# Building highly reliable data pipelines @ Datadog

Quentin FRANCOIS

Team Lead, Data Engineering

# Building highly <u>reliable</u> data pipelines @ Datadog

Quentin FRANCOIS

Team Lead, Data Engineering

DataEng Barcelona '18

**Reliability** is the probability that a system will produce correct outputs up to some given time **t.**

Source:  E.J. McClusky & S. Mitra (2004). "Fault Tolerance" in Computer Science Handbook 2ed. ed. A.B. Tucker. CRC Press.

# Highly reliable data pipelines

1. Architecture

# Highly reliable data pipelines

1. Architecture
2. Monitoring

# Highly reliable data pipelines

1. Architecture
2. Monitoring
3. Failures handling

# Historical metric queries

## Avg of system.load.1



## Time series data

| metric | system.load.1 |
|--------|---------------|
| timestamp | 1526382440 |
| value | 0.92 |
| tags | host:i-xyz,env:dev,... |

# Historical metric queries



1 point/second

# Historical metric queries



1 point/day

# Historical metric queries

High resolution data — 1pt /sec

AWS S3

Rollups pipeline

Low resolution data — 1pt /min, 1pt /hour, 1pt /day

- Runs once a day.
- Dozens of TBs of input data.
- Trillions of points processed.

# Highly reliable data pipelines

1. Architecture
2. Monitoring
3. Failures handling

# Our big data platform architecture

| USERS | Web | CLI | Scheduler | Datadog monitoring |
|---|---|---|---|---|
| WORKERS | Luigi | | Spark | |
| CLUSTERS | EMR | EMR | EMR | EMR |
| DATA | S3 | | | |

# Our big data platform architecture

| | | | | Datadog monitoring |
|---|---|---|---|---|
| **USERS** | Web | CLI | Scheduler | |
| **WORKERS** | Luigi | | Spark | |
| **CLUSTERS** | EMR aws | EMR aws | EMR aws | EMR aws |
| **DATA** | S3 aws | | | |

17

# Many ephemeral clusters

- New cluster for every pipeline.
- Dozens of clusters at a time.
- Median lifetime of ~3 hours.

# Total isolation



We know what is happening and why.

# Pick the best hardware for each job

**c3**

**r3**

For CPU-bound jobs

For memory-bound jobs

# Scale up/down clusters

- If we are behind.
- Scale as we grow.
- No more waiting on loaded clusters.

# Safer upgrades of EMR/Hadoop/Spark

**5.13**

**5.12**

**5.12**

**5.12**

**5.12**

# Spot-instance clusters

**+** Ridiculous savings
(up to 80% off the on-demand price)

**-** Nodes can die at any time

# Spot-instance clusters

**(+)** Ridiculous savings
(up to 80% off the on-demand price)

**(-)** Nodes can die at any time

How can we build **highly reliable** data pipelines with **instances killed randomly** all the time?

# No long running jobs

- The longer the job, the more work you lose on average.

- The longer the job, the longer it takes to recover.

# No long running jobs

# No long running jobs

# Break down jobs into smaller pieces

**Vertically** - persist intermediate data between transformations.

**Horizontally** - partition the input data.

# Example
## Rollups pipeline

Input
data

Raw time series
data

Output
data

Aggregated time
series data
(custom file format)

# Example
## Rollups pipeline

**(1)** **Aggregate** high resolution data.

**(2)** **Store** the aggregated data in our custom file format.

Input data — Raw time series data

Output data — Aggregated time series data (custom file format)

# Example
## Vertical split

**1** **Aggregate** high resolution data.

**2** **Store** the aggregated data in our custom file format.

**1**

**Input data** — Raw time series data

**1**

**Checkpoint data** — Aggregated time series data (Parquet format)

**2**

**Output data** — Aggregated time series data (custom file format)

32

# Example
## Horizontal split

**1** **Aggregate** high resolution data.

**2** **Store** the aggregated data in our custom file format.



A  B

C  D

**1**  Raw time series data

Checkpoint data

Aggregated time series data (Parquet format)

**2**

Output data

Aggregated time series data (custom file format)

33

# Example
## Horizontal split

**1** **Aggregate** high resolution data.

**2** **Store** the aggregated data in our custom file format.

Raw time series data

**1**

Aggregated time series data (Parquet format)

**2**

Aggregated time series data (custom file format)

34

# Break down jobs into smaller pieces

Performance

Fault tolerance

# Lessons

- Many clusters for better isolation.
- Break down jobs into smaller pieces.
- Trade-off between performance and fault tolerance.

# Highly reliable data pipelines

1. Architecture
2. Monitoring
3. Failures handling

# Cluster tagging

#anomaly
-detection

#rollups

# Monitor cluster metrics

# Monitor cluster metrics

# Monitor cluster metrics

# Monitor cluster metrics



Decommissioned EMR Nodes - by AZ

cluster_az:us-east-1e
2.85

Sep 9    Mon 10    Tue 11

# Monitor work metrics



More details: datadoghq.com/blog/monitoring-spark/

# Monitor work metrics

```scala
log.info("5. Merged records for all aggregations")
sc.setJobDescription(s"Get records, merge, and encode for all aggregations ($jobIdentifier)")
val numberOfProcessedPoints = sc.longAccumulator("Points")
val recordsForAllAggregations = RawlsBaseRdds.mergeRecordsForAllAggrs(
  filteredRecords, numberOfProcessedPoints, doublePointsCounter, mergeTime)
DatadogMetricsClient.gauge("points_processed", numberOfProcessedPoints.value, tags: _*)
```

# Monitor work metrics

# Monitor work metrics

# Monitor data lag



47

# Lessons

- Measure, measure and measure!
- Alert on meaningful and actionable metrics.
- High level dashboards.

# Highly reliable data pipelines

1. Architecture
2. Monitoring
3. Failures handling

BRACE YOURSELF

FAILURES ARE COMING

# Data pipelines will break

**Hardware failures**

**Increasing volume of data**

**Upstream delays**

**Bad code changes**

# Data pipelines will break

1. Recover fast.
2. Degrade gracefully.

# Recover fast

- No long running job.
- Switch from spot to on-demand clusters.
- Increase cluster size.

# Recover fast: easy way to rerun jobs

- Needed when jobs run but produce some bad data.
- Not always trivial.

# Example: rerun the rollups pipeline

s3://bucket/

📁  2018-01

📁  2018-02

📁  2018-03

📁  2018-04

📁  2018-05

# Example: rerun the rollups pipeline

s3://bucket/2018-05/

📁  as-of_2018-05-01

📁  as-of_2018-05-02

📁  ...

📁  as-of_2018-05-21

# Example: rerun the rollups pipeline

s3://bucket/2018-05/

📂    as-of_2018-05-01

📂    as-of_2018-05-02

📂    ...

📂    as-of_2018-05-21      ⟵    Active location

# Example: rerun the rollups pipeline

s3://bucket/2018-05/

---

📁 as-of_2018-05-01

📁 as-of_2018-05-02

📁 ...

📁 as-of_2018-05-21  ⟵ Active location

📁 as-of_2018-05-22

# Example: rerun the rollups pipeline

s3://bucket/2018-05/

📂 as-of_2018-05-01

📂 as-of_2018-05-02

📂 ...

📂 as-of_2018-05-21

📂 as-of_2018-05-22  ⟵ Active location

# Example: rerun the rollups pipeline

s3://bucket/2018-05/

📂 as-of_2018-05-01

📂 as-of_2018-05-02

📂 ...

✅ 📂 as-of_2018-05-21 ← Active location

❌ 📂 as-of_2018-05-22

# Example: rerun the rollups pipeline

s3://bucket/2018-05/

📁 as-of_2018-05-01

📁 as-of_2018-05-02

📁 ...

📁 as-of_2018-05-21

❌ 📁 as-of_2018-05-22

📁 as-of_2018-05-22_run-2     ← Active location

# Degrade gracefully

- Isolate issues to a limited number of customers.

- Keep the functionalities operational at the cost of performance/accuracy.

# Degrade gracefully: skip corrupted files



- Job failure caused by limited corrupted input data.
- Don't ignore real widespread issues.

# Lessons

- Think about potential issues ahead of time.
- Have knobs ready to recover fast.
- Have knobs ready to limit the customer facing impact.

# Conclusion

Building highly reliable data pipelines

# Conclusion

Building highly reliable data pipelines

- Know your time constraints.

# Conclusion

Building highly reliable data pipelines

- Know your time constraints.
- Break down jobs into small survivable pieces.

# Conclusion

Building highly reliable data pipelines

- Know your time constraints.

- Break down jobs into small survivable pieces.

- Monitor cluster metrics, job metrics and data lags.

# Conclusion

Building highly reliable data pipelines

- Know your time constraints.
- Break down jobs into small survivable pieces.
- Monitor cluster metrics, job metrics and data lags.
- Think about failures ahead of time and get prepared.

# Thanks!

We're hiring!

qf@datadoghq.com
https://jobs.datadoghq.com