# Kudos

Akamai

aruba CLOUD

beSharp.
we make IT run.

DATADOG

adesso.it

aws

EssilorLuxottica
Our Vision. Your Future.

software one

Cleafy

CoNDENSE

# Unprecedented Growth of Data

There is more data and more diversity
of data than people think

**Data growth**

**>10x**

every 5 years

**Data
platforms needs**

To live for

**15+**

years

To scale

**1,000x**

*IDC, "Data Age 2025"*

# Traditional Data Architecture

In the past, decision-making revolved around the **enterprise data warehouse**.



OLTP

LOB

ERP

CRM

Enterprise data warehouse

Business intelligence

# Diversity of Roles and applications

Data scientists

Business users

Analysts

Applications

Machine learning
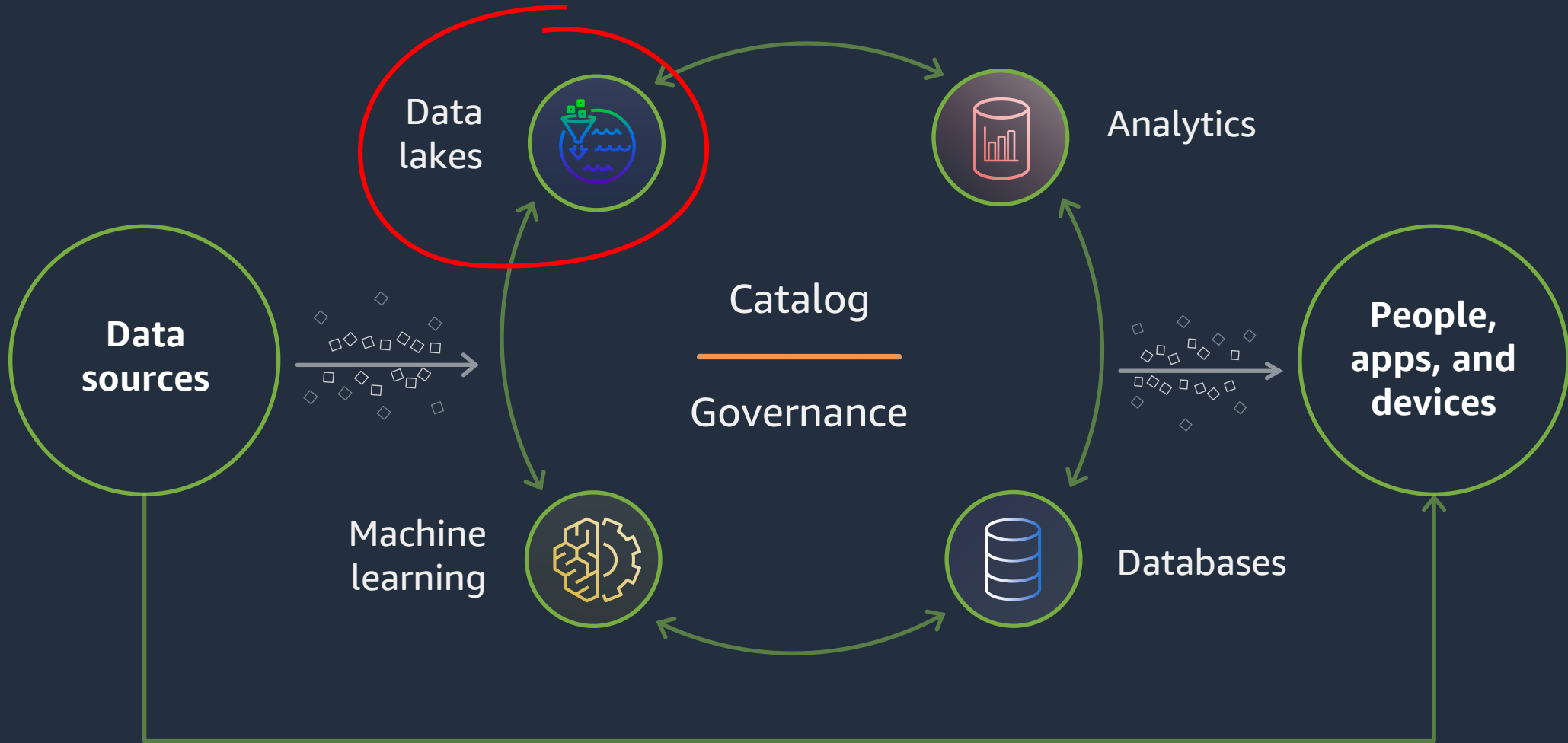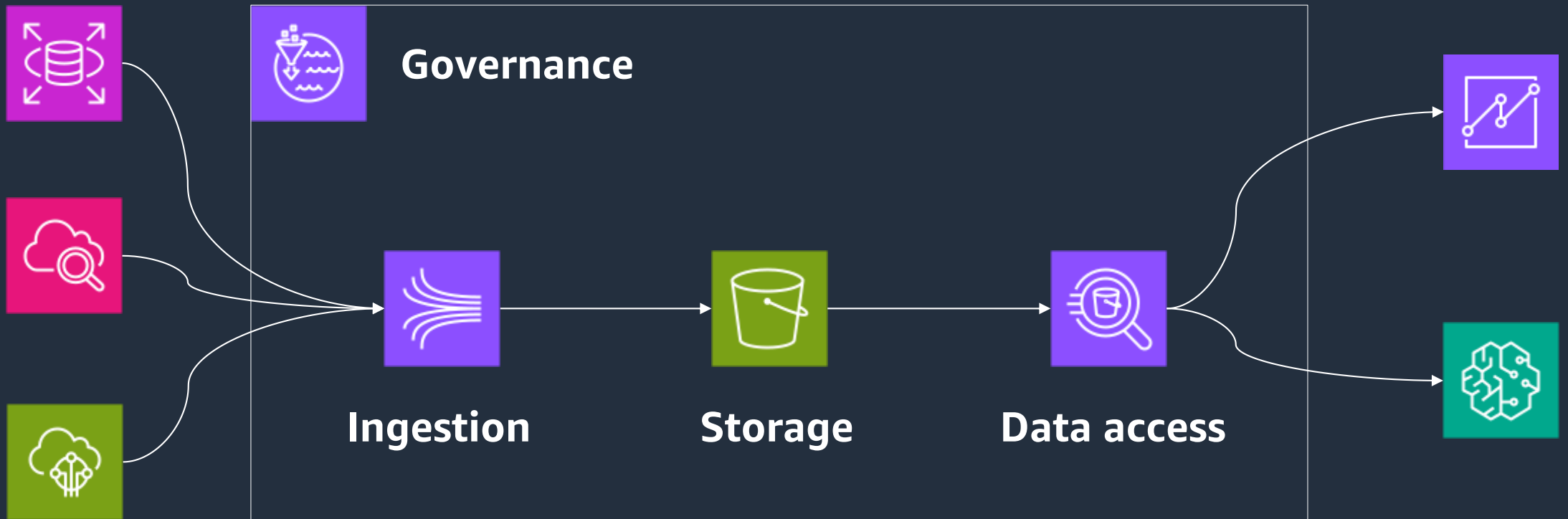
SQL analytics

Scientific

Real-time, streaming

There are **more people** accessing data

And in **different ways**
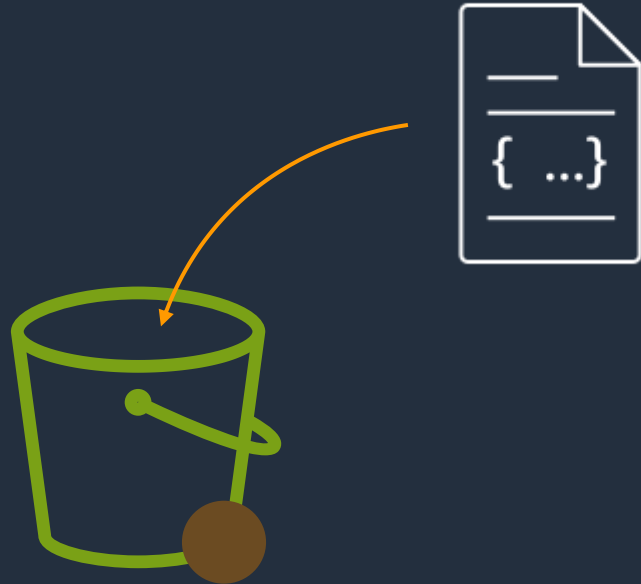
# Modern Data Architecture



Data lakes

Analytics

Data sources

Catalog

Governance

People, apps, and devices

Machine learning

Databases

# Data Lake Fundamentals

**Governance**

**Ingestion**    **Storage**    **Data access**

**Data sources**                                        **Consumers**

# Ingestion Bucket



## Bronze

Misc formats
(e.g. JSONL, CSV, PDF)

Raw data

# Format Transformation



**Bronze**

Misc formats
(e.g. JSONL, CSV, PDF)

Raw data

**Silver**

Write Optimised
(e.g. Avro)

Curated data

# Data Lake Conventions

s3://my-data-lake/customers/year=2023/month=03/day=26/hour09.xyz

*base name*   *table*   *partition 1*   *partition 2*   *partition 3*   *data file*

# Data Lake Formats

**Avro**
*row-based*

**Parquet**
*column-based*

**ICEBERG**
*transactional*

*Not in scope for this approach*

# Three-tier Data Lake



**Bronze**

Misc formats
(e.g. JSONL, CSV, PDF)

Raw data

**Silver**

Write Optimised
(e.g. Avro)

Curated data

**Gold**

Read Optimised
(e.g. Parquet)

Consumable data

# Transformations

{ "id": "123", "user": { "id": "987", "name": "Jack Johnson", "email": "jj@example.com", "city": "Aachen" }, "items": [ {"p_id": "456", "qty": 1, "price": 37.81}, {"p_id": "567", "qty": 2, "price": 42.35} ], "datetime": "2024-05-23T18:23:18Z" }

{ "id": "123", "user_id": "987", "user_city": "Aachen", "order_total": 122.51, "items": [ {"p_id": "456", "qty": 1, "price": 37.81}, {"p_id": "567", "qty": 2, "price": 42.35} ], "datetime": "2024-05-23T18:23:18Z" }

Apache Spark

AWS Glue

# Cost Drivers

With **Terabytes** of total volumes and **Gigabytes** per transaction, be aware of typical cost drivers …

**Compute**

**Storage**

**Data Transfer**

# Partitioning & Compression Examples

```
select count(*) from datalake where dt >= '20170515' and dt < '20170516'
```

| Partitioning | Size on S3 | Run Time | Data Scanned | Cost |
|---|---|---|---|---|
| NO | 74 GB | 10.41 sec | 74.1 GB | $0.36 |
| YES | 74 GB | 2.73 sec | 871.39 MB | $0.004 |
| **Result** | **same** | **4x faster** | **85% less** | **98% cheaper** |

Process all dataset

| Compression | Size on S3 | Run Time | Data Scanned | Cost |
|---|---|---|---|---|
| Text | 1.15 TB | 3m 56s | 1.15 TB | $5.75 |
| Parquet | 130 GB | 6.78s | 2.51 GB | $0.013 |
| **Result** | **87% less** | **34x faster** | **99% less** | **99.7% savings** |

# Metadata Catalog



Table properties

Data statistics

Table schema
& Partitions

Nested fields

# APIs and Query Semantics for Big Data

{ "id": "123", "user": { "id": "987", "name": "Jack Johnson", "email": "jj@example.com", "city": "Aachen" }, "items": [ {"p_id": "456", "qty": 1, "price": 37.81}, {"p_id": "567", "qty": 2, "price": 42.35} ], "datetime": "2024-05-23T18:23:18Z" }

{ "id": "124", "user": { "id": "988", "name": "John Jackson", "email": "jj@example.net", "city": "Milan" }, "items": [ { "p_id": "678", "qty": 3, "price": 43.19 }, { "p_id": "789", "qty": 25, "price": 2.88 } ], "time": "2024-11-20T09:43:10Z" }

…

*Task: "Get orders from Milan in Nov. 2024 or later"*

# Query Semantics: REST

/orders?customerCity=Berlin&orderDateFrom=2023-03-01

# Query Semantics: OData

/orders?$filter=City eq 'Berlin'&$expand=Orders($filter=OrderDate ge 2023-03-01;$expand=OrderItems)

# Query Semantics: GraphQL

```graphql
query {
  customers(filter: { city: "Berlin" }) {
    id
    name
    orders(filter: { datetime: { gte: "2023-03-01" } }) {
      id
      date
      orderItems {
        p_id
        qty
        price
      }
    }
  }
}
```

# Query Semantics: SQL

```sql
SELECT c.*, o.*, oi.*
FROM customers c
JOIN orders o ON c.customer_id = o.customer_id
JOIN order_items oi ON o.order_id = oi.order_id
WHERE c.city = 'Berlin'
    AND o.order_date >= DATE '2023-03-01'
```

trino

Amazon
**Athena**

# Data API Security

# Column-level Access

# Row-level Access

Row-level access

Choose whether this filter should have row-level restrictions.

○ Access to all rows
● Filter rows

Row filter expression

Enter the rest of the following query statement SELECT * FROM nested-table WHERE...
Please see the documentation for examples of filter expressions.

customer.customerName <> 'John'

# Workloads Control - Athena Workgroups



Unique query output location per Workgroup

Encrypt results with unique AWS KMS key per Workgroup

Collect and publish aggregated metrics per Workgroup to AWS CloudWatch

Use Workgroup settings eliminating need to configure individual users

**Workgroup name*** `adhoc-users`
Use 1 - 128 characters. (A-Z,a-z,0-9,_,-,.)

**Description** `Workgroup for ad-hoc analytics users`
Use up to 1024 characters.

**Query result location** `s3://bucket/adhocusers/` 📁 Select
Enter a path to an S3 bucket or prefix.

**Encrypt query results** ☑ Encrypt results stored in S3

**Encryption type** `SSE-KMS` ▼ ⓘ

**Encryption key** `aws/s3` ▼ ⓘ ⧉ Create KMS key

**Metrics** ☑ Publish query metrics to AWS CloudWatch ⓘ

**Override user settings** ☑ ⓘ

# Example 1: Sync Query API



AWS IAM
(Identity & Access)

Amazon API
Gateway

AWS Lambda
(FaaS)

Amazon Athena
(Query Engine)

AWS Glue
(Data Catalog)

Amazon S3
(Data Lake)

# Example 2: Async Query API

# Example 3: Data Temperature Routing



Amazon Aurora
(Relational DB)

AWS Glue
(Data Catalog)

Hot data

Warm &
cold data

Query data

Amazon API
Gateway

AWS Lambda
(FaaS)

Amazon Athena
(Query Engine)

Amazon S3
(Data Lake)

# Do it Yourself!

# Additional resources



All resources can be found here:
https://de3n2axe8vokl.cloudfront.net

- 📑 Presentation Slides
- 🔗 Blog Links
- 📚 Suggested Training Material
- 📊 Survey

Linkedin @venturir

**Thanks!**

Riccardo Venturi

Senior Solution Architect @ AWS

improove

**CL▶UD DAY 2024**

**Milano, Nov 20**