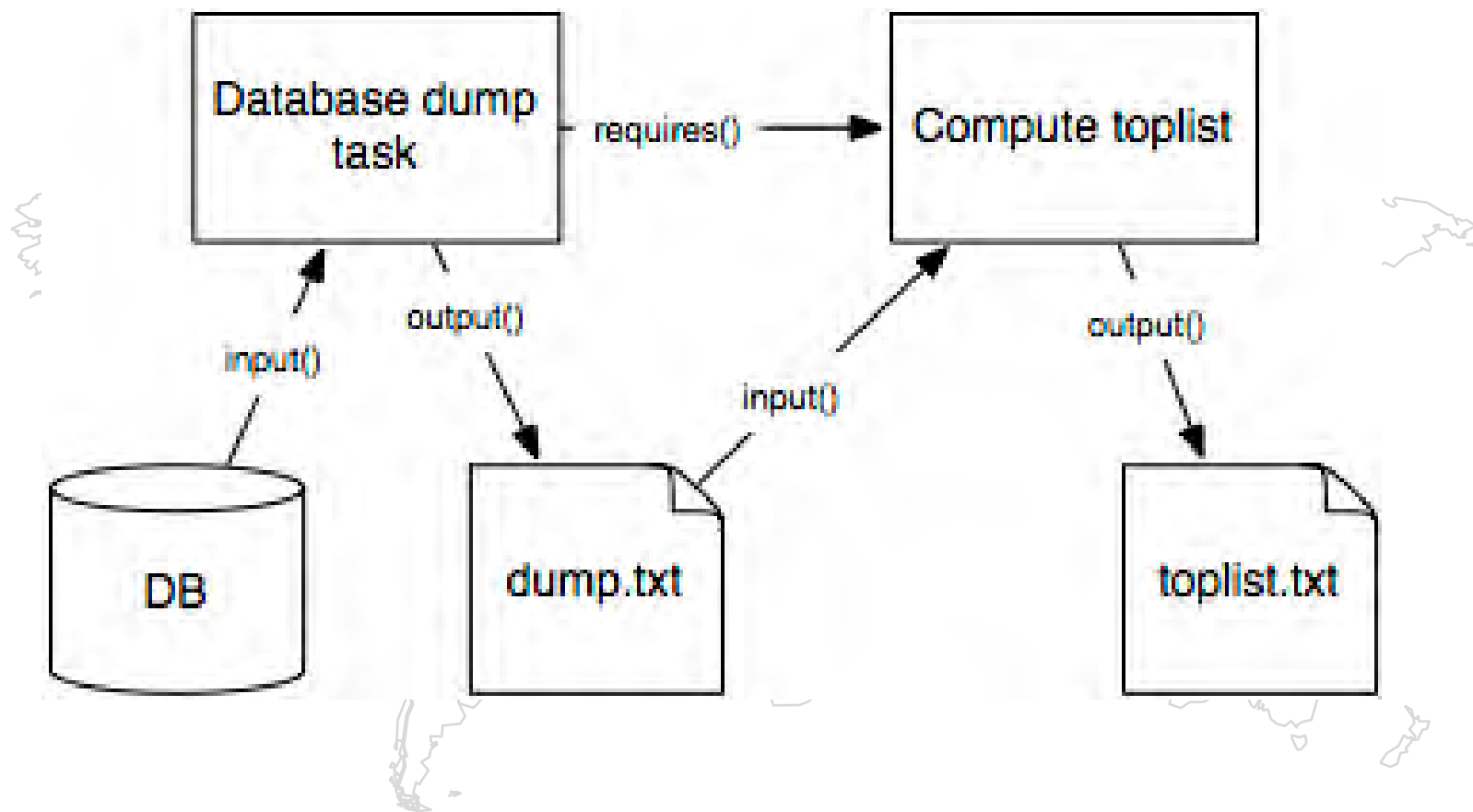
The business logic of the task

Where it writes output

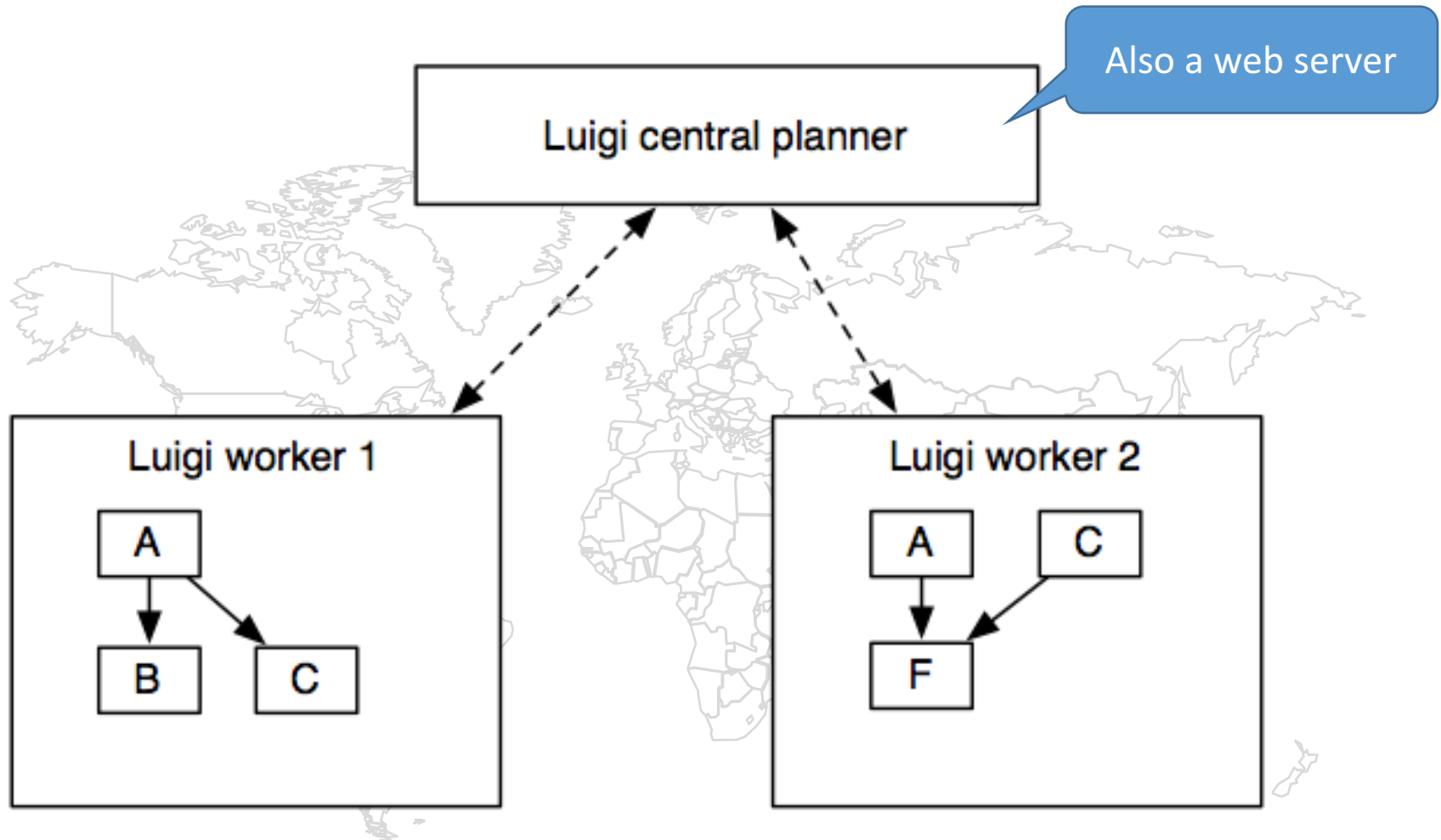
What other tasks it depends on

Parameters for this task

- Task Execution



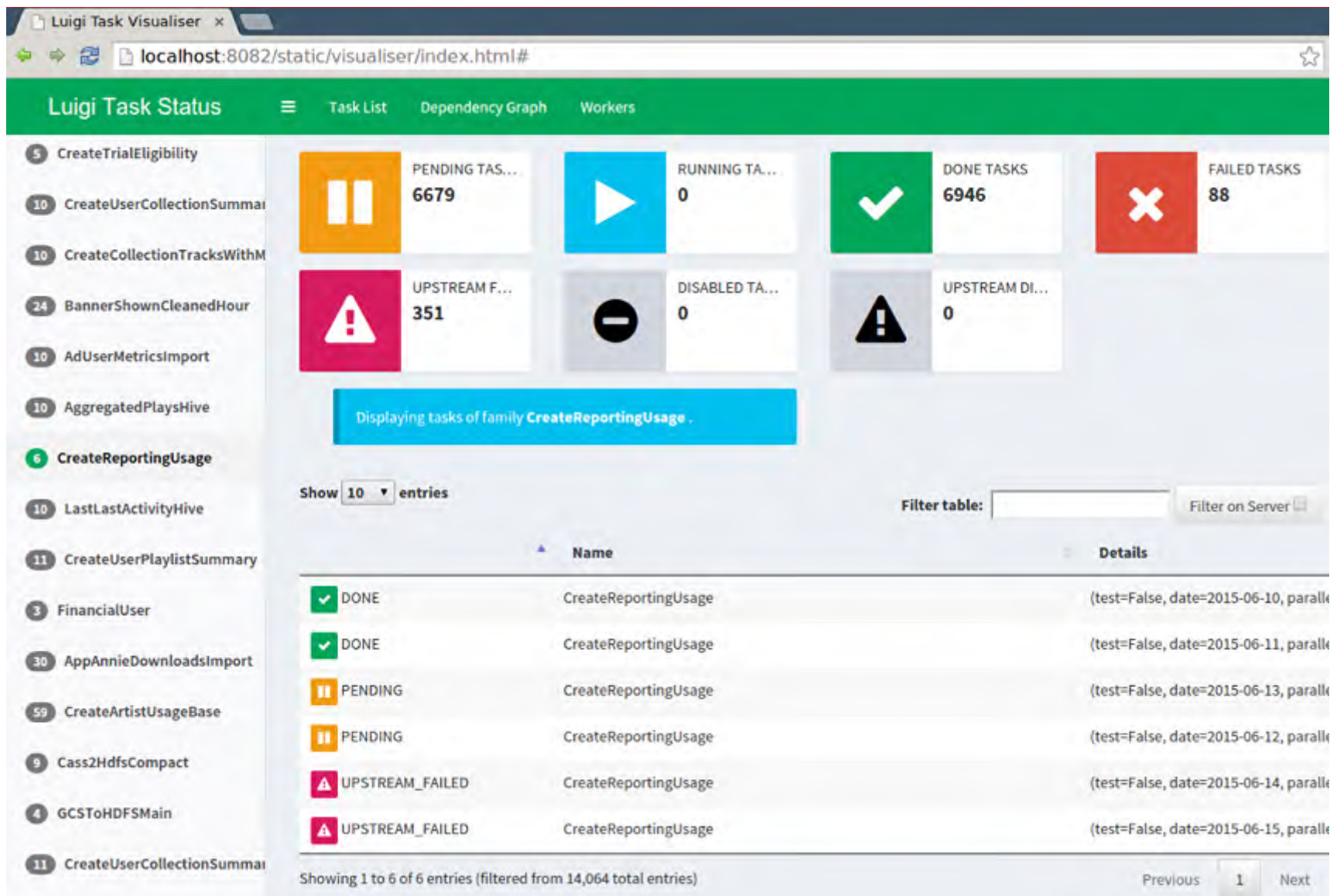
- Architecture



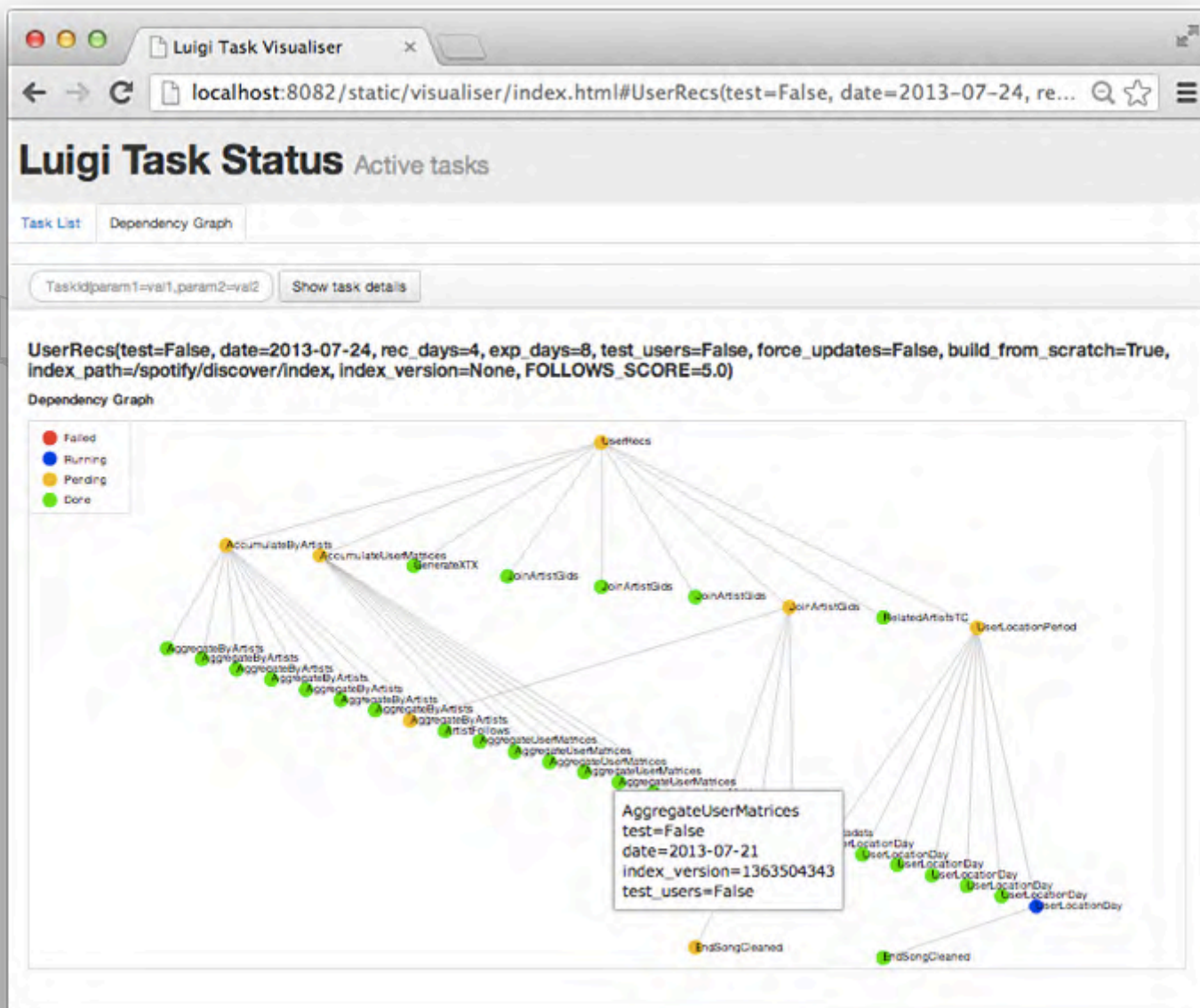
● Architecture Notes

- 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



• 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...



● 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: Jinja2



Airflow Web



Work-flow & Dependency




Log Access & Monitoring



Data Profiling

- Web UI – Overall status

 AirFlow

DAGs






















Tools ▾

Browse ▾

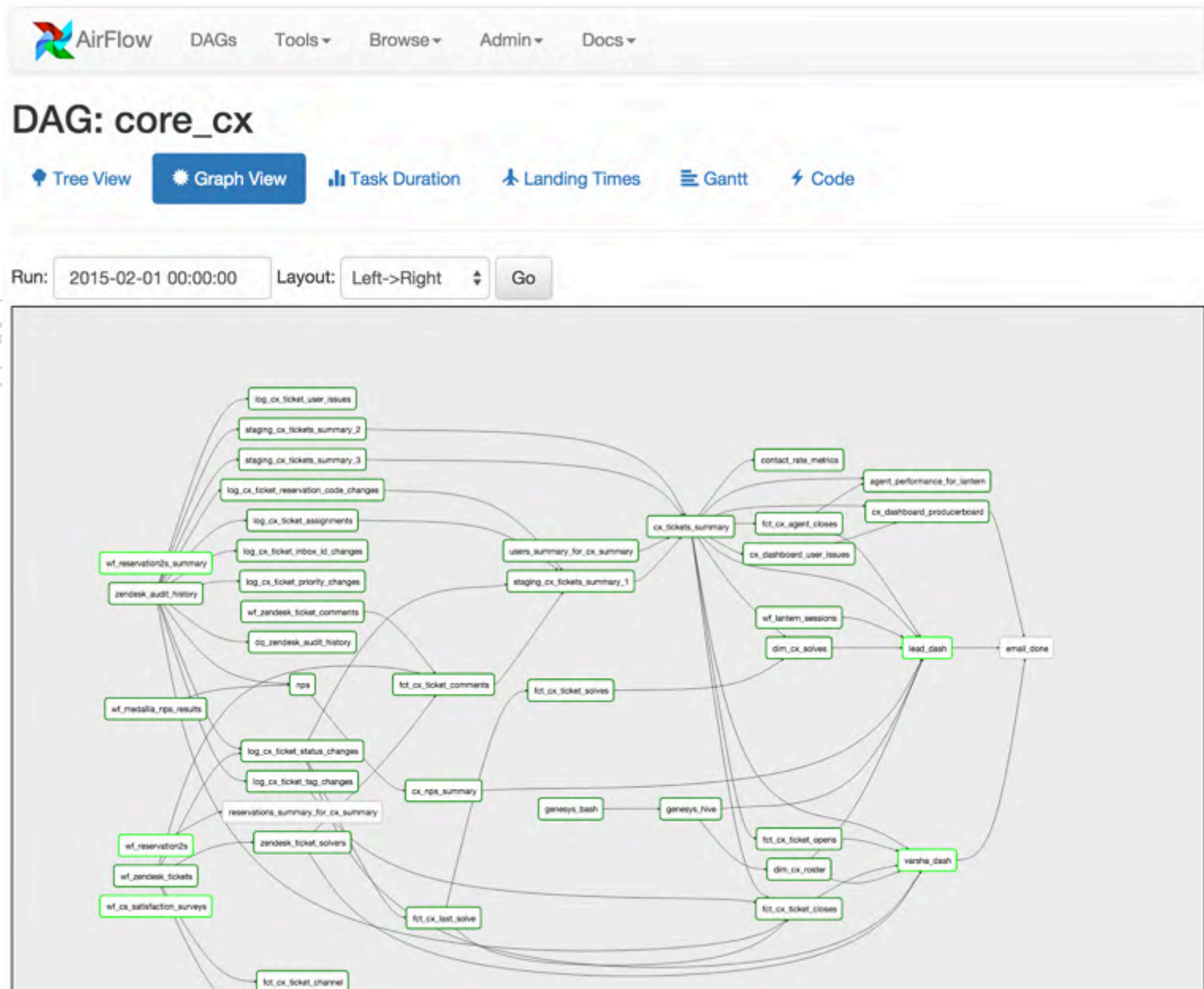
Admin ▾

Docs ▾

DAGs

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

Web UI – workflow visualization



● Web UI – execution history



AirFlow

DAGs

Tools ▾

Browse ▾

Admin ▾

Docs ▾

DAG: example2

Tree View

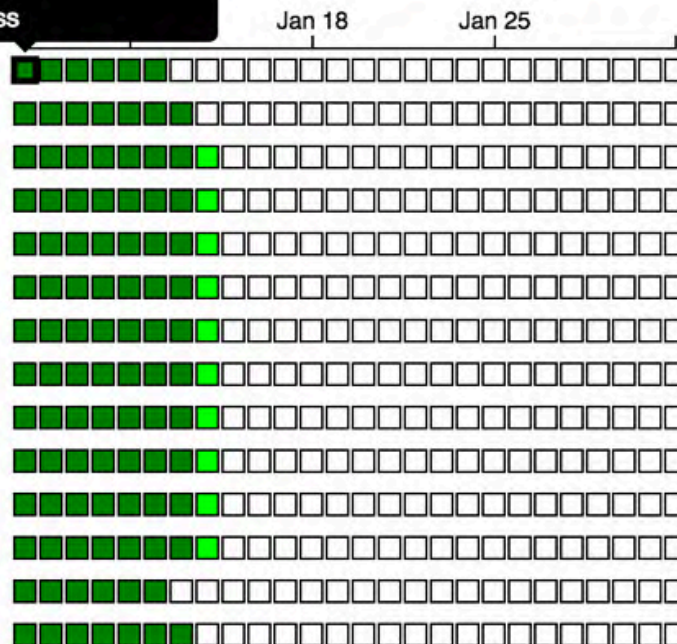
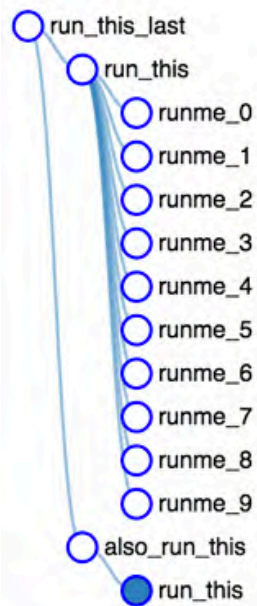
Graph View

Run: 2015-01-07T00:00:00
Started: 2015-02-01T20:22:22
Ended: 2015-02-01T20:22:22
Duration: 0
State: success

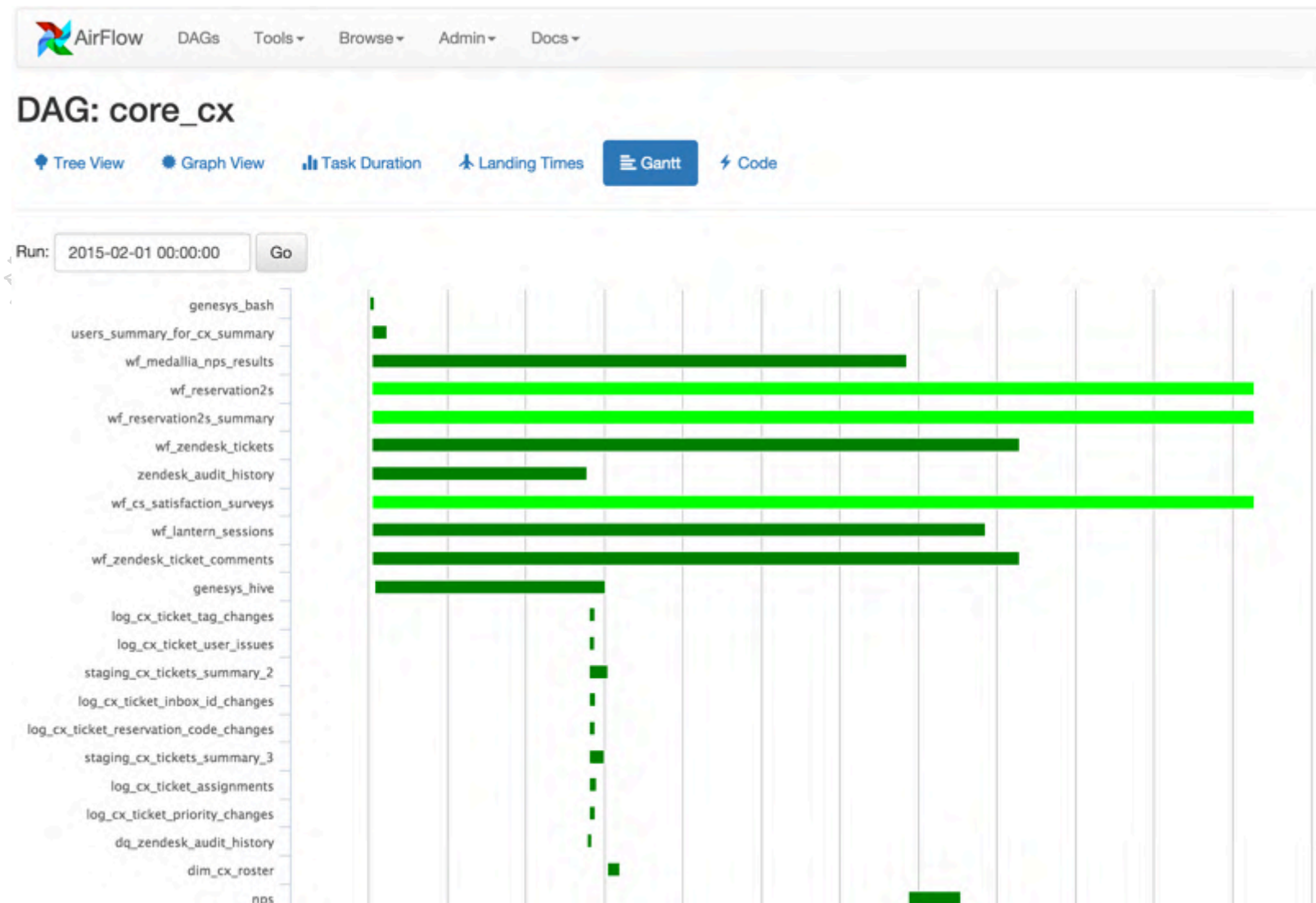
Landing Times

Gantt

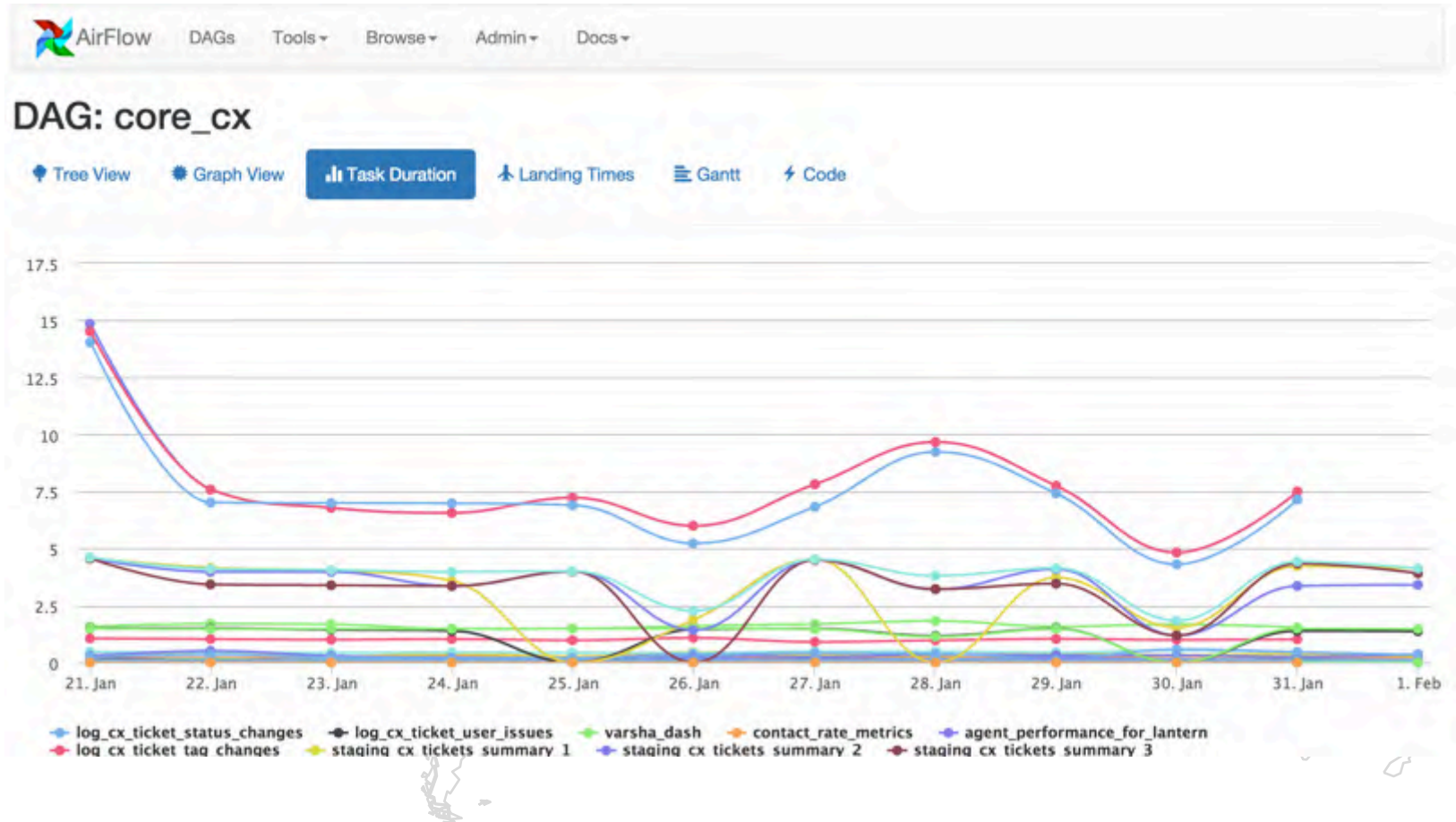
Code



● Web UI – performance profile



Web UI – Performance stats over time



- Web UI – Deep dive for task execution

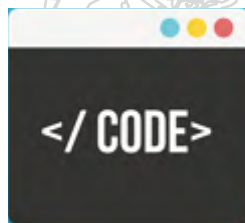
The screenshot displays the AirFlow web interface. At the top, the 'AirFlow' logo and navigation tabs (DAGs, Tools, Admin, etc.) are visible. The main heading is 'DAG: core_cx'. Below it, there are tabs for 'Tree View' and 'Graph View', with 'Graph View' being the active view. A 'Run:' dropdown is set to '2015-02-01 00:00:00'. A modal window is open in the center, titled 'Task Instance: staging_cx_tickets_summary_3 on 2015-02-01T00:00:00'. The modal contains several buttons: 'Task Details', 'Rendered', 'Task Instances', and 'View Log' in the top row; 'Clear', 'Past', 'Future', 'Upstream', and 'Downstream' in the middle row; and 'Graph View' and 'Gantt Chart' in the bottom row. A 'Close' button is located at the bottom right of the modal. The background shows a DAG graph with various task nodes, including 'wf_reservation2s_summary', 'log_cx_ticket_reservation_code_changes', 'log_cx_ticket_assignments', 'log_cx_ticket_inbox_id_changes', 'users_summary_for_cx_summary', and 'cx_tickets_summary'.



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
```

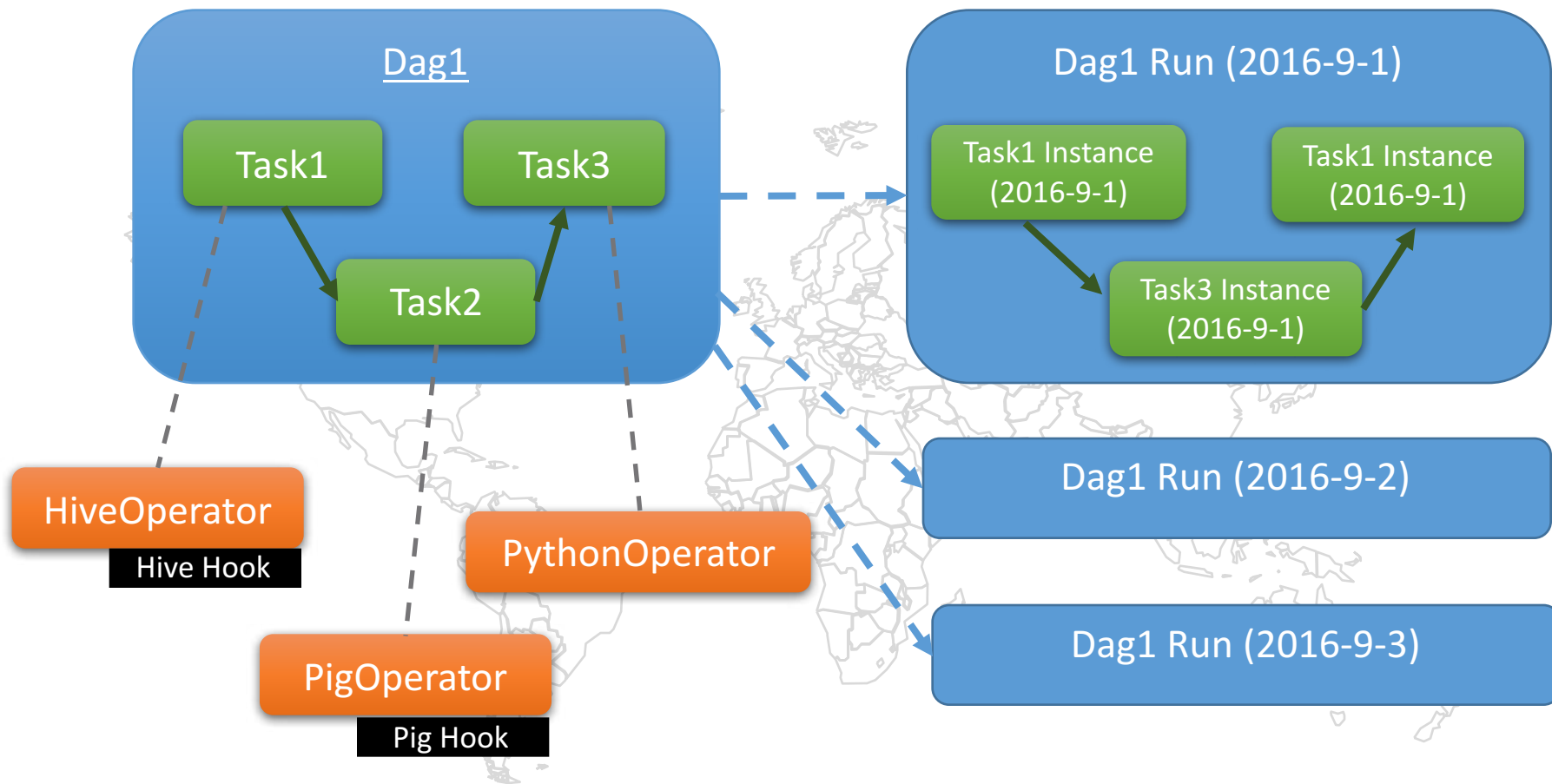
```
task2 = PythonOperator(
    task_id='mydag_task2',
    python_callable=python_logic,
    op_kwargs={'param1': 10},
    dag=dag)
```

An task of another kind of
PythonOperator

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

Setup the dependencies

● 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, Cloud 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 Jinja for batch or any places that fit
- Xcom
 - Share data between Tasks



Architecture Insight



Scalability

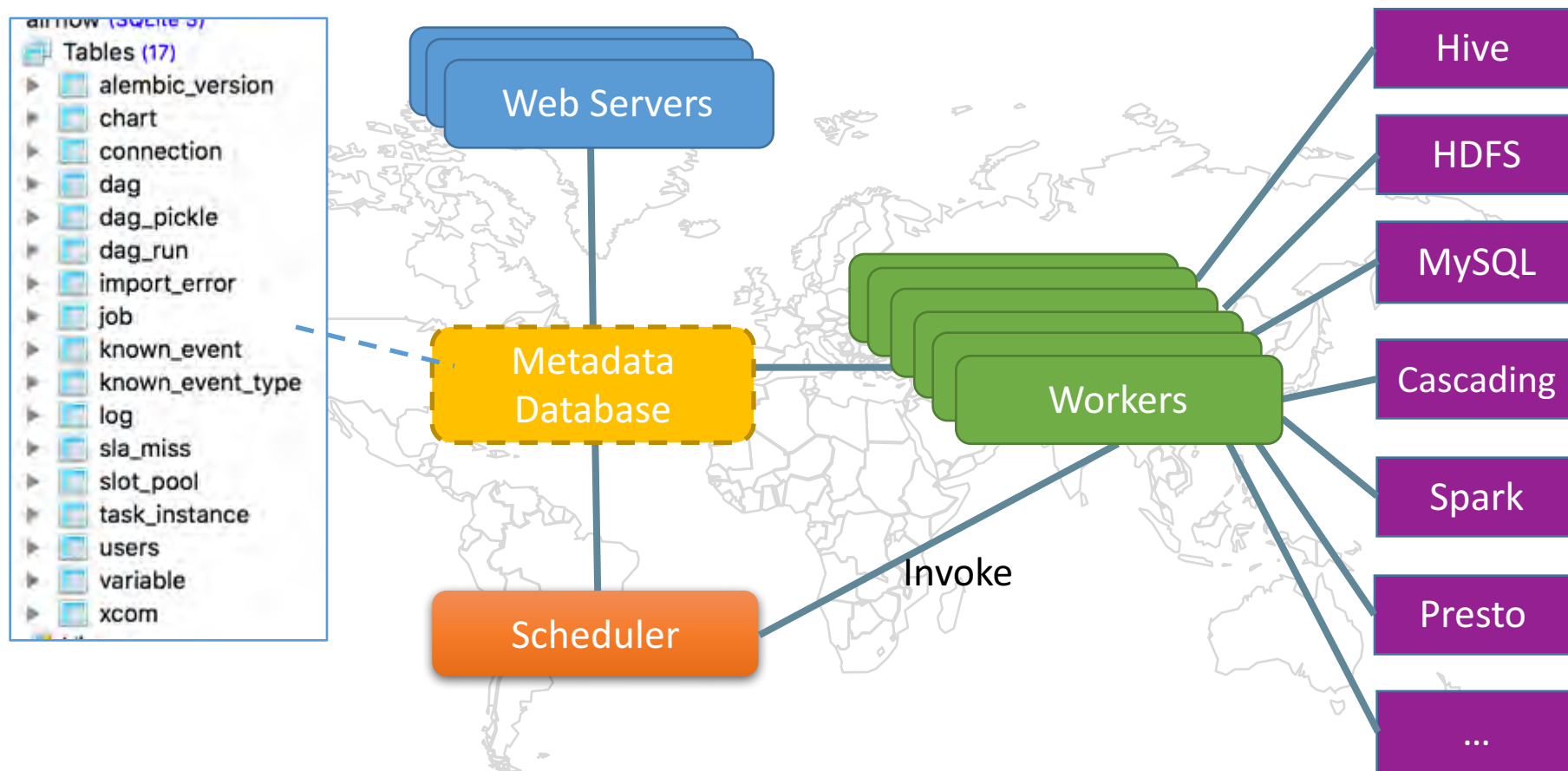


High Availability

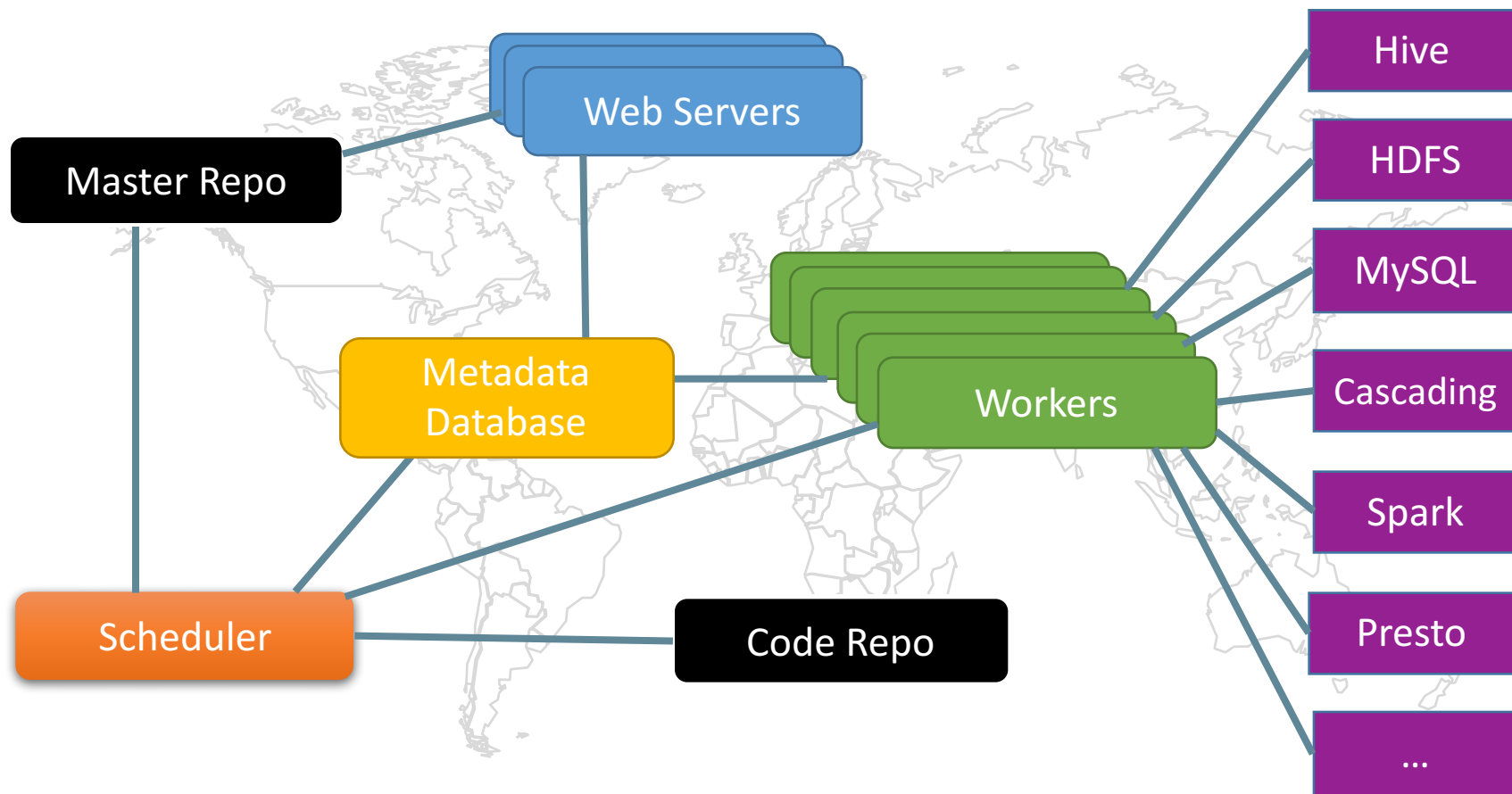


Versioning

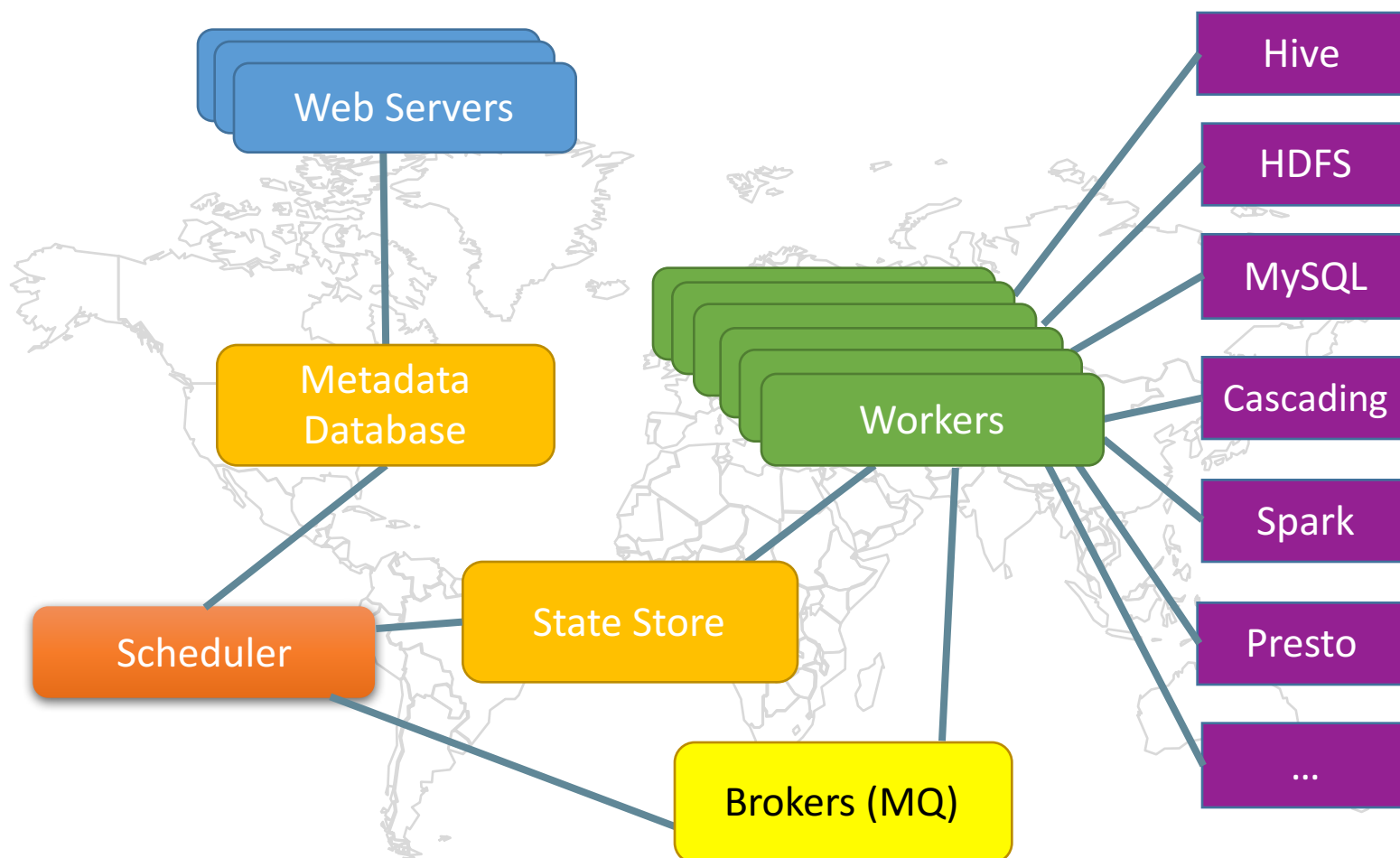
● Airflow Architecture (Local Scheduler)



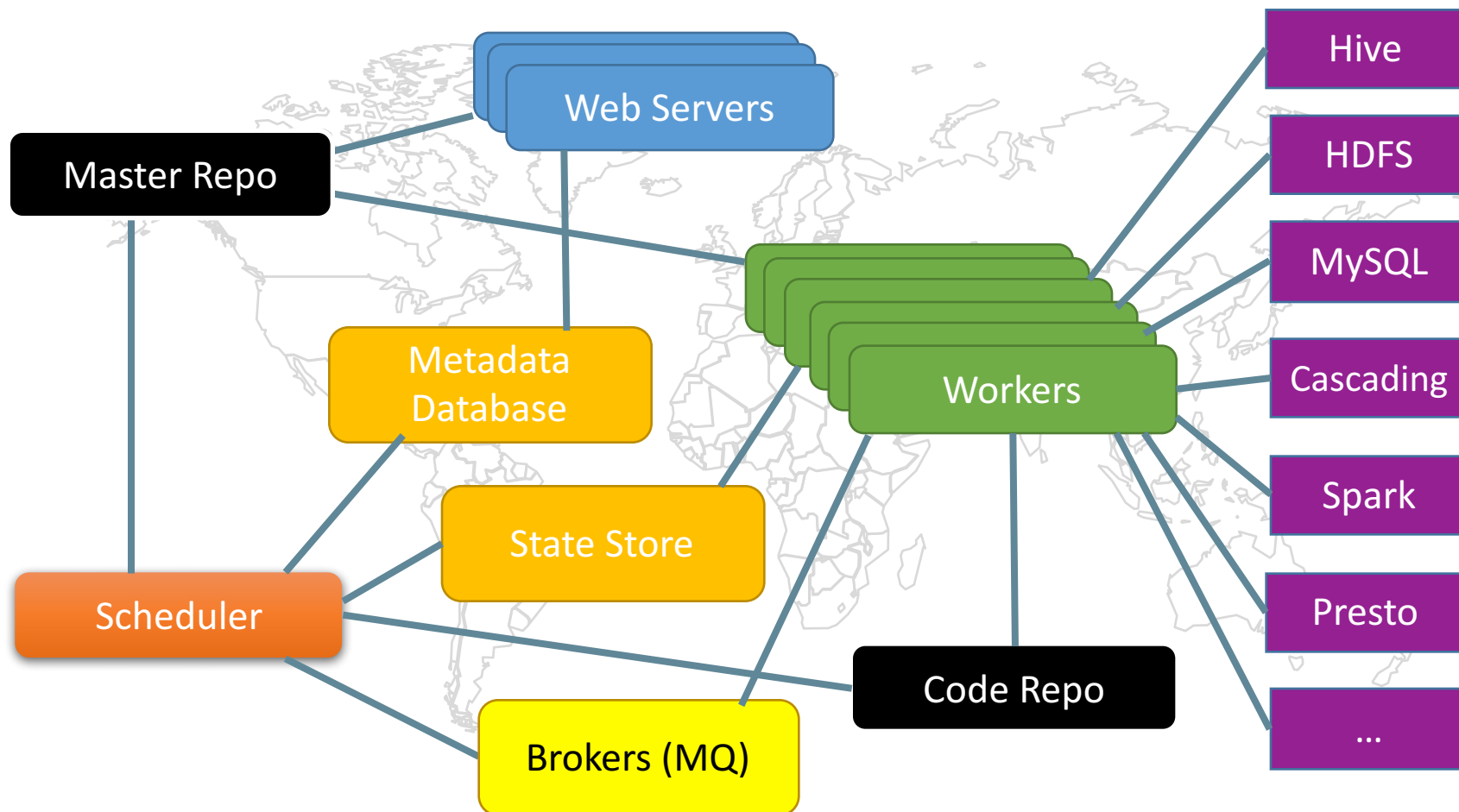
- Local Scheduler – w/ version control



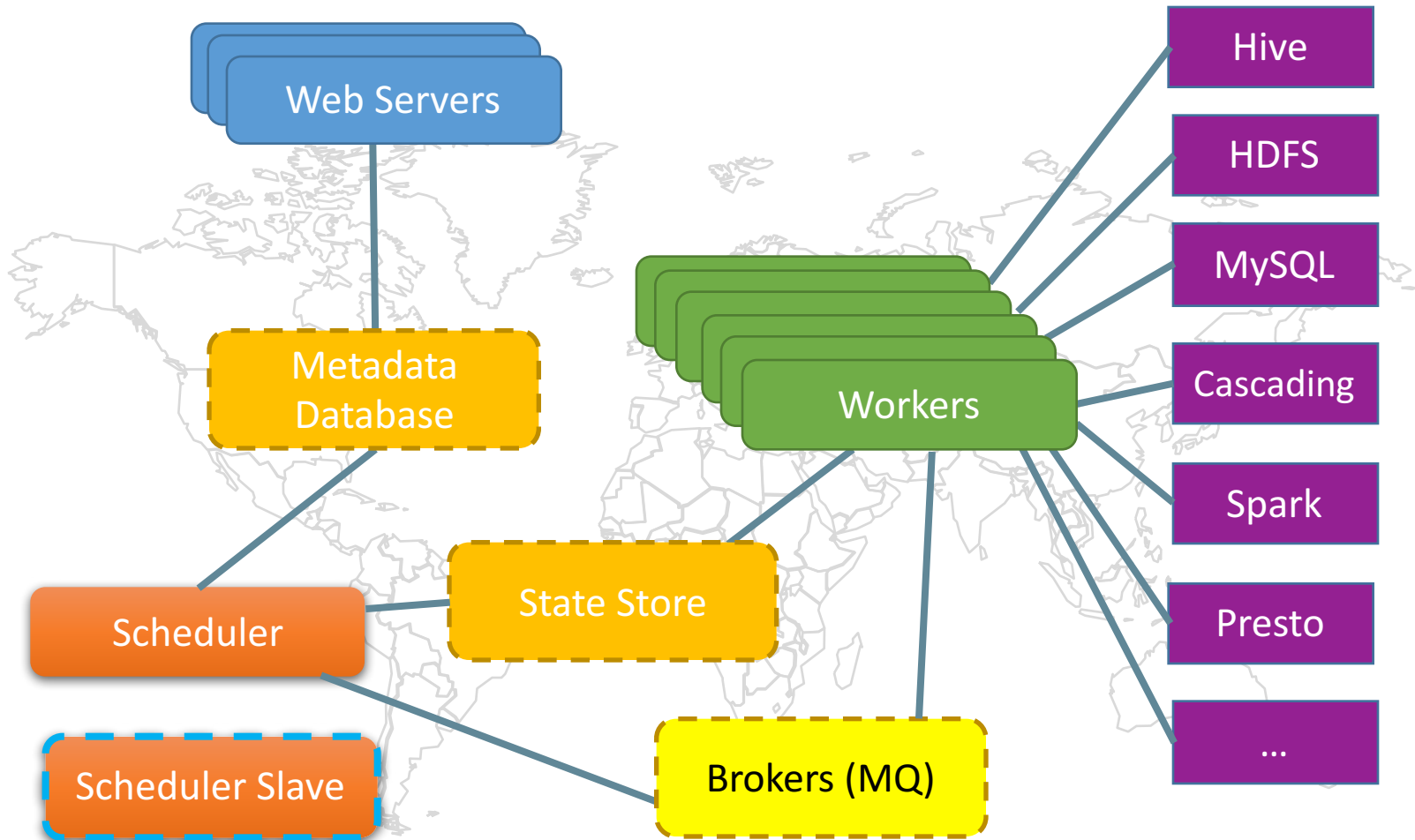
● Airflow Architecture (Celery Scheduler)



- Celery Scheduler – w/ version control



- Airflow Architecture - HA





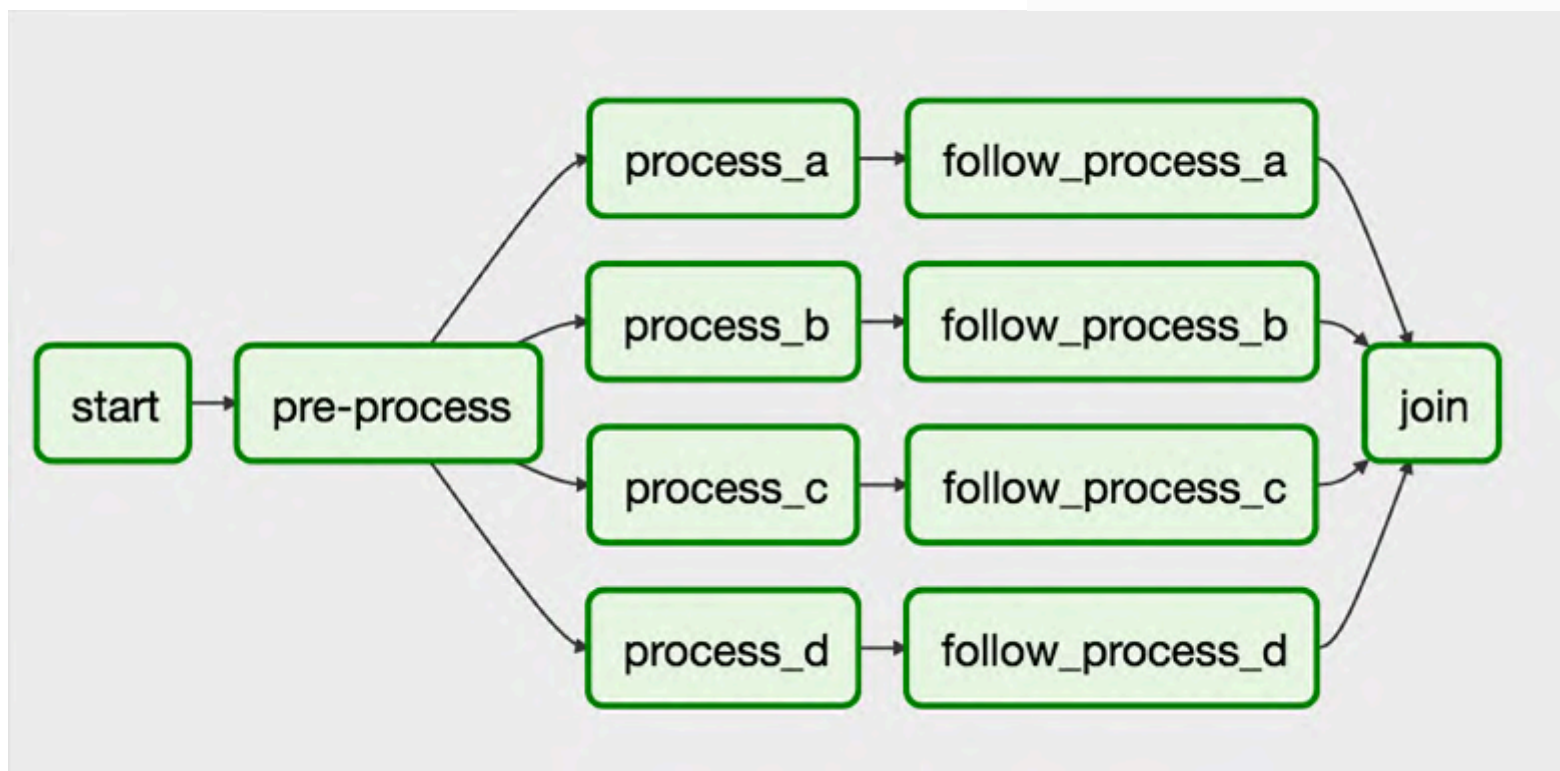
Workflow Patterns



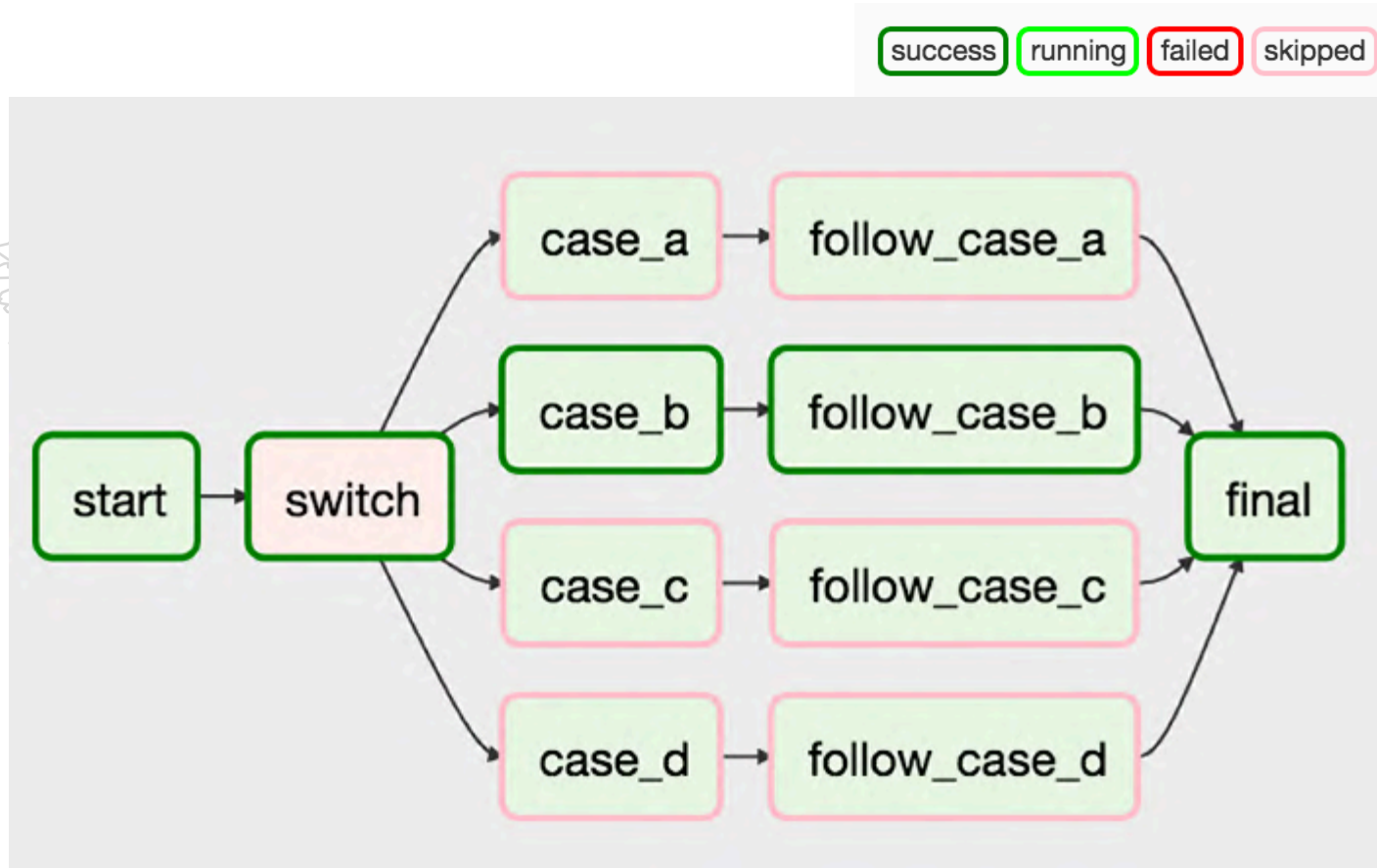
Complex Rule

- Process in parallel

success running failed skipped

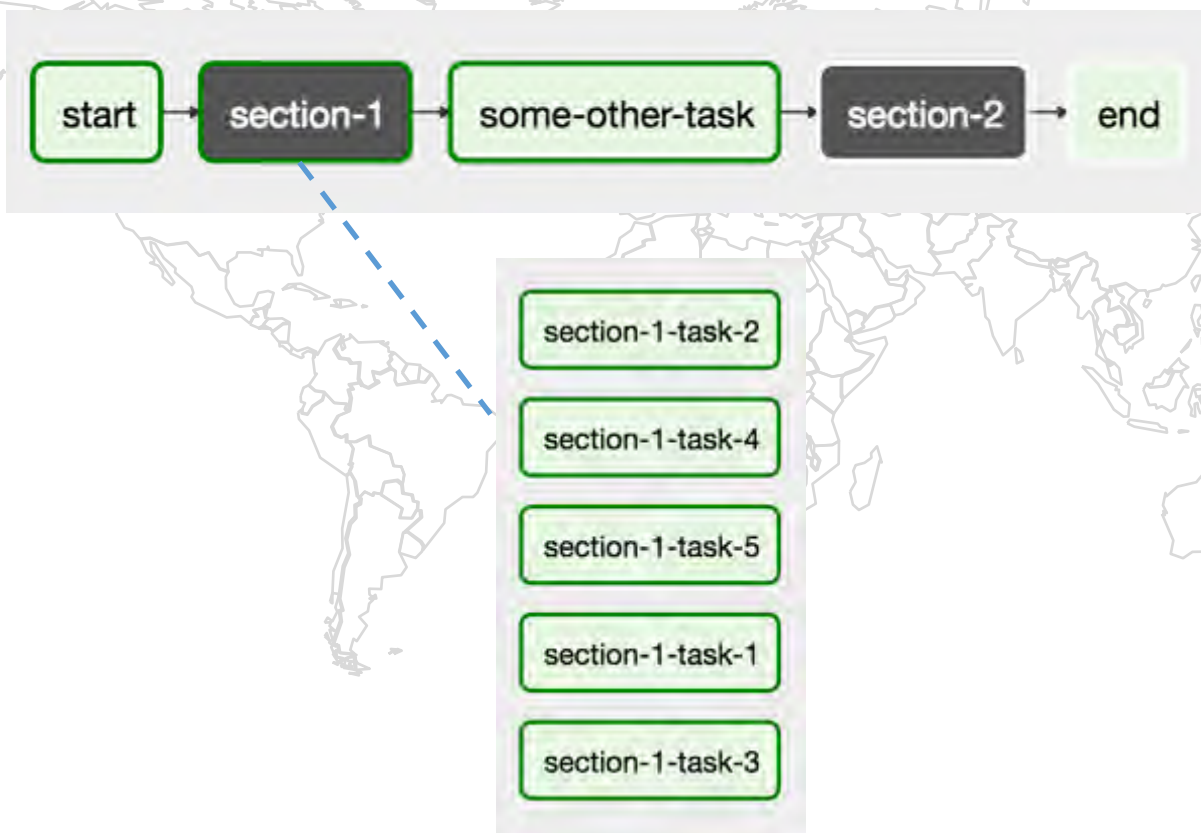


● Switch

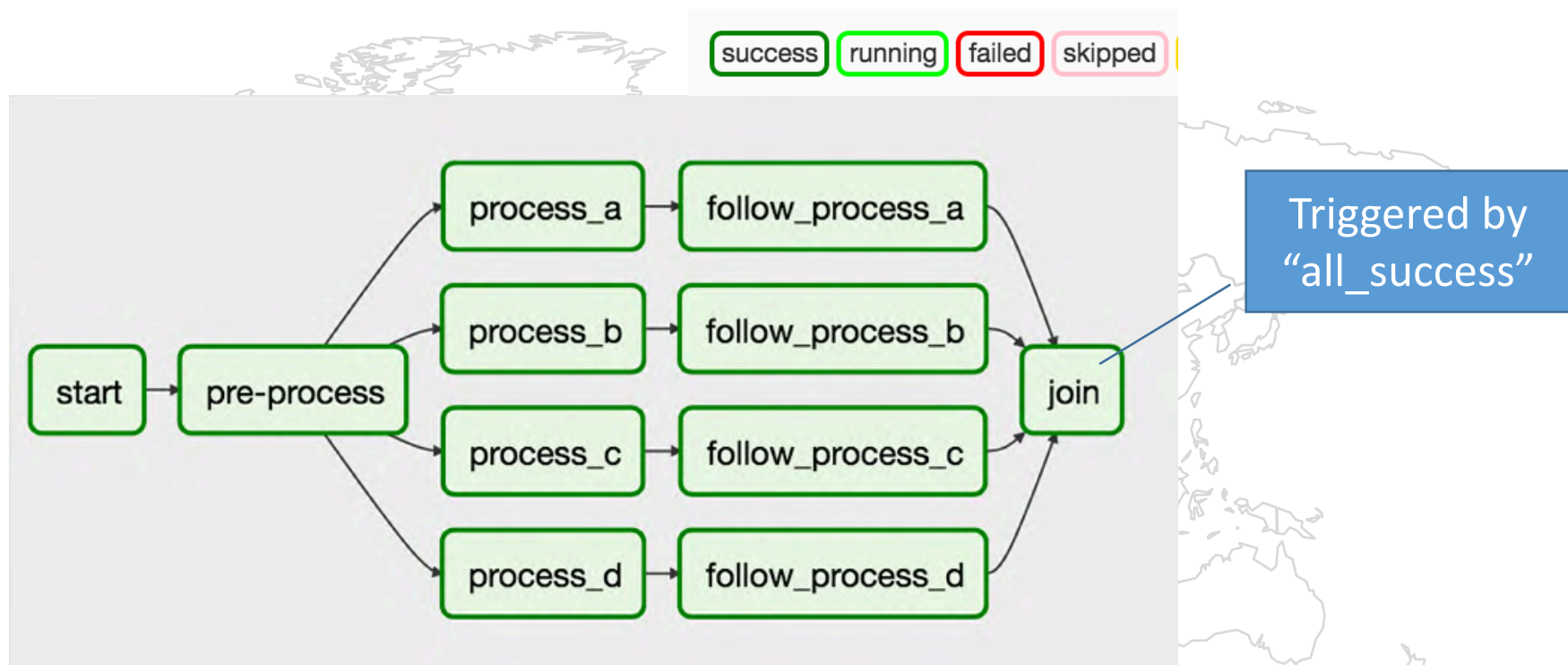


- Sub Dag

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

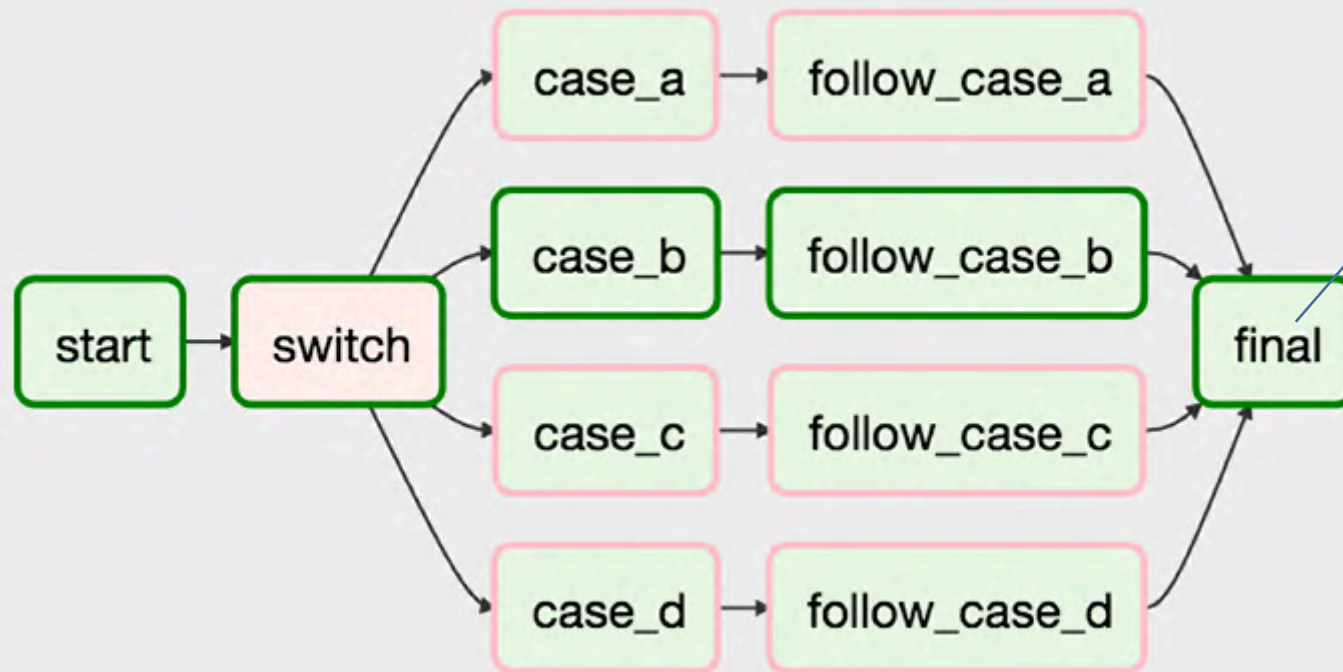


● Trigger Rule – all success



● Trigger Rule – one success

success running failed skipped



Triggered by
“one_success”

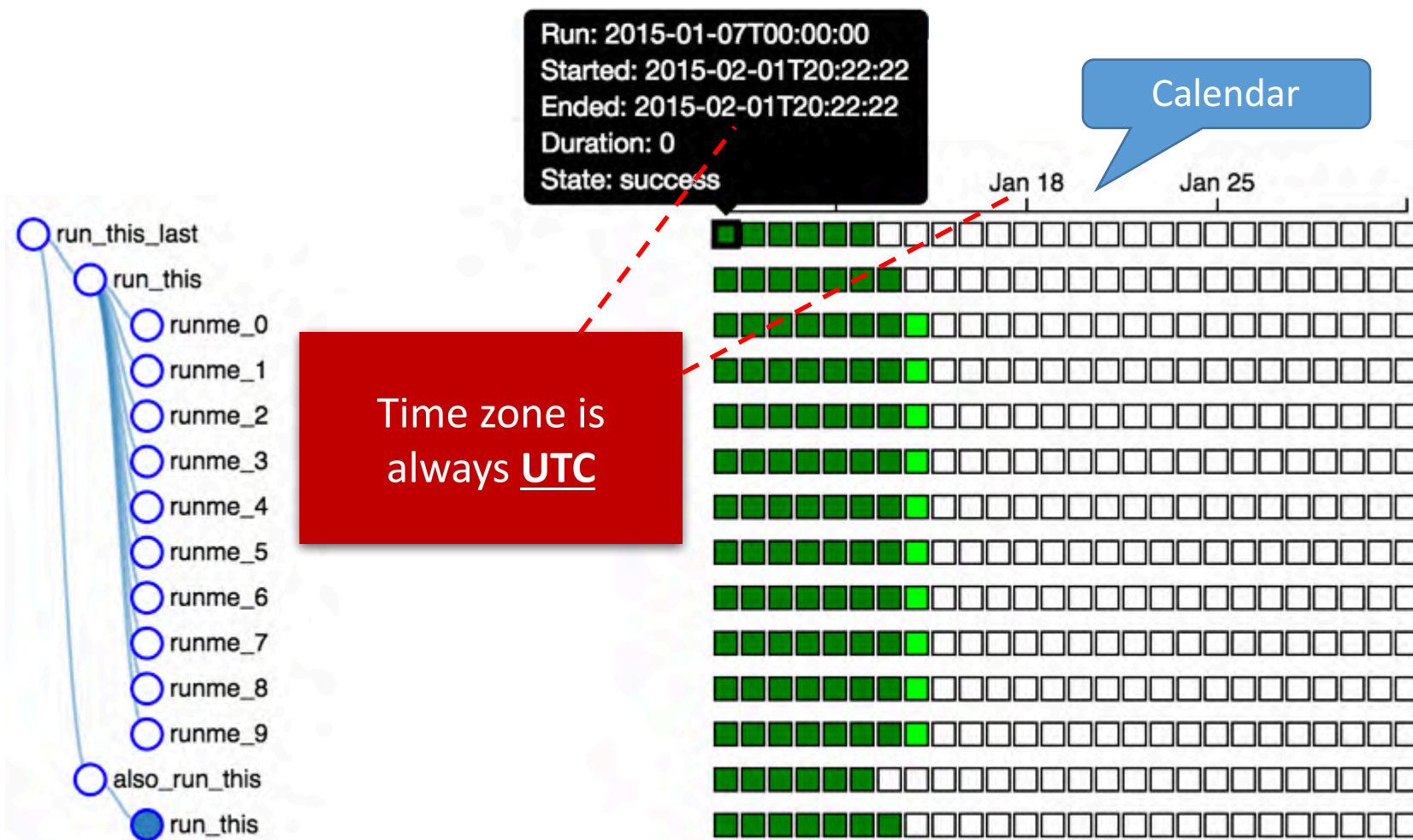


Scheduling Practice



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

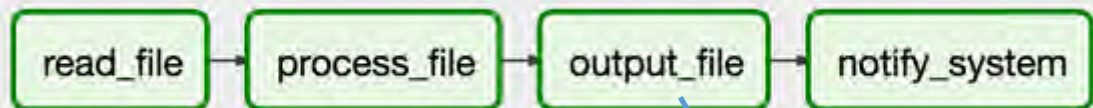
Note

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

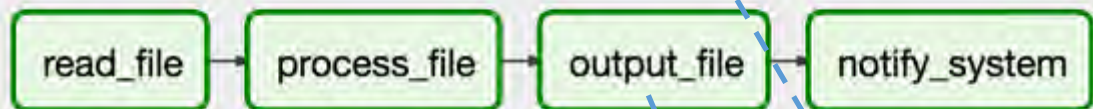
Run at 04:00:00

- 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

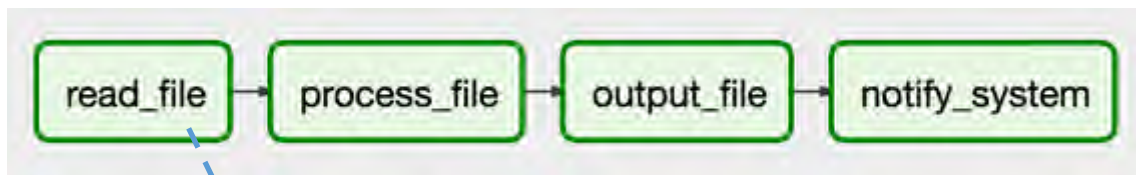
It's ideal case, in real cases...



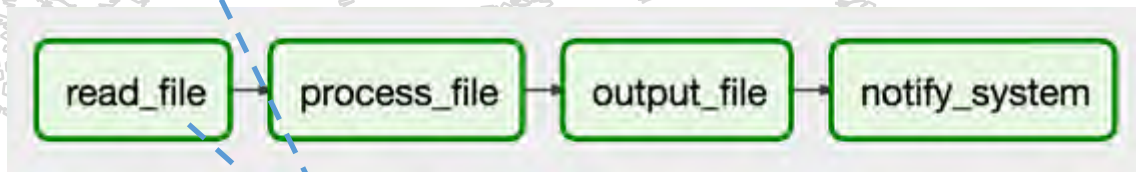
图样图森破

● 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

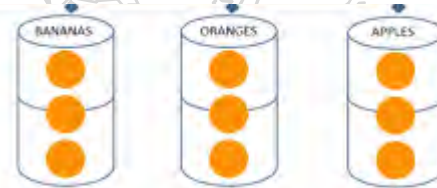
Turn on option “**depends_on_past**” for task **read_file**

● Resource Control

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



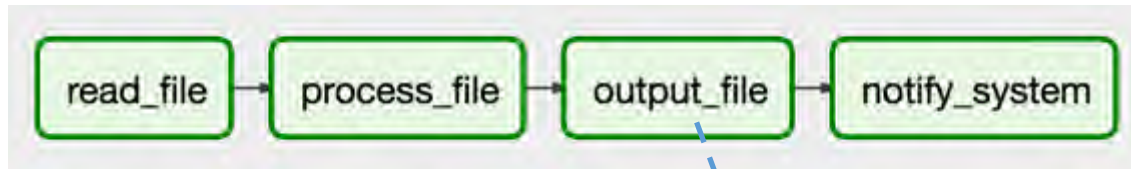
Pool (Limit concurrency + priority)



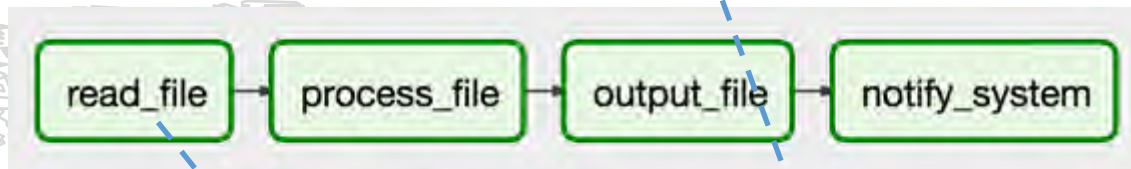
Queue (Affinity)

- Scheduler – more recursive dependency

Run1 (09-01 00:00)



Run2 (09-01 04:00)



What if read_file in Run2 relies 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 *',

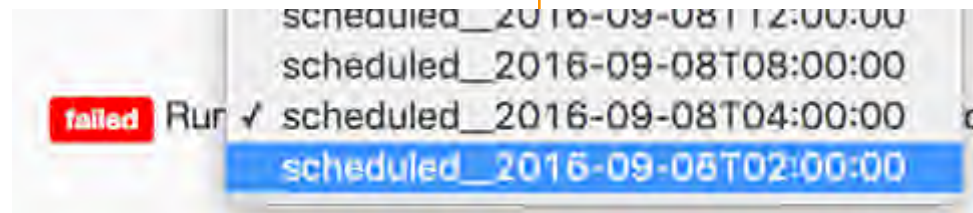
Note

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



● 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

