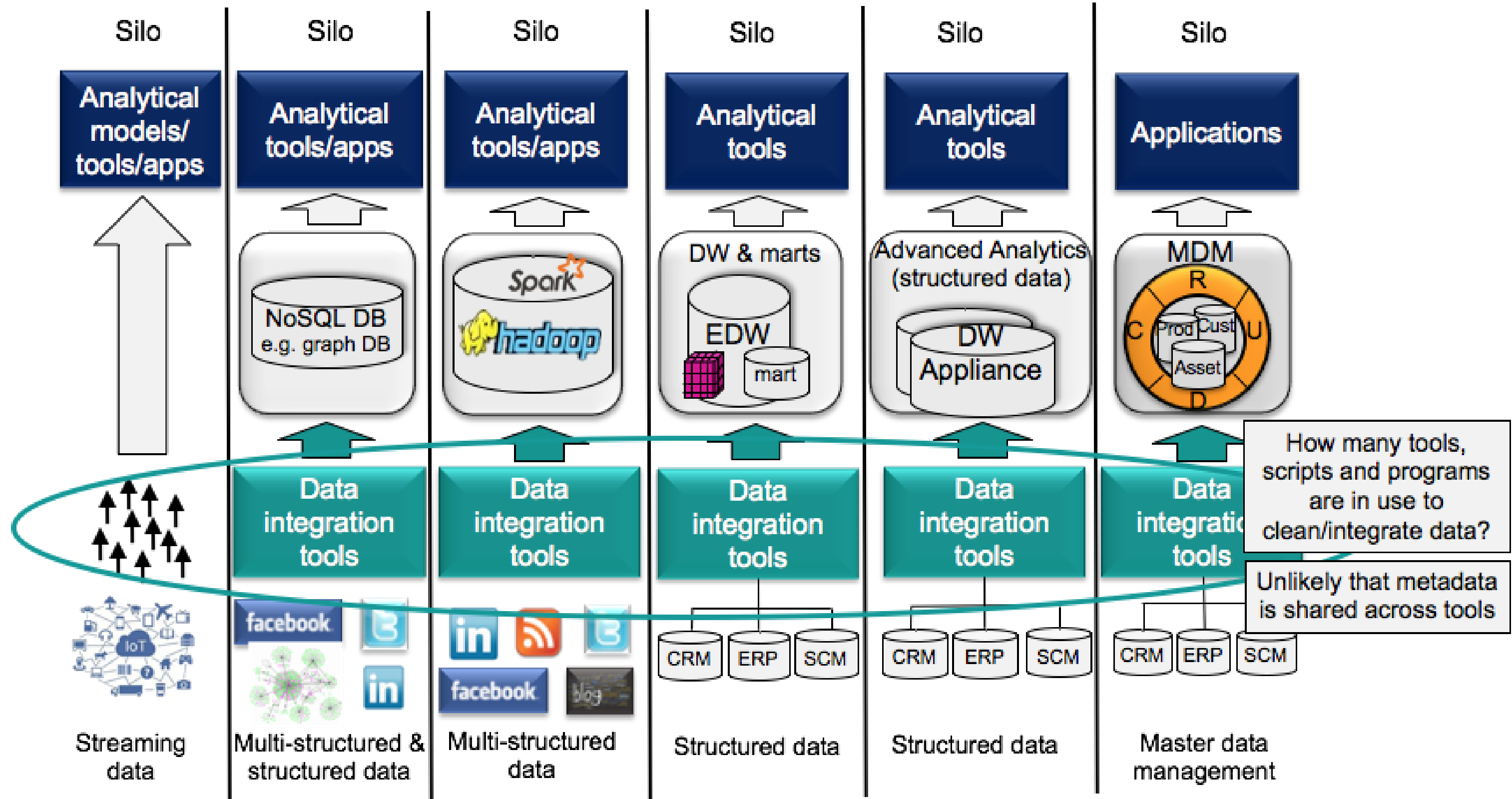# Using a Data Lake to Accelerating ETL Development

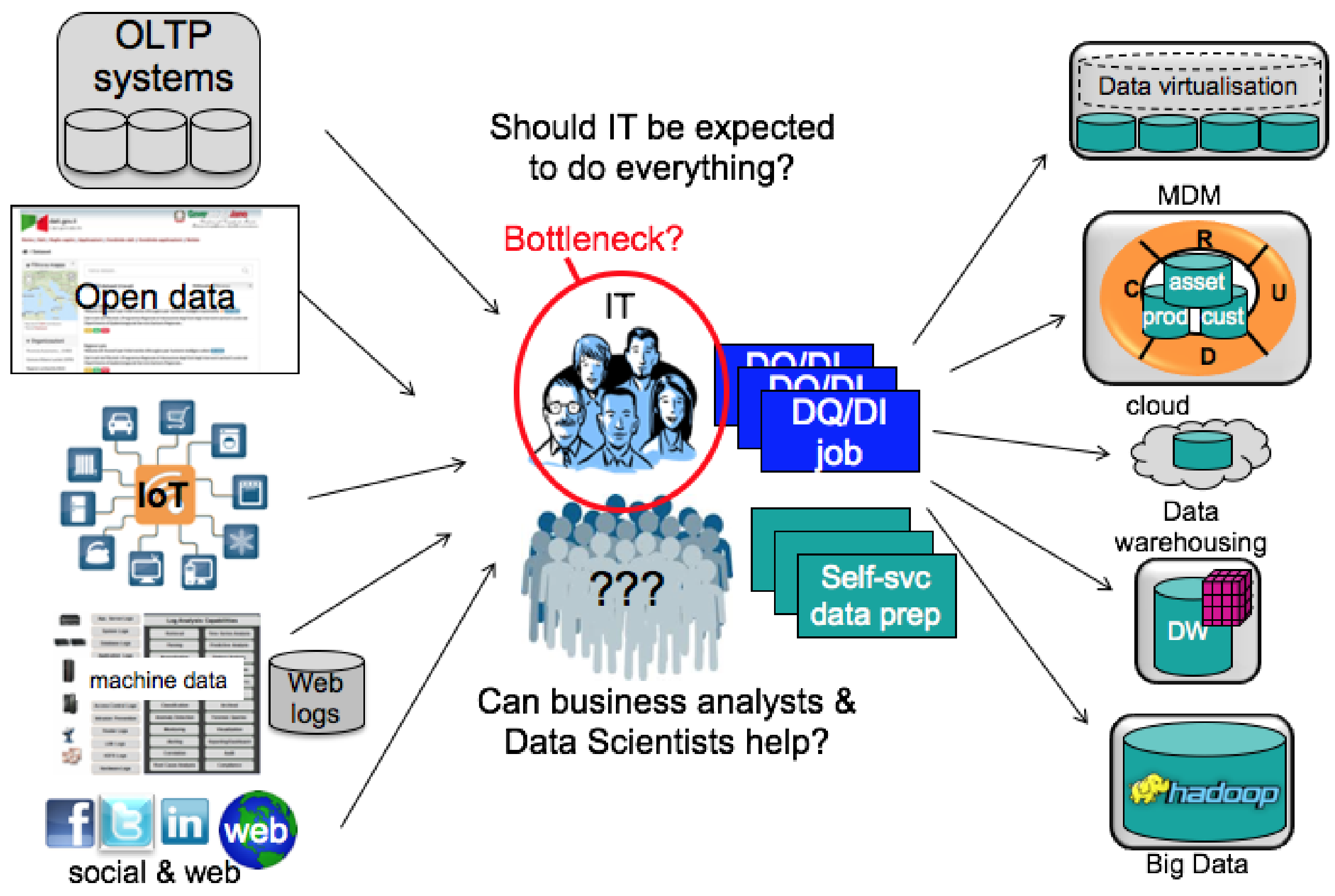Part 2 of the Modern Data Series

Mike Ferguson

quest for knowledge®

q4k.com

Many companies today have a siloed approach to data and analytics, with many tools, scripts and code in use to clean, transform and integrate data

—

But with 000's of data sources that business want to analyse, IT will likely become a bottleneck unless they can work with business to integrate data
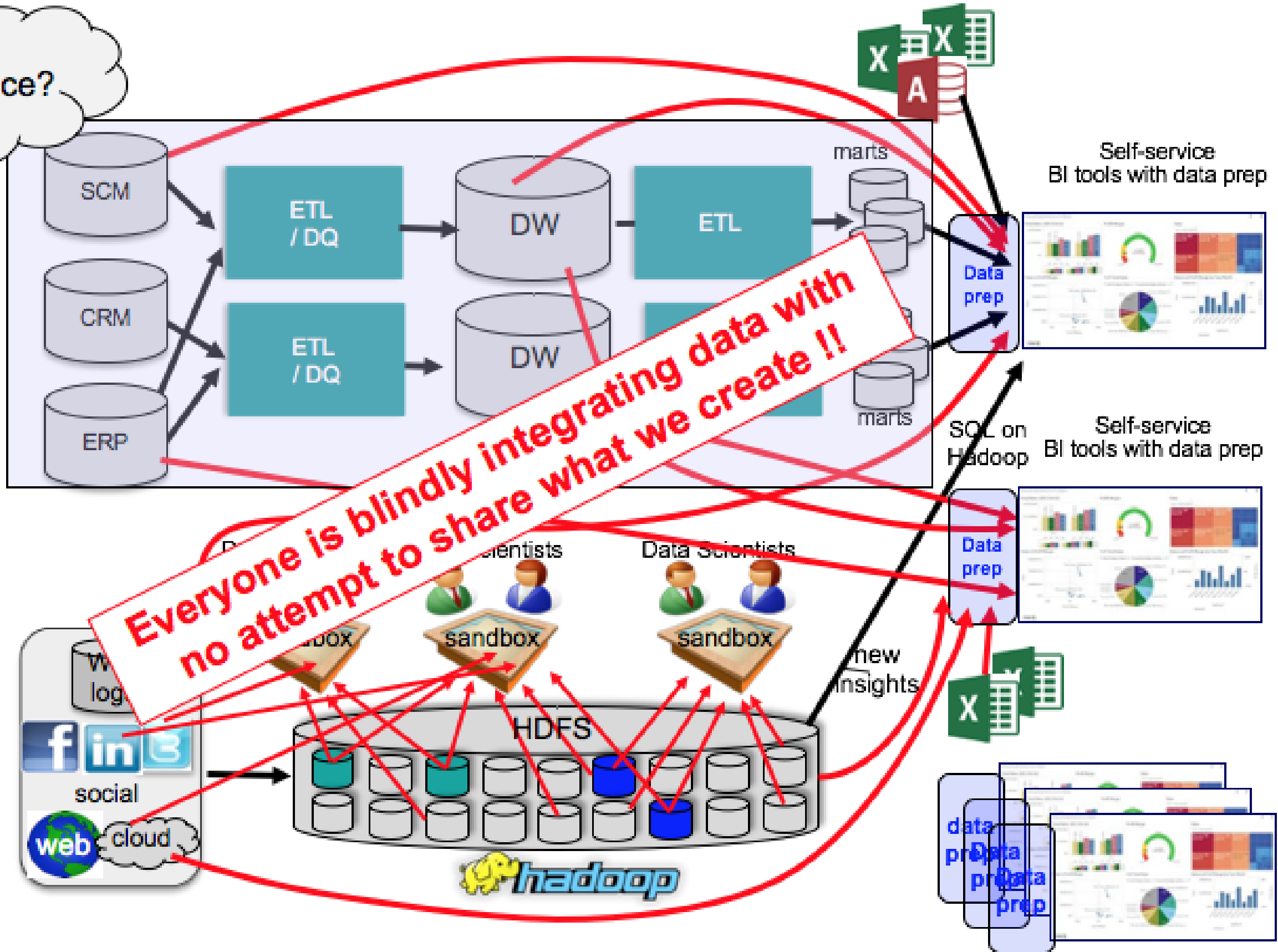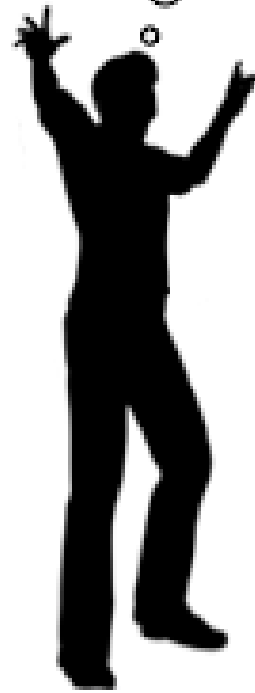
—

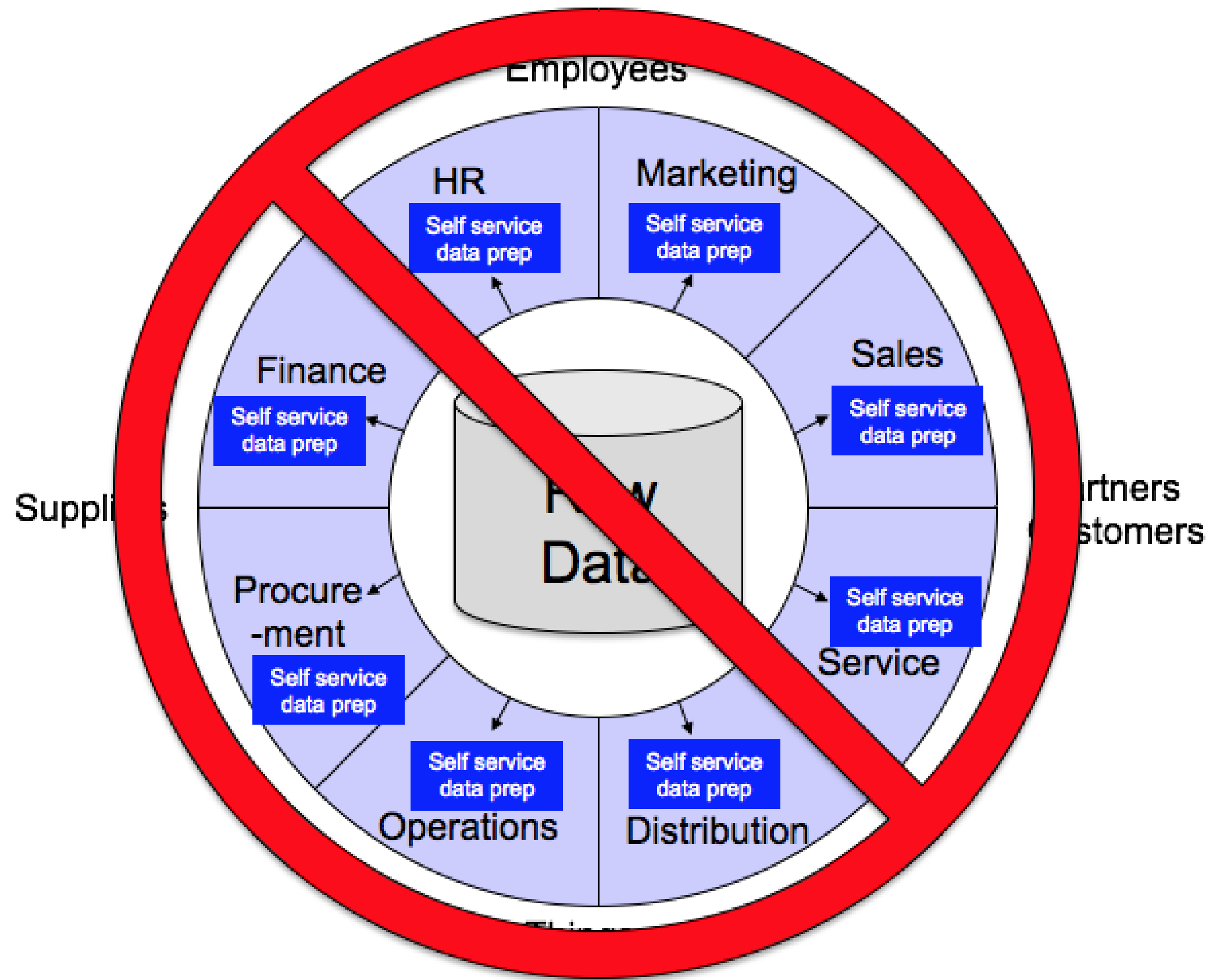# The danger of self-service data preparation
## An explosions of personal silos!

—

# Challenges – How do you govern self-service data preparation to avoid chaos in the enterprise when departments are buying their own tools?

—

INTELLIGENT
BUSINESS
STRATEGIES
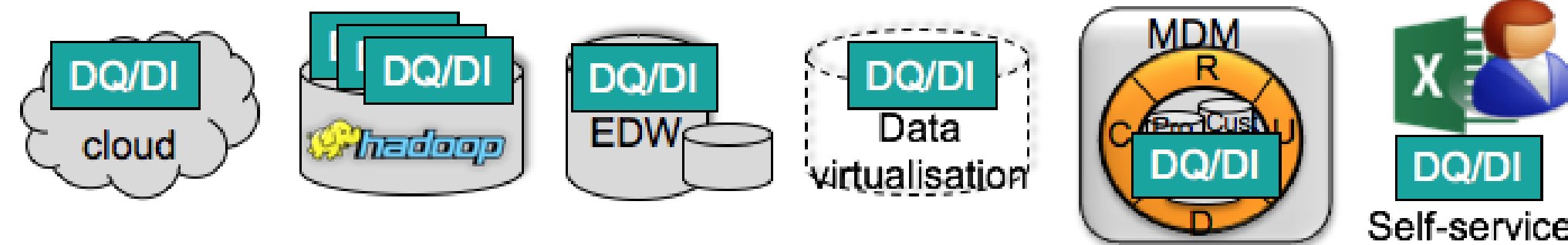
quest for
knowledge
q4k.com

# We do **NOT** want self-service chaos with many people creating inconsistent data

—

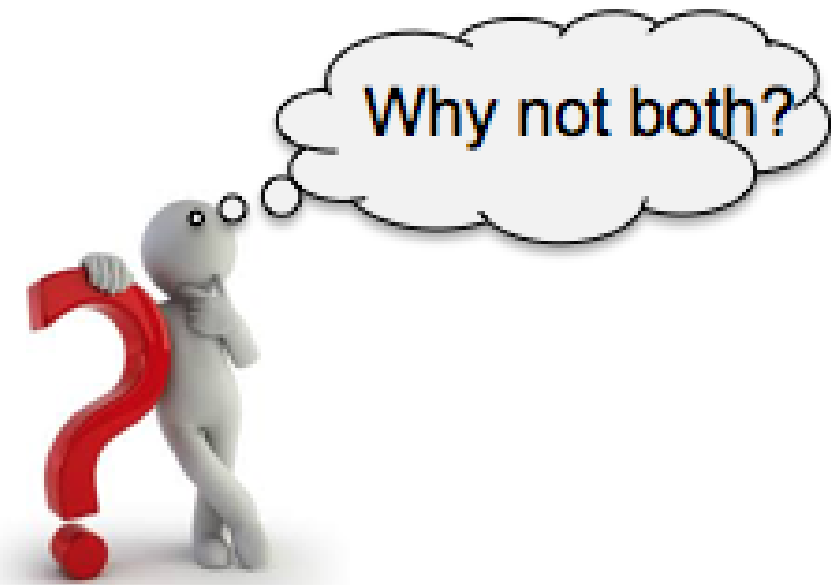# Problems with the current approach of many different tools and scripts being used by IT and business to integrate data

—

INTELLIGENT BUSINESS STRATEGIES

quest for knowledge
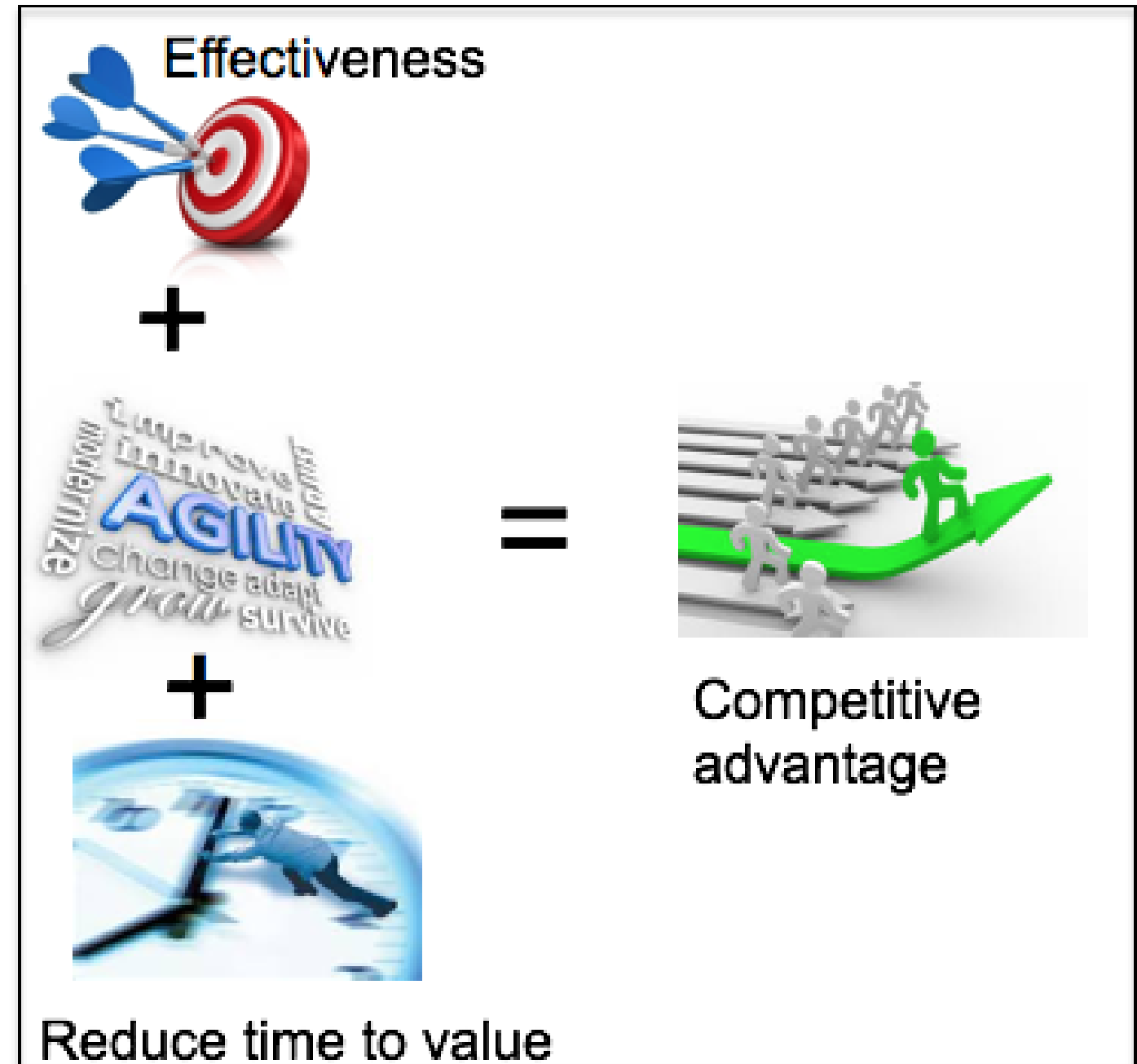q4k.com

DQ/DI = Data quality / data integration

- Cost of data integration is much too high
- Slow speed of development
- Project oriented siloed approach to DQ/DI with limited collaboration
- Multiple DQ/DI technologies and techniques in use that are not integrated
- Multiple skill sets fractured across different tools and projects
- Lots of re-invention rather than re-use
- Fractured metadata across multiple tools or no metadata at all
- Risk of duplicate inconsistent DQ/DI rules for same data
- Metadata lineage is unavailable in many places
- Repetition of mistakes

We need a way to accelerate ETL processing to be more agile, reduce time to value and improve effectiveness while also governing data

—

# How can you accelerate ETL processing?

**Build once, re-use everywhere**
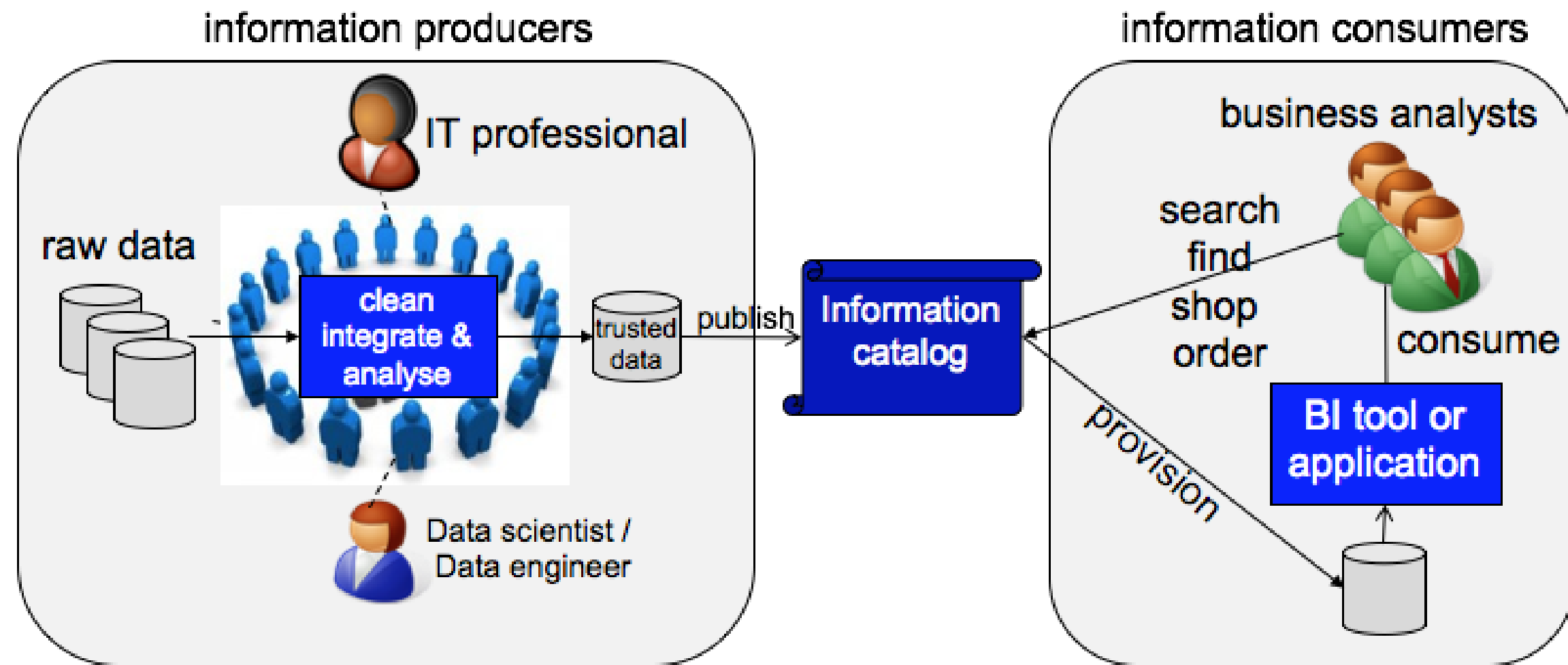
—

INTELLIGENT
BUSINESS
STRATEGIES

quest for
knowledge

q4k.com

# We want trusted data assets and analytics built once, reused everywhere in a data driven-enterprise

—

INTELLIGENT BUSINESS STRATEGIES

The Intelligent Business

# Organize to become data driven – information producers and information consumers

_

Information producers and information consumers need to make use of

- A business glossary and information catalog
- Role-based data management tools aimed at IT AND business
- A collaborative approach by business and IT to produce data and analytical assets
- A catalog to quickly find reusable trusted assets to drive business value

We need to create trusted, business ready data for users to easily find, consume and use to drive value **faster**

_

INTELLIGENT BUSINESS STRATEGIES

quest for knowledge
q4k.com

# Data available as a Service



**BUSINESS READY**
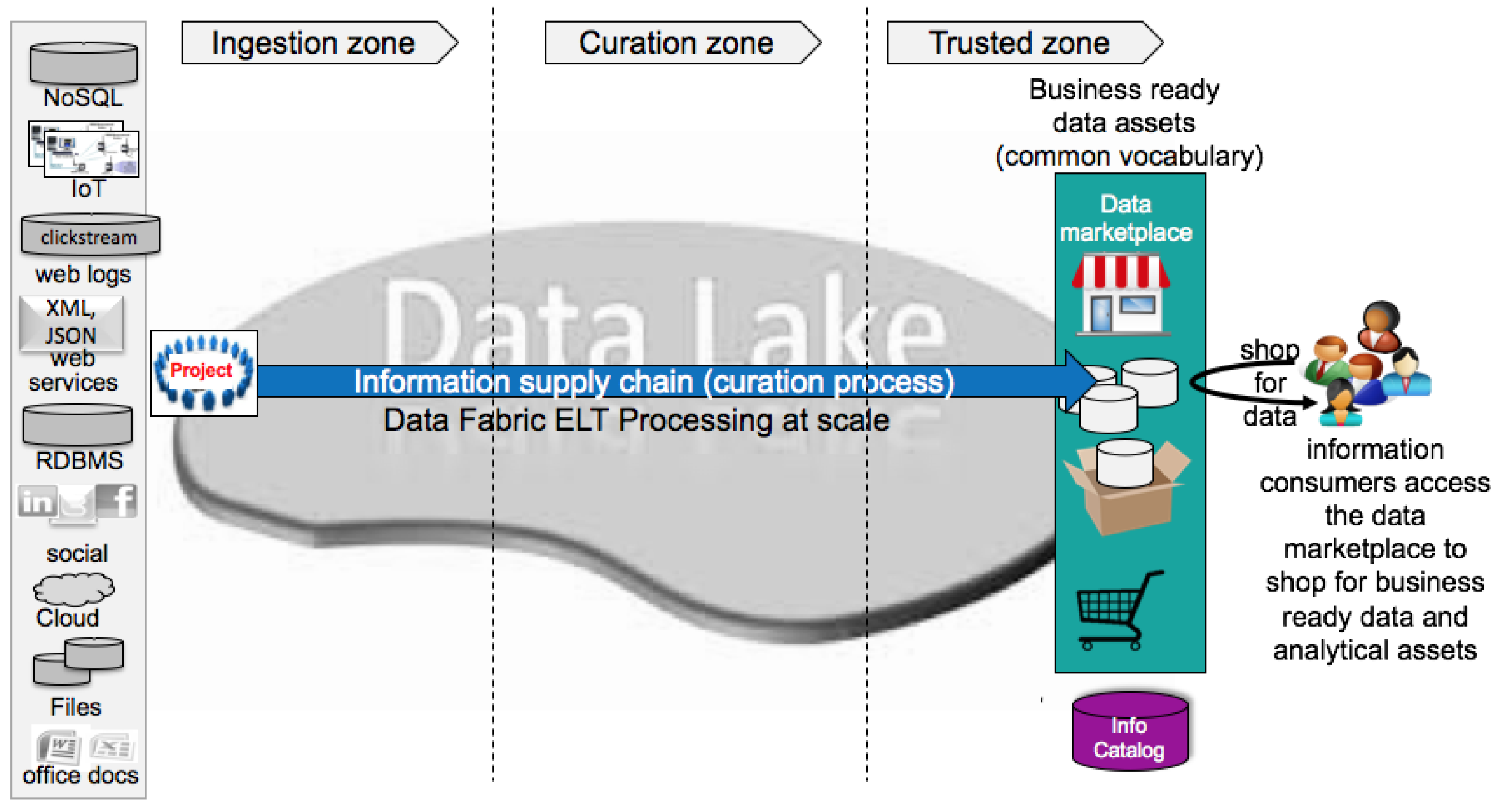
Business ready data products can be logical entities

**Master Data**
- Customers
- Products
- Suppliers
- Assets
- Employees
- Materials

**Transaction Data**
- Orders
- Shipments
- Payments
- Adjustments
- Returns

Create a **data lake** to organize data so 'business ready' data and analytical assets can be produced and published in a marketplace for users to consume
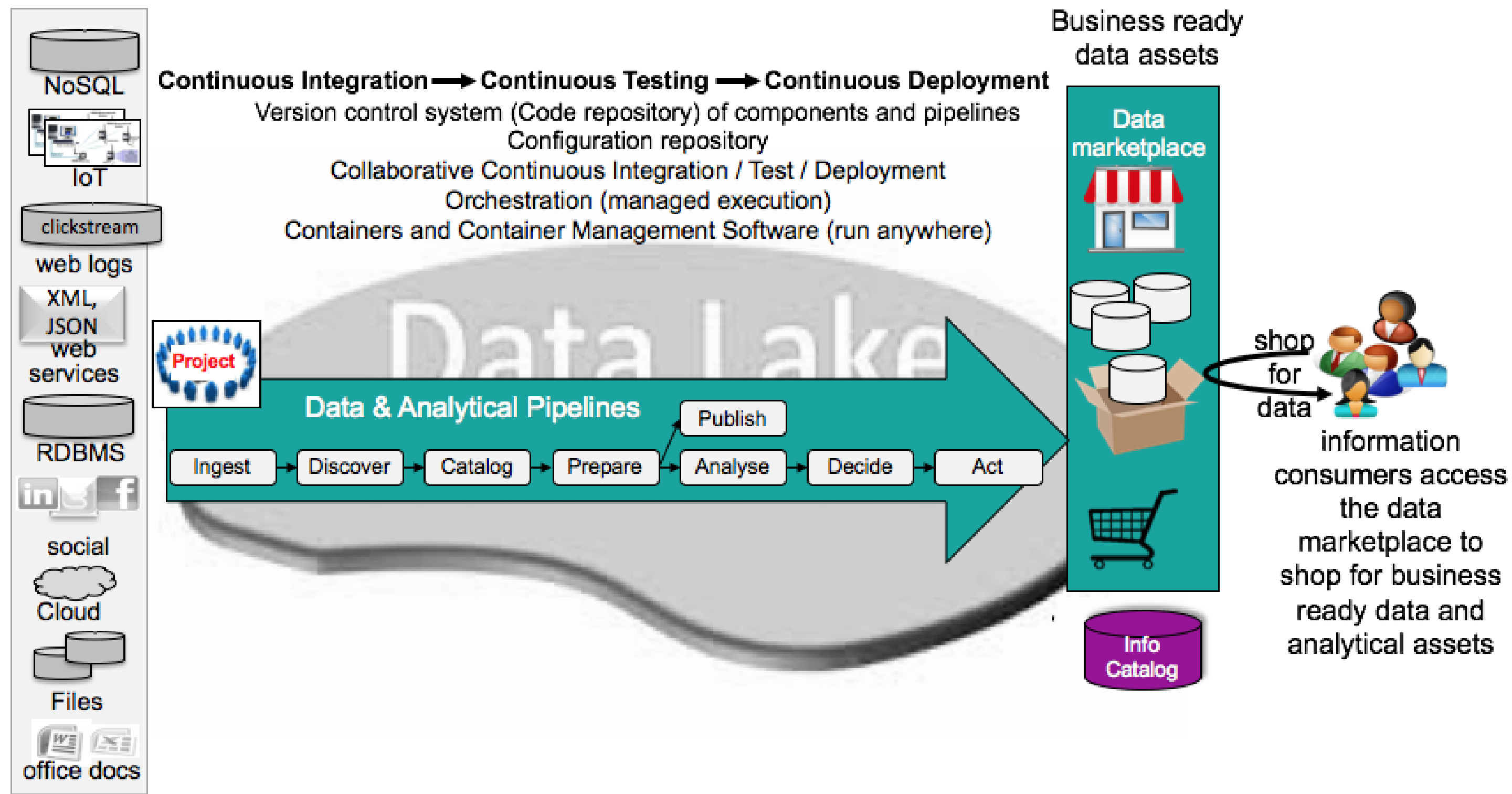
—

# Use trusted data assets to consistently build MDM, DW, data sets for data science and virtual data services

—

**Dataops** – The continuous integration / continuous deployment (CI/CD) of component based pipelines to produce 'business ready' data & analytical assets –

INTELLIGENT BUSINESS STRATEGIES

quest for knowledge

q4k.com

# Upcoming Courses

Big Data Architecture and Technology for Analytics

Cloud Data Warehouse Migration

Data Warehouse ETL: The Kimball Approach

Data Warehouse Lifecycle: The Kimball Approach

Data Warehouse Modernization

Designing, Operating and Managing an Enterprise Data Lake

Hands-on Data Science for BI Professionals and Data Analysts

Dimensional Modeling: The Kimball Approach

Enterprise Data Governance & Master Data Management

**STAY TUNED FOR PART 3**

# Modern Data Pipelines

quest for
knowledge

q4k.com