# The Problem

# What we hear from IT …

"ETL is causing our EDW project to fail"

"I wish we could do predictive analytics with our data"

"We failed at our last EDW project"

"We have lots of data today, but we don't know how to do analytics on it"

"We'd like to integrate social media but don't know how"

# What we hear from business users …

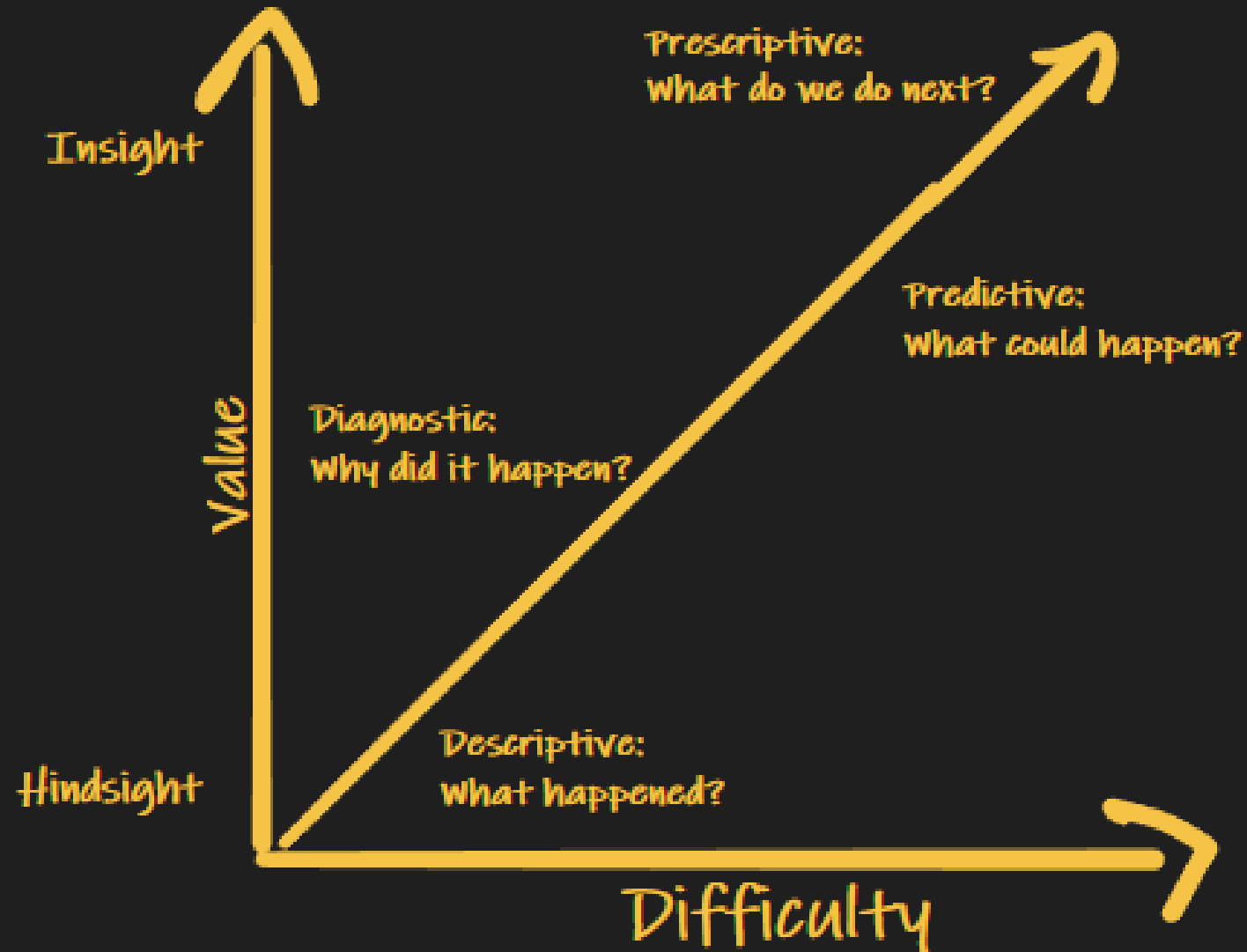"The data warehouse data is nightly, I need real-time data"

"I want to control my own data"

"Even with the data warehouse and reports, I do analysis in spreadmarts"

"I just want to get my job done"
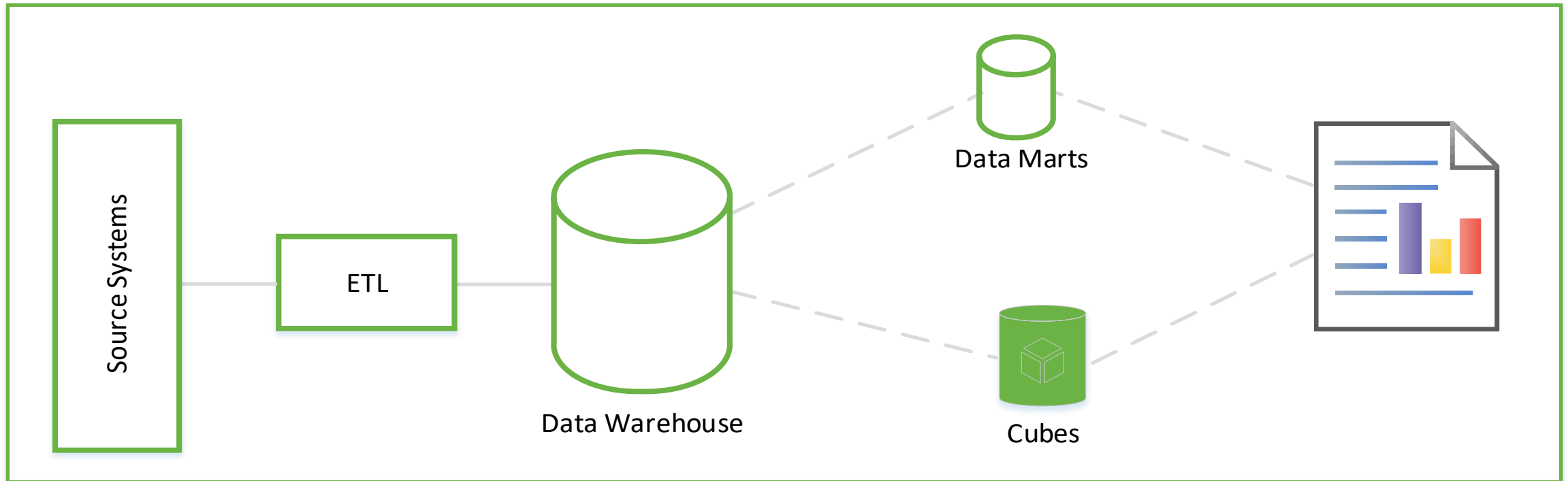
"The data warehouse doesn't answer my questions"

# Past Solutions

# Legacy Thinking

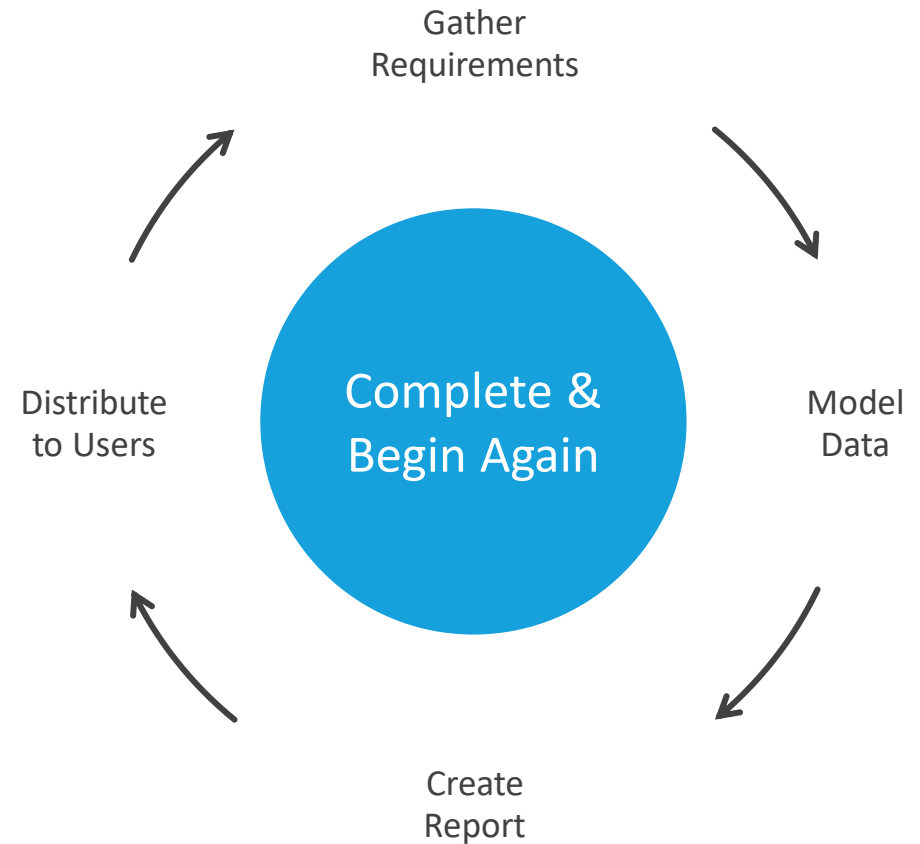**The Philosophy:** Model data » Transform data » Load data » *Understand* data

# Data Projects have a high fail rate

Too much time is spent in:

- Requirements gathering
- Data modeling
- ETL

Users only see the fruits of the endeavor after the reports are created

Gather
Requirements

Model
Data

Create
Report
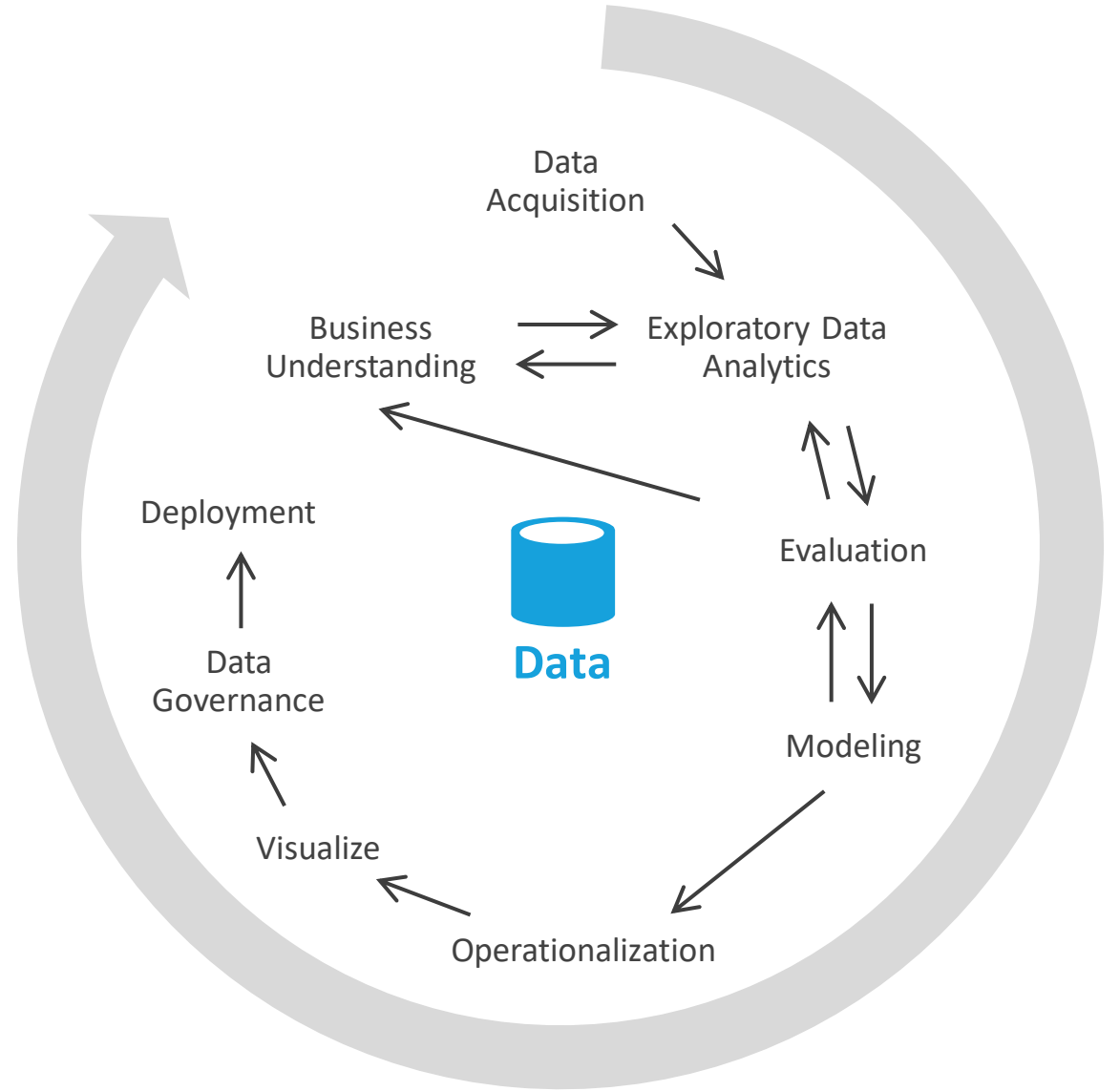
Distribute
to Users

Complete &
Begin Again

# Our Modern Approach

# Data Sandboxing

- A robust and well-proven methodology.

- Data science-like.

- Iterative.

- Stresses up-front understanding of data.

- Modeling is done later in the process (schema-on-read).

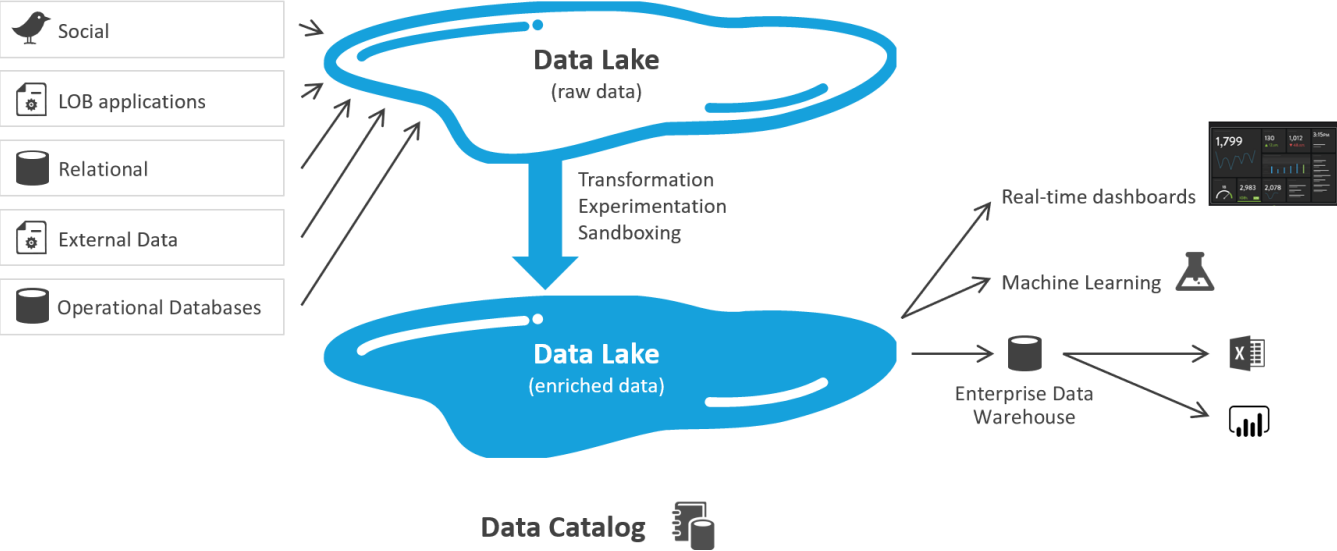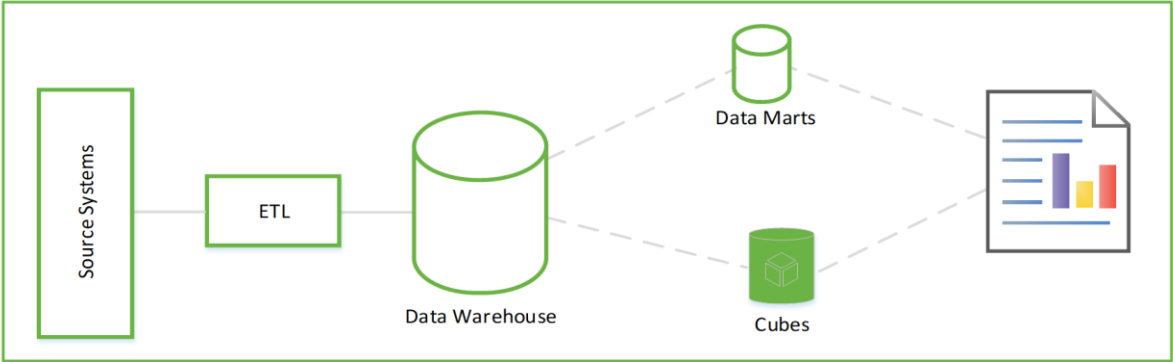- ETL might not be needed

# Self-Service Enabler

**A Data Lake solves 80% of analytical needs**

It is not meant to provide operational reporting

- Data ingestion is more real-time, enabling prediction
- The Data Lake, as a source of all data, is built to efficiently feed a data warehouse.
- Fetches all data, no longer have to go back to source systems for minor changes

# Real-World Example – Customer 360



| Ingest all Data | Store all data | Do analysis | Operationalize |
|---|---|---|---|
| Extract and Load, NO Transform | In native format | Using almost any tool | Create schemas and pipelines |

Social

LOB applications

Relational

External Data

Operational Databases

**Data Lake** (raw data)

Transformation
Experimentation
Sandboxing

**Data Lake** (enriched data)

Real-time dashboards

Machine Learning

Enterprise Data Warehouse

**Data Catalog**

# Data Lake Design–Folder Structures

**A Data Lake Is Just a Folder Structure with Smart Organization and significant processing power**
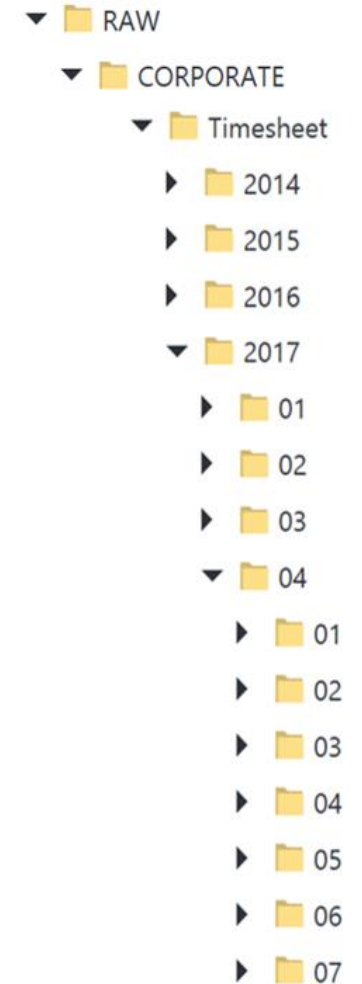
**Authentication/Governance**
- Accounts/Folders/Files

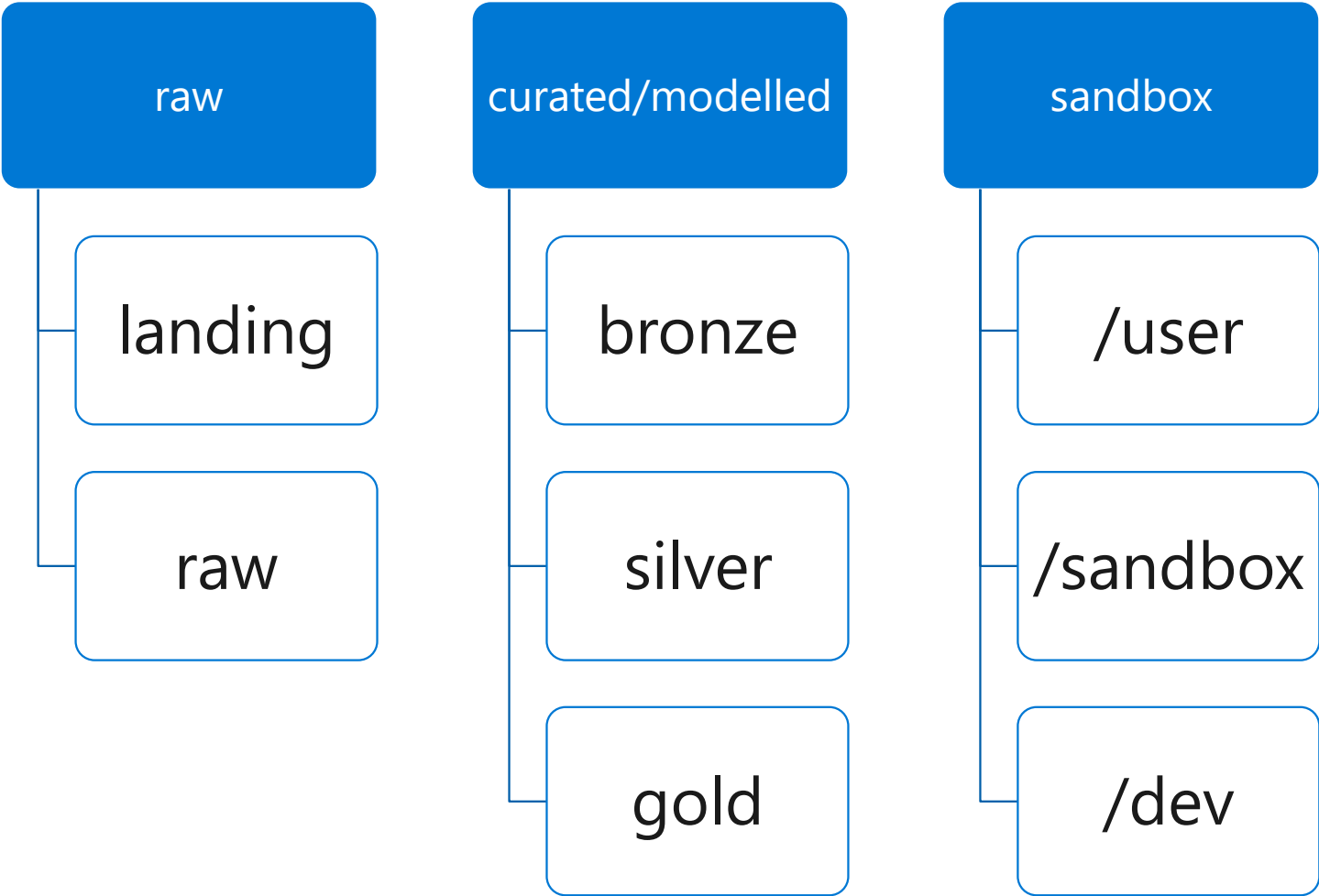**Obvious, Self-documenting Paths**
- dev/raw/{datasource}/{object}/YYYYMMDD/
- dev/reject/{datasource}/{object}/YYYYMMDD/
- prod/snapshot/{datasource}/{object}/YYYYMMDD/
- laboratory/jsmith

**Time partitioning schemes are important**
- AVRO/TPFS format
- .../YYYY/MM/DD/HH/MM

▶ 📁 ENRICHED

▶ 📁 LABORATORY

▶ 📁 RAW

▼ 📁 RAW
  ▼ 📁 CORPORATE
    ▼ 📁 Timesheet
      ▶ 📁 2014
      ▶ 📁 2015
      ▶ 📁 2016
      ▼ 📁 2017
        ▶ 📁 01
        ▶ 📁 02
        ▶ 📁 03
        ▼ 📁 04
          ▶ 📁 01
          ▶ 📁 02
          ▶ 📁 03
          ▶ 📁 04
          ▶ 📁 05
          ▶ 📁 06
          ▶ 📁 07

# Physical Structure

| raw | curated/modelled | sandbox |
|---|---|---|
| landing | bronze | /user |
| raw | silver | /sandbox |
| | gold | /dev |

# Analytics workflow

Trusted

Staging

RAW

Sandbox

① Development / Data Ingestion

② Production

③ Data Scientists / Business Analysts/ Developers

④ Development

⑤ Production

# Data Lake vs. Data Warehouse

| Data Lake | Data Warehouse |
|---|---|
| Complementary to the EDW | Can be sourced from the Data Lake |
| Load first, understand later | Understand first, load later |
| Schema-on-read | Schema-on-write |
| System of Insight | System of Record |
| Detailed Data | Refined Data |
| Supports varied data formats | Structured data |
| Adapts to changing requirements | Difficult to change structure |
| Optimized for Cost | Optimized for Performance |

Home | Discover | Publish | Settings | User

Search 🔍

Results Per Page: 10 ⌄

Highlight ●

Open In ... ⌄ | 🗑 Delete

## Searches

**▼ Current Search**

Save | Clear All

Source Type:
Azure Data Lake Store ⊗

**▼ Filters**

**Tags**
☐ earthquakes (4)
☐ datalaketest (3)
☐ Demo (2)
☐ datalaketest3 (2)

**Object Type**
☐ Directory (3)
☐ File (3)
☐ Data Lake (2)

**Source Type**
☐ SQL Server (4423)
☐ Filesystem (1551)
☐ SQL Server Analysis Services Tabular (396)
☐ SQL Server Analysis Services (140)

see more

### Daily Earthquake Records... ☑

Directory represents logical data set. Contains one file per day.

Experts:
amitkul@microsoft.com

earthquakes

Contained In Data Lake:
adlsadc.azuredatalakes...

📁 **DATA LAKE DIRECTORY**

Open In ... ⌄ | Explore Data Lake ❯

### region_names.asc

click tile to add a description...

Experts:
amitkul@microsoft.com

earthquakes

Contained In Data Lake:
adlsadc.azuredatalakes...

📄 **DATA LAKE FILE**

Open In ... ⌄ | Explore Data Lake ❯

### AzureDataCatalogErrorEve...

click tile to add a description...

Experts:
linm@test.com

Demo    datalaketest

### test1

click tile to add a description...

Experts:

datalaketest3

## Properties

| Preview | Columns | Docs |

**Name:**
DailyEarthquakes

**Friendly Name:**

Daily Earthquake Records from USGS

**Description:**

Directory represents logical data set. Contains one file per day.

**Experts:**

amitkul@microsoft.com

Add...

**Tags:**

earthquakes

Preview    Columns    **Documentation**

We have a winning idea, help us operationalize it

# Data Ingestion Factory
## (Kappa Architecture)

# Handling Streaming Data

**Twitter**
Tweets

**WebJob**
Twitter Feed
Consumer

**Event Hubs**
Streaming
Data

**Stream
Analytics**
Process Data

**ADLS Gen 2**
Storage Account
Data Lake

**Synapse SQL**
(Provisioned)

**Power BI**
Dashboards
& Reports

**IoT Sensors**
In-store Traffic
Data

**IoT Hub**
Streaming
Data