# Beyond Ad-hoc Data Science

Oscar Boykin
@posco
oscar@twitter.com

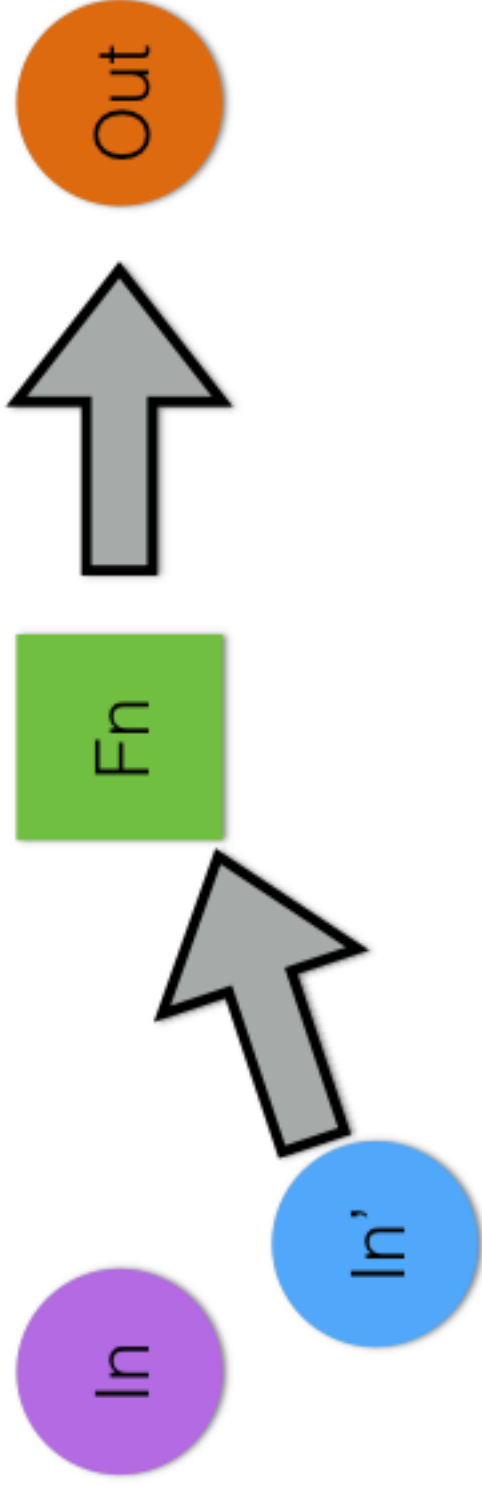…towards a solution for reliable data pipelines

# Ad-hoc Data Science

# ML/Data Science in Prod

- Dashboard of unique active users by geograph, client/version.

- Process to train and evaluate click prediction models (joins impression features + user features and click events)

- Evaluate a classification model on an event to support search (show me news tweets).

- Eventually 100s - 1000s of jobs with complex dependencies.
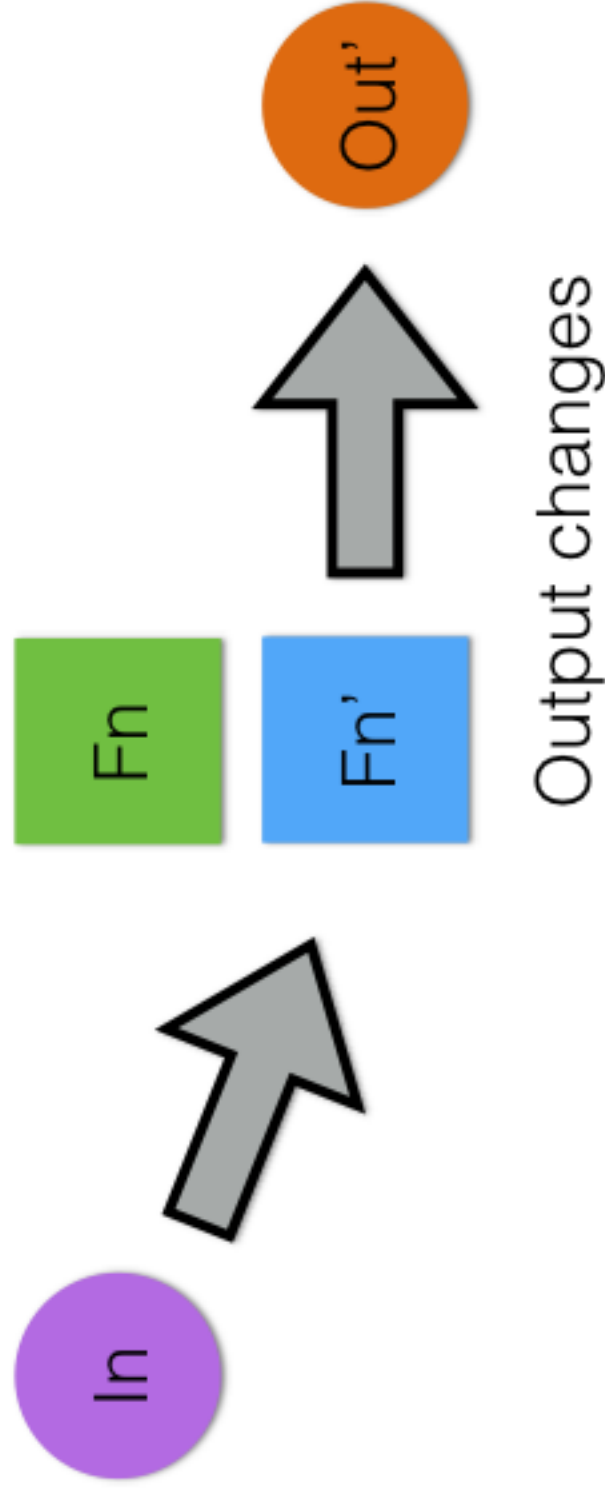
# How do things change?

# How do things change?



Out

Fn

In

In'

Input changes

# How do things change?

In → Fn → Out

# How do things change?

In → Fn / Fn' → Out'

Output changes

# Basic Job Safety

- Tests! At Twitter, all production jobs should have tests. Virtually all jobs are in scalding (compiler checks and standard testing works).

- All jobs should be pure functions: output only depends on input, implies idempotency (safe to re-run the job).
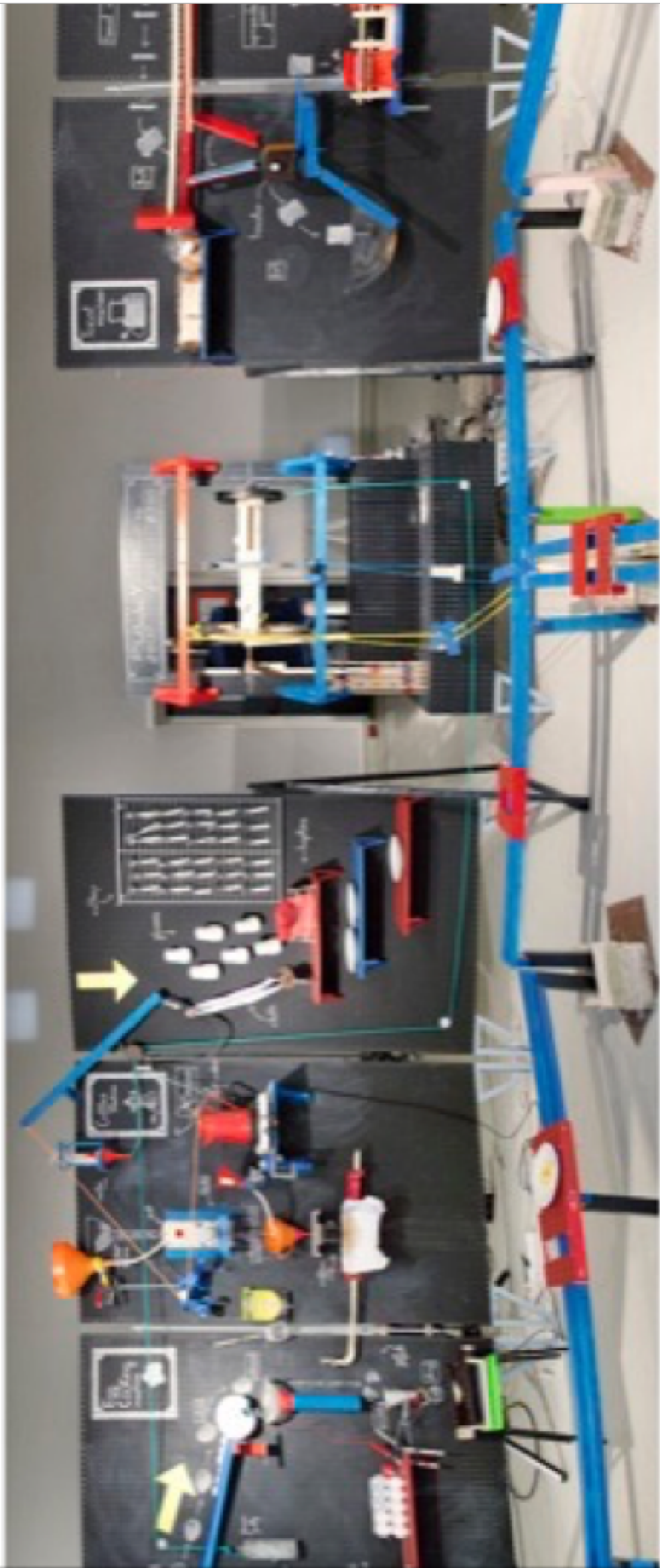
# How Inputs Change

- Data moves: new HDFS cluster, Kinesis to Kafka, private cluster to EMR,

- Product changes => logging changes.

- Data dependency is removed: some team stops producing a model/aggregate/feature.

# How Output Changes

- Want to produce more detail/better model:

  - e.g. make a dashboard of unique users by geo/ client+login frequency. <= new data schema

  - consume more features to improve prediction

  - change modeling technique => rollout Deep Learning!!!

  - bugs are removed (and added)

O(1000) batch jobs, O(100) streaming jobs => things break all the time without some rigor.

# It gets complicated

# The Problem

- How to manage changing data and changing jobs.

- Somewhat unique to data science/ML which is interested in data-artifacts.

- Somewhat similar to the problem of compiling large code bases from complex source dependencies (but that generally has a simpler notion of time).

# 3 Solutions

- Restrict how schemas to evolve

- Decouple data locations from job logic

- Build a system to track data and job dependencies

# Schemas

- a Good Thing!

- some folks like unstructured data.

- logic always requires some (subset of) a schema, *better to be explicit.*

# Schemas

- At Twitter, all data has a schema which we can express using thrift. Written to HDFS/Kafka/etc...

- Thrift has boolean, numbers, strings, lists, sets, maps, structs, and unions. Structs members (columns) can be optional (nullable) or required (not nullable).

- Columns: (add|remove|switch) * (optional|required)

# Thrift Example

```thrift
struct UniqueUserCount {
  0: optional string geo;
  1: optional string client;
  2: required long timestamp;
  3: required long uniqueCount;
}
```

# Thrift Example

```
struct UniqueUserCount {
    0: optional string geo;
    1: optional string client;
    2: required long timestamp;
    3: required long uniqueCount;
}
```

Can generalize this to other schema systems (SQL, protobuf, etc...)

| | Old Reader / New Data | New Reader / Old Data |
| --- | --- | --- |
| add optional | Ok (ignored) | Ok (always absent) |
| remove optional | Ok (always absent) | Ok (ignored) |
| add required | Ok (ignored) | Fail (missing data) |
| remove required | Fail (missing data) | Ok (ignored) |
| opt -> req | Ok (always present) | Fail (missing data) |
| req -> opt | Fail (missing data) | Ok (always present) |

| | Old Reader / New Data | New Reader / Old Data |
|---|---|---|
| add optional | Ok (ignored) | Ok (always absent) |
| remove optional | Ok (always absent) | Ok (ignored) |
| add required | Ok (ignored) | Fail (missing data) |
| remove required | Fail (missing data) | Ok (ignored) |
| opt -> req | Ok (always present) | Fail (missing data) |
| req -> opt | Fail (missing data) | Ok (always present) |

# How to make this work?

- Schemas live in the same repo as the code: jobs won't compile if schema change out of sync with logic. Does not solve the problem of deployed jobs (can't redeploy everything synchronously).

- Tool checks the previous version, ensures that all (public) schemas only add or remove optional fields. Otherwise, the check fails, and the code is not merged to master.

# Decoupling Logic and Data

# Storage

- At Twitter, production data pipelines are written in with scalding.

- Problem: overly tight coupling of paths to jobs.

# Example Job

```scala
 1 case class Impression(id: Long, contentId: Long, timestamp: Long)
 2 case class CountRecord(start: Long, end: Long, count: Long)
 3
 4 object MegaCount {
 5
 6   case class ImpressionSource(dateRange: DateRange) extends
 7     HourlyDataSource[Impression]("/logs/impression/", dateRange)
 8
 9   case class CountSink(dateRange: DateRange) extends
10     HourlyDataSink[CountRecord]("/aggregates/unique_users/", dateRange)
11
12   val OneDay = 1000 * 86400 /* (millis / sec) * (secs / day) */
13
14   def job(dr: DateRange) = {
15     TypedPipe.from(ImpressionSource(dr))
16       .map { imp ⇒ (dayOf(imp.timestamp), Set(imp.id)) }
17       .sumByKey
18       .mapValues { set ⇒ set.size }
19       .map { case (start, size) ⇒ CountRecord(start, start + OneDay, size) }
20       .write(CountSink(dr))
21   }
22 }
```

# Example Job

```scala
1  case class Impression(id: Long, contentId: Long, timestamp: Long)
2  case class CountRecord(start: Long, end: Long, count: Long)
3
4  object MegaCount {
5
6    case class ImpressionSource(dateRange: DateRange) extends
7      HourlyDataSource[Impression]("/logs/impression/", dateRange)
8
9    case class CountSink(dateRange: DateRange) extends
10     HourlyDataSink[CountRecord]("/aggregates/unique_users/", dateRange)
11
12   val OneDay = 1000 * 86400 /* (millis / sec) * (secs / day) */
13
14   def job(dr: DateRange) = {
15     TypedPipe.from(ImpressionSource(dr))
16       .map { imp ⇒ (dayOf(imp.timestamp), Set(imp.id)) }
17       .sumByKey
18       .mapValues { set ⇒ set.size }
19       .map { case (start, size) ⇒ CountRecord(start, start + OneDay, size) }
20       .write(CountSink(dr))
21   }
22 }
```

Physical paths

move data => breakage

Logical Names

```scala
1  case class Impression(id: Long, contentId: Long, timestamp: Long)
2  case class CountRecord(start: Long, end: Long, count: Long)
3
4  object MegaCount {
5    case class ImpressionSource(dr: DateRange) extends
6      Source[Impression]("impressions", dr)
7    case class CountSink(dr: DateRange) extends
8      Sink[CountRecord]("counts", dr)
9
10   val OneDay = 1000 * 86400 /* (millis / sec) * (secs / day) */
11
12   def job(dr: DateRange) = {
13     TypedPipe.from(ImpressionSource(dr))
14       .map { imp ⇒ (dayOf(imp.timestamp), Set(imp.id)) }
15       .sumByKey
16       .mapValues { set ⇒ set.size }
17       .map { case (start, size) ⇒ CountRecord(start, start + OneDay, size) }
18       .write(CountSink(dr))
19   }
20 }
```

Logical names

# Logical Names

```
 1 case class Impression(id: Long, contentId: Long, timestamp: Long)
 2 case class CountRecord(start: Long, end: Long, count: Long)
 3
 4 object MegaCount {
 5   case class ImpressionSource(dr: DateRange) extends
 6     Source[Impression]("impressions", dr)
 7   case class CountSink(dr: DateRange) extends
 8     Sink[CountRecord]("counts", dr)
 9
10   val OneDay = 1000 * 86400 /* (millis / sec) * (secs / day) */
11
12   def job(dr: DateRange) = {
13     TypedPipe.from(ImpressionSource(dr))
14       .map { imp ⇒ (dayOf(imp.timestamp), Set(imp.id)) }
15       .sumByKey
16       .mapValues { set ⇒ set.size }
17       .map { case (start, size) ⇒ CountRecord(start, start + OneDay, size) }
18       .write(CountSink(dr))
19   }
20 }
```

Logical names    Makes a service

call at runtime,

match schema,

resolve paths

Tracking Logical to Physical

# Mappings

- Now we need a service that understands this mapping.

- Hive Metastore is a common choice (works well with Hive, Presto,....)

- Twitter has complex schemas, and needed better scalding support. Wanted a service that tracks not only data, but jobs.

# Tracking Job Evolution

# Tracking Job Evolution

- Wanted a system that understands that jobs change, that schemas change.

- Wanted a full change log: append only structure (immutable records).

- How to handle bugfixes when you have an immutable log?

# \<design\>

logical plane (independant of location, format)

**Application**
role, env
domain
name

owns

composed of

**AppComponent**
role, env
domain
name

produces

consumes

owns

owns

Application

owns

**Logical Dataset**
role, env
name
schema
isPublic

**Deployment**
deployment info
zone (aka dc)
sha
mesos role
client type

**Physical Dataset**
rootLocation (URL)
dalV1URL
type:Partitioned|Snapshot|Stream

LDAP

**group**

controls

**Service
account
(role)**

owns

**user**

**Application**
role, env
domain
name

composed of

**AppComponent**
role, env
domain
name

Application

owns

produces

consumes

**Logical Dataset**
role, env
name
schema
isPublic

owns

logical plane (independant of location, format)

owns

**Deployment**
deployment info
zone (aka dc)
sha
mesos role
client type

**Physical Dataset**
rootLocation (URL)
dalV1URL
type:Partitioned|Snapshot|Stream

**AppComponent**
role, env
domain
name

consumes

consumes

**Logical Dataset**
role, env
name
schema
isPublic

**Deployment**
deployment info
zone (aka dc)
sha
mesos role
client type

**Run Event**
type (BatchRun |
StreamingStarted |
StreamingEnded)
BatchRange (start, end)

**Physical Dataset**
rootLocation (URL)
dalV1URL
type:Partitioned|Snapshot|Stream

**Segment**
physical location = (URL)
physical schema
format
SegmentRange(start, end)
revision

**AppComponent**
role, env
domain
name

**Logical Dataset**
role, env
name
schema
isPublic

produced

consumed

consumes

replicated from

Physical plane (immutable)

**Deployment**
deployment info
zone (aka dc)
sha
mesos role
client type

**Run Event**
type (BatchRun |
StreamingStarted |
StreamingEnded)
BatchRange (start, end)

produced

consumed

**Physical Dataset**
rootLocation (URL)
dalV1URL
type:Partitioned|Snapshot|Stream

**Segment**
physical location = (URL)
physical schema
format
SegmentRange(start, end)
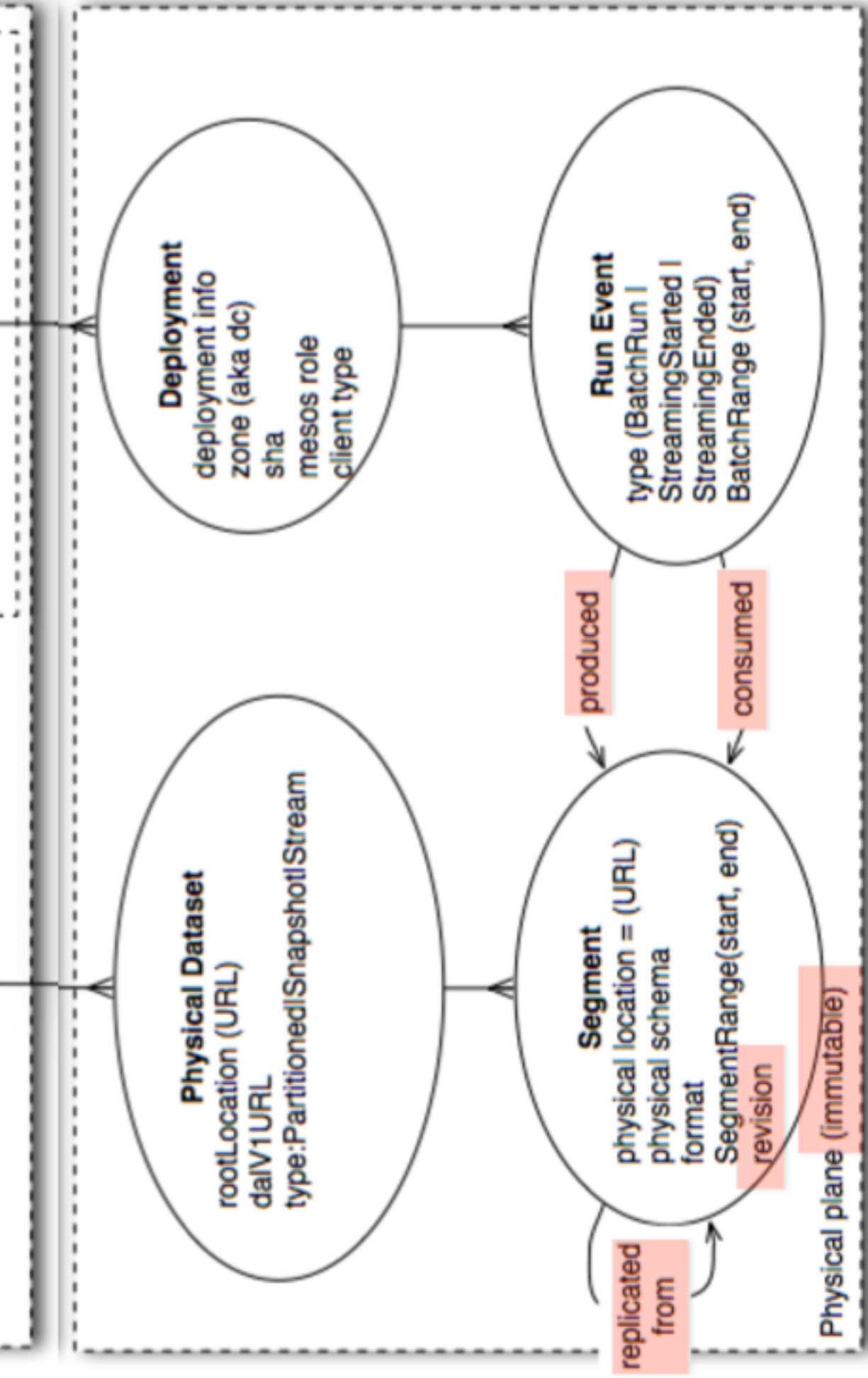revision

replicated
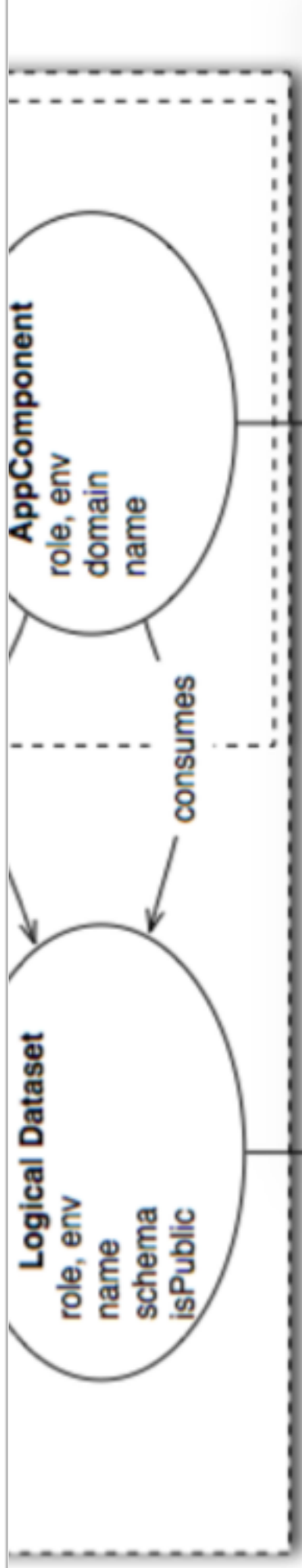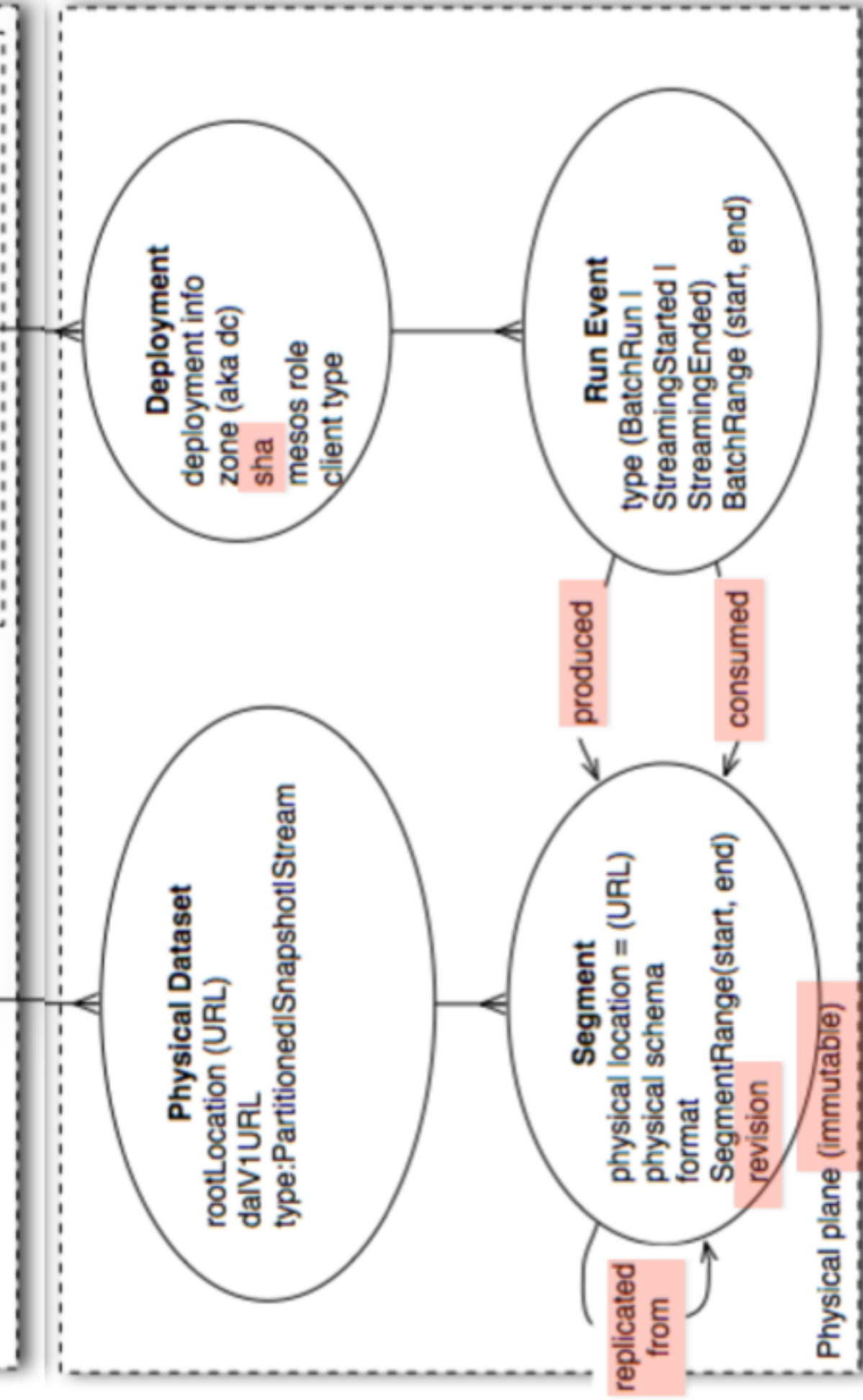from

Physical plane (immutable)

- Each segment has a physical schema that was actually used to write it.

- All the schemas for a given dataset are always compatible (never allow incompatible schema changes).

- When jobs are moved to new clusters, we can replicate needed segments automatically, and record their existence.

**AppComponent**
role, env
domain
name

consumes

**Logical Dataset**
role, env
name
schema
isPublic

**Deployment**
deployment info
zone (aka dc)
sha
mesos role
client type

**Run Event**
type (BatchRun |
StreamingStarted |
StreamingEnded)
BatchRange (start, end)

produced

consumed

**Physical Dataset**
rootLocation (URL)
dalV1URL
type:Partitioned|Snapshot|Stream

**Segment**
physical location = (URL)
physical schema
format
SegmentRange(start, end)
revision

replicated from

Physical plane (immutable)

# Other Issues

- If we discover a bug in a given git sha, we can query to see which data was produced. Can re-run the job and downstream jobs to produce new revisions of segments.

- Currently, jobs are run by a mesos cron scheduler. Read the job state to see if dependencies are ready. Can move to a non-polling model where jobs are run as soon as dependencies are available: *open question how to best express (time) dependencies in jobs.*

- Alerting is external to this: alert if a given job has not had a successful run in the last time interval. Can explain why: which Jobs/Segments are blocking.

The Future

- We want to get have an artifact build system that works for thousands of jobs, running 10-100k instances per day, with 100-1000s of users.

- Logic and data dependencies in the same place, independent of physical storage/format/etc… details.

- Understands version control, can easily (automatically?) schedule backfills when a job's code changes.

- Has a manageable model of time, understands goal deadlines and priorities. Can predict when a goal deadline can't be met, good alerts (actionable, early)

Thank You!
@posco
oscar@twitter.com