



Securely explore your data

# BIG DATA INFRASTRUCTURE: LESSONS LEARNED

Adam Fuchs, CTO  
Sqrri Data, Inc.  
October 29, 2013



# RAPID APPLICATION DEVELOPMENT REQUIRES...

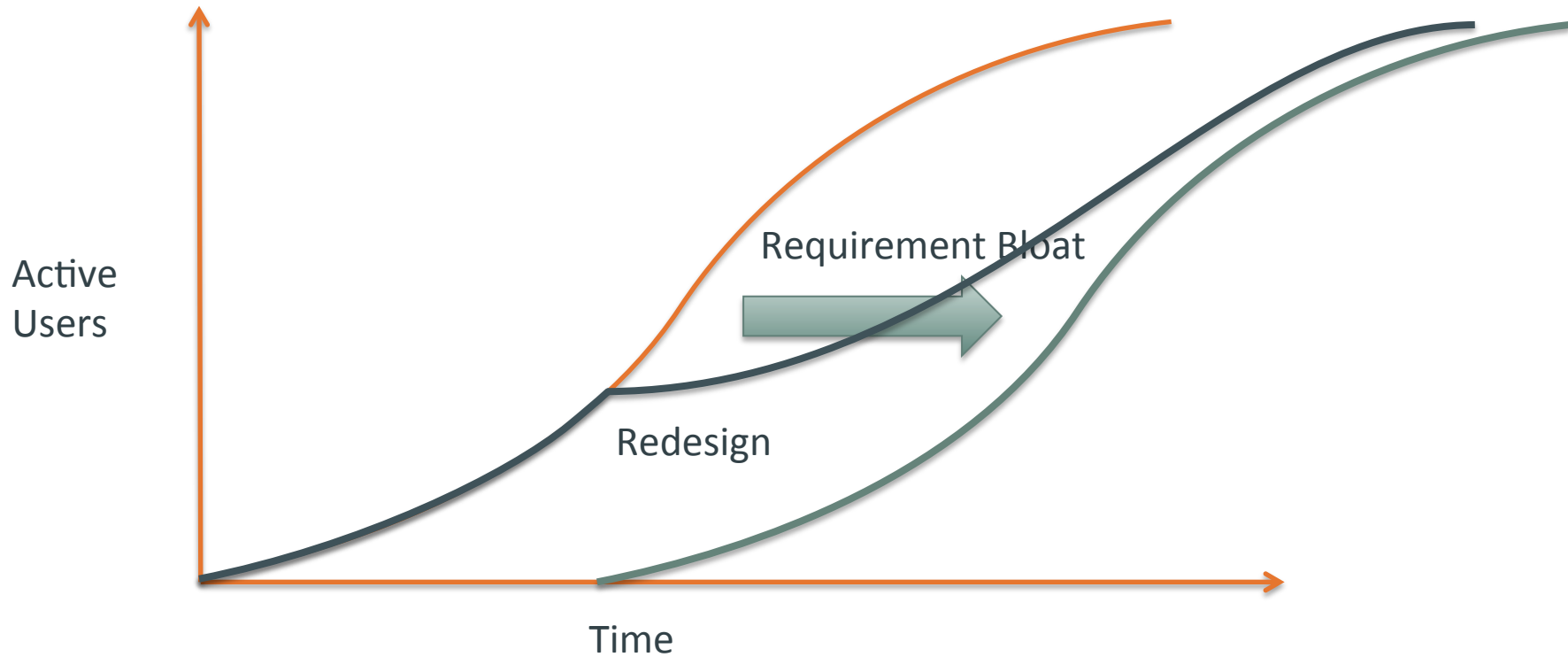
- Prototyping with operational data



- Smooth transition from dev to production
  - Scaling users, data
  - Achieving durability, consistency, availability
  - Security and compliance

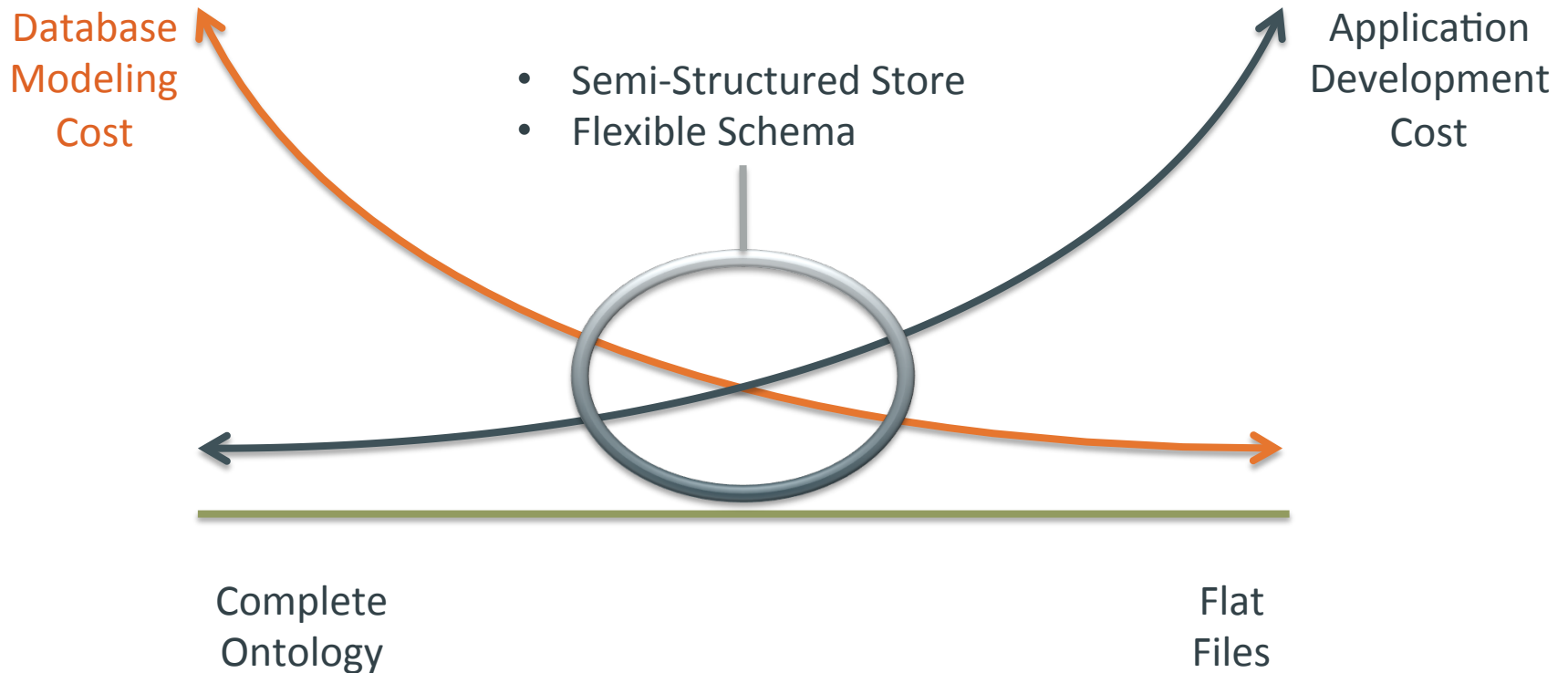
# ADOPTION CURVES

Start small, but design for scalability



# DATA MODELING COSTS

Seek a happy medium...



# DATA-CENTRIC SECURITY

## Simplifying complex security models

The screenshot displays a medical software interface with several key sections:

- Header:** "Handy patients enterprise edition" with standard window controls.
- Left Panel:** Patient information for David Anderson, born 5 January 2009. It includes a list of appointments (e.g., "2 month checkup", "Respiration problem") and a "Diagnosis" section with checkboxes for "General", "My Diagnosis", and "Social".
- Center Panel:** A "Digestive" examination report for Thursday, 22 Jan 2009. It includes sections for "Digestive inspection" (Normal), "Digestive auscultation" (Normal abdomen noises), and "Digestive palpation" (Little pain on the right lower area).
- Right Panel:** A diagram of the human digestive system with labels for the Esophagus, Liver, Stomach, Colon, Gall bladder, Rectum, and Anus. A red arrow points to the lower right quadrant of the abdomen.
- Bottom Panel:** A "Notes" section containing the text: "Father ask many questions, add 10 minutes to consuikation".

Annotations on the image highlight sensitive data:

- A blue circle highlights the patient's name "Anderson" and "David".
- A blue circle highlights the birth date "5 January 2009".
- A blue circle highlights the "Diagnosis" section.
- A blue circle highlights the "Notes" section.
- The text "PII" (Personally Identifiable Information) is overlaid in large blue letters.
- The text "Sensitive Diagnoses" is overlaid in large blue letters with red arrows pointing to the diagnosis and notes sections.
- The text "Doctor's Notes" is overlaid in large blue letters.

# DISCOVERY ANALYTICS

Building blocks for secure, scalable, operational apps.

Common solutions for operational requirements, like:

- Scalability
- Security
- Schema Adaptation
- Multi-tenancy

In domains like:

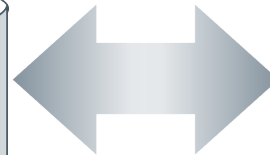
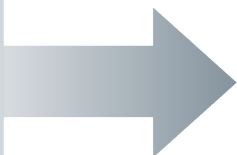
- “Data Lake”
  - Indexed Search
  - Discovery
  - Big-Picture Views
  - Bulk slicing
- “Operational Apps”
  - High concurrency, low-latency
  - Real-time
  - Denormalized

# DATA LAKE

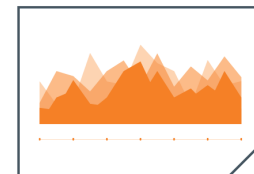
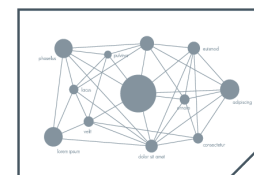
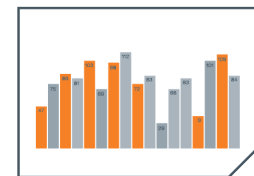
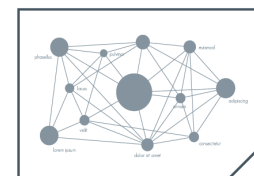
+Data-Centric Security

+Operational, Real-Time

Non-sensitive data  
Sensitive data  
Highly sensitive data  
Highly sensitive data  
Sensitive data  
Non-sensitive data  
Non-sensitive data



Real-Time  
Big Apps

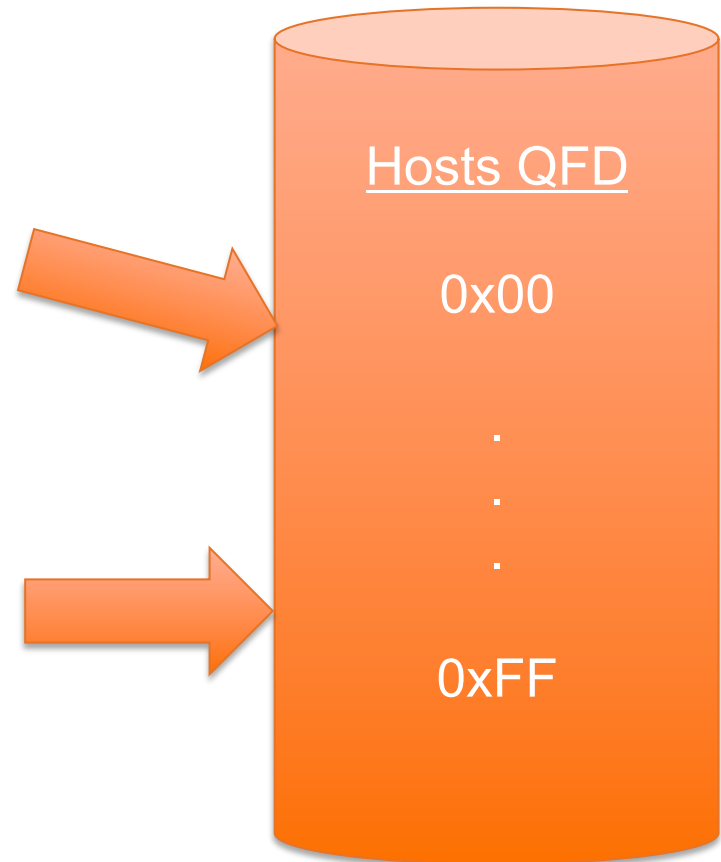




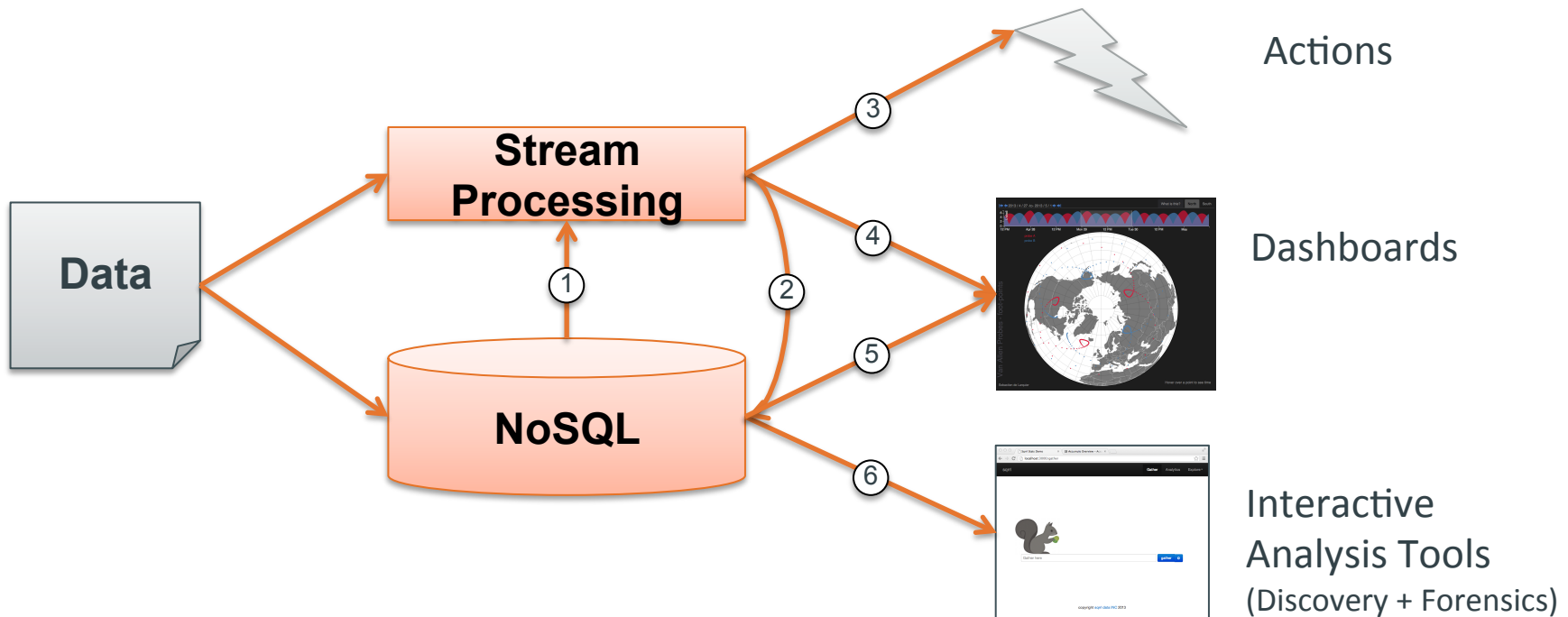
# QUESTION-FOCUSED DATASETS

Source	Destination	Port	Bytes In	Bytes Out	Protocol
10.1.2.3	10.9.8.7	80	73,824	15,632	http

<u>Key</u>	<u>-&gt; Value</u>
10.1.2.3, Bytes In	-> +73,824
10.1.2.3, Bytes Out	-> +15,632
10.1.2.3, Ports Used	-> +{80}
10.1.2.3, Protos Used	-> +{http}
10.9.8.7, Bytes In	-> +15,632
10.9.8.7, Bytes Out	-> +73,824
10.9.8.7, Ports Hosted	-> +{80}
10.9.8.7, Protos Hosted	-> +{http}



# NOSQL: REAL-TIME Real-Time Ecosystem



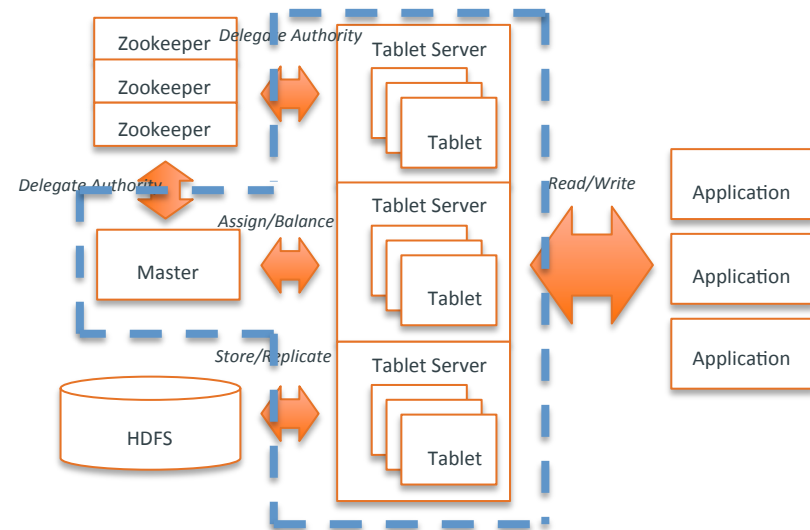
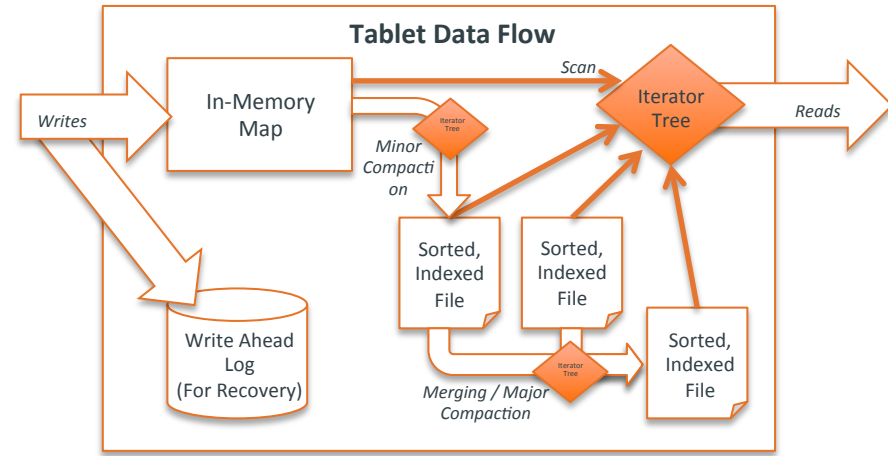
1. SPE queries NoSQL to enrich streaming data
2. SPE persists results in NoSQL for future query
3. SPE takes action automatically
4. SPE issues data-driven alerts
5. Sqrri provides context for dashboards
6. Analysis tools query use Sqrri to search and manipulate historical data

# APACHE ACCUMULO

- Founded in 2008 (internal NSA)
- Modeled after Bigtable
- Open-sourced to ASF in 2011

## Strengths

- **Shared-Nothing** => Scalability
- Micro-Batching for Efficient **Random I/O**
- **Iterators** for local, server-side analytics on sorted data
- **High Concurrency, Low Latency** for Denormalized Data
- Sparse, **Flexible Schema** supports dynamic and diverse data models
- **Cell-level Security** promotes sharing



# SECURITY AT ACCUMULO'S CORE

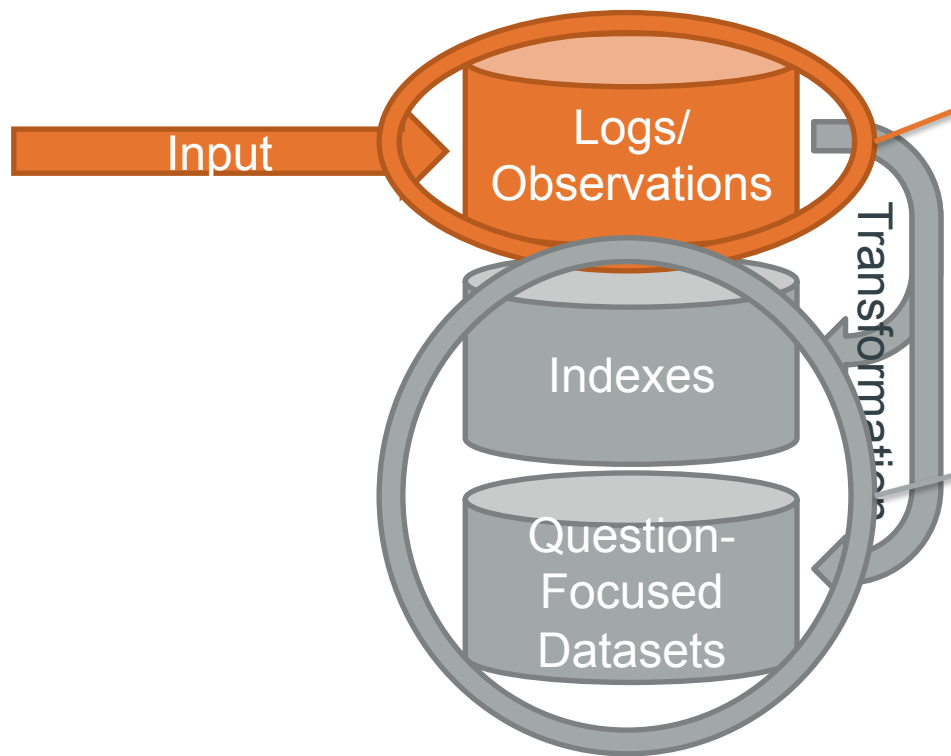
An Accumulo key is a 5-tuple, consisting of:

- **Row:** Controls Atomicity
- **Column Family:** Controls Locality
- **Column Qualifier:** Controls Uniqueness
- **Visibility Label:** Controls Access
- **Timestamp:** Controls Versioning

Row	Col. Fam.	Col. Qual.	Visibility	Timestamp	Value
John Doe	Notes	PCP	PCP_JD	20120912	Patient suffers from an acute ...
John Doe	Test Results	Cholesterol	JD PCP_JD	20120912	183
John Doe	Test Results	Mental Health	JD PSYCH_JD	20120801	Pass
John Doe	Test Results	X-Ray	JD PHYS_JD	20120513	1010110110100...

Accumulo Key/Value Example

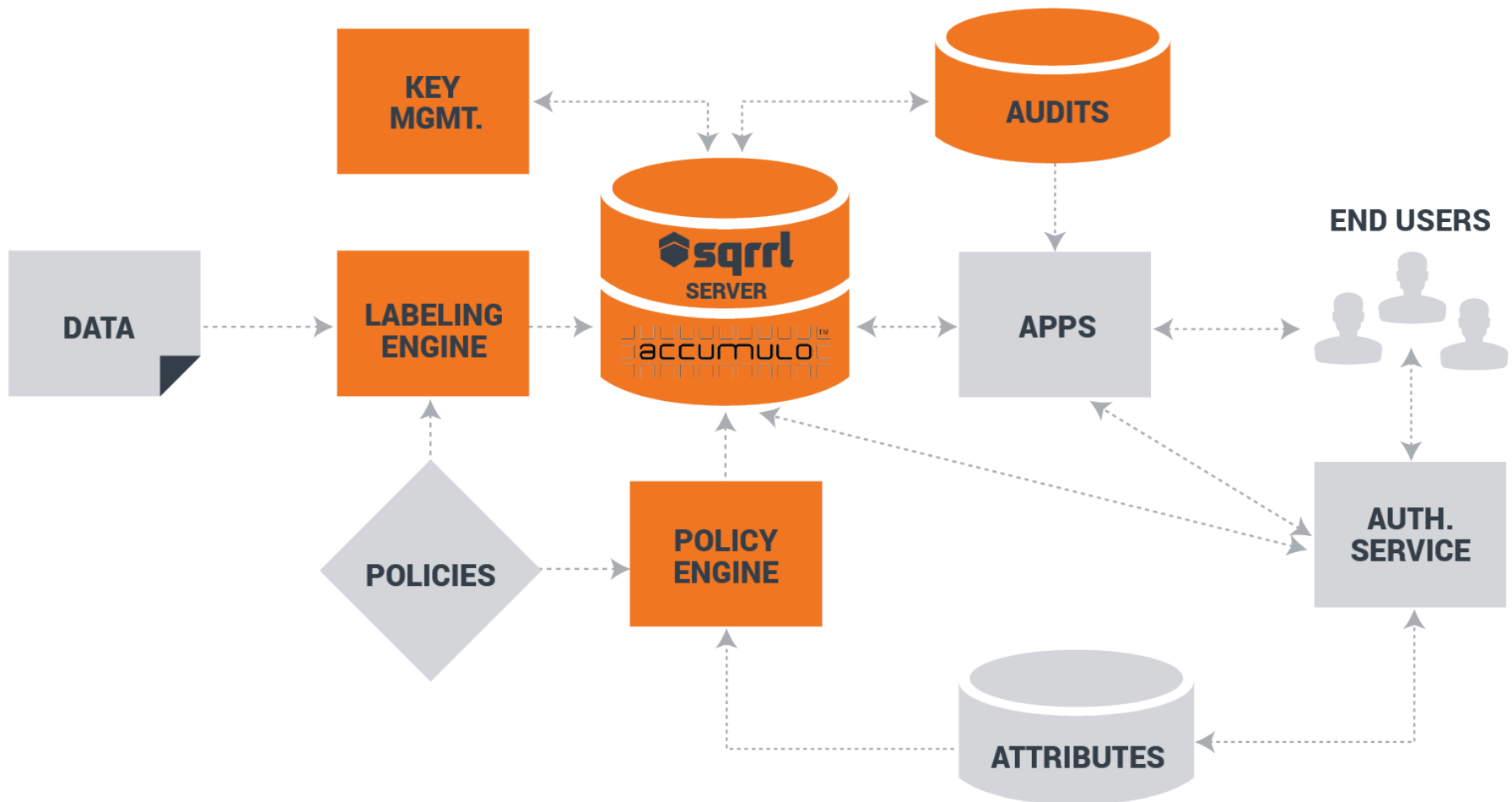
# SECURITY IN INDEXES AND QFDS



- **Simple Provenance**
- *Simple Security Model*
  
- **Complex Provenance**
- *Complex Security Model*

# DATA-CENTRIC SECURITY

**Definition:** Data carries with it information that is required to make policy decisions on its releasability.



# LAYERED ARCHITECTURE

Turtles all the way down...



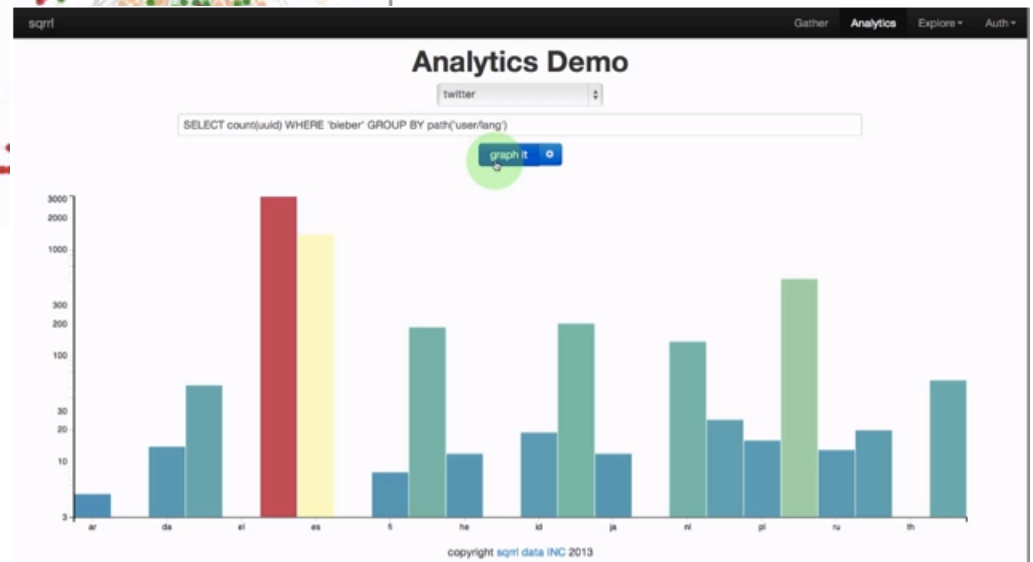
*Sqrrel API over Apache Thrift RPC*  
(JSON, Graph, Aggregation,  
Search, etc.)

*Accumulo RPC*  
(Sorted Key/Value I/O)

*Hadoop RPC*  
(File I/O)

# REAL-TIME OPERATIONAL APPS

## Innovation through rapid development







Securely explore your data

**QUESTIONS?**

**Adam Fuchs, CTO**  
**Sqrri Data, Inc.**  
**October 29, 2013**