# From batch to streaming to both

Herman Schaaf, Senior Software Engineer

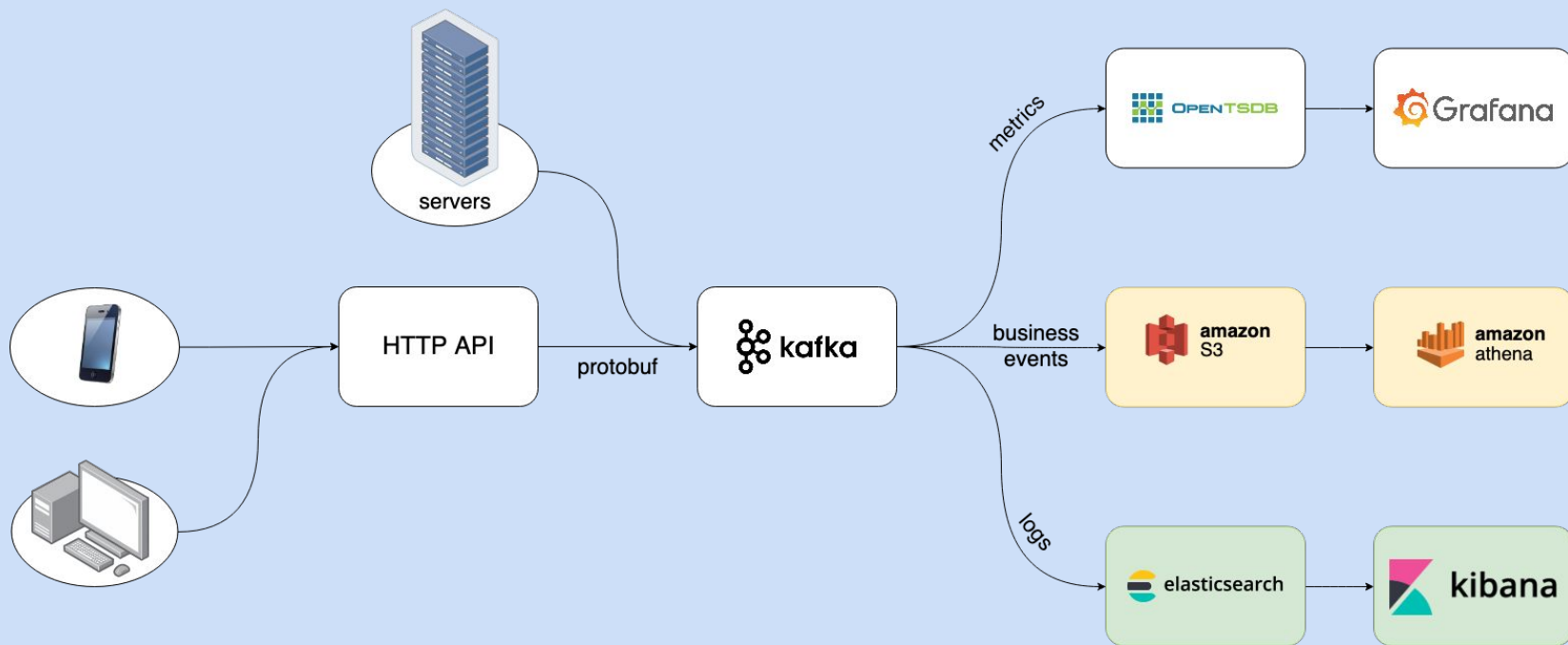Skyscanner

@ironzeb

# A Story

# About me

Herman Schaaf, Senior Software Engineer
Data Platform Tribe

Skyscanner

@ironzeb

# "The Cube"

# From batch to streaming

@ironzeb

# The Single Unified Log

@ironzeb

# The Single Unified Log

@ironzeb

Lesson 1:

Conway's Law is true for data platforms

"Organizations which design **data platforms** are constrained to produce designs which are copies of their communication structures"

Being self-serve is good

…but then metadata is critical

@ironzeb

So let's talk about metadata

prod.identity-service.AuditLog.identity.AuditMessage

prod.flyingcircus.applog.applog.Message

prod.raccoon_bandit.experiment.bandit.Metric

# A simple convention

<prod|sandbox|local>.<service-name>.<event-name>.<schema>.<message>

**Descriptive**

What does it mean?
Who owns it?
What does it contain?
Where does it come from?

*we had some of this*

**Structural**

How does it relate to other data sets?
How is it organized?
How is it sorted / partitioned?

*Some, from using protobuf schemas*

**Administrative**

How far does it date back?
How frequently is it updated?
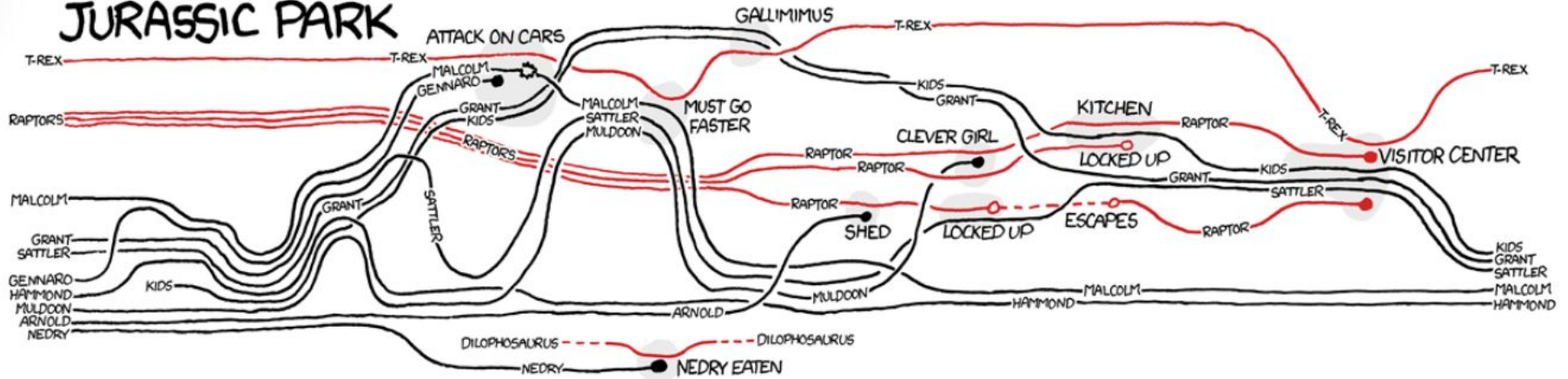How large is it?
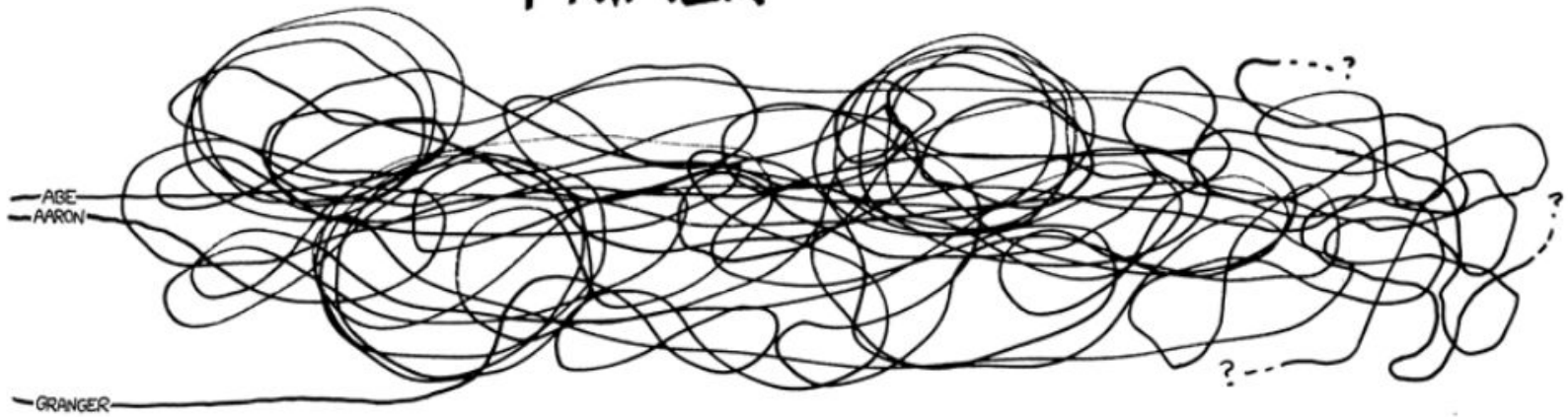How complete is it?

*nope.*

Lesson 2:
Metadata is Critical

- Especially relationships

- Ideally automated

- Ideally from the start

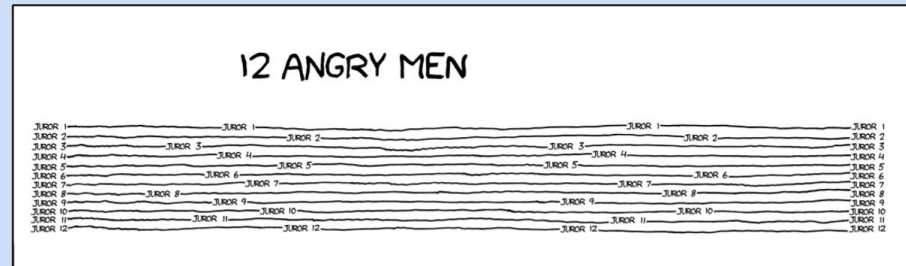- Tools like Schema Registry are a start, but not the full solution

@ironzeb

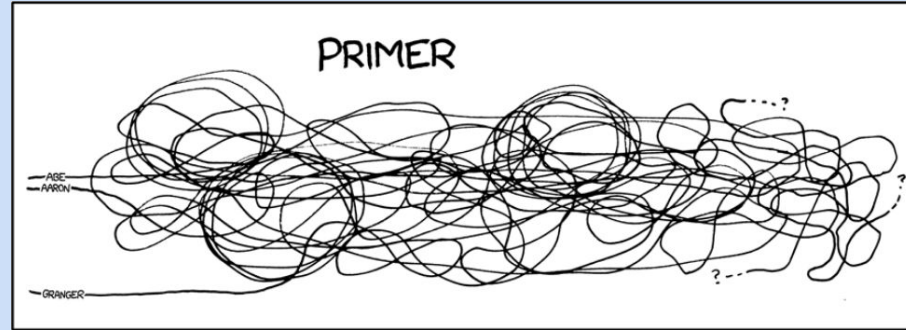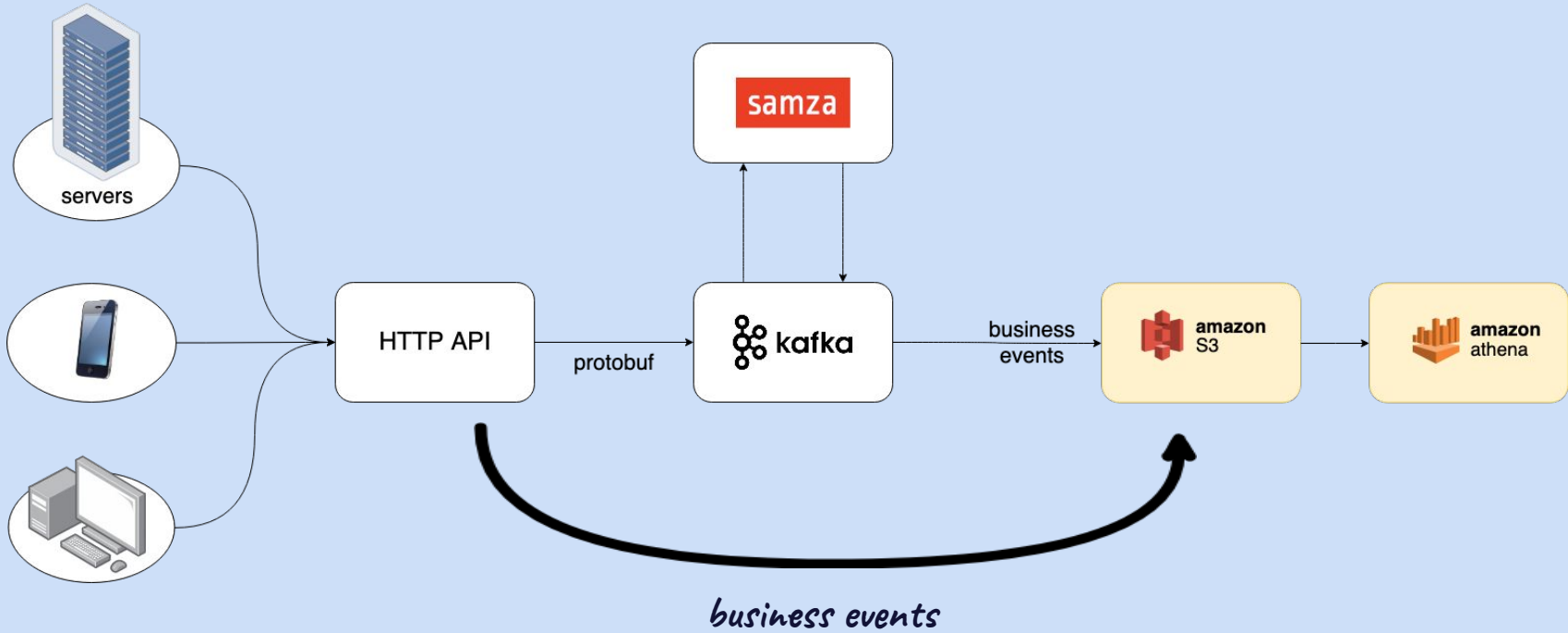https://xkcd.com/657/

https://xkcd.com/657/

# Lesson 3:

## Data Engineers Control the Plot Line

servers

HTTP API

protobuf

samza

kafka

business
events

amazon S3

amazon athena

business events

@ironzeb

# From streaming to both

@ironzeb

Kafka

logs, metrics,
events

HTTP API

business
events

amazon
Kinesis

Flink

amazon
S3

measure    completeness

verify ownership
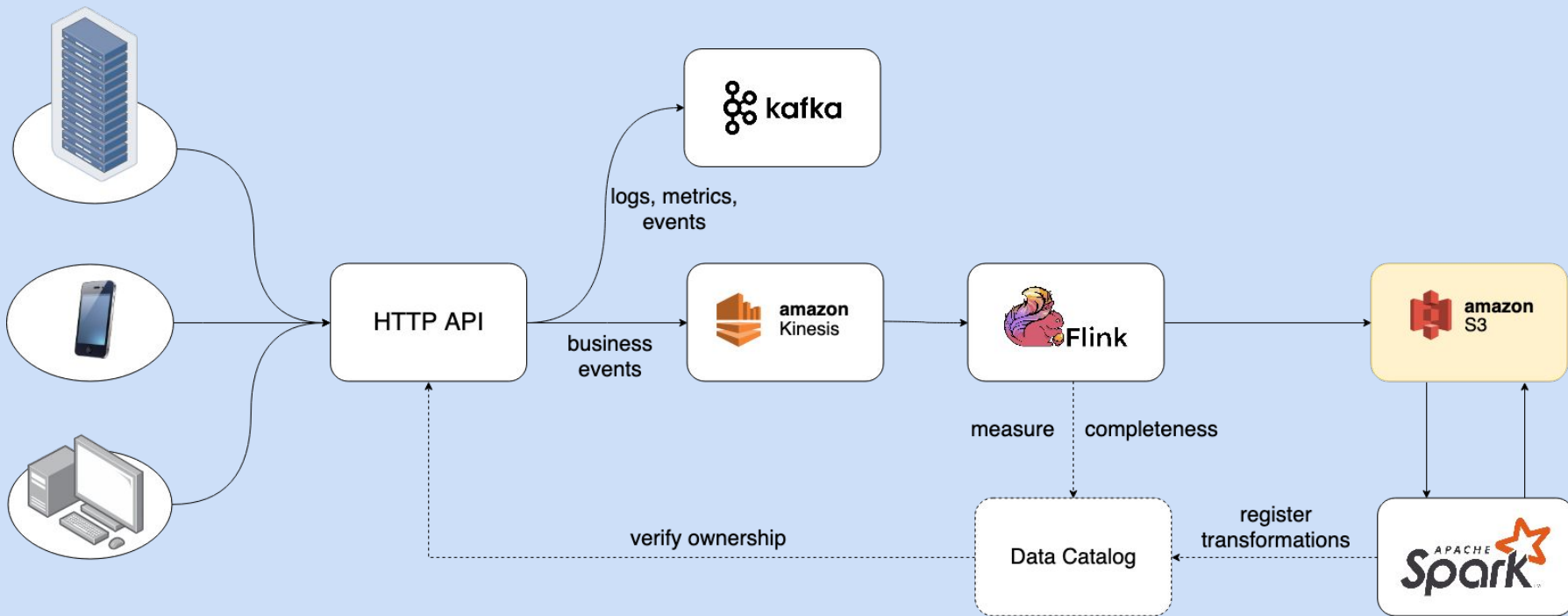
Data Catalog

register
transformations

APACHE
Spark

@ironzeb

## Lesson 4:
## **Repeatability is important**

- Streams have to choose between replays and accepting errors as permanent

- Batch processing can be done again any time

- Going straight to the archive in small batches gets the benefits of both.

# Key Takeaways

- Conway's Law is true for data platforms

- Metadata is Critical

- Data Engineers Control the Plot Line

- Repeatability is important

@ironzeb

# Thanks

**Contact**
If you have any questions regarding Skyscanner please contact:

**Herman Schaaf**
herman.schaaf@skyscanner.net

**Herman Schaaf**
@ironzeb

Skyscanner