

# *Introduction to*



*Apache Streams™*

ApacheCon Europe, September 2015

[sblackmon@apache.org](mailto:sblackmon@apache.org)

# *Agenda*

- **Problem: Proliferation**
- **Activity Streams**
- **Apache Streams**
- **Compatibility**
- **Schemas**
- **Resources**

# *Problem: Proliferation!*

- S **Silos**
- S **Standards**
- S **Schemas**
- S **SDKs**
- S **Databases**
- S **Frameworks**
- S **Runtimes**

# *Silos*

- S It's challenging to get a composite picture of a person or organization because data resides in many systems that are not easily integrated.

# *Standards*

- S We have no universally adopted standard for structuring social profiles, or for transmitting activities across data silos.
- S This is true across web sites, as well across enterprise applications.

# *Schemas*

- S Most silos make minimal if any effort to to promote interoperability by publishing machine-readable schemas for their APIs, or supporting standardized data formats.

# *SDKs*

- S Many data silos recommend usage of one of their SDKs to use their data services, however:
- S These SDKs impose their preferred libraries (such as HTTP clients and json libraries) on us without actually making development easier.

# *Databases*

- S We have an unprecedented range of choices for how and where we store data.
- S Developers often have a handful they prefer to use, and aren't eager to learn the protocols and assumptions of a new DB.
- S Many applications require a polyglot architecture to scale.



# *Frameworks*

- S Frameworks can be very helpful when building scalable systems, but they all enforce conventions and have constraints.
- S Frameworks lead to lock-in, unless your team is extra-ordinarily vigilant.

# *Runtimes*

- S Running code in the cloud may be cheaper, but runtime-specific variation impacts the way we:
  - S Package
  - S Deploy
  - S Configure
  - S Monitor
- S Runtimes lead to lock-in, unless your team is extra-ordinarily vigilant.

# *Activity Streams*

- S A public specification for describing digital activities and identities in JSON format
- S 1.0 – 2011
- S 2.0 – WIP

# *Activity Streams Objectives*

- S Language agnostic**
- S Cross-application interoperability**
- S Support for multiple schemas**
- S Stream Federation**
- S Stream Filtering**

# *Activity Streams Basics*

**S Normalized form for entities and events**

**S** <actor> did <verb> with <object> (to <target>) at <published>

**S objectType**

**S** Person, Organization, Image, Video, etc...

**S Verbs**

**S** Post, Share, Like, etc...

# *Implementation Challenges*

## S Adoption

S Industry support has been tepid at best

## S Ambiguity

S The spec itself is open to interpretation

## S Extensions

S The spec rightly allows for arbitrary extensions

## S Flexibility

S As a result, activities from any two providers are just barely interoperable

## S Validation

S Data correctness or coherence is not covered by spec

# *Apache Streams*

- S A lightweight (yet scalable) framework for Activity Streams
- S An SDK for building data-centric JVM applications
- S A set of patterns for building reliable, adaptable, data processing pipelines

# *Philosophies*

- S Be Database agnostic
- S Be Runtime agnostic
- S Enforce task and document serializability
- S Documents as the core unit of processing
- S Support any type of documents and arbitrary metadata
- S Encourage explicit specification of documents via json schema and xml schema
- S Assist with conversion to and from activitystrea.ms



# Interfaces

## S Provider

S Task running within Activity Streams deployment that sources documents for the stream, likely in their original data format.

## S Processor

S Task running within Activity Streams deployment that transforms documents, perhaps with a synchronous call to an external system.

## S Persist Reader

S Task running within an Activity Streams deployment that sources documents from a file system or database.

## S Persist Writer

S Task running within an Activity Streams deployment that saves documents to a file system or database.

# *Compatibility Dimensions*

- S Providers
- S Persistence
- S Pipelines
- S Runtimes
- S Schemas

# *Compatibility: Providers*

- S Datasift
- S Facebook
- S GMail
- S Gnip
- S Google Plus
- S Instagram
- S Moreover
- S RSS
- S Sysomos
- S Twitter
- S YouTube

# *Compatibility: Persistence*

- S Buffer (file system)
- S Cassandra
- S Elasticsearch
- S Graph (neo4j)
- S HBase
- S HDFS
- S MongoDB
- S Kafka
- S Kinesis
- S S3

# *Compatibility: Runtime Frameworks*

S Docker

S Dropwizard

S Pig

S Spark

S Storm

# *Compatibility: Runtime Roadmap*

- S Crunch
- S Flink
- S Logstash
- S NiFi
- S Samza
- S Spark Streaming
- S Twill

# *Compatibility: Schemas*

- S Schemata are:
  - S The presence and absence of fields and structure
  - S Different from class and from format
- S Strategies for Schema Management
  - S Many-to-Many
  - S Many-to-Mine
  - S Many-to-One
- S Schema-related Challenges

# *Schema Management: Many-To-Many*

- S** For every provider and type, map and convert to compatible types from all other providers
- S** This is the default modality for data and it sucks



# *Schema Management: Many-To-Mine*

- S Specify internal types, then for every provider and type: assess, align and convert to preferred internal representation*
- S This is better, but it fails as soon as we want to interoperate with other departments or organizations who are all using their own internal schemas*
- S Expect to change your internal spec relatively often in early stages, meaning you probably have to either
  - S upgrade your data or*
  - S guarantee backward compatibility in-application**

# *Schema Management: Many-To-One*

- S** For every provider and type, a community dedicated to the inter-operability of that dataset sorts out a reasonable mapping to a relatively static public specification
- S** Where the existing public specs are inadequate, the community can find a way to establish compatibility via convention
- S** Open-source communities and standards bodies can collaborate for benefit of all

# *Schema Challenges: Sharing*

## **S** Business-as-usual:

**S** Schemas are often implicit, shared via unstructured web documentation and language specific sdks

## **S** Streams:

**S** Streams source code contains json and xml schemas for many supported providers

**S** Anyone can import or extend these schemas (via HTTP!)

# *Schema Challenges: Date-Times*

**S** Business-as-usual:

**S** Here's a string, have fun!

**S** Streams:

**S** Every library on the classpath declares its preferred format(s)

**S** Framework resolves any known format and uses Joda to convert to RFC3339

# *Schema Challenges: Versioning*

## **S** Business-as-usual:

**S** Schemas change as product and API features evolve, and everyone just muddles through.

## **S** Streams:

**S** Schemas get published with every release and every snapshot for benefit of those responsible for dependent libraries

**S** Changes get described in release notes

**S** Updates to unit and integration tests

# *Schema Challenges: IDE Support*

**S** Business-as-usual:

**S** Import our SDK or GTFO

**S** Streams:

**S** All streams types have a Serializable POJO representation

**S** Importable with maven to specific version

**S** Convertible to ancestor, sibling, and child types with a cast

**S** Convertible to other types with a one-liner

# *Schema Challenges: Imports*

**S** Business-as-usual:

**S** Every service is an island

**S** Streams:

**S** 'Extends' capability of json schema allows for emergence of a web of related types

**S** Describe your objects as a delta to base schemas or a mashup of several

**S** Undeclared fields propagate by default

# Schema Challenges: Conversion

## S Business-as-usual:

- S Either get too much type safety or none, take your pick
- S If you're lucky, framework helps with serialization and compression

## S Streams:

- S Includes multiple type conversion options, available as processors for your streams or singleton utility classes to embed in your code
  - S jackson conversion
  - S hocon conversion
  - S via java/scala



# Resources

S Website

S <http://streams.incubator.apache.org/>

S Source Code

S <https://github.com/apache/incubator-streams>

S Documentation

S <http://streams.incubator.apache.org/site/0.2-incubating/streams-project/index.html>

S Examples

S <https://github.com/apache/incubator-streams-examples>

S Examples Documentation

S <http://streams.incubator.apache.org/site/0.2-incubating-SNAPSHOT/streams-examples/index.html>