

DATA LAKE

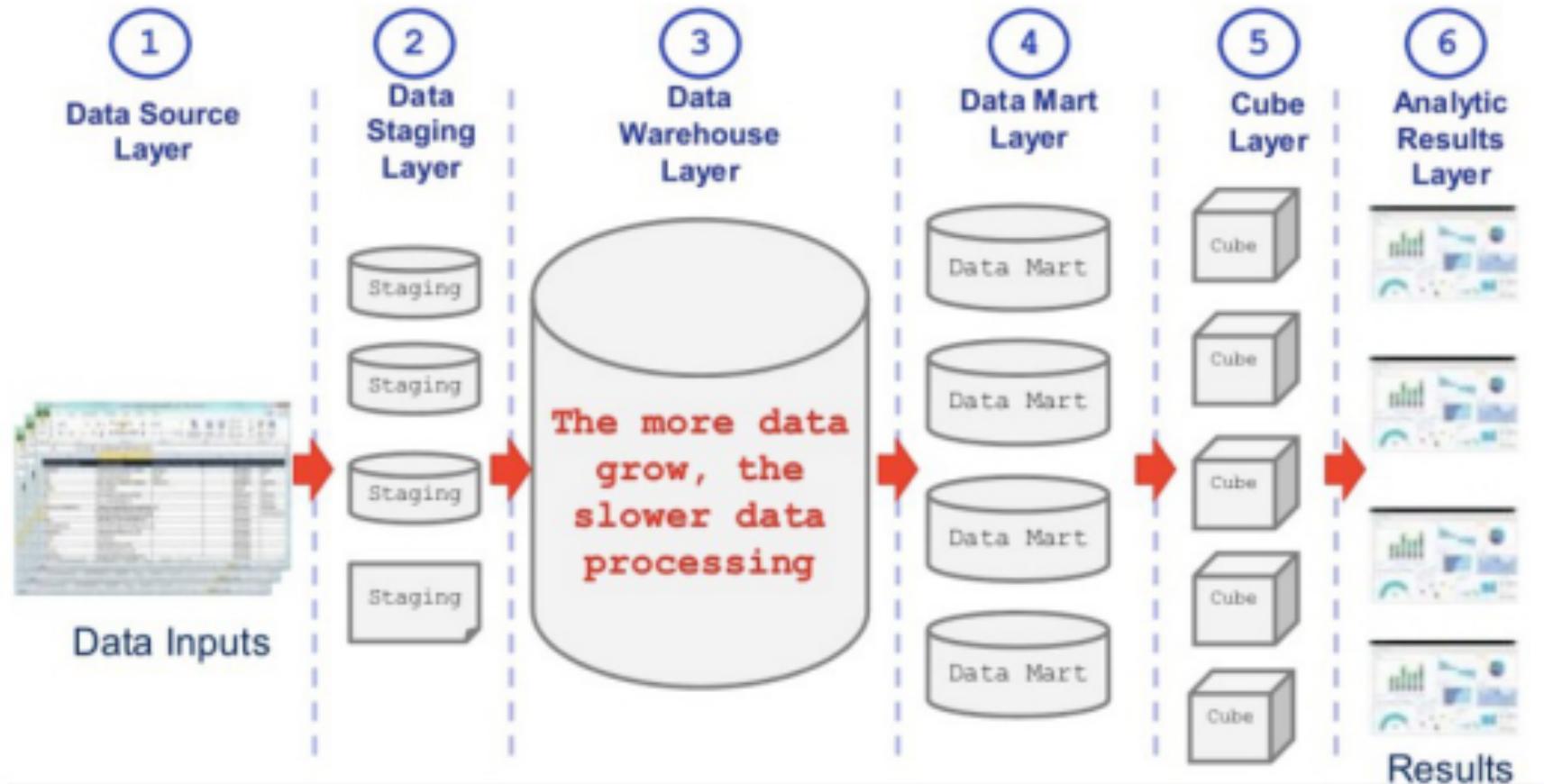
DATA LAKE



Data Warehouse



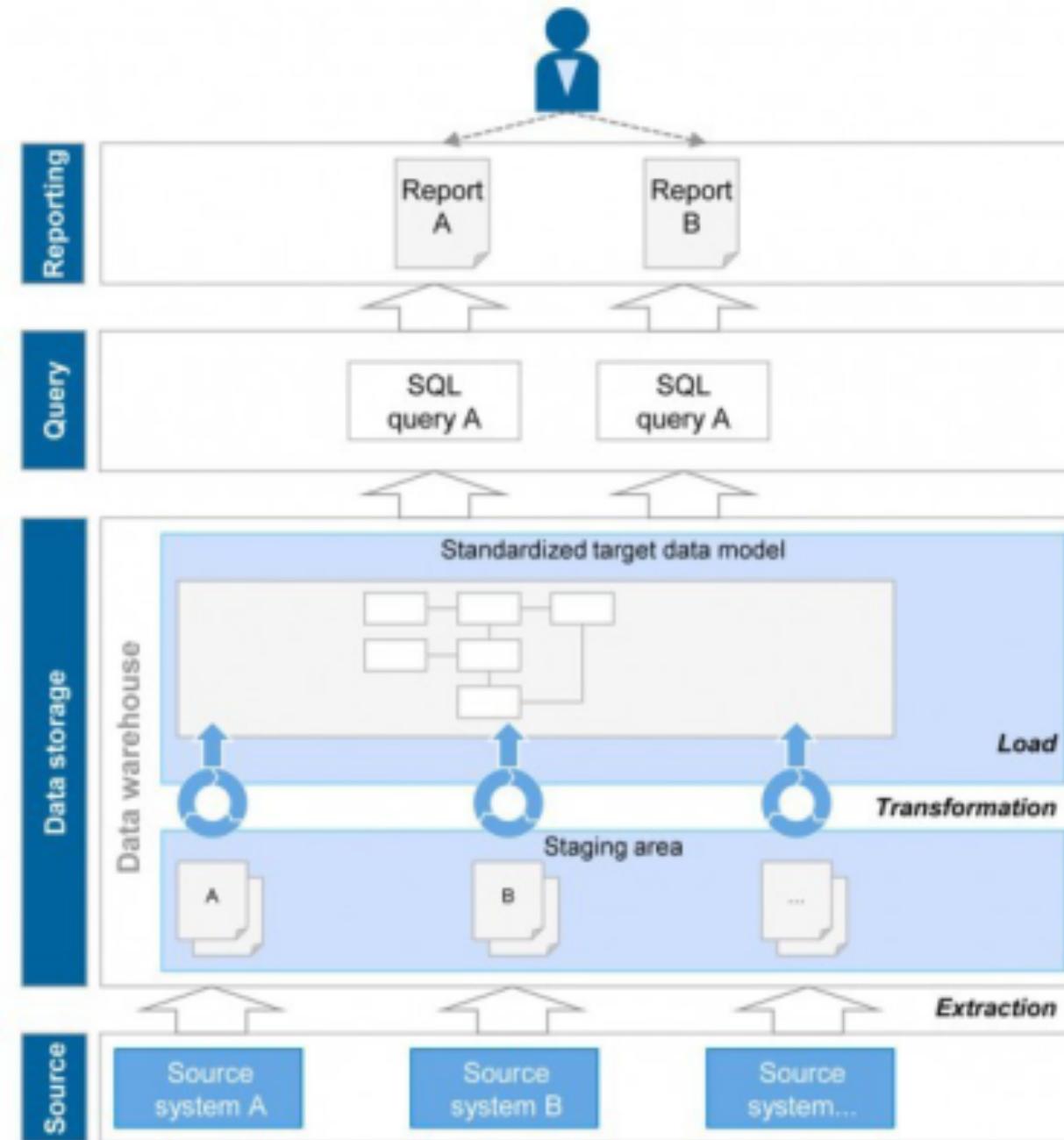
FLAGSHIP FOR LIFE



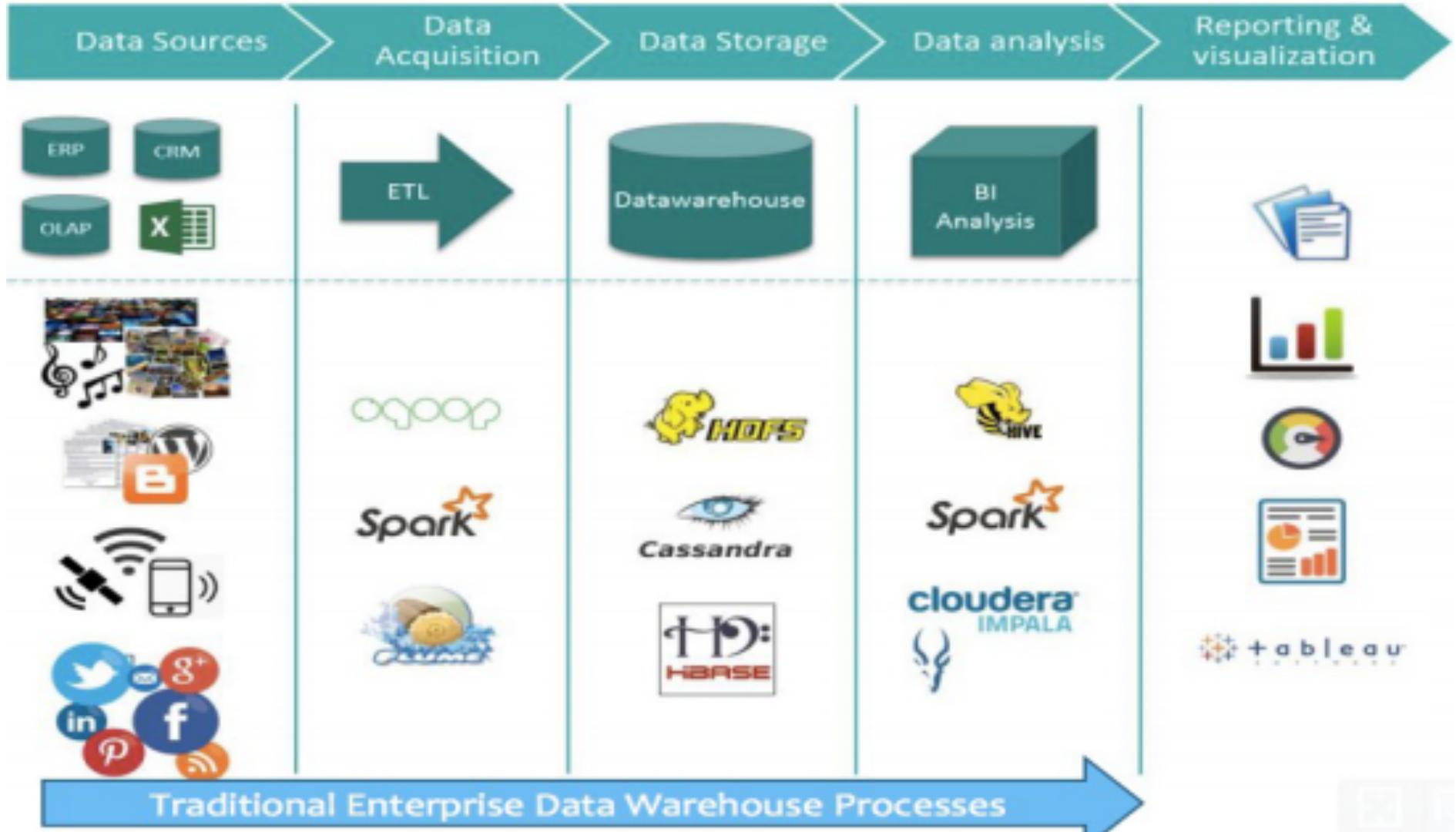
Top Concerns from Traditional Data Warehouse Architecture

1. A lot of data duplication lead to cost of data store/storage issue
2. Very slow of data processing and need to restart/roll back the job if any failed
3. Data security issue due to keep data too many copies and various formats

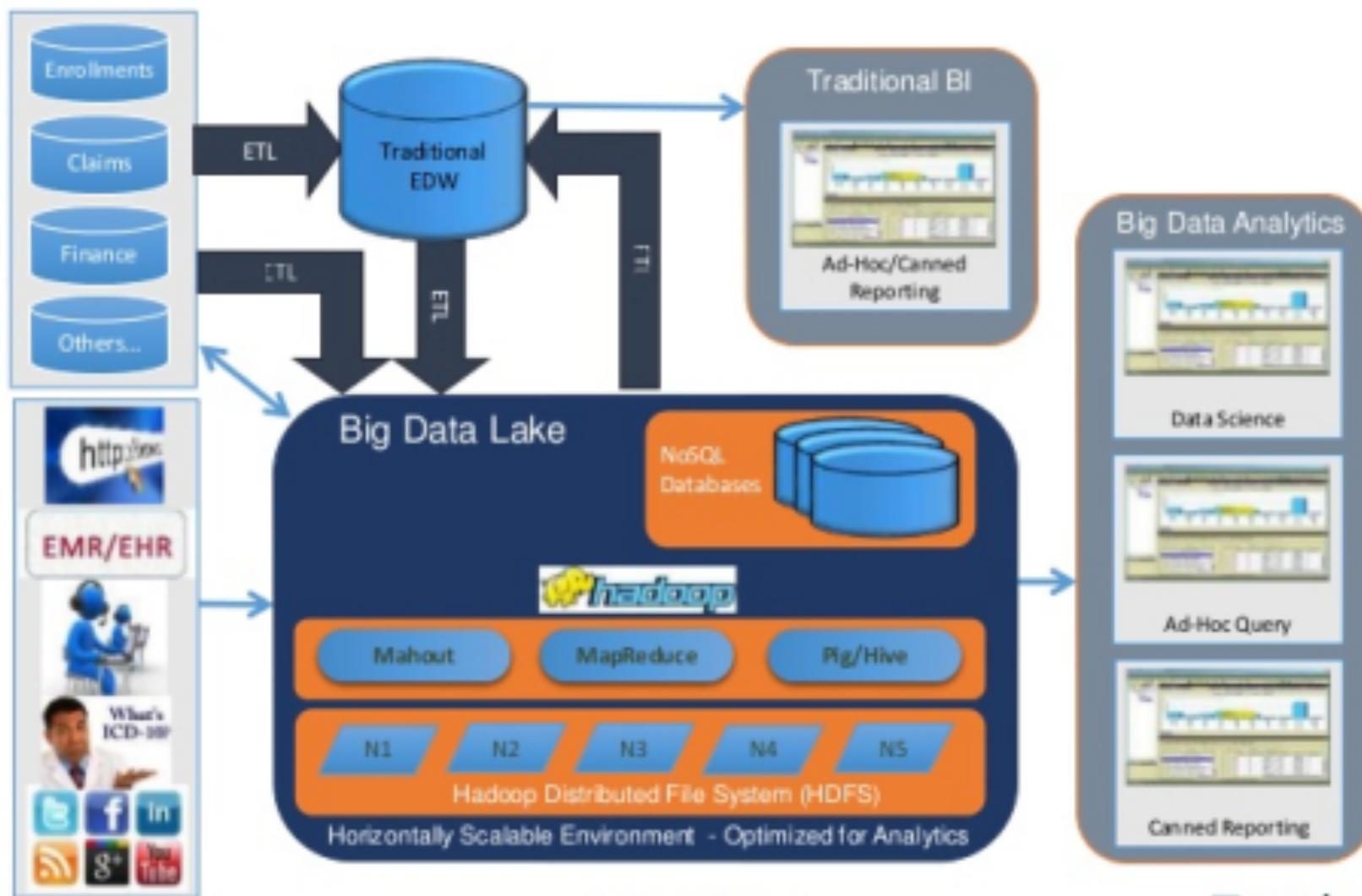
Data Warehouse : ETL



How Data Lake Works?



Today's business environment requires Big Data



#edwdc15

@joe_Caserta

 caserta

maveewat Khanan

2603615 Organization Data Management

Data Lake isn't just a technology
It is an architecture

Data Lake: Definition

Is a huge repository

Is not Hadoop

is not a database, it uses various NoSQL and In-Memory databases.

Stores large volumes of both unstructured and structured data.

Is a data scientist's favorite hunting ground.

Data Lake cannot be implemented in isolation.

Data Lake: Key Benefits

- Scale as much as you can**
- Plug in disparate data sources**
- Acquire high-velocity data: Store in native format**
- Don't worry about schema**
- Unleash your favorite SQL**
- Advanced algorithms**
- Administrative resources**

Data Lake v.s. Data Warehouse

Complementary to EDW (not replacement)	Data lake can be source for EDW
Schema on read (no predefined schemas)	Schema on write (predefined schemas)
Structured/semi-structured/Unstructured data	Structured data only
Fast ingestion of new data/content	Time consuming to introduce new content
Data Science + Prediction/Advanced Analytics + BI use cases	BI use cases only (no prediction/advanced analytics)
Data at low level of detail/granularity	Data at summary/aggregated level of detail
Loosely defined SLAs	Tight SLAs (production schedules)
Flexibility in tools (open source/tools for advanced analytics)	Limited flexibility in tools (SQL only)

Data Lake Risks

More Data Sources

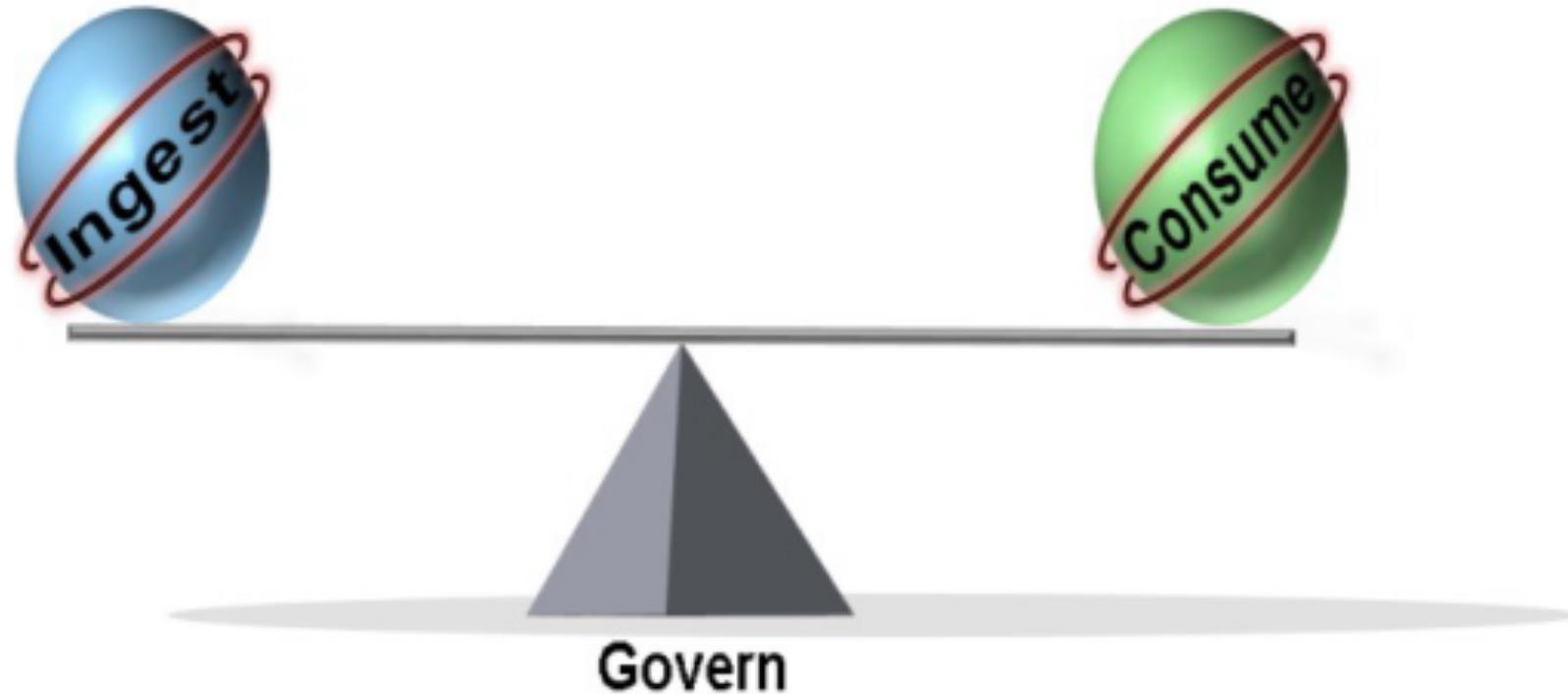
More Applications

More Business Units

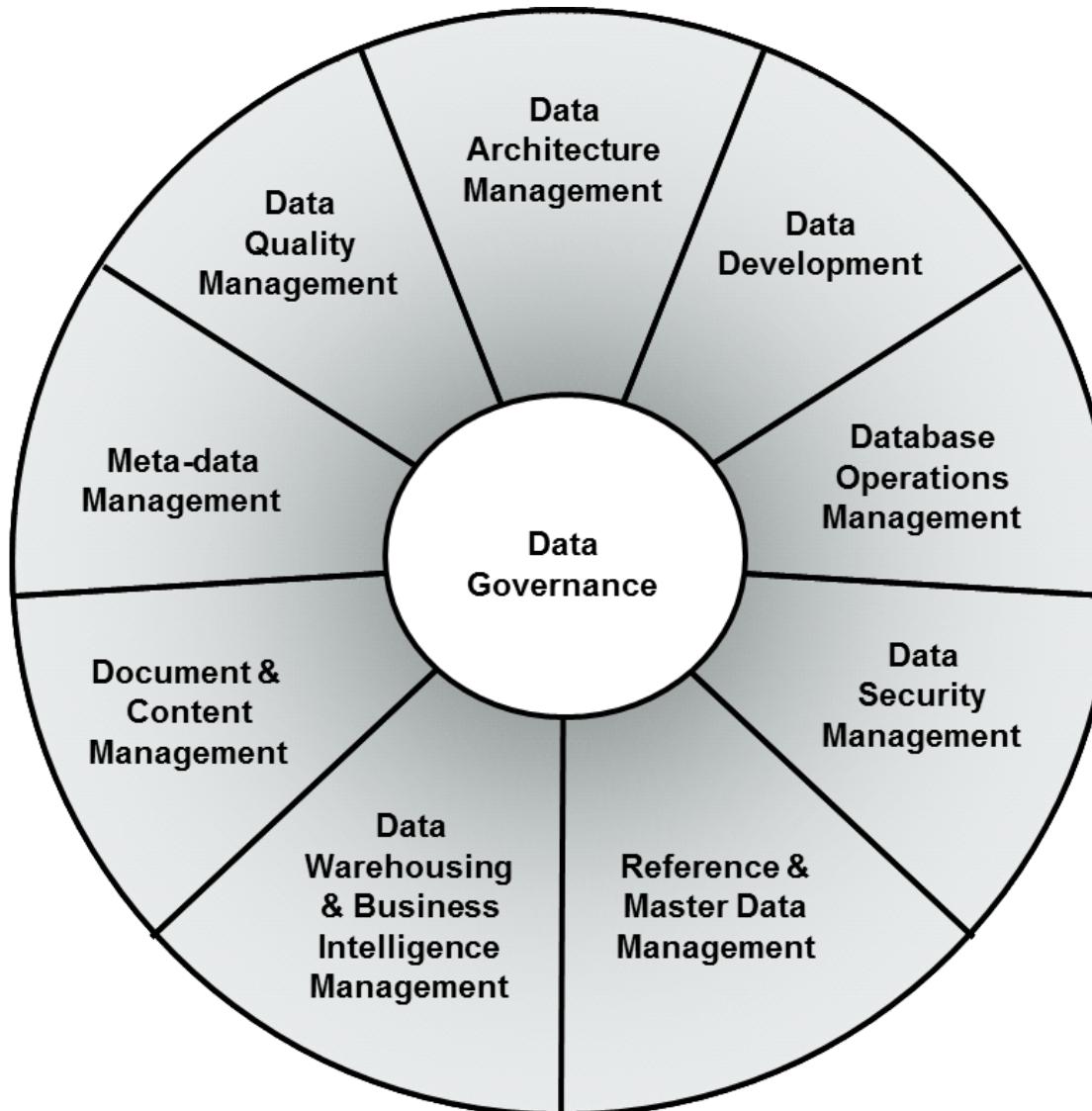
More Users

Without proper governance mechanisms
Data lakes risk turning data swamps

Data Lake Governance



Data Governance

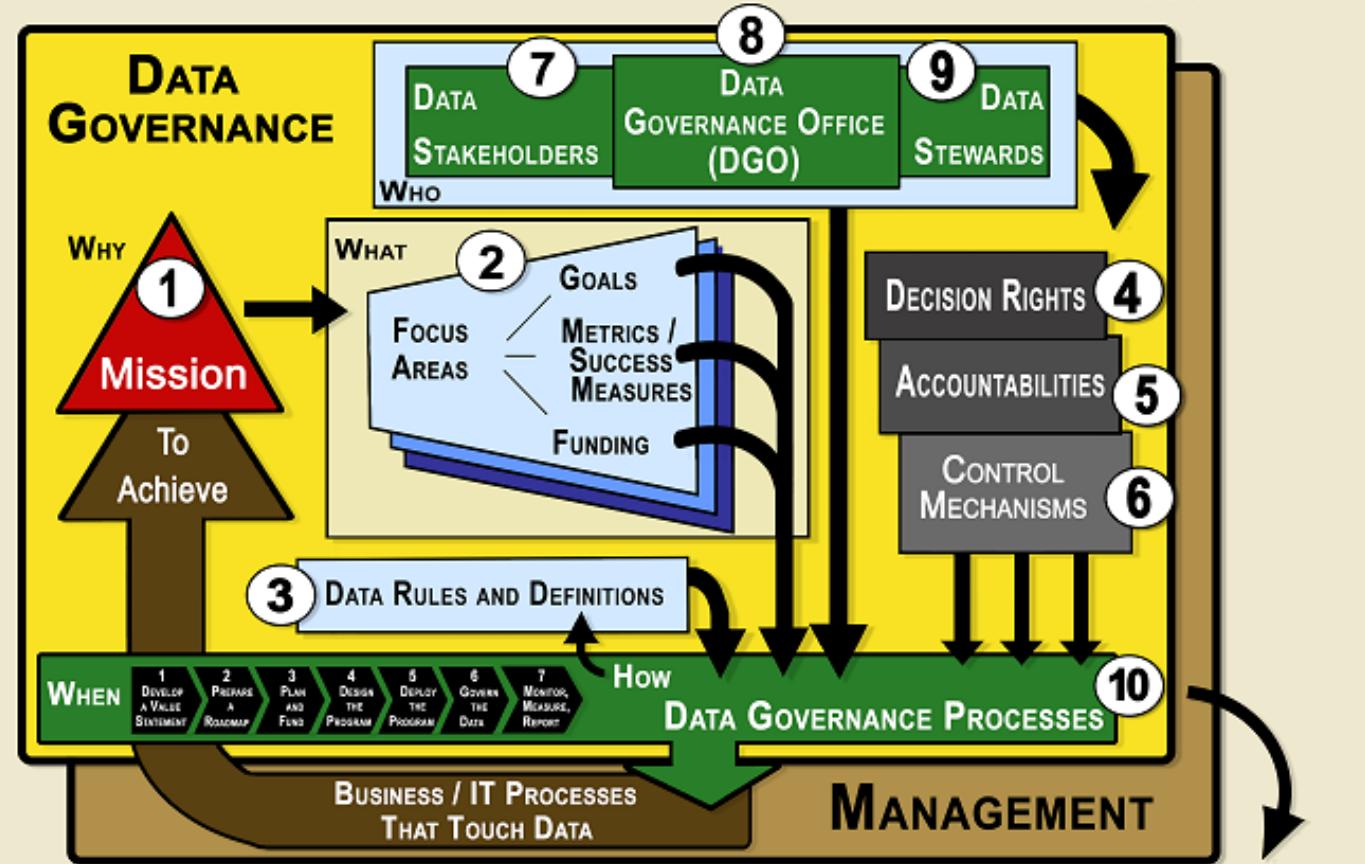


Copyright © by DAMA International

PEOPLE
&
ORGANIZATIONAL
BODIES

RULES
&
RULES OF
ENGAGEMENT

PROCESSES



Definition:

Data Governance is the exercise of decision making and authority for data-related matters.

It's a system of decision rights and accountabilities for information-related processes, executed according to agreed upon models which describe who can take what actions with what information and under what circumstances, using what methods.

Processes for governing how data is used, and when, and by whom

- 1. Aligning Policies, Requirements & Controls
- 2. Establishing Decision Rights
- 3. Establishing Accountability
- 4. Performing Stewardship
- 5. Managing Change
- 6. Defining Data
- 7. Issue Resolution
- 8. Specifying Data Quality Requirements
- 9. Building Governance into Technology
- 10. Stakeholder Care and Support
- 11. Stakeholder Communications
- 12. Measuring and Reporting Value

Data Lake Governance

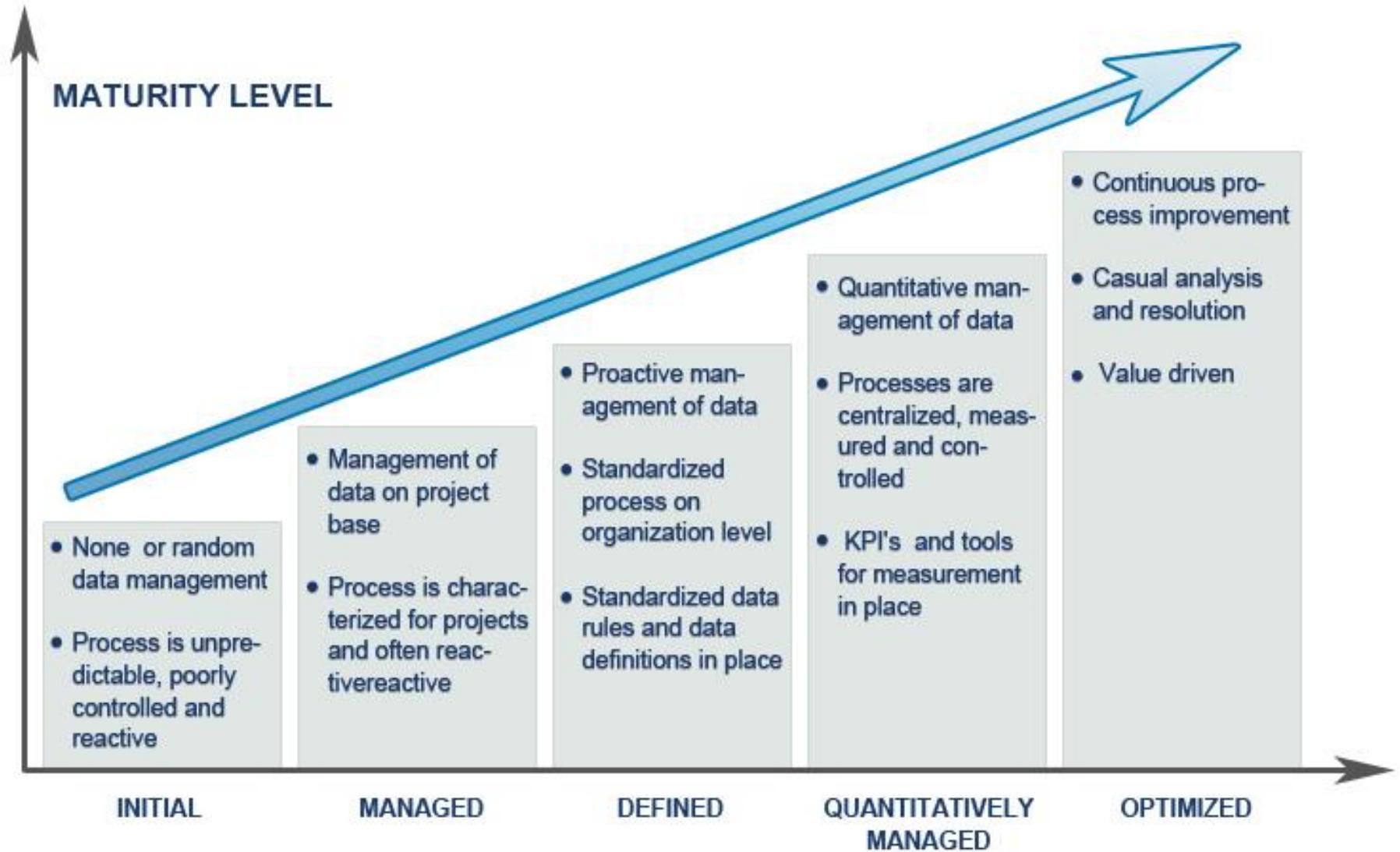
Fundamental Capabilities

- **The definition of the incoming data from a Business use perspective**
- **Documentation of the context, lineage, and frequency of the incoming data;**
- **Security level classification of the incoming data;**
- **Documentation of creation, usage, privacy, regulatory, and encryption business rules which apply to the incoming data.**

What can it do for my Data Lake

- Where did my data come from ? How is it being transformed ?
- Track usage, resolve anomalies, visualize, optimize and clarify data lineage Search and access data
- Assess data quality and fitness for purpose
- Govern who can/cannot access the data
- Data life cycle management, archiving and retention policies
- Auditing, compliance

Data Governance Maturity



Summary

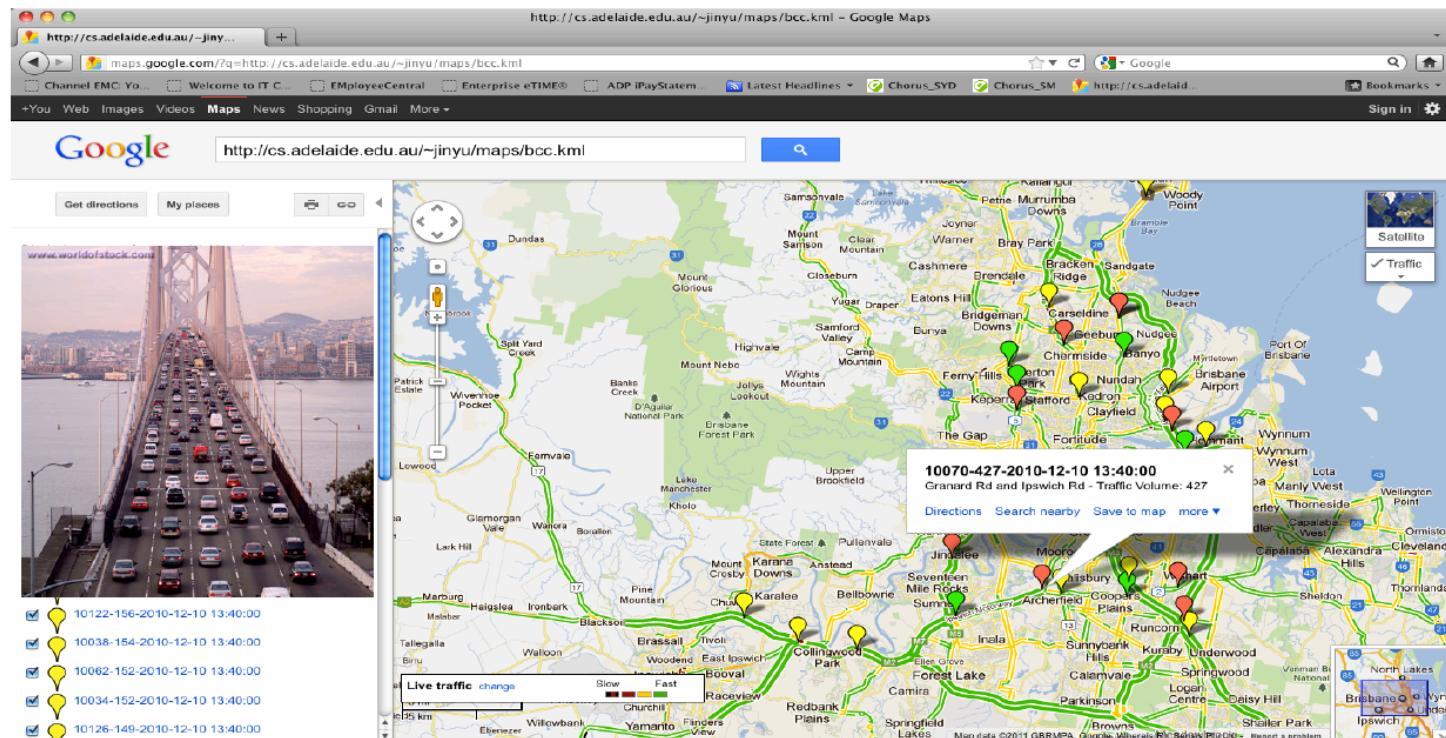
- **Big Data: Data lake instead of data warehouse?**
- **Data Lake is not only a technology, it is an architecture**
- **Data Lake components: Data acquisition (intake), Data management, Data Storage, Data consumption (Discovery)**
- **Data Lake governance is very important**

Customer Case Studies on Big Data Lake

Customer Example: Analytics

Municipal Traffic Analysis to Simulate Traffic Velocity Patterns and Reduce Delays

- Correlate multiple types of data (GPS, weather, sensor, video, social media)
- Simulation techniques to model traffic transition points
- Signal retiming to minimize stops and delays
- Peak delays reduced by 16% and stops reduced by 22%



Pivotal™

Customer Case Studies on Big Data Lake

China Railways scales online sales for the largest rail way in the world with Pivotal Gemfire

- Reliable, High Performance and Continuous uptime with thousands of transactions per second
- On demand scaling for growth
- Cost effective operations

