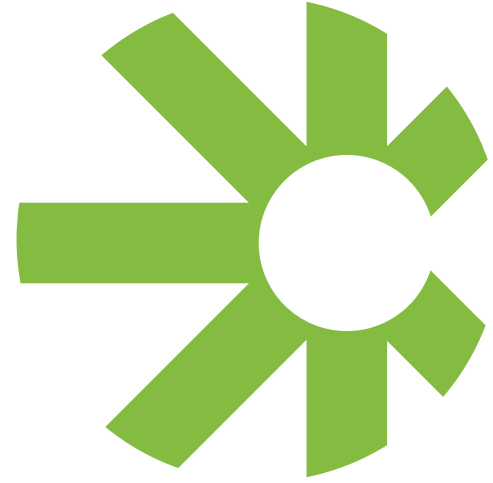DRIVING INNOVATION THROUGH DATA

# ACCELERATING BIG DATA APPLICATION DEVELOPMENT WITH CASCADING

Supreet Oberoi
VP Field Engineering, Concurrent Inc

CONCURRENT

Founded: *2008*
HQ: *San Francisco, CA*

CEO: *Gary Nakamura*
CTO, Founder: *Chris Wensel*

www.concurrentinc.com

**Leader in Application Infrastructure for Big Data**

- Building enterprise software to simplify Big Data application development and management

**Products and Technology**

- *CASCADING*

  *Open Source -* The most widely used application infrastructure for building Big Data apps with over 175,000 downloads each month

- *DRIVEN*

  Enterprise data application management for Big Data apps

**Proven — Simple, Reliable, Robust**

- Thousands of enterprises rely on Concurrent to provide their data application infrastructure.

# ENTERPRISE NEEDS FOR DATA APP INFRASTRUCTURE

- Need reliable, reusable tooling to quickly build and consistently deliver data products

- Need the degrees of freedom to solve problems ranging from simple to complex with existing skill sets

- Need the flexibility to easily adapt an application to meet business needs (latency, scale, SLA), without having to rewrite the application

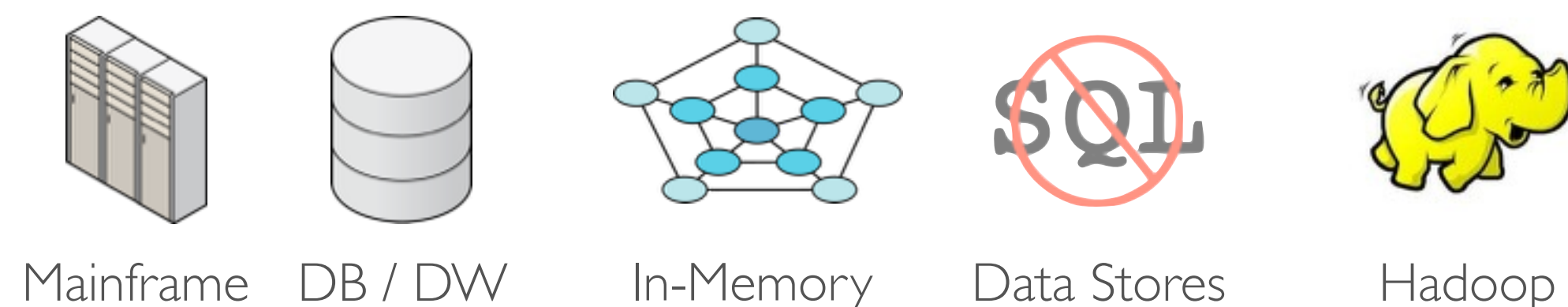- Need operational visibility for entire data application lifecycle

CONCURRENT

**Cascading Apps**

SQL · PMML Predictive Model Markup Language · Scala · Clojure · Ruby · python

**CASCADING**

**System Integration**

Mainframe · DB / DW · In-Memory · ~~SQL~~ Data Stores · Hadoop

**New Fabrics**

MapReduce · Spark · Tez · Storm

- Standard for enterprise data app development

- Your programming language of choice

- Cascading applications that run on MapReduce will also run on Apache Tez, Spark, Storm, and …

**CONCURRENT**

```
String docPath = args[ 0 ];
String wcPath = args[ 1 ];
Properties properties = new Properties();
AppProps.setApplicationJarClass( properties, Main.class );
HadoopFlowConnector flowConnector = new HadoopFlowConnector( properties );
```

**configuration**

```
// create source and sink taps
Tap docTap = new Hfs( new TextDelimited( true, "\t" ), docPath );
Tap wcTap = new Hfs( new TextDelimited( true, "\t" ), wcPath );
```

**integration**

```
// specify a regex to split "document" text lines into token stream
Fields token = new Fields( "token" );
Fields text = new Fields( "text" );
RegexSplitGenerator splitter = new RegexSplitGenerator( token, "[ \\[\\]\\(\\),.]" );
// only returns "token"
Pipe docPipe = new Each( "token", text, splitter, Fields.RESULTS );
// determine the word counts
Pipe wcPipe = new Pipe( "wc", docPipe );
wcPipe = new GroupBy( wcPipe, token );
wcPipe = new Every( wcPipe, Fields.ALL, new Count(), Fields.ALL );
```
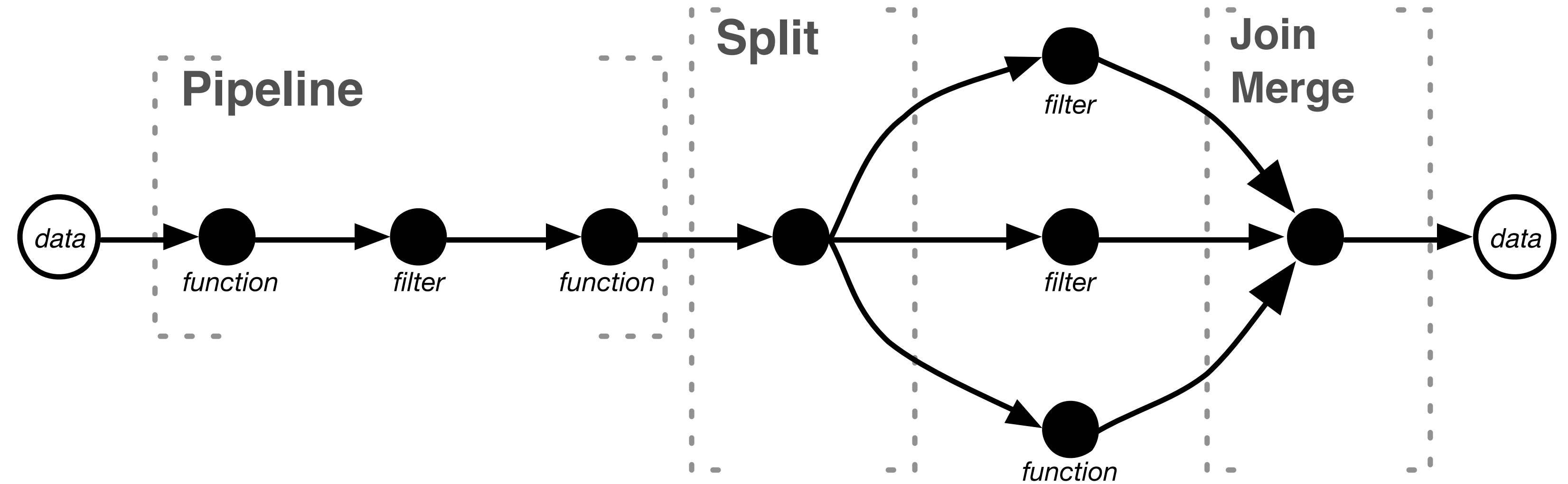
**processing**

```
// connect the taps, pipes, etc., into a flow definition
FlowDef flowDef = FlowDef.flowDef().setName( "wc" )
 .addSource( docPipe, docTap )
 .addTailSink( wcPipe, wcTap );
// create the Flow
Flow wcFlow = flowConnector.connect( flowDef ); // <<-- Unit of Work
wcFlow.complete();                              // <<-- Runs jobs on Cluster
```

**scheduling**

CONCURRENT

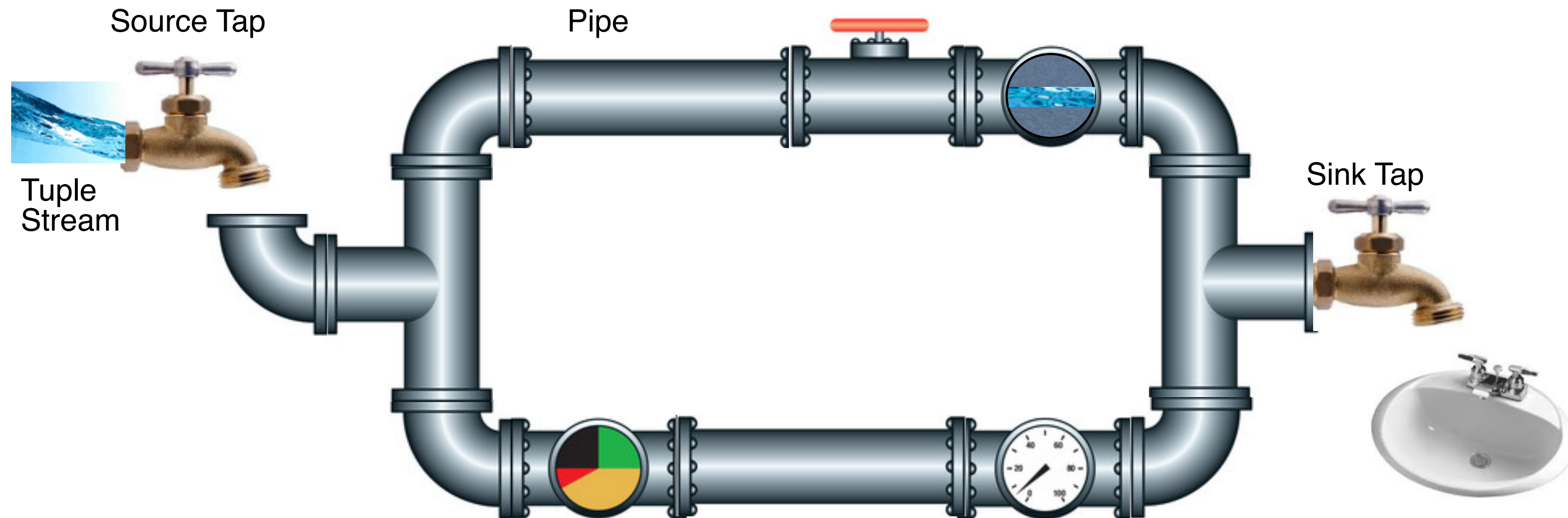# SOME COMMON PATTERNS

- Functions
- Filters
- Joins
  - Inner / Outer / Mixed
  - Asymmetrical / Symmetrical
- Merge (Union)
- Grouping
  - Secondary Sorting
  - Unique (Distinct)
- Aggregations
  - Count, Average, etc



**Topology**

**CONCURRENT**

- The Cascading processing model is based on a metaphor of flows based on patterns

# CASCADING PROCESSING MODEL TERMINOLOGY

| | |
|---|---|
| Tuple Stream | Series of tuples (data record) |
| Fields | Representation of the Tuple Stream, used in operations |
| Pipe | Applies operations to tuples or groups of tuples |
| Branch | Pipes linked together under a common Pipe name |
| Pipe Assembly | An interconnected set of pipe branches |
| Tap | Source or sink for data |
| Flow | Pipe assembly with taps |
| Cascade | Multiple flows grouped together & executed as a single process |

CONCURRENT

# TUPLE STREAM

- A Tuple represents a set of values.

- Consider a Tuple the same as a database record where every value is a column in that table.

- A "tuple stream" is a set of Tuple instances passed consecutively through a Pipe assembly.
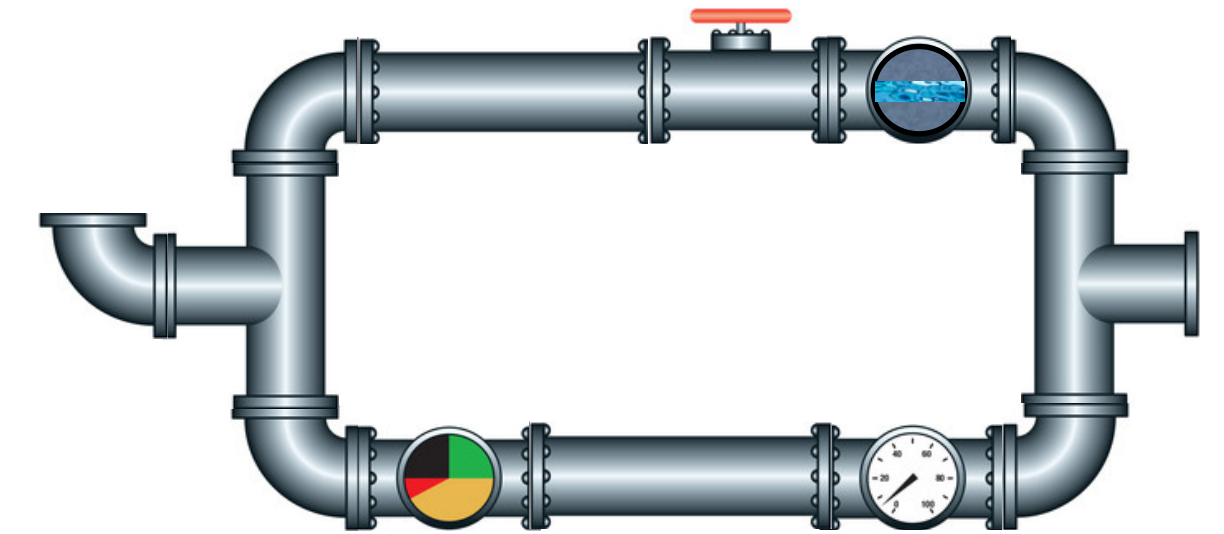
CONCURRENT

# PIPES CAN BE CHAINED TO PERFORM COMPLEX OPERATIONS

- Pipes control the flow of data applying operations to each Tuple or groups of Tuples.

- Pipes work on fields of one or more tuples.

- Pipes allow you to manage a data flow such as doing:

  - Grouping
  - Joining
  - Filtering
  - Buffering
  - Aggregating

CONCURRENT

# PIPES CAN BE BRANCHED AND MERGED

- Pipe Assemblies are an interconnected set of pipe branches modeled as a DAG (Directed Acyclic Graph)

- Pipe Assemblies can consist of splits and/or merges.

- Pipe assemblies are specified independently of the data source they are to process.

- For a pipe assembly to be executed, it must be bound to data sources and sinks (which becomes a flow)

*DAG: collection of vertices and directed edges, each edge connecting one vertex to another, such that there is no way to start at some vertex v and follow a sequence of edges that eventually loops back to v again.*

CONCURRENT

# TAPS ABSTRACT INTEGRATION TO THIRD-PARTY SYSTEMS

- Taps provide the ability to read and write data.

- Taps can be shared between flows and can be restricted to being either sources or sinks.

- Taps can be set up to have the actual file identifiers determined when they run.

- Examples of Taps are:

  - File on the local file system

  - File on a Hadoop distributed file system
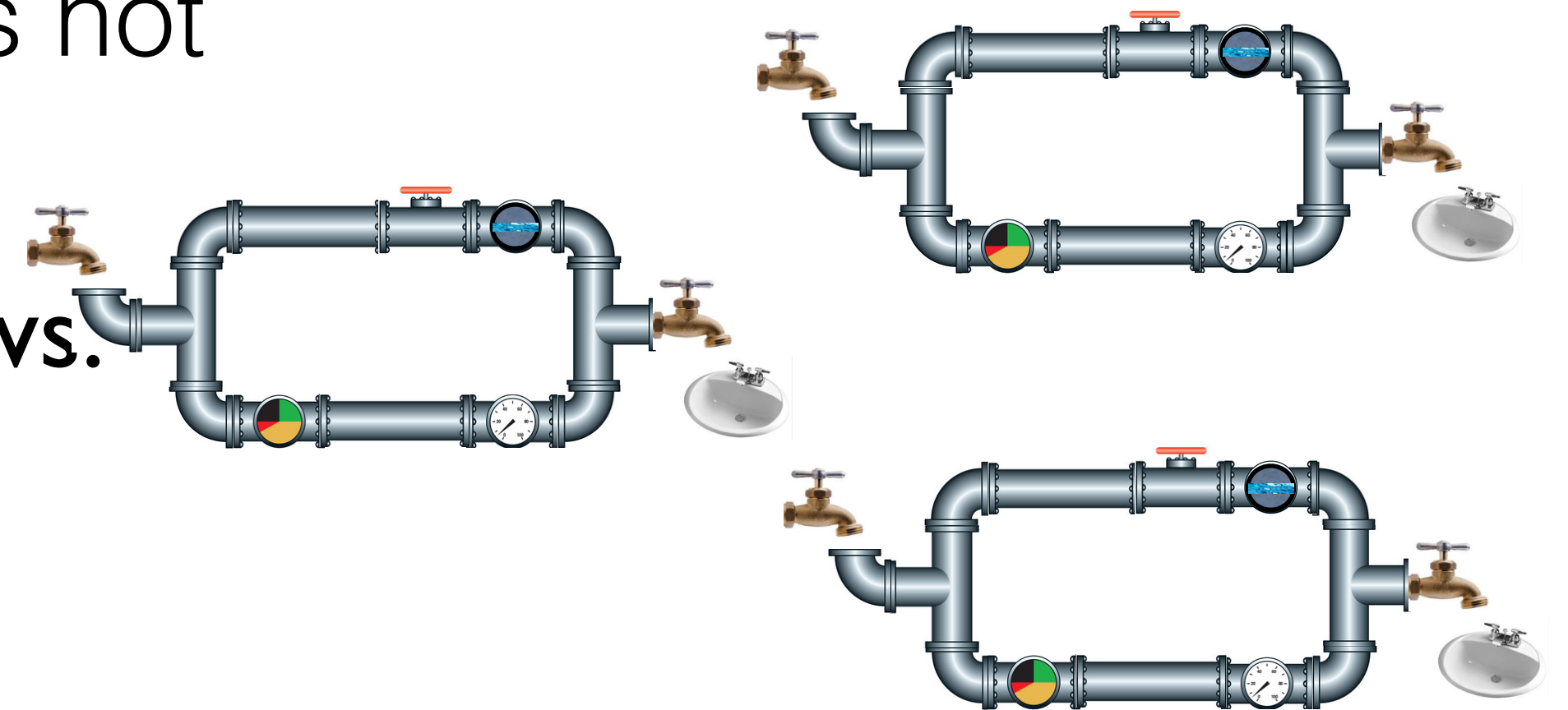
  - File on Amazon S3

CONCURRENT

- Flows consist of pipe assemblies with data sources and sinks

- Flows contain one or more data sources, a DAG (Directed Acyclic Graph) of pipes, and one or more data sinks.

- Flows are designed to be re-useable units of work.

- Flows show the business and programming process.

- A flow is a basic unit of work of arbitrary size.

CONCURRENT

- Cascade joins together multiple flows.

- Use Cascade if there are dependencies among the Flows:

  - Cascade will cause a flow to not be executed until all of its data dependencies are satisfied.

  - A cascade can determine that a Flow does not need to run.

- A CascadeConnector makes a Cascade from Flows.

CONCURRENT

- Java API

- Separates business logic from integration

- Testable at every lifecycle stage
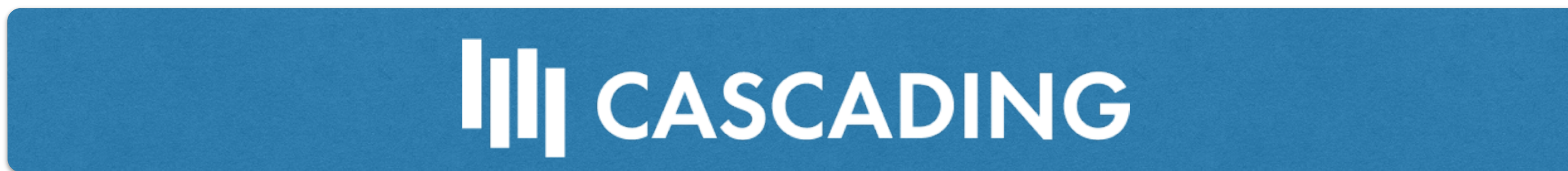
- Works with any JVM language

- Many integration adapters

**Scripting**
Scala, Clojure, JRuby, Jython, Groovy

**Enterprise Java**

**Cascading**

**Scheduler API**

**Processing API**

**Integration API**

*Process Planner*

*Scheduler*

**Apache Hadoop**

**Data Stores**

**CONCURRENT**

**Third-party Systems**

cassandra · splunk> · Amazon **Redshift** · elasticsearch. · **memcached** · **MySQL** · H·BASE · **TERADATA** · HIVE · mongoDB · ORACLE · Neo4j · Apache Solr · **COBOL PROGRAMMING** · JDBC Driver

Source

Sink

**CASCADING**

http://www.cascading.org/extensions/

CONCURRENT

**"Write once and deploy on your fabric of choice."**

- The Innovation — Cascading allows for data apps to execute on existing and emerging fabrics through its new customizable query planner.

- Cascading 3.0 supports — Local In-Memory, Apache MapReduce and Apache Tez.  1H 2015  - Apache Spark and Apache Storm

- Flexibility to meet changing business needs

# THE STANDARD FOR DATA APPLICATION DEVELOPMENT

**CASCADING**

Proven application development framework for building data apps

www.cascading.org

## Application platform that addresses:

**Build data apps that are scale-free**

Design principals ensure best practices at any scale

**Systems Integration**

Hadoop never lives alone. Easily integrate to existing systems

**Application Portability**

Write once, then run on different computation fabrics

**Staffing Bottleneck**

Use existing Java, Scala, SQL, modeling skill sets

**Test-Driven Development**

Efficiently test code and process local files before deploying on a cluster

**Operational Complexity**

Simple - Package up into one jar and hand to operations

**CONCURRENT**

# CASCADING DATA APPLICATIONS

### Enterprise IT
Extract Transform Load
Log File Analysis
Systems Integration
Operations Analysis

### Corporate Apps
HR Analytics
Employee Behavioral Analysis
Customer Support | eCRM
Business Reporting

### Telecom
Data processing of Open Data
Geospatial Indexing
Consumer Mobile Apps
Location based services

### Marketing / Retail
Mobile, Social, Search Analytics
Funnel Analysis
Revenue Attribution
Customer Experiments
Ad Optimization
Retail Recommenders

### Consumer / Entertainment
Music Recommendation
Comparison Shopping
Restaurant Rankings
Real Estate
Rental Listings
Travel Search & Forecast

### Finance
Fraud and Anomaly Detection
Fraud Experiments
Customer Analytics
Insurance Risk Metric

### Health / Biotech
Aggregate Metrics For Govt
Person Biometrics
Veterinary Diagnostics
Next-Gen Genomics
Argonomics
Environmental Maps

CONCURRENT

- Cascading Java API

- Data normalization and cleansing of search and click-through logs for use by analytics tools, Hive analysts

- Easy to operationalize heavy lifting of data in one framework

**THE CLIMATE CORPORATION**

- Cascalog (Clojure)

- Weather pattern modeling to protect growers against loss

- ETL against 20+ datasets daily

- Machine learning to create models

- Purchased by Monsanto for $930M US

CONCURRENT

**TWITTER**

- Scalding (Scala)

- Makes complex analysis of very large data sets simple

- Machine learning, linear algebra to improve

- 30,000 jobs a day — this works @ scale

- Ad quality (matching users and ad effectiveness)

CONCURRENT

*Hadoop ecosystem supports Cascading*

**DRIVEN**

# Visibility from Development to Production



## Development — Building and Testing

- Design & Development
- Debugging
- Tuning

## Production — Monitoring and Tracking

- Maintain Business SLAs
- Balance & Controls
- Application and Data Quality
- Operational Health
- Real-time Insights

## Operational Meta-data

- Automatically Collected
- Business critical meta-data
- Scalable & searchable store
- Programmatically accessible

**CONCURRENT**

# DRIVEN ARCHITECTURE



Accessibility

UI          CLI

Driven Cluster

Driven Server    Driven Server    Driven Server

App Cluster

Web App Server    Web App Server    Web App Server

Storage Cluster

WAR file

High Availability & Scale

Hadoop Cluster

Driven Plugin

Cascading Application (JAR file)

Telemetry (SSL)

*Debug and optimize your Hadoop applications more effectively with Driven*



- Easily comprehend, debug, and tune your data applications

- Get rich insights on your application performance

- Monitor *applications* in real-time

- Compare app performance with historical (previous) iterations

CONCURRENT

# GET OPERATIONAL INSIGHTS WITH DRIVEN

DRIVEN

*Visualize the activity of your applications to help maintain SLAs*



- Quickly breakdown how often applications execute based on their tags, teams, or names

- Immediately identify if any application is monopolizing cluster resources

- Understand the utilization of your cluster with a timeline of all applications running

CONCURRENT

*Segment your applications for greater insights across all your applications*





- Easily keep track of all your applications by segmenting them with user-defined tags

- Segment your applications for trending analysis, cluster analysis, and developing chargeback models

- Quickly breakdown how often applications execute based on their tags, teams, or names

# COLLABORATE WITH TEAMS

*Utilize teams to collaborate and gain visibility over your set of applications*



- Invite others to view and collaborate on a specific application
- Gain visibility to all the apps and their owners associated with each team
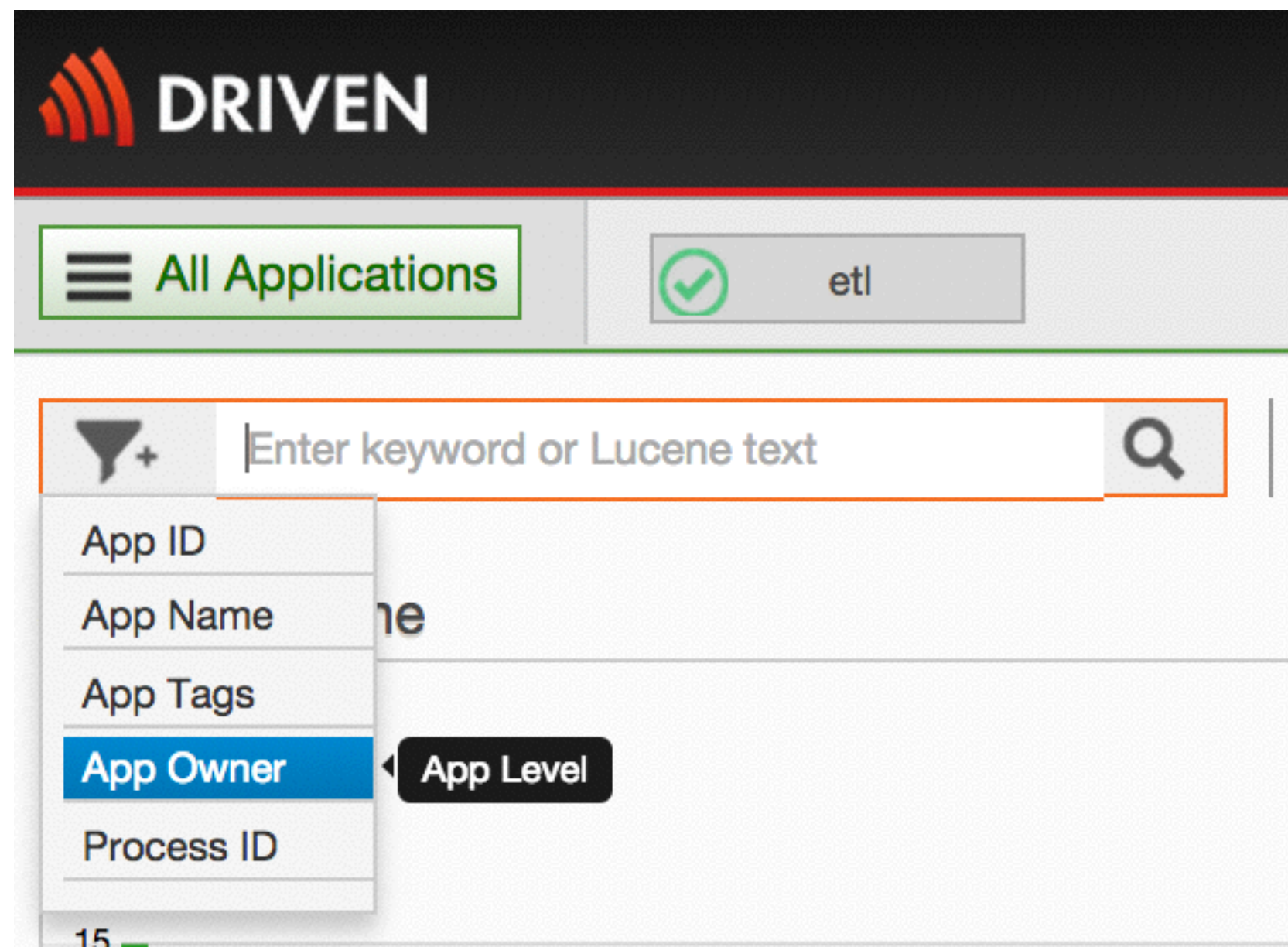- Simply manage your teams and the users assigned to them

*Fast, powerful, rich search capabilities enable you to easily find the exact set of applications that you're looking for*
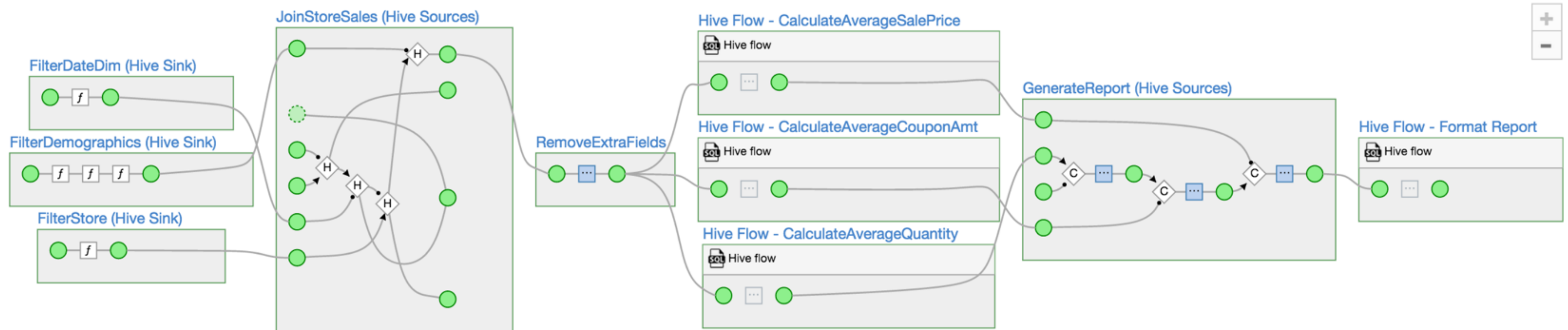


- Identify problematic apps with their owners and teams
- Search for groups of applications segmented by user-defined tags
- Compare specific applications with their previous iterations to ensure that your application can meet its SL
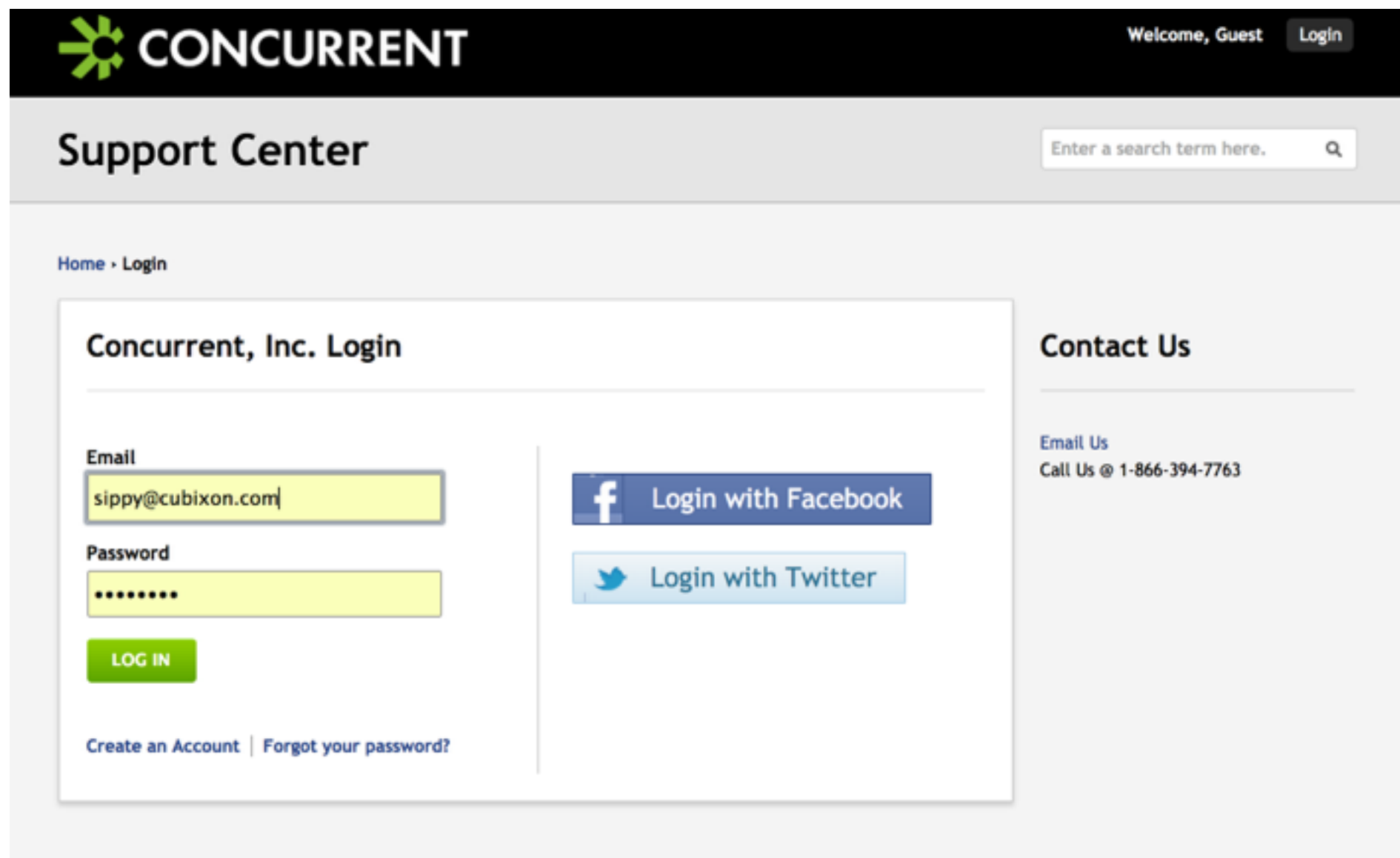
- Understand the anatomy of your Hive app
- Track execution of queries as single business process
- Identify outlier behavior by comparison with historical runs
- Analyze rich operational meta-data
- Correlate Hive app behavior with other events on cluster

- Support for Cascading over email, phone, support portal and web forums that meet your operational SLAs
- Availability of on-site and public training classes for Cascading & Scalding
- Services of experienced technical resources provide custom design solutions
- Presence of thriving community building mission-critical applications for data-driven businesses

DRIVING INNOVATION
THROUGH DATA

# THANK YOU

Supreet Oberoi

**☀ CONCURRENT**