

Effective Data Pipelines: Data Management from Chaos

Katharine Jarmul (@kjam)

QCon - London - March 6, 2017

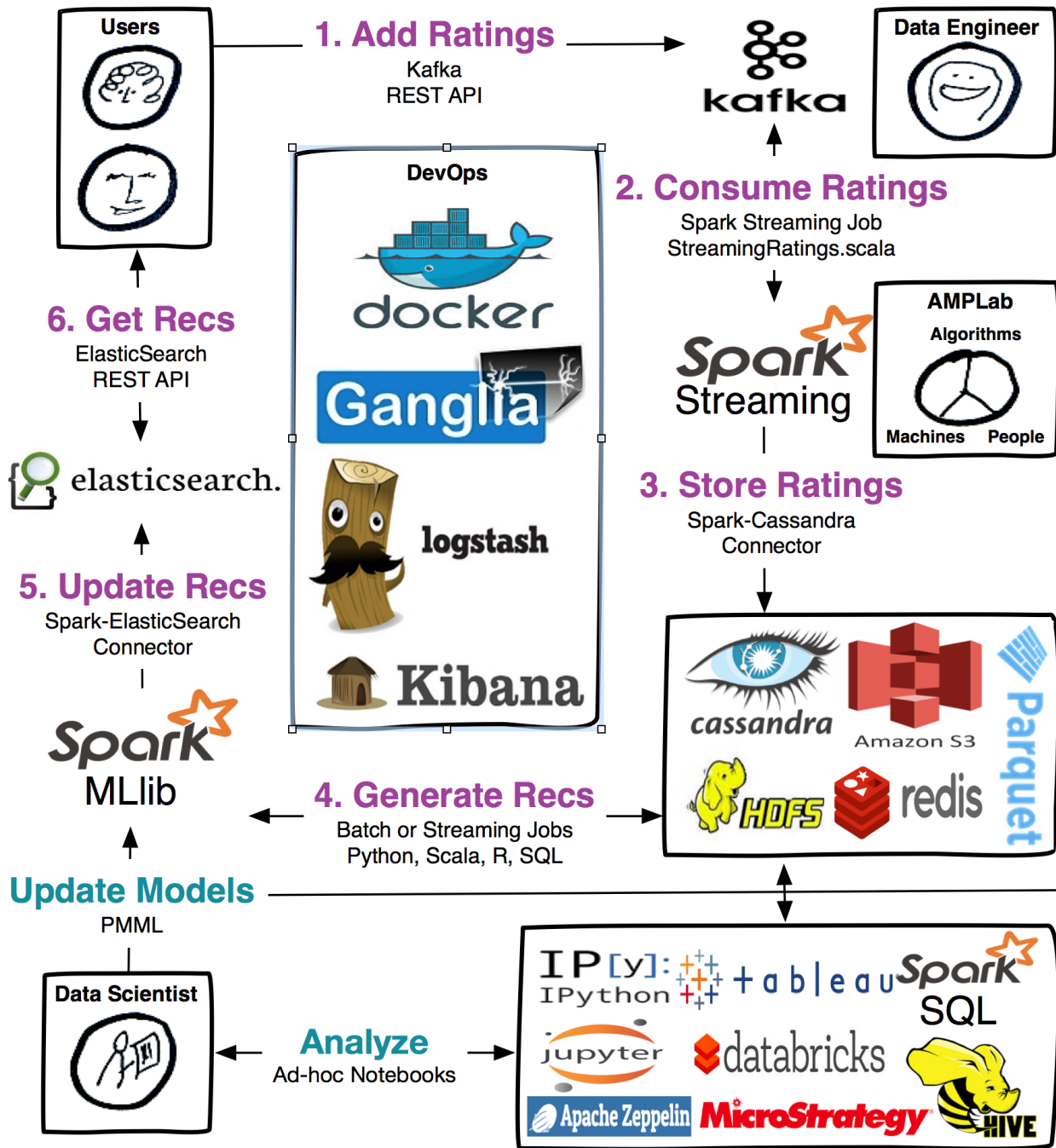
About Katharine

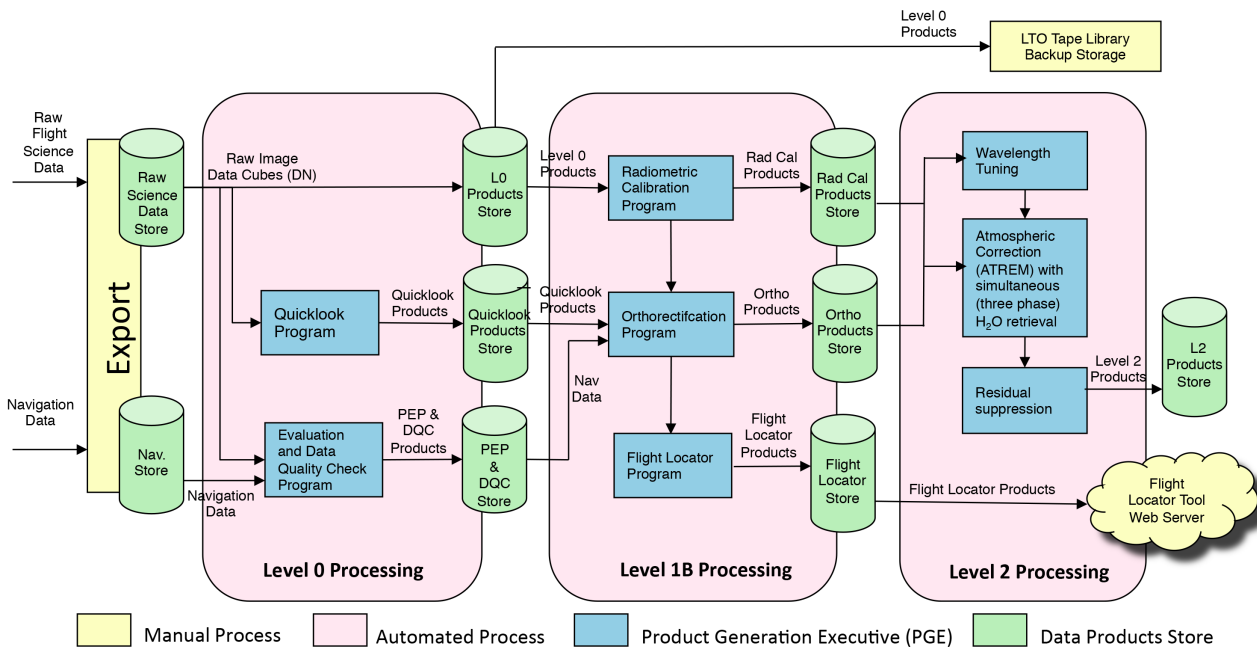
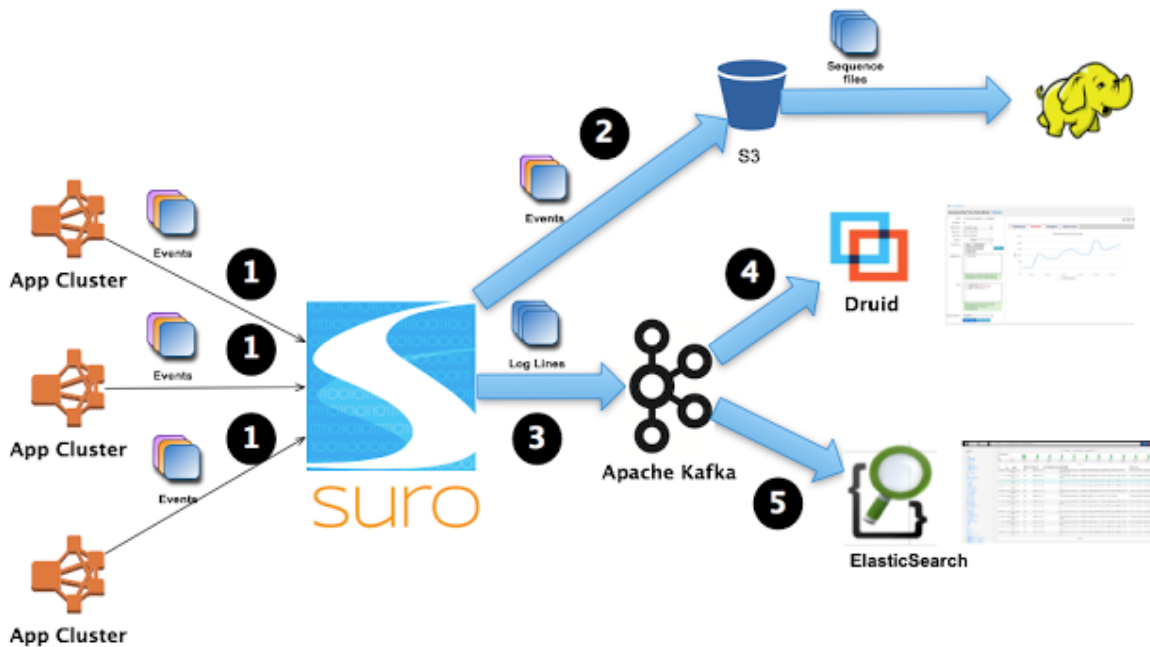


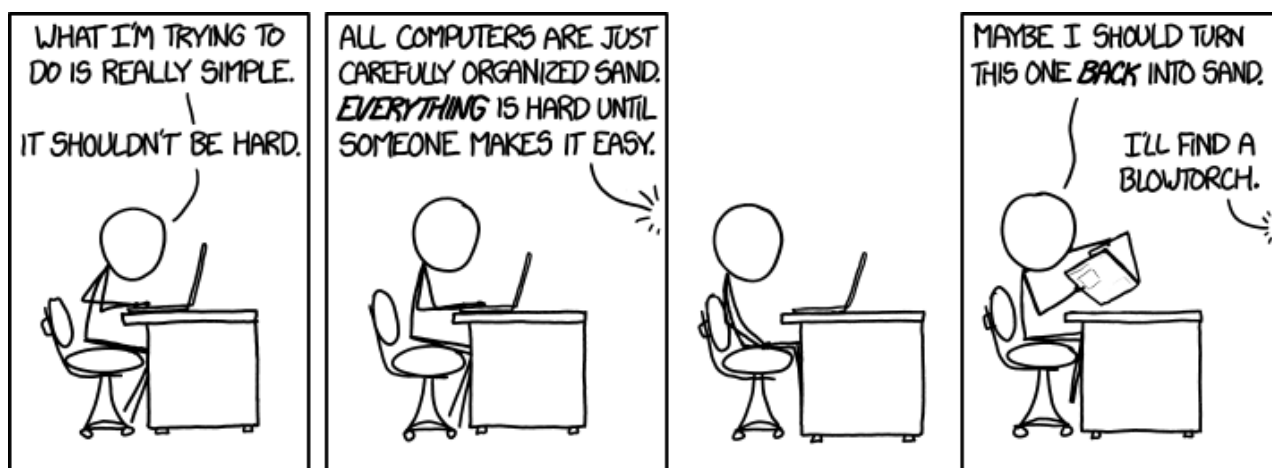
Data Scientist, Engineer, Author, Pythonista

Founder @ kjamistan UG: data science consulting & engineering

Find me at: kjamistan.com - katharine@kjamistan.com - @kjam







Three Questions when Building Data Workflows

1. Who is the producer? Who is the consumer?
2. Where, What, When is the data?
3. What are the constraints? When might they change?

(sorry, that was more like seven.)

Three Tips when Building Data Pipelines

1. **Premature [architecture | optimization | infrastructure] is a bad idea.**
2. **Untested == Unreliable**
3. **Security today, not tomorrow.**

Three Practical Steps for Pipelines

1. **Automate the easy stuff, testing and deployment. Slowly automate the difficult things.**
2. **It *is* infrastructure. Treat it as such.**
3. **Monitoring, alerting and debugging are meaningless without a chain of responsibility.**

Qualities of an Ideal Data Pipeline

- **Idempotent with State Handling**

-- You *will* need to interrupt and rerun tasks (due to bugs, upstream errors, data validation issues).

-- State management is a core part of most pipeline / streaming frameworks. When you can, rely on the framework to do it.

Qualities of an Ideal Data Pipeline

- Scalable and Resilient

-- You may face bursty periods and slow ones. Is autoscaling or provisioning an option?

-- The fallacies of distributed computing often apply to pipelines.

Qualities of an Ideal Data Pipeline

- Replacable or Programmable

-- It's very difficult to foresee where and how your pipeline might grow and change. Be adaptable.

-- Open-source or clear programmability allows for transparent and easy additions.

Qualities of an Ideal Data Pipeline

- Testable and Traceable

-- Upstream, instream, downstream bugs will happen. Make them

easier to find.

-- Find good ways to mock, mirror and replay production data for integration and regression testing.

Qualities of an Ideal Data Pipeline

- Documented and Automated

-- A pipeline without proper documentation is legacy code.

-- Use automated deploys with continuous integration.

Qualities of an Ideal Data Pipeline

- Idempotent with State Handling

- Scalable and Resilient

- Replacable or Programmable

- Testable and Traceable

- Documented and Automated

Pipeline Testaments

- My pipeline is easy to test, debug and monitor.

- There are clear solutions for replaying, rerunning and interrupting tasks or dataflow in my pipeline.**
- There are several teams involved in my pipeline (for security, maintainability and development); however, there is a clear chain of responsibility and protocol for when things go wrong.**
- We have reviewed business and stakeholder use cases. We chose a pipeline structure fitting our current constraints with a straightforward path for growth and change.**

Thank you for listening!

Questions?

Now?

Later? @kjam / katharine@kjamistan.com

Want to talk about pipelines? Data unit testing? Data wrangling? (come find me!)