

Prefect for data pipelines in Python

Franklin Ventura
ResBaz 20210520

<https://github.com/VenturaFranklin/resbazTucson2021-Prefect>

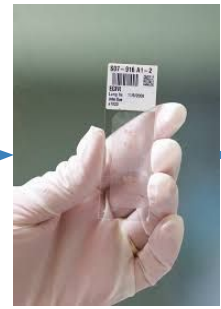
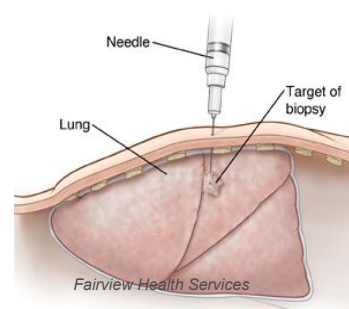


PREFECT

Franklin Ventura

Roche Tissue Diagnostics

- BS Biomedical Engineering
- MS Software/Systems Engr
- 9 years
 - Senior Software Engineer
 - Research and Early Development (tRED)
 - Digital Pathology
 - Cancer Diagnostics
 - RTD Programming Club
 - Beginners & More Experienced
 - Microsoft TEALS
 - Software Carpentry Instructor
 - Tucson Python Language Enthusiasts (Local Meetup)

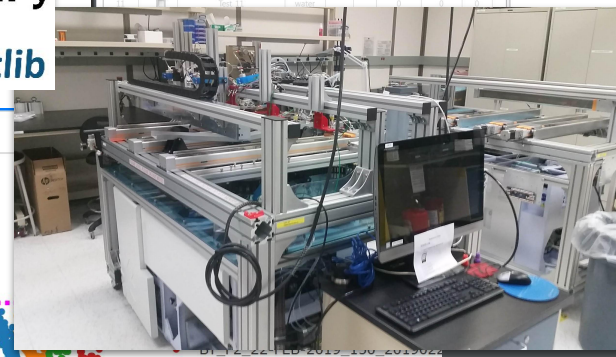
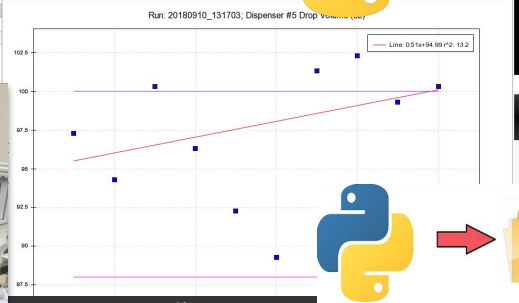
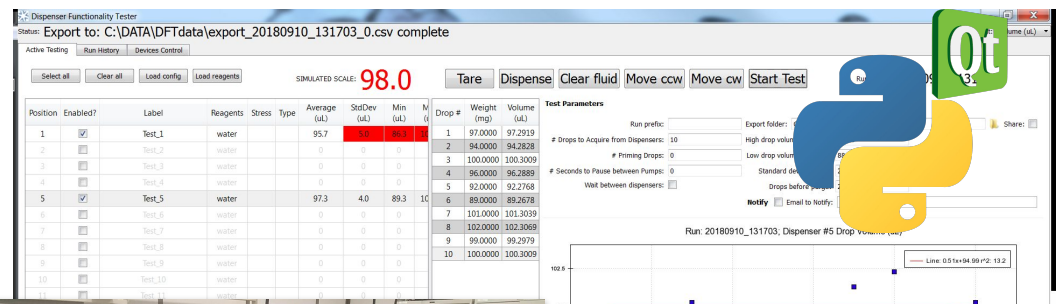
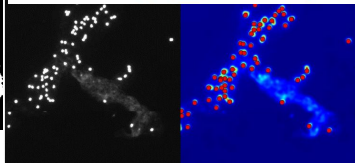
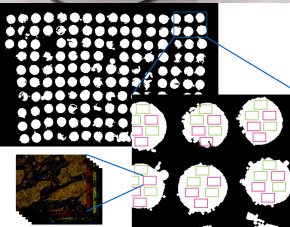


µManager

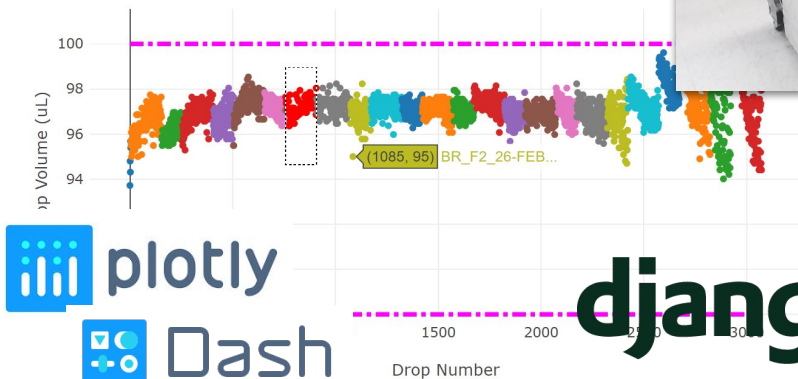
THE OPEN SOURCE MICROSCOPY SOFTWARE

Variety of Python Applications at RTD

Roche



- BT_F2_25-FEB-2019_100_2019022
- BT_F2_25-FEB-2019_150_Runn 3
- BT_F2_25-FEB-2019_100_Runn 3
- BR_F2_26-FEB-2019_Run_4_150_2
- BR_F2_26-FEB-2019_Run_5_150_2
- BT_F2_27-FEB-2019_100_Run_5_2
- BT_F2_27-FEB-2019_150_Run_6_2
- BT_F2_27-FEB-2019_100_Run_6_2
- BT_F2_27-FEB-2019_100_Run_7_150_2



Dispenser Functionality Tester Data Analysis

Report Name: C:\DATA\DFTdata\export_20180910_131703_0.csv complete

Version: 0.1.0

Analysis Version: 0.1.0

Dispenser Details

Label	Reagent	Type	Reagent	Average (uL)	StdDev (uL)	Min (uL)	Max (uL)	Drop #	Weight (mg)	Volume (uL)
BR_F2_25-FEB-2019_100_2019022	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768
BR_F2_25-FEB-2019_150_Runn 3	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768
BR_F2_25-FEB-2019_100_Runn 3	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768
BR_F2_26-FEB-2019_Run_4_150_2	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768
BR_F2_26-FEB-2019_Run_5_150_2	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768
BT_F2_27-FEB-2019_100_Run_5_2	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768
BT_F2_27-FEB-2019_150_Run_6_2	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768
BT_F2_27-FEB-2019_100_Run_6_2	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768
BT_F2_27-FEB-2019_100_Run_7_150_2	water	CONTROL	water	95.7	4.0	89.3	10.0	5	92.0000	92.2768

Document Generated: 04-Mar-2019

Page 1 of 2

ReportLab
Open Source





PREFECT

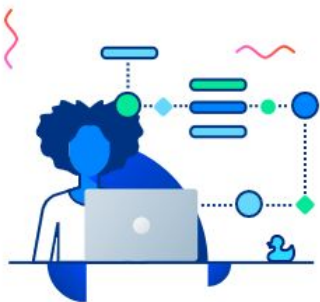
<https://www.prefect.io/>



- programmatically author, schedule, and monitor workflows
- Extract, Transform, and Load data processing expressed as code

organize and kick off machine learning jobs running on external Dask clusters

“Prefect takes care of scheduling, infrastructure, error handling, retries, logs, triggers, data serialization, parameterization, dynamic mapping, caching, concurrency, and more...”



STEP 01 - CORE



STEP 02 - CLOUD



STEP 03 - CORE

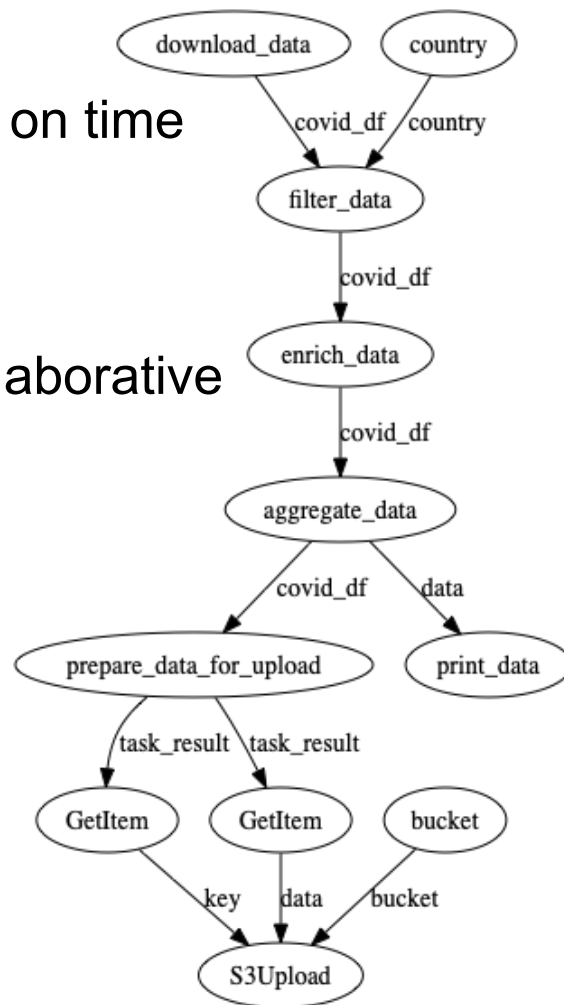
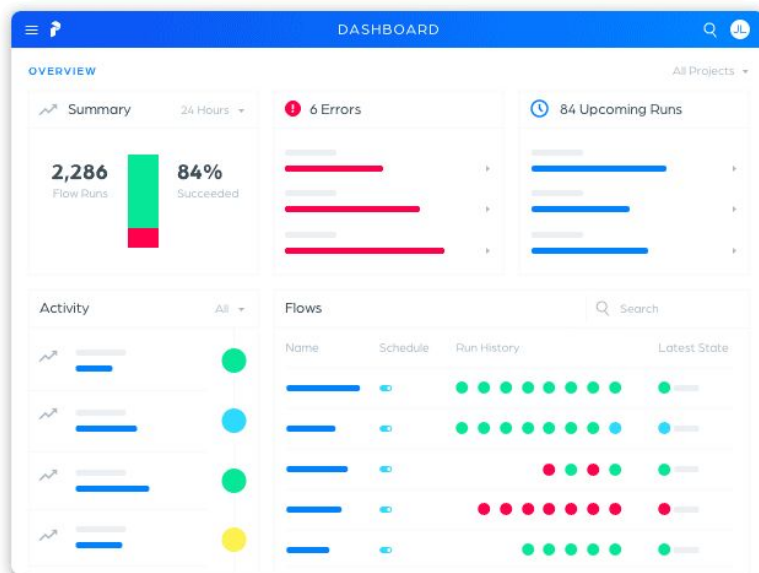


STEP 04 - CLOUD

Why use Prefect at all?

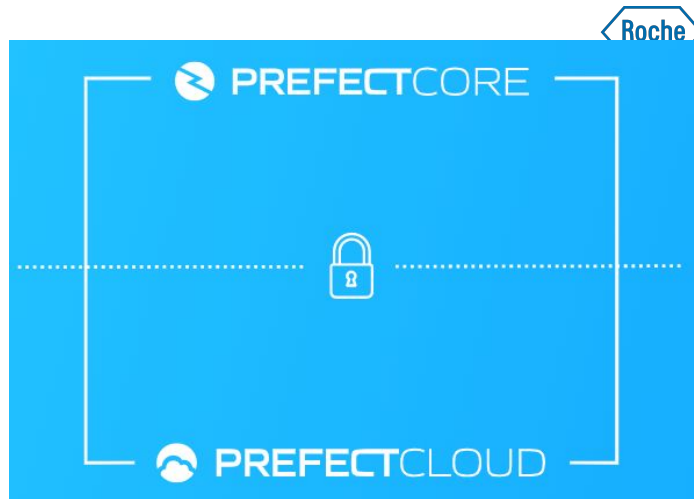


- making sure everything is running smoothly on time
- scheduling and coordination
- recovering from failures
- workflows defined as code, become more maintainable, versionable, testable, and collaborative



Will Prefect Get My Data?

- Prefect only gets what you consciously decide to send them
- By default only function names, structure of DAG and log info is sent
- Cloud: “hybrid execution keeps code and data completely private”



OR

- Run your own prefect Server which comes with the UI

See:

[Prefect Server vs. Prefect Cloud - which should I choose?](#)



Prefect - Alternatives



Apache
Airflow

































PREFECT

Class Based	API	Function Based
strict dependency on a specific time	Scheduling	any time, with any concurrency, for any reason
Flow not meant to change	Parameters	Allow slightly different flow
Second class citizen	Data	First class operation
Internal function	Dynamic Workflows	Task mapping

Prefect - Building a Flow

It's just Python functions



 Airtable 	 AWS 	 Azure 	 Azure ML 	 Databricks 
 DBT 	 Docker 	 Dropbox 	 Email 	 GitHub 
 Google Cloud 	 Google Sheets 	 Great Expectations 	 Jira 	 Kubernetes 

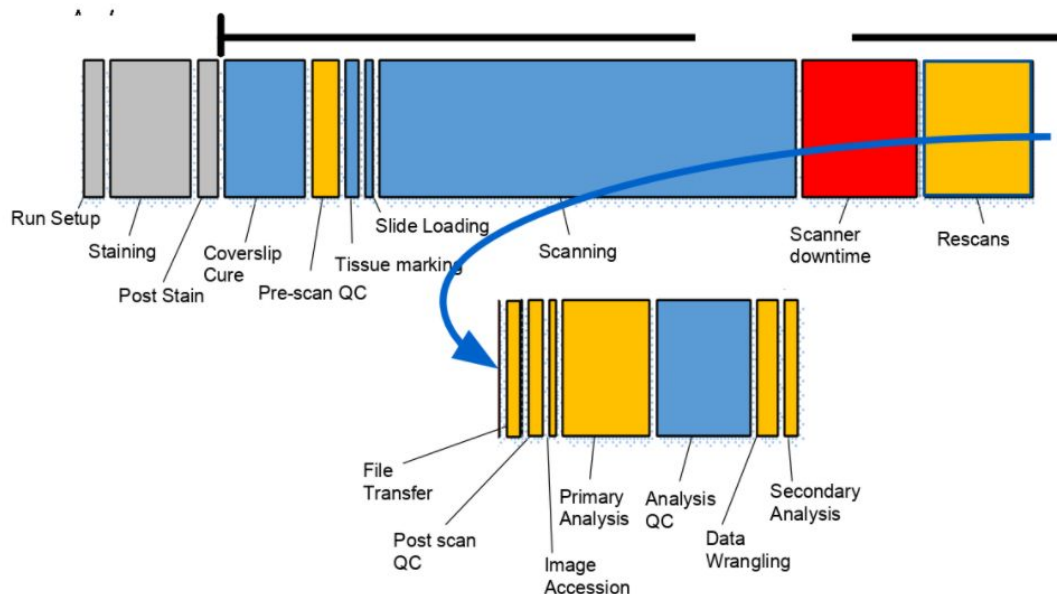
```
def test_flow():  
    with Flow() as flow:  
        StartContainer(**)  
        WriteGsheetRow(**)  
        status = WaitOnContainer(**)  
        EmailTask(status, **)  
        SlackTask(**)  
        custom_cleanup(**)
```


Data

Orchestration

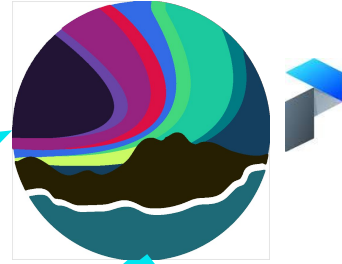
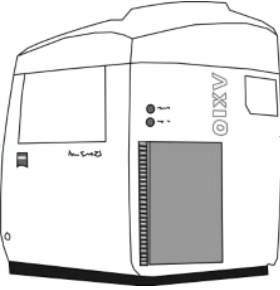
Transfer images from scanner; Track metadata for images
Parallelization relieves the bottleneck

- Removing scanner bottle necks
- Reduction of touch points
- Flexible Metadata collection
- Data pipeline analytics, monitoring, and error tracking



Solution

Process Improvement



PREFECT

Summary of Slides ☆

File	Edit	View	Insert	Format	Data	Tools	Add-ons	Help	All c
100%	\$	%	.0	.00	123	Arial	10		
A	B	C	D	E	F				
Include in Name?	REQUIRED	REQUIRED	REQUIRED	REQUIRED	REQUIRED				
Barcode	Study	Experiment	Panel	Scanner	Sca				
620995	Prec01	Run 2	P1	Cardinal					
645750	Prec01	Run 2	P1	Cardinal					
317335	Prec01	Run 2	P1	Cardinal					
407871	Prec01	Run 2	P1	Cardinal					
465551	Prec01	Run 2	P1	Cardinal					
754788	Prec01	Run 2	P1	Cardinal					



Data Pipeline

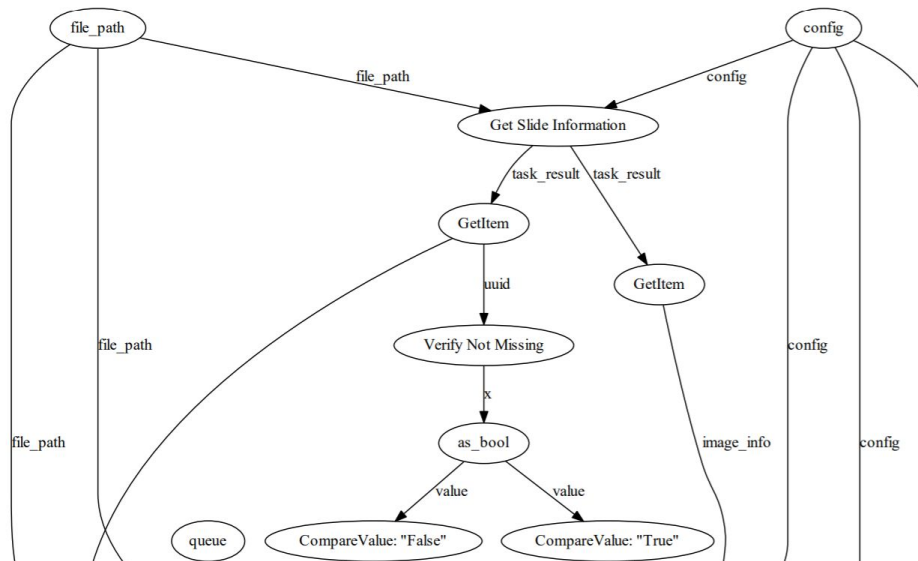


Clearly organize the various steps and the associated order

- Configure retry policies into individual tasks
- Set up alerting in the case of failures, retries, as well as tasks running longer than expected
- Comes with an intuitive UI with powerful tools for monitoring and managing jobs

Version 2.5.0

Version 2.10.5



Data Flow Orchestration

Track images through system - Gain Analytics



DASHBOARD

All Projects

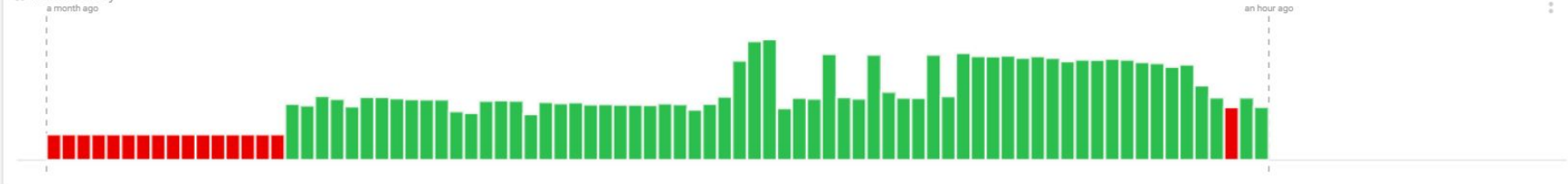
📁 All Projects ▼

OVERVIEW

FLOWs

AGENTS

Run History



Summary

🕒 24 Hours ▼

11
flow runs

In the last day



90.9 %
succeeded



1 Failed Flows

🕒 24 Hours ▼

Scan Image Flow: Raven
2:19pm MST



0 Upcoming Runs

UPCOMING ⌚

LATE ⌚

✓ No upcoming runs.

Data Flow Monitoring

Monitor systems



DASHBOARD

Aurora

OVERVIEW

FLOWs

AGENTS

Cardinal Agent

LocalLocalAgent

LAST QUERY

4:37pm | 6 seconds ago

CORE VERSION0.12.0

TOKEN ID27ce4d63-595f-4...

LABELS

azure-flow-storage

Cardinal

gcs-flow-storage

RTUMWXCOEXT0065

s3-flow-storage

Raven Agent

LocalLocalAgent

LAST QUERY

4:37pm | 10 seconds ago

CORE VERSION0.12.0

TOKEN ID83f162a6-5690-4...

LABELS

azure-flow-storage

gcs-flow-storage

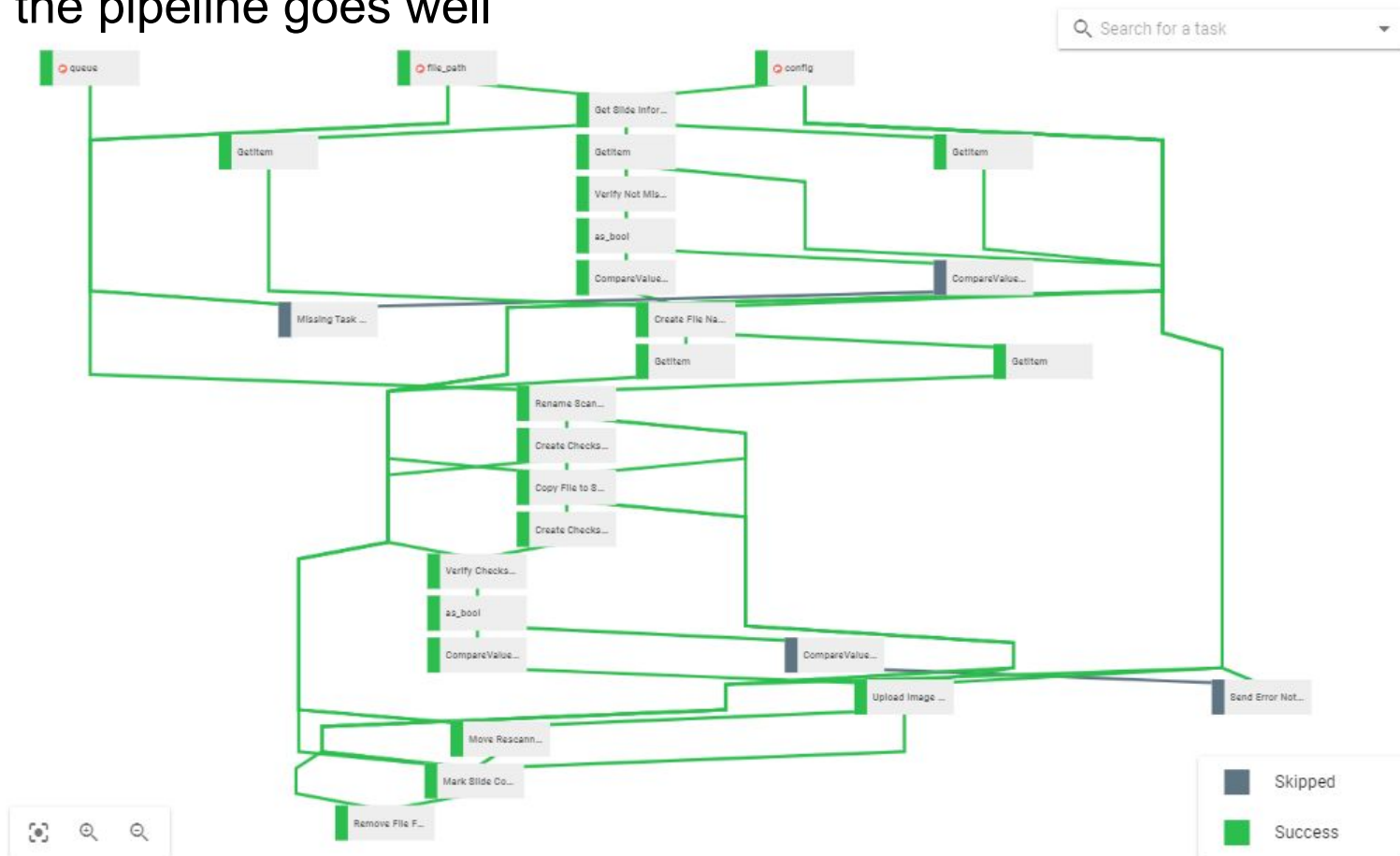
Raven

RTUMWXCOEXT0050

s3-flow-storage

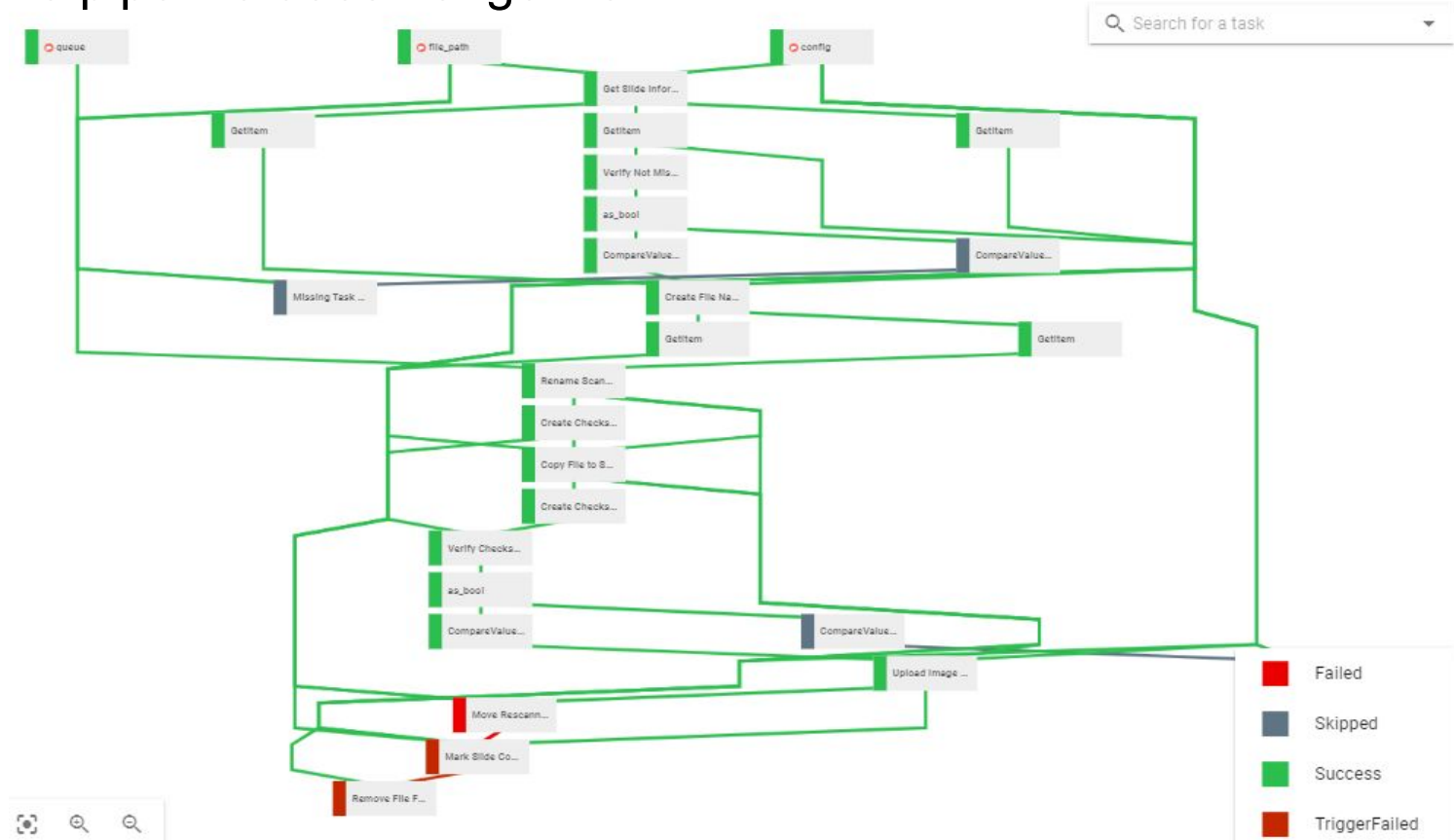
Data Pipeline Visualization

When the pipeline goes well



Data Pipeline Visualization

When the pipeline does not go well



Error Investigation

Dive into errors remotely, no combing through logs



545891 - Move Rescanned Images

[Aurora](#) > [Scan Image Flow: Raven](#) > [545891](#) >

RESTART



MARK AS



OVERVIEW

LOGS



Move Rescanned Images

Failed



Flow Run

545891

Failed

Last State Message

[6 Aug 2020 2:19pm]: Unexpected error:
FileNotFoundException('No image found for
feb90bff5f5f73db')

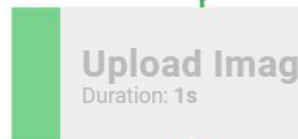
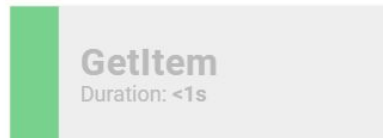
Started	6 Aug 2020 2:19pm
Ended	6 Aug 2020 2:19pm
Duration	1 second
Updated	6 Aug 2020 2:19pm
Class	FunctionTask
Result Type	NoResultType
Result Location	None



Dependencies

3 Upstream • 2 Downstream

☒ Show Cards



Search for a task

Task Run

545891 - Move Rescan...

Failed - 6 Aug 2020 2:1...

Task

Move Rescanned Images

Task Run State Message:

Unexpected error:
FileNotFoundException('No image found
for feb90bff5f5f73db')

Max retries: 0

Class: FunctionTask

Activity

All

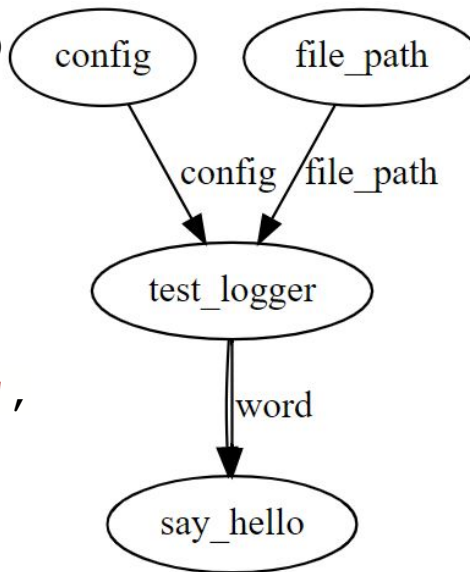
Prefect - Hello World Simple Example



```
def define_test_flow():  
    with Flow() as flow:  
        file = Parameter("file_path")  
        config = Parameter("config")  
        word = test_logger(file, config)  
        hello(word,  
                upstream_tasks=[word])  
    return flow
```

```
flow = define_test_flow()  
flow.visualize()  
flow.run(file_path="TEST_PATH",  
         config={1: "Test"})
```

```
@task  
def test_logger(file,  
               config):  
    logger.debug(  
        f"Received {file}")  
    return "Hello"  
  
@task(log_stdout=True)  
def hello(word):  
    print(f"Say, {word}.")
```





PREFECT

**Questions?
Prefect or Other**

***Doing now what patients need
next***