CSCI E-103

*Data Engineering for Analytics to Solve Business Challenges*

# BI Analytics & Data Visualization

## *Lecture 06*

Anindita Mahapatra & Eric Gieseke

Harvard Extension, Fall 2025

# Agenda

- History of Data warehouses and why they are still popular for BI use cases
- Business Intelligence(BI) & Business Analytics(BA)
- JDBC connection to retrieve data
- KPIs: Concurrency & Latency Requirements
- Data Visualization
- Using the Lakehouse architecture for facilitating BI


- Lab
  - BI Reporting Dashboard

# Review

- Main differences between Lakes & Warehouses?
- What are some ways of hydrating a Data Lake?
- Data Silo Vs Data Swamp?
- What are the 3 phases of the medallion architecture? What is the significance and why is it important?
- 'Data as a product' by decentralized domain centered teams is an example of a _____ architecture?
- Access data wherever it resides - cloud, on-prem, edge etc is an example of a _____ architecture?

# Questions that we'll look at tonight

- What is BI?
- Name a few popular Warehouses
- Name a few popular BI Tools
- Who is the primary data persona for consuming BI data?
- What is the primary skill set of a BI persona?
- How is BI different from AI?
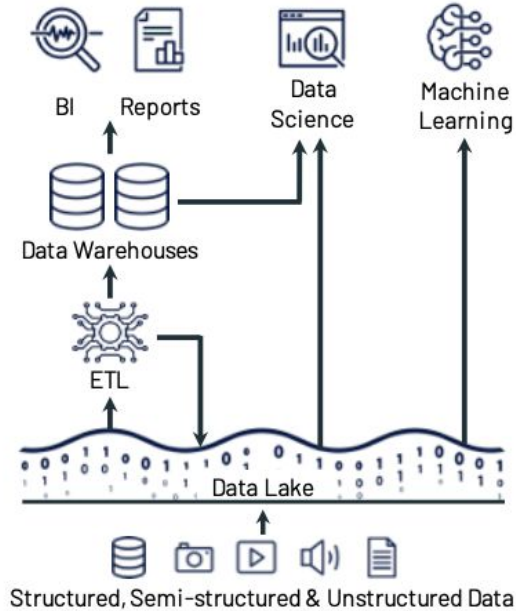
# Data Lake Vs Data Warehouse

| (Dimension) | Data Lake | | Data Warehouse | |
|---|---|---|---|---|
| | **Pro** | **Con** | **Pro** | **Con** |
| **Storage** | Open-format<br>All File Types | Lower quality<br>Coarse file-level access to data | More reliable<br>Fine-grained access control | Mostly structured<br>Proprietary format |
| **Compute** | More economical especially for larger datasets | Operational complexity | Easy to Use<br>High concurrency, low latency | Expensive to scale<br>Limiting historical datasets |
| **Consumption** | Rich ecosystem of tools/frameworks | BI use cases are not first class | Pro-sql | Limited to no ML & streaming use cases |

*Lakehouse gives you the best characteristics of Lakes & Warehouses*
*Performance of a Warehouse with the economics of a Lake*
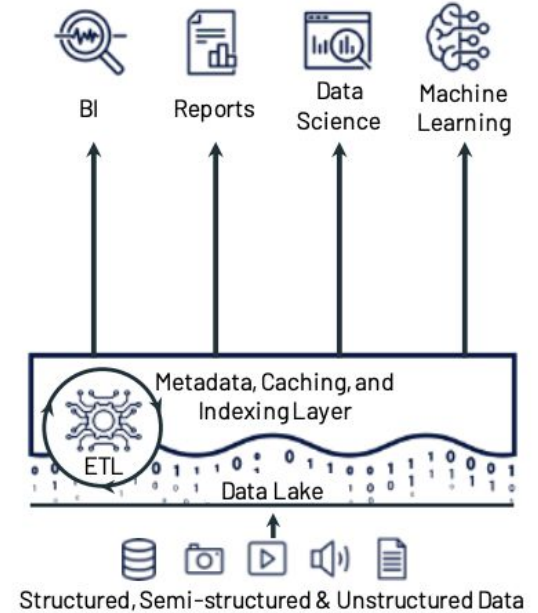*Leading to the important metric **price-performance***

# Lakehouse architecture



(a) First-generation platforms.

(b) Current two-tier architectures.
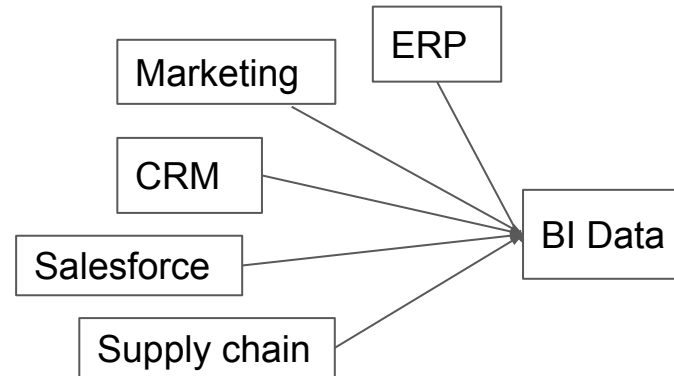
(c) Lakehouse platforms.

# Business Intelligence (BI)

"Data is what you need to do Analytics.
Information is what you need to do Business"
*John Owen*, a theologian.

Refers to technologies, applications and practices for the collection, integration, analysis, and presentation of business information to support better business decision making.
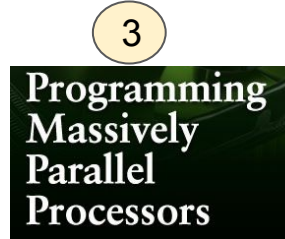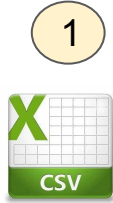
Examples

- Contact and Interaction Analytics
- Closed Deal Analysis
- Website Traffic



Marketing

ERP

CRM

Salesforce

Supply chain

BI Data

real time reporting, dashboards, and analysis.

# Data Store  Evolution



1

2

**Data Warehouse**

3

**Programming Massively Parallel Processors**

4

**N**ot **O**nly **SQL**

5

hadoop

**DATA LAKE**

**DATA SWAMP**

Data As a Service
Data As a Product

6

THE LAKEHOUSE

**Data Mesh/Fabric**: Distributed data architecture, under centralized governance and standardization
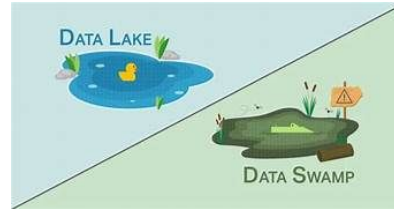
**Bill Inmon**: Father of DW uses ER model in enterprise data warehouse and dimensional model for data marts only

Doug Cutting and Mike Cafarella
Creators of Hadoop

**Ralph Kimball**: Father of DW Proposed dimensional model such as star schemas or snowflakes to organize the data

8

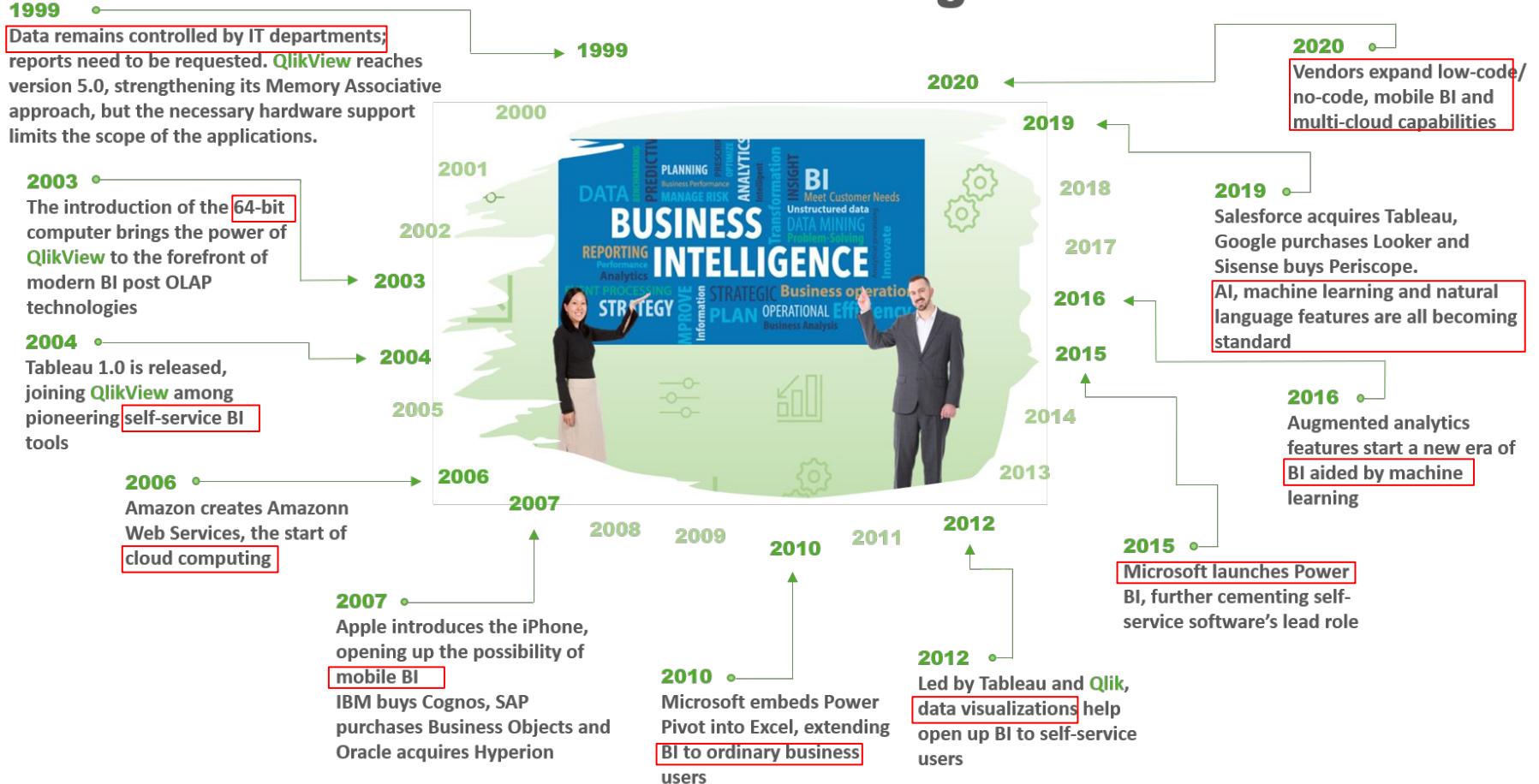# The Evolving Landscape of Business intelligence

**1999**
Data remains controlled by IT departments; reports need to be requested. **QlikView** reaches version 5.0, strengthening its Memory Associative approach, but the necessary hardware support limits the scope of the applications.

**2003**
The introduction of the 64-bit computer brings the power of **QlikView** to the forefront of modern BI post OLAP technologies

**2004**
Tableau 1.0 is released, joining **QlikView** among pioneering self-service BI tools

**2006**
Amazon creates Amazonn Web Services, the start of cloud computing

**2007**
Apple introduces the iPhone, opening up the possibility of mobile BI
IBM buys Cognos, SAP purchases Business Objects and Oracle acquires Hyperion

**2010**
Microsoft embeds Power Pivot into Excel, extending BI to ordinary business users

**2012**
Led by Tableau and **Qlik**, data visualizations help open up BI to self-service users

**2015**
Microsoft launches Power BI, further cementing self-service software's lead role

**2016**
Augmented analytics features start a new era of BI aided by machine learning

**2016**
AI, machine learning and natural language features are all becoming standard

**2019**
Salesforce acquires Tableau, Google purchases Looker and Sisense buys Periscope.

**2020**
Vendors expand low-code/no-code, mobile BI and multi-cloud capabilities

1999  2000  2001  2002  2003  2004  2005  2006  2007  2008  2009  2010  2011  2012  2013  2014  2015  2016  2017  2018  2019  2020

BUSINESS INTELLIGENCE

9

# BI Process for primarily <u>descriptive</u> analysis, trend towards <u>prescriptive</u>

Analyze data & present <u>actionable insights</u> to business stakeholders for <u>decision making</u>

- Self-Service Capabilities,
- Usage monitoring,
- Performance optimization,
- Security controls

Data storytelling features show data in an easy-to-grasp way

**STEP 1** Data from source systems is integrated and loaded into a data warehouse of other analytics repository.

**STEP 2** Data sets are organized into analytics data models or OLAP cubes to prepare them for anlysis.

**STEP 3** BI analysts, other analytics professionals and business users run analytical queries against the date.

**STEP 4** The query results are built into data visualizations, dashboards, reports and online portals.

**STEP 5** Business executives and workers use the information for decision-making and strategic planning.

Data Mining

Predictive Analysis

Text Mining

Statistical Analysis

Big Data

Data Visualization

KPI

Performance Benchmarking

Querying

10

# Business Intelligence (BI) vs Business Analytics (BA)

BI uses past+current data to address the <u>what</u> & <u>how</u>

BA uses past data to explain present and predict future addressing the <u>why</u> & <u>what next</u>

**Answers the questions:**
- → What happened?
- → When?
- → Who?
- → How many?

- → Why did it happen?
- → Will it happen again?
- → What will happen if we change X?
- → What else does the data tell us that we never thought to ask?

**Includes:**
- → Reporting (KPIs, metrics)
- → Automated monitoring and alerting (thresholds)
- → Dashboards
- → Scorecards
- → OLAP* (cubes, slice and dice, drilling)
- → Ad hoc query
- → Operational and-real time BI

- → Statistical or quantitative analysis
- → Data mining
- → Predictive modeling
- → Multivariate testing
- → Big data analytics
- → Text analytics

# Where is the BI data?

- **BI Analyst** is a different persona as compared to Data Engineer & Data Scientist
  - Primarily skilled in sql
- BI data is **curated**
- BI data is typically stored in enterprise Data Warehouses or in specialized Data marts
- In recent years, Data Lakes have also been added to that list
  - Initial curation of data before it is pushed to a Warehouse
  - BI tools can now directly tap into all the data in the Data Lake (Lakehouse)
- Improvement in data democratization efforts is allowing for better
  - Self-service BI
  - Data Discovery
  - Data Mining for better what-if predictive scenarios

# Warehouse Terminology

| Catalog | Stores Metadata |
|---|---|
| Database/Schema | Namespace |
| Table | Data with storage<br>● Managed<br>● UnManaged/External |
| Keys/Indexes/Constraints | ● PK, FK (not enforced, just for relationship understanding)<br>● Identity columns<br>● Surrogate Keys Vs natural keys<br>● Constraints for data quality |
| View | Virtual table (hydrated by a query) |
| Federated Query | Cross Data-Store Boundaries<br>Push down predicates |
| Materialized View | Pre-computed to facilitate faster access (more frequent access) |
| Stored Procedure/UDF | Logic Encapsulation |
| Semantic Data Model | Relationships captured using business terminology Vs referential constraints |

# Components of good BI

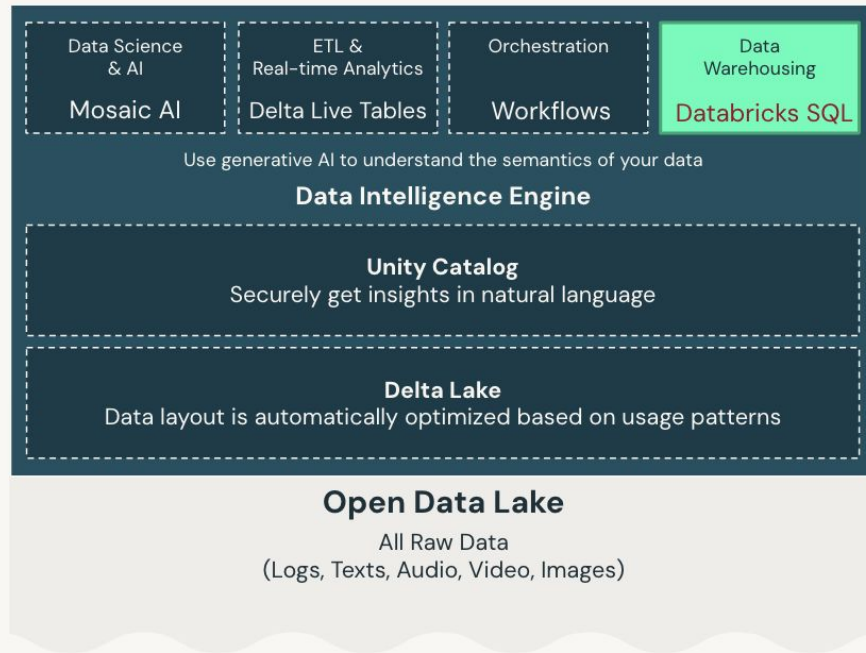| | |
|---|---|
| **Data Collection** | <u>Business data</u> of any nature, that lies scattered across flat files, feeds, databases, cloud storage and business applications is <u>gathered</u> for further analysis and reporting. |
| **Data Preparation** | Data collected from different sources goes through a sequence of steps : integration, modelling, cleansing, preparation and enrichment, before organizing it into an analytics-ready format. |
| **Intelligent Analytics** | Derive maximum value out of the available data, by doing analysis to uncover insights about - 1.what had happened 2. why and how did it happen and even go ahead to predict 3. What might happen. |
| **Data Visualisation** | Analytical insights can be made <u>easily consumable through dashboards and reports</u>, that shall be built with an easy-to-use drag-and-drop interface. |
| **Sharing and Collaboration** | The insightful reports and dashboards can be <u>shared</u> with each other, for collaborative analytics and informed decision-making. |
| **Data Governance** | Who has access to what information |
| **Strategy Documentation** | Centralized Data Catalogs |
| **Ease of Use, Implementation & Integration** | Data democratization; How long to implement a BI solution; Integration with existing technology stack |

14

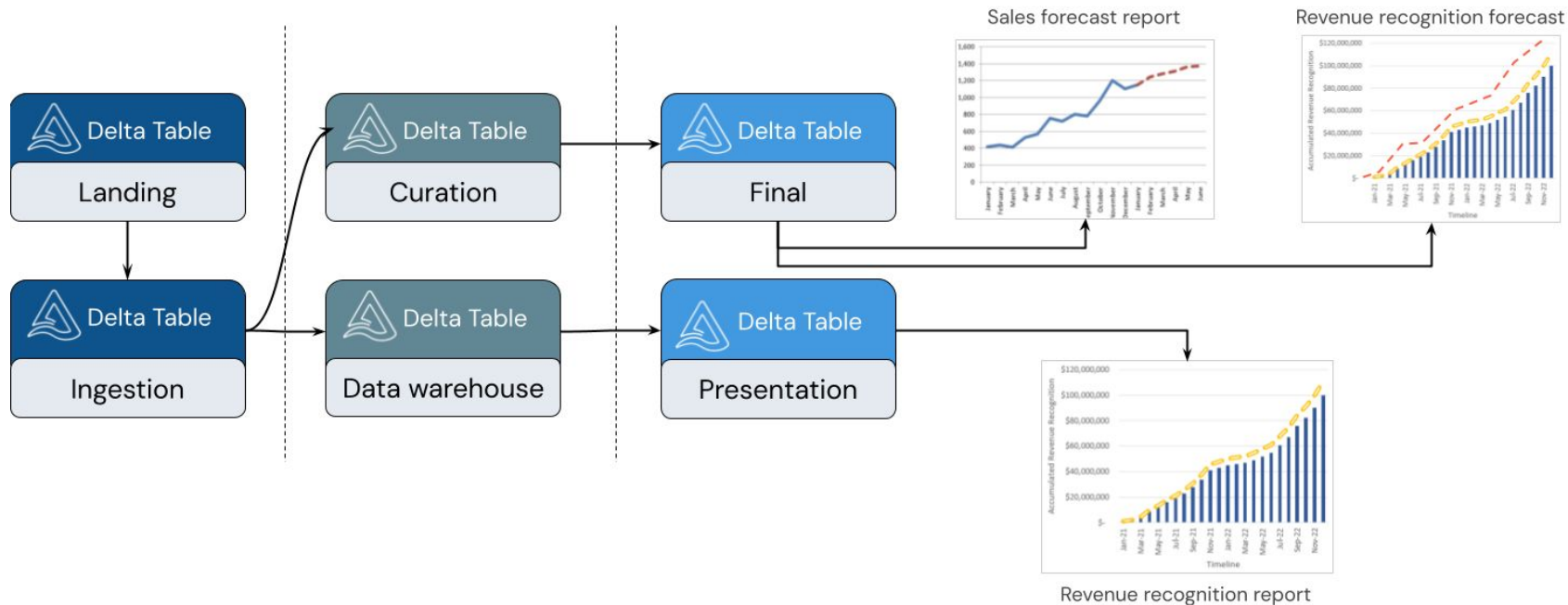# Databricks Data Intelligence Platform on AWS

| Source | Ingest | Transform | Query and Process | Serve | Analysis |
|---|---|---|---|---|---|

**Integration**

## Data Intelligence Platform on AWS

### Orchestration
Workflows

**Federation**

Amazon Redshift …

**ETL**

Sensors and IoT (unstructured)

RDBMS (structured)

Files / Logs (semi-structured)

Media (unstructured)

Business Apps (structured)

Other clouds

### Batch & Streaming
Auto loader

### Batch & Streaming
Amazon AppFlow

AWS Glue

AWS IoT Core

Amazon Kinesis

AWS DMS

### Data Science & Gen AI
ML Modeling

MLOps

Gen AI

Mosaic AI

Feature Serving

Model Serving

Vector Search

### Processing, ETL, Real-time Analytics
Delta Live Tables

Spark / Photon

### Data Warehousing
Databricks SQL

### Data Intelligence Engine
Assistant

DatabricksIQ

Data Rooms

### Data and AI Governance
Catalog & Lineage

Access Control

Unity Catalog

Models / Features

Lakehouse Monitoring

### Data
Delta Lake

bronze  silver  gold

### Collaboration
Delta Sharing

Market place

**AI Apps**
AI App

**Dashboards**
Lakeview

**BI Tools**
Amazon QuickSight

**Operat. DBs**
Amazon RDS

Amazon Dynamo DB

**Business Apps**
Biz App

**3rd party**
Data Consumer

15

**ID Provider**
ID Provider

**Governance**
Enterprise Catalog

**AI Services**
Hugging Face

Amazon Bedrock
…

**Key**
Domain
Key capability

Amazon S3

## Storage

# Databricks SQL

Delivering analytics on the freshest data with data warehouse performance and data lake economics

- Home for data analysts
- Use ANSI SQL to query data
- Built on open foundation
- Broad Integration with other BI tools like Tableau or Power BI
- Partner connect to aid data hydration
- Better price / performance than other cloud data warehouses
- Simplify discovery and sharing of new insights
- Simplified administration and governance – setup catalog, discover data, view lineage
- Lakeview Dashboards
- AI/BI Genie space

| Data Science & AI | ETL & Real-time Analytics | Orchestration | Data Warehousing |
|---|---|---|---|
| Mosaic AI | Delta Live Tables | Workflows | Databricks SQL |

Use generative AI to understand the semantics of your data

**Data Intelligence Engine**

**Unity Catalog**
Securely get insights in natural language

**Delta Lake**
Data layout is automatically optimized based on usage patterns

**Open Data Lake**
All Raw Data
(Logs, Texts, Audio, Video, Images)

# Data modeling for Data Warehouse

# Modeling for the Warehouse

Link

- Approach
  - Understand relations aka OOP - Is a,Has a … Nouns & Verbs
  - Semantic -> Logical -> Physical
- Patterns
  - 3 NF (normalization)
    - Tight referential integrity
  - Dimensional Modeling Star Schema Snowflake Schema
    - Facts
    - Dimensions
    - *Query-optimized to support BI*
  - Data Vault (more flexible/adaptable)
    - Hub (core business concepts,  eg. ids)
    - Links (PK/FK between Hubs)
    - Satellite (descriptive attributes)
    - Dimensional Model on top
    - *Adaptable to change, supports data integration and governance*

A subscription business integrating customer, product, and transaction data from disparate sources. This involves creating **Hub tables** for core entities (e.g., customers, products), **Link tables** for relationships (e.g., customers purchasing products), and **Satellite tables** for detailed, time-varying attributes of those entities.

This model enables rapid integration of new data sources by adding new Satellite tables without disrupting the existing structure and allows for robust historical tracking and an ELT process for greater agility.

# A few different approaches

| Bronze | Silver | Gold |
|--------|--------|------|



Raw Ingestion          Filtered, Cleaned          Curated: Business-level
                            Augmented                    & aggregates

**Expose & Query Gold Tables**
- Data engineering or "analytics engineering" curates & provides access to curated / gold level tables to the rest of the organization.
- Usually follows best-practices with proper modeling. (e.g. Kimball, denormalized reporting or mixed)
- Works well for less-technical users, as well as serving external users (e.g. companies selling data / insights)
- In this model, end-users typically do less self-service / "last-mile" ETL and rather rely on curated assets.

# Persona Handoffs



**Data Steward**
**Setup & Administration**

- Central Query Log
- Usage attribution
- Debug & Troubleshoot

**Data**

Metastore

**Data Integrators**
**DE/ML**

1

2

**SQL Endpoint**
**(Compute)**

Concurrency
and scaling

Performance

4

CDO

3

- SQL Editor (auto-complete)
- Parameterized Queries
- Built-in Visualizations
- Dashboard
- Data Refresh (Scheduling)
- Alerting (Notification)
- Built-in connectors for BI tools

**BI Analysts**

Data Integrators bring in the data

Data Engineers ETL and curate it

DS add ML insights

Administrators provide compute and data access

BI Analysts work off curated data

Executive business users consume the reports

20

# Newer Capabilities of Warehouses

- Federated Data Warehouse
  - Join multiple sources jdbc/odbc or use true federation
- Virtual Data Warehouse
  - Views and materialized views (MV) on the data with ACL and caching
- Realtime Warehouse
  - Caching, Streaming tables, union of hot data and some cached immutable historical data
  - Eg. Create a real-time orders table that combines historical and hot data
- ETL despite support for federate & virtualize
  - Latency, low throughput from source, CDC
- Time series Data Warehouse
  - Create view that uses window function, order by (desc) key to return the first one
- Data Lake
  - Capture as much data, use metastore for definitions(schema, data loc, format, partitions)
- Schema Evolution
- Intelligent data warehousing
  - Access for everyone to ask questions of their data using natural language
  - Intelligent, automated management and tuning
  - Optimal TCO

# Databricks SQL Features

| | |
|---|---|
| **EXPLORATORY SQL** | SQL Editor with intelligent auto complete, ANSI SQL |
| **MANAGEMENT & GOVERNANCE** | Query History & Profile, Data Explorer (Unity Catalog), Managed Data Sharing |
| **CONNECTIVITY** | SQL Rest API, Python, Node.js, Go*, Partner Connect |
| **PERFORMANCE** | Photon Engine (Massively Parallel Processing) |
| | Predictive I/O |
| **SQL ETL/ELT** | Query Federation*, Materialized Views*, Workflows Integration* |
| **DATA SCIENCE & ML** | Python UDFs*, Notebooks Integration*, Geospatial* |
| **SERVERLESS DATA WAREHOUSE** | Instant, Elastic, Fully Managed Compute* |
| **HIGH CONCURRENCY BI** | Intelligent Workload Management* |
| | Serverless Query Result Caching* |

# Governed and secured by Unity Catalog

## Governance for all your data and AI assets



Simplified **data discovery**, **governance**, **federation**, **lineage**, and **compliance** with enhanced **security** and **auditing** with Unity Catalog and Databricks SQL

# Simple and fast performance

## Accelerating federated workloads

**Federation 💛 Materialized views:**

- Consistent latency & concurrency for data outside of the Lakehouse

- Accelerate cross-source joins and complicated transformation logic

- Offload access to underlying databases via materialized views to avoid high/concurrent loads on operational databases

# Simple streaming with SQL

```
CREATE STREAMING TABLE  web_clicks
AS
SELECT *
FROM STREAM
  read_files('s3://mybucket')
```

```
CREATE STREAMING TABLE  server_logs
AS
SELECT from_json(...) data
FROM STREAM
  read_kafka(...)
```

Data stream

Cloud Storage
(S3, ADLS, GCS)

Message Queues
(Kafka, Pub/Sub, Kinesis, etc)

## Benefits:

1. **Unlock real-time use cases.** Ability to support real-time analytics/BI, machine learning and operational use cases with streaming data.

2. **Better scalability.** More efficiently handle high volumes of data via incremental processing vs. large batches.

3. **Enable more practitioners.** Simple SQL syntax makes data streaming accessible to all data engineers and analysts.

# Simple and fast BI with Materialized Views

```
CREATE MATERIALIZED VIEW  customer_orders
AS
SELECT
  customers.name,
  sum(orders.amount),
  orders.orderdate
FROM orders
  LEFT JOIN customers ON
    orders.custkey = customers.c_custkey
GROUP BY
  name,
  orderdate;
```

Results are
pre-computed and
incrementally
refreshed

```
customers
(Table)
```

```
orders
(Table)
```

## Benefits:

1. **Accelerate BI dashboards.** Much faster to query data that is pre-computed vs querying base tables.

2. **Reduce data processing costs.** MV results are refreshed incrementally avoiding the need to completely rebuild the view when new data arrives.

3. **Improve data access control.** More tightly govern what data can be seen by consumers by controlling access to base tables.

# Simple orchestration for your SQL and more

## Queries, notebooks dashboards, alerts, and more!



Automate and schedule Databricks SQL workloads with advanced workflow orchestration, reliable monitoring and observability

# Simply access any LLMs directly in Databricks SQL

```sql
SELECT
 sku_id,
 product_name,
 ai_query (
    "my-external-openai-chat",
    "You are a marketing expert for a winter holiday promotion
targeting GenZ. Generate a promotional text in 30 words mentioning a
50% discount for product: " || product_name
 )
FROM
 uc_catalog.schema.retail_products
WHERE
 inventory > 2 * forecasted_sales
```

**Integrate any LLMs in SQL** to enrich data
and empower analysts to extract actionable insights

# Simple help from AI in your SQL

Write SQL to get insight from unstructured text data via LLMs

## SQL AI ANALYZE SENTIMENT

```
> SELECT ai_analyze_sentiment('I am happy');
  positive

> SELECT ai_analyze_sentiment('I am sad');
  negative
```

## AI SQL CLASSIFY

```
SELECT ai_classify("My password is leaked.", ARRAY("urgent", "not urgent"));
urgent

SELECT
  description,
  ai_classify(description, ARRAY('clothing', 'shoes', 'accessories', 'furniture')) AS category
FROM
  products
```

## SQL AI EXTRACT

```
> SELECT ai_extract(
    'John Doe lives in New York and works for Acme Corp.',
    array('person', 'location', 'organization')
  );
{"person": "John Doe", "location": "New York", "organization": "Acme Corp."}

> SELECT ai_extract(
    'Send an email to jane.doe@example.com about the meeting at 10am.',
    array('email', 'time')
  );
{"email": "jane.doe@example.com", "time": "10am"}
```

## SQL AI FIX GRAMMAR

```
SELECT ai_fix_grammar('This sentence have some mistake');
"This sentence has some mistakes"

SELECT ai_fix_grammar('She dont know what to did.');
"She doesn't know what to do."
```

## SQL AI MASK

```
SELECT ai_mask(
  'John Doe lives in New York. His email is john.doe@example.com.',
  array('person', 'email')
);
[MASKED] lives in New York. His email is [MASKED]."

SELECT ai_mask(
  'Contact me at 555-1234 or visit us at 123 Main St.',
  array('phone', 'address')
);
Contact me at [MASKED] or visit us at [MASKED]"
```

## SQL AI SIMILARITY

```
SELECT ai_similarity('Apache Spark', 'Apache Spark');
1.0

SELECT
  company_name
FROM
  customers
ORDER BY ai_similarity(company_name, 'Databricks') DESC
LIMIT 1

Databricks Inc.
```
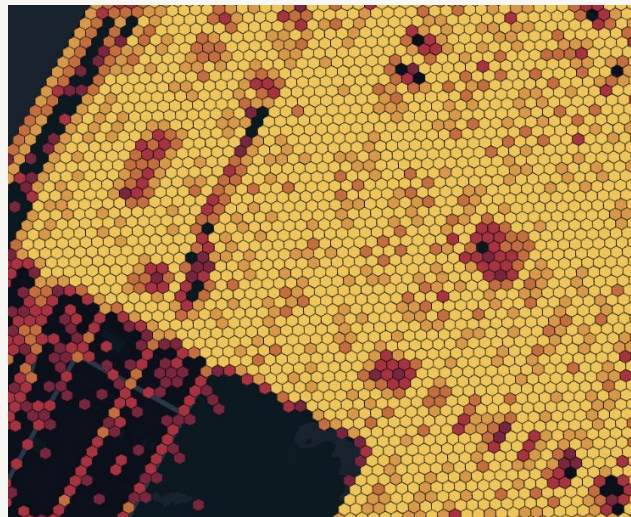
# Geospatial support

## Supercharge your geospatial processing

**Efficient storage** for spatial data in both large and small sizes

**Fast SPATIAL JOINs** and binning support

**Easy to visualize and integrate with ML** you don't need to switch betweens tools to maximize geospatial data value



Rideshare pick-up locations in New York City visualized in a Databricks Notebook using Kepler.gl

# Databricks SQL is a complete data warehouse

**Data engineering, ETL**

Auto-loader
Materialized views
Streaming tables
Data lineage in UC
Lakehouse federation
PK/FKs, ERD in Catalog ANSI
SQL by default
Rich, Tabbed SQL Editor
Notebooks on SQL WH
SQL Execution API
Python UDFs
SQL session variables
Row level security
Column masking
Dark mode
Schema browser

**Enterprise scale and perf**

Serverless
Intelligent autoscaling
Adaptive routing
Predictive optimization
Predictive I/O
Results caching
Liquid clustering
100K+ user support
Query scheduler
SQL tasks in workflows
Statement history
System tables: WH events
System tables: Billing
System tables: Audit log

**Native BI + DW ecosystem**

Lakeview
Databricks assistant for LV
Delta sharing
Data marketplace
Data rooms
Clean rooms
100+ Integrations
Partner connect (25+)
Power BI, Tableau
Publish to PowerBI Online
dbt: incremental models
dbt: materialized views
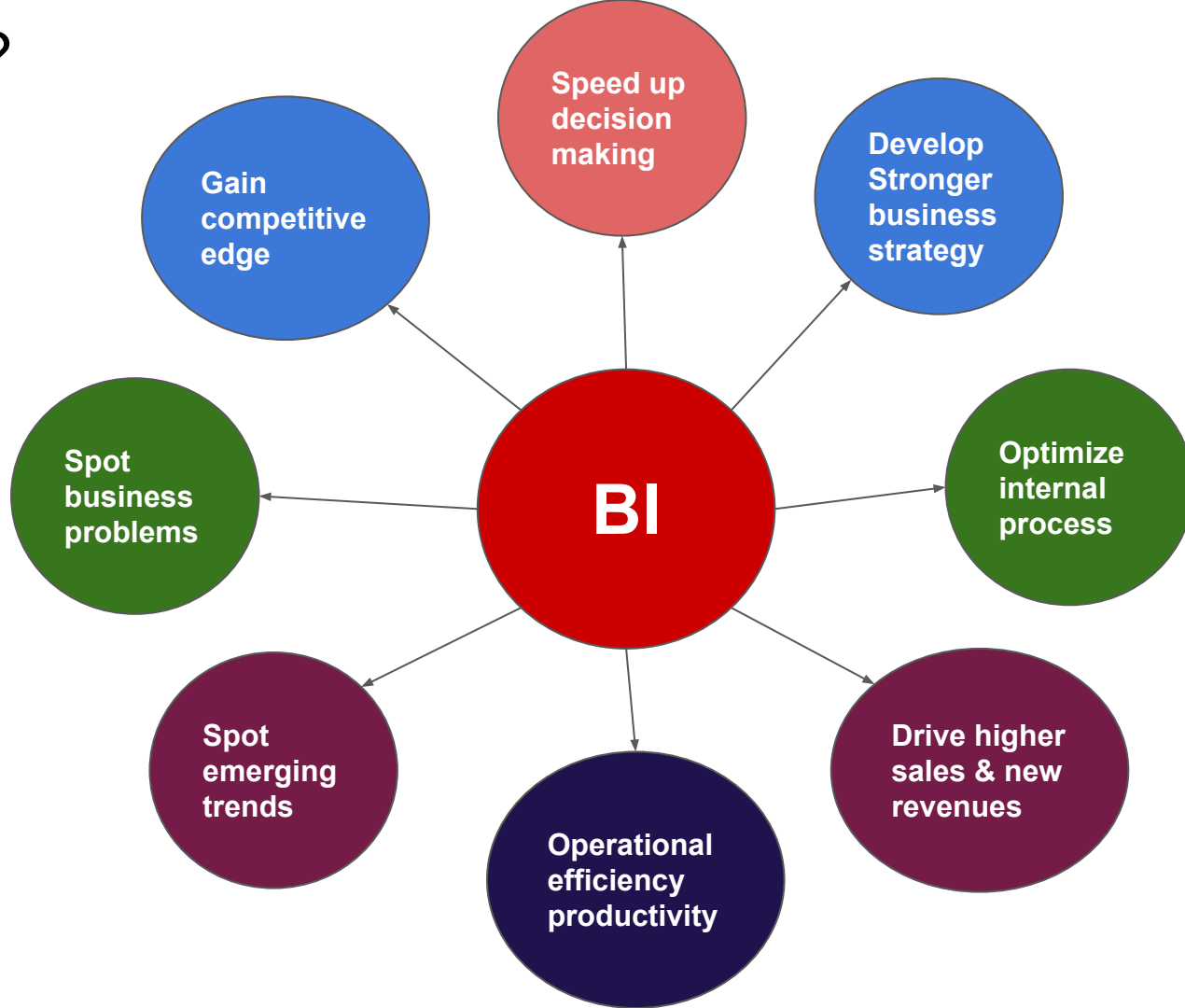Fivetran
OAuth
Cloud Fetch fast results

# Lab

- Use the Databricks UI to create a catalog, schema, table, and view.
- Use the Databricks UI to upload data and create a managed table.
- Use the SQL Editor to complete multiple data analytics tasks.
- Create a data visualization associated with a query.
- Create an interactive dashboard.
- Create a refresh schedule and alert.
- Share data based assets in the Databricks DI Platform with others.

# Appendix

# Why BI?

# Simple migrations at your pace



## CRAWL
**with Lakehouse Federation**

–Get your data in one place, but don't migrate… FEDERATE
–Based on which datasets your teams are moving and where business value is created, THEN migrate those over to Databricks

**IN PRODUCTION TODAY**

## WALK
**with Materialized Views**

–Set up materialized views on top of the federated source data
–This will create a copy of the meta data in Delta Lake, relieving the pressure and cost on the source data system!

**IN PRODUCTION TODAY**

## RUN
**with Change Data Capture (CDC)**

– CDC will only process the pieces of the data that CHANGED, making it SIMPLER and CHEAPER
–This is made possible by the Arcion acquisition and will replace the Materialized Views from WALK stage.

**COMING SOON**

**ALL OF THIS IS POWERED BY UNITY CATALOG**