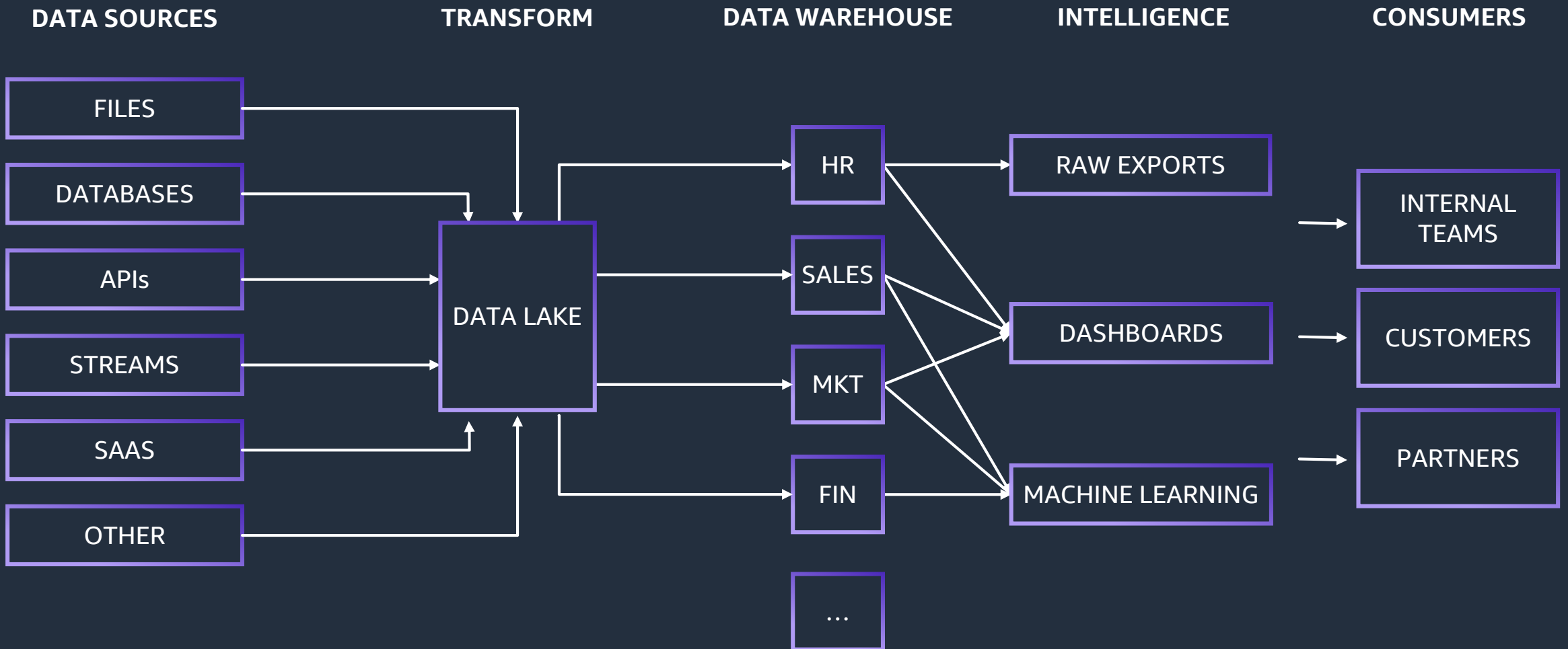aws

**ADC-C3**

# The Zero-ETL future of analytics
## Solving data integration challenges

Sandipan Bhaumik (he/him)
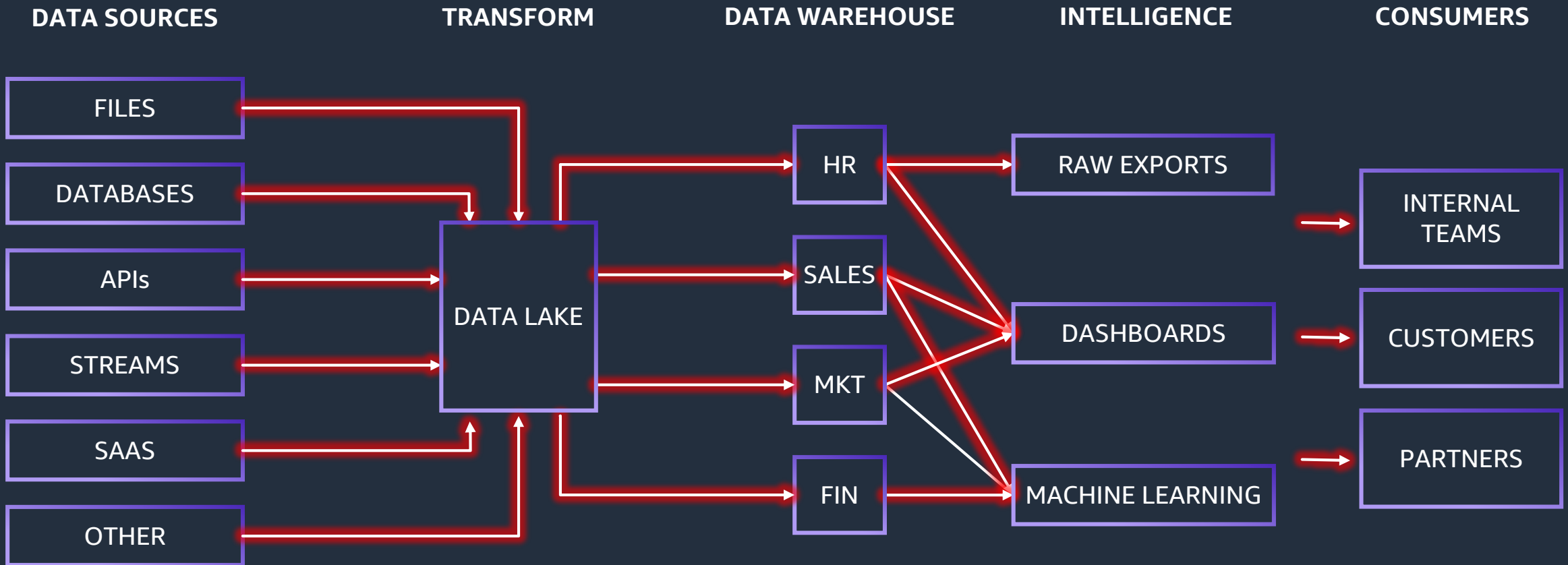
Sr. Analytics Solutions Architect
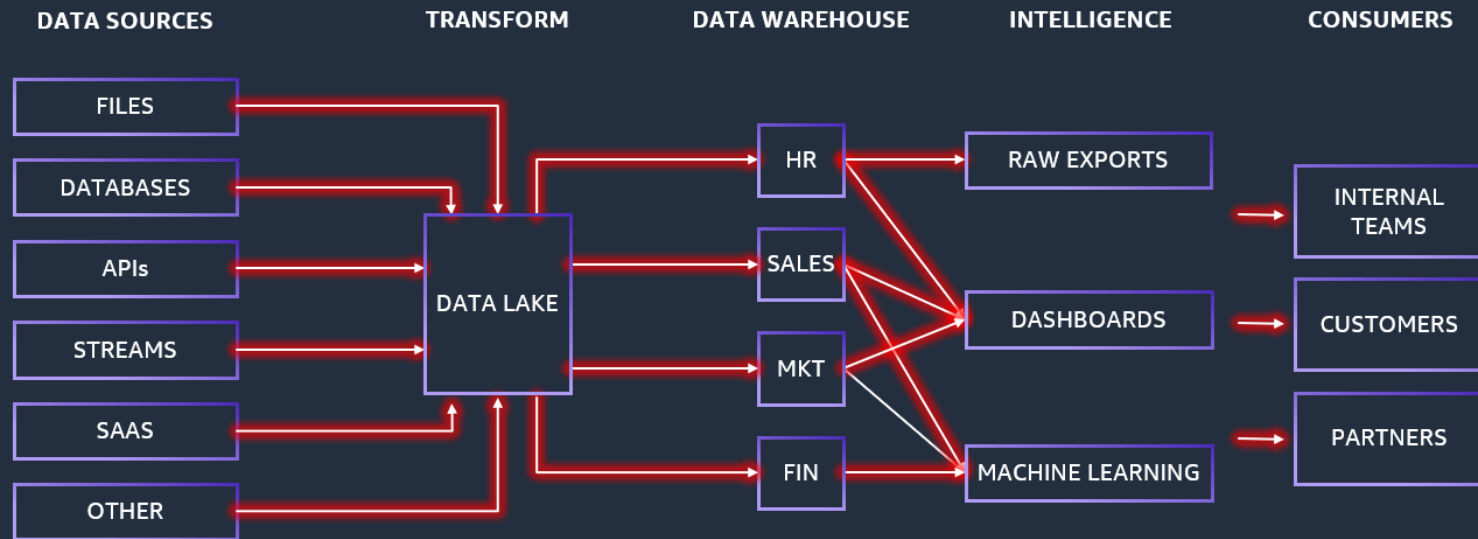/in/sandipanbhaumik

# Data lifecycle in an organization

# You need to build data pipelines, aka ETL

**DATA SOURCES**  **TRANSFORM**  **DATA WAREHOUSE**  **INTELLIGENCE**  **CONSUMERS**

FILES

DATABASES

APIs

STREAMS

SAAS

OTHER

DATA LAKE

HR

SALES

MKT

FIN

RAW EXPORTS

DASHBOARDS

MACHINE LEARNING

INTERNAL TEAMS

CUSTOMERS

PARTNERS

# Building ETL is complex



DATA SOURCES     TRANSFORM     DATA WAREHOUSE     INTELLIGENCE     CONSUMERS

FILES

DATABASES

APIs

STREAMS

SAAS

OTHER

DATA LAKE

HR

SALES

MKT

FIN

RAW EXPORTS

DASHBOARDS

MACHINE LEARNING

INTERNAL TEAMS

CUSTOMERS

PARTNERS

Manage connections, secure them

Handle failure scenarios

Manage ETL infrastructure

Additional monitoring processes
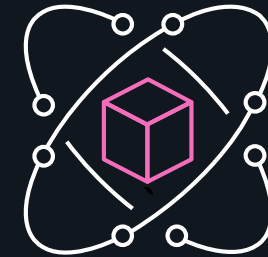
Write complex code

Spend time on overheads

# So, to make it easier for you we are building...



**Zero ETL
future** enabled
by service integrations

**AWS Glue** for
value-add data
transformations and more

Services for
**connecting to 100s
of data sources**

# AWS Glue: Key capabilities

## SERVERLESS DATA INTEGRATION SERVICE

### Scalable data integration engine

Built-in data transforms

Execution engine

Monitor

### Centralized and unified data governance

AWS Glue Data Catalog

AWS Glue Data Quality

AWS Glue crawlers

AWS Lake Formation

### Connect and ingest data

AWS Glue connectors

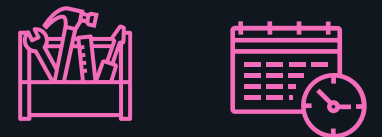AWS Glue connector marketplace

Various interfaces

### User productivity and data ops
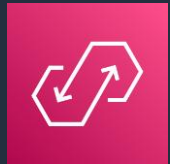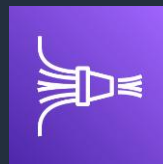
Persona specific tools
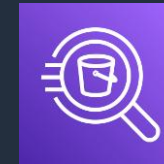
Productivity tools

Data ops tools

# Services for connecting to 100s of data sources

Connect to **50+** SaaS applications with **Amazon AppFlow**
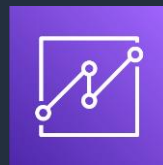
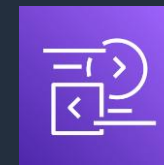Stream data in real time from **30+** sources with **Amazon Kinesis Data Firehose**

Query **25+** data sources in place with **Amazon Athena**

Build models on **Amazon SageMaker** using data from **40+** sources

**30+** sources to build interactive dashboards in **Amazon QuickSight**

Access third-party data from **300+** data providers with **AWS Data Exchange**

# Connect to 100s of data sources

Adobe Analytics · ahana · AMPLITUDE · ATLASSIAN Confluence · box · Cloudant · CLOUDERA · Cloudera IMPALA · Collibra · Coralogix

databricks · DATADOG · DELTA LAKE · dremio · Dropbox · dynatrace · Exasol · freshdesk · honeycomb.io · HubSpot

Informatica · LogicMonitor · MariaDB · MarkLogic · MongoDB · new relic · OKĒRA · privacera · redis · Sage Intacct

salesforce · SAP · TWILIO SendGrid · servicenow · singular · slack · snowflake · splunk> · Starburst · stripe

sumo logic · +ableau · teradata. · TREND MICRO · upsolver · Veeva · VERTICA · workday. · zaloni · zendesk

INGESTION & INTEGRATION | ANALYTICS | AI & MACHINE LEARNING | BUSINESS INTELLIGENCE | GOVERNANCE

aws

# Zero ETL future enabled by service integrations

**DATA SOURCES**



FILES

DATABASES

APIs

STREAMS

SAAS

OTHER

Amazon Simple Storage Service
(Amazon S3)

**DATA LAKE**

AWS Glue
Data Catalog

*crawlers*

*SQL/Spark*

Amazon Athena

# Zero ETL future enabled by service integrations

**DATA SOURCES**

**DATA WAREHOUSE**

FILES

DATABASES

APIs

STREAMS

SAAS

OTHER

*Data sharing*

Amazon Redshift

*Data sharing*

Amazon Redshift

*Data sharing*

Amazon Redshift
(Serverless)

*Auto copy (preview)*

Amazon Simple Storage Service
(Amazon S3)

**DATA LAKE**

# Continuous file ingestion from Amazon S3 (preview)

DATA SOURCES

DATA WAREHOUSE

**AWS** **VPC**

## Automatic data ingestion from Amazon S3 using Copy Jobs

### External sources

Data Source 1(RDBMS)
(SQL-Oracle-PG, My SQL)

Data Source 2 (Flat Files)
(CSV, JSON)

**Extract**

Data Source N (Applications)

Bucket -1  Folder 1 -
File1... N **(CSV)**

**Copy Job1**

Bucket -1  Folder 2 -
File1... N **(Json)**

**Copy-job 2**

Bucket -2  Folder 1 -
File1... N **(CSV)**

**Copy Job 3**

**Table 1**

**Table 2**

**Amazon S3- Landing Zone**

**Amazon Redshift-
Data Warehouse**

Amazon Simple Storage
(Amazon S3)

OTHER

**DATA LAKE**

```
COPY public.target_table
FROM 's3://mybucket-bucket/staging-folder'
IAM_ROLE 'arn:aws:iam::123456789012:role/MyLoadRoleName'
JOB CREATE my_copy_job_name
AUTO ON;
```

# Zero ETL future enabled by service integrations

**DATA SOURCES**


Amazon Relational Database
Service (Amazon RDS)

*federated querying*

**DATA W**


Amazon Redshift

*Auto copy (previe*


Amazon Simple Storage Service
(Amazon S3)

**DATA LAKE**

1. **Create credentials for accessing database**

- DB instance with user name and password authentication
- Amazon Redshift cluster with a cluster maintenance version that supports federated queries.

**2. Create External Schema**

```
CREATE EXTERNAL SCHEMA amysql
FROM MYSQL
DATABASE 'functional'
URI 'endpoint to remote hostname'
IAM_ROLE 'arn:aws:iam::              .role/Redshift-SecretsManager-RO'
SECRET_ARN 'arn:aws:secretsmanager:          ..        .    :secret:fede
```

**3. Query table**

```
SELECT level FROM amysql.employees LIMIT 1;

  level
-------
     8
```

*Documentation: Querying data with*
*federated queries in Amazon Redshift*

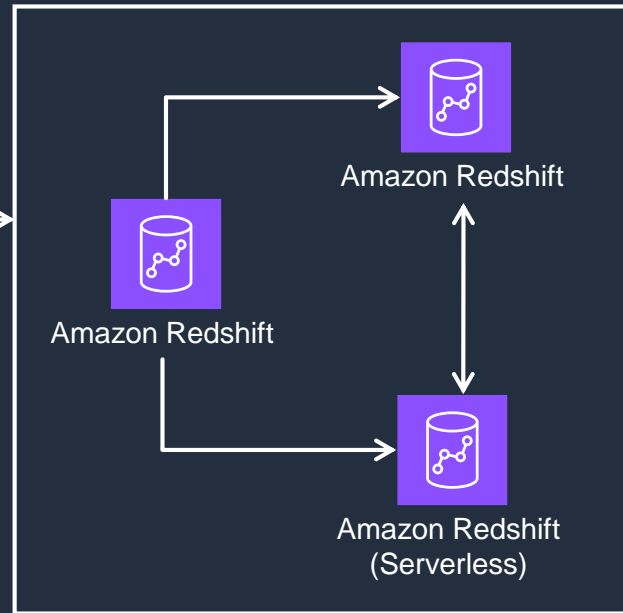# Zero ETL future enabled by service integrations

**DATA SOURCES**

**DATA WAREHOUSE**

Amazon Aurora

*Zero-ETL integration*

Amazon Redshift

Amazon Redshift

Amazon Redshift
(Serverless)

Amazon Simple Storage Service
(Amazon S3)

**DATA LAKE**

# Transaction analytics at scale



Requires building and managing complex data pipelines

Transactional Applications

Analyst

Analytics Applications

Amazon Aurora

DMS

Amazon S3

AWS Glue

Amazon EMR

Amazon S3

Amazon Redshift

Analytics

Analyst

Data engineer

# Amazon Aurora zero-ETL integration with Amazon Redshift



Transactional Applications

Analyst

Analytics Applications

Amazon Aurora

Zero-ETL integration

Amazon Redshift

Analytics

Analyst

Aurora storage

Redshift storage

Seed once

CDC streaming
(continuous change data capture)

# Zero ETL future enabled by service integrations

# Near-real-time analytics using Amazon Redshift streaming ingestion

**Amazon DynamoDB** → **Amazon Kinesis** → *Streaming ingestion* →

1. **Create External Schema**

```sql
CREATE EXTERNAL SCHEMA demo_schema
FROM KINESIS
IAM_ROLE { default | 'iam-role-arn' };
```

2. **Create Materialized View**

```sql
CREATE MATERIALIZED VIEW demo_stream_vw AS
    SELECT approximate_arrival_timestamp,
    partition_key,
    shard_id,
    sequence_number,
    json_parse(kinesis_data) as payload
    FROM demo_schema."demo-data-stream";
```
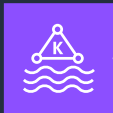
3. **Refresh Materialised View**

```sql
REFRESH MATERIALIZED VIEW demo_stream_vw;
```

*AWS  Blog: Near-real-time analytics using Amazon Redshift streaming ingestion with Amazon Kinesis Data Streams and Amazon DynamoDB*

# Zero ETL future enabled by service integrations

**DATA SOURCES**

**DATA WAREHOUSE**

1. **Create External Schema**

```
CREATE EXTERNAL SCHEMA MySchema
FROM MSK
IAM_ROLE { default | 'iam-role-arn' }
AUTHENTICATION { none | iam }
CLUSTER_ARN 'msk-cluster-arn';
```

2. **Create Materialised View**

```
CREATE MATERIALIZED VIEW MyView AUTO REFRESH YES AS
SELECT kafka_partition,
 kafka_offset,
 kafka_timestamp_type,
 kafka_timestamp,
 kafka_key,
 JSON_PARSE(kafka_value) as Data,
 kafka_headers
FROM MySchema."mytopic"
WHERE CAN_JSON_PARSE(kafka_value);
```
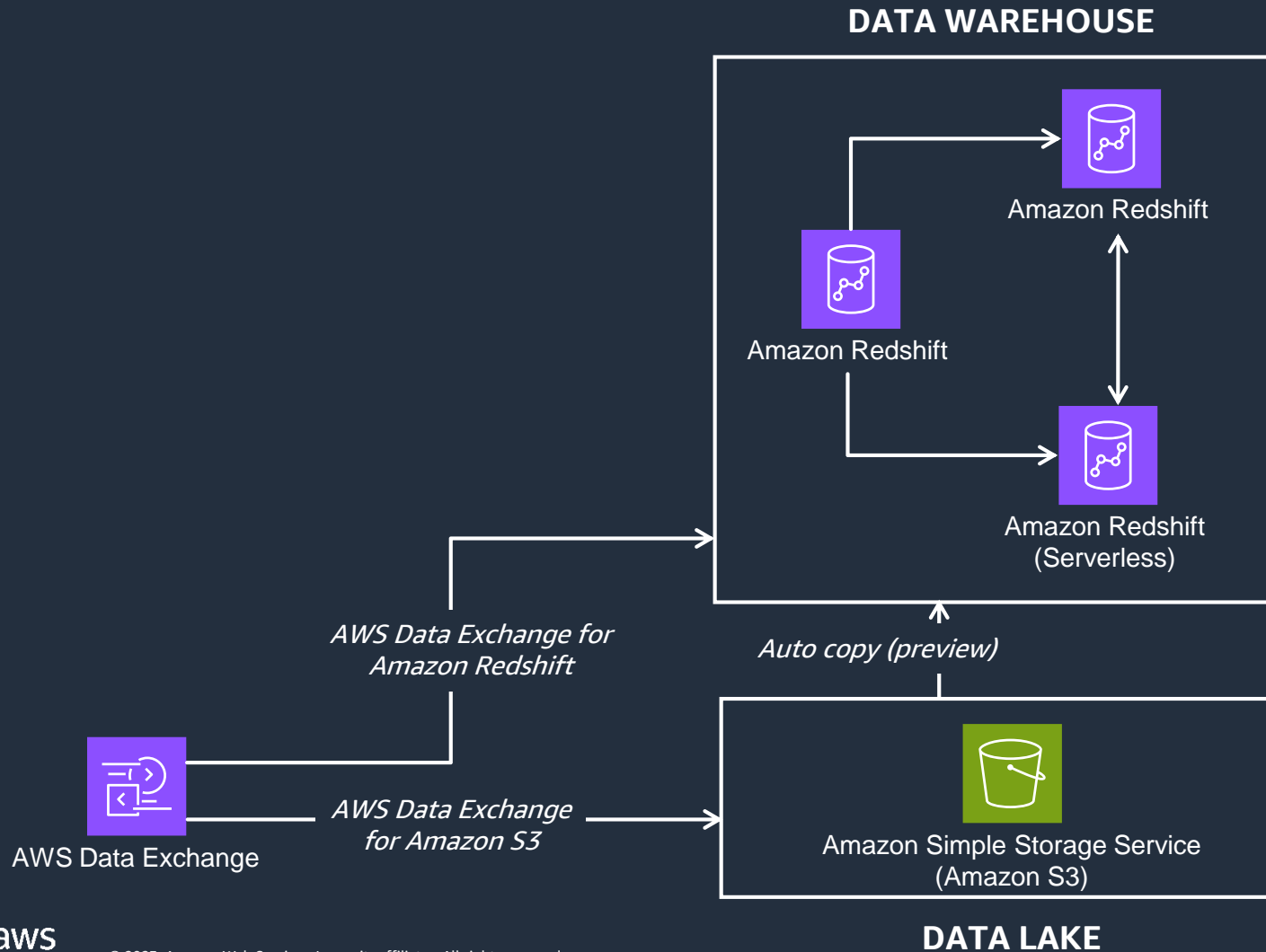
Amazon Managed Streaming for
Apache Kafka

—Streaming ingestion—

*Documentation: Getting started with streaming ingestion from Amazon Managed Streaming for Apache Kafka*

(Amazon S3)

**DATA LAKE**

aws

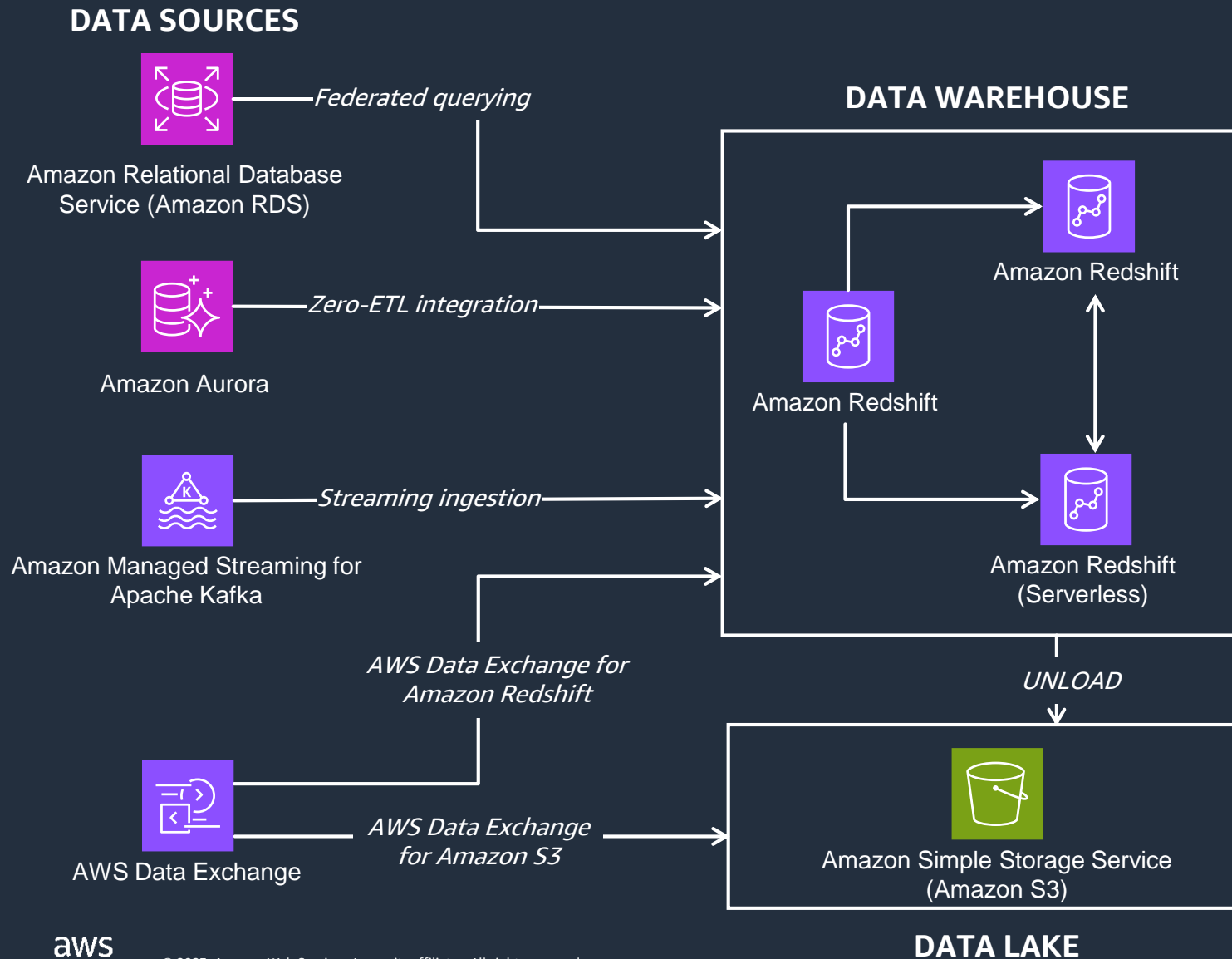# Find, subscribe to, and query third-party data without any ETL



With **AWS Data Exchange for Amazon S3**, subscribers directly query provider's S3 bucket.

Data subscribers can subscribe to and directly query **Amazon Redshift data sets** without writing ETL processes.
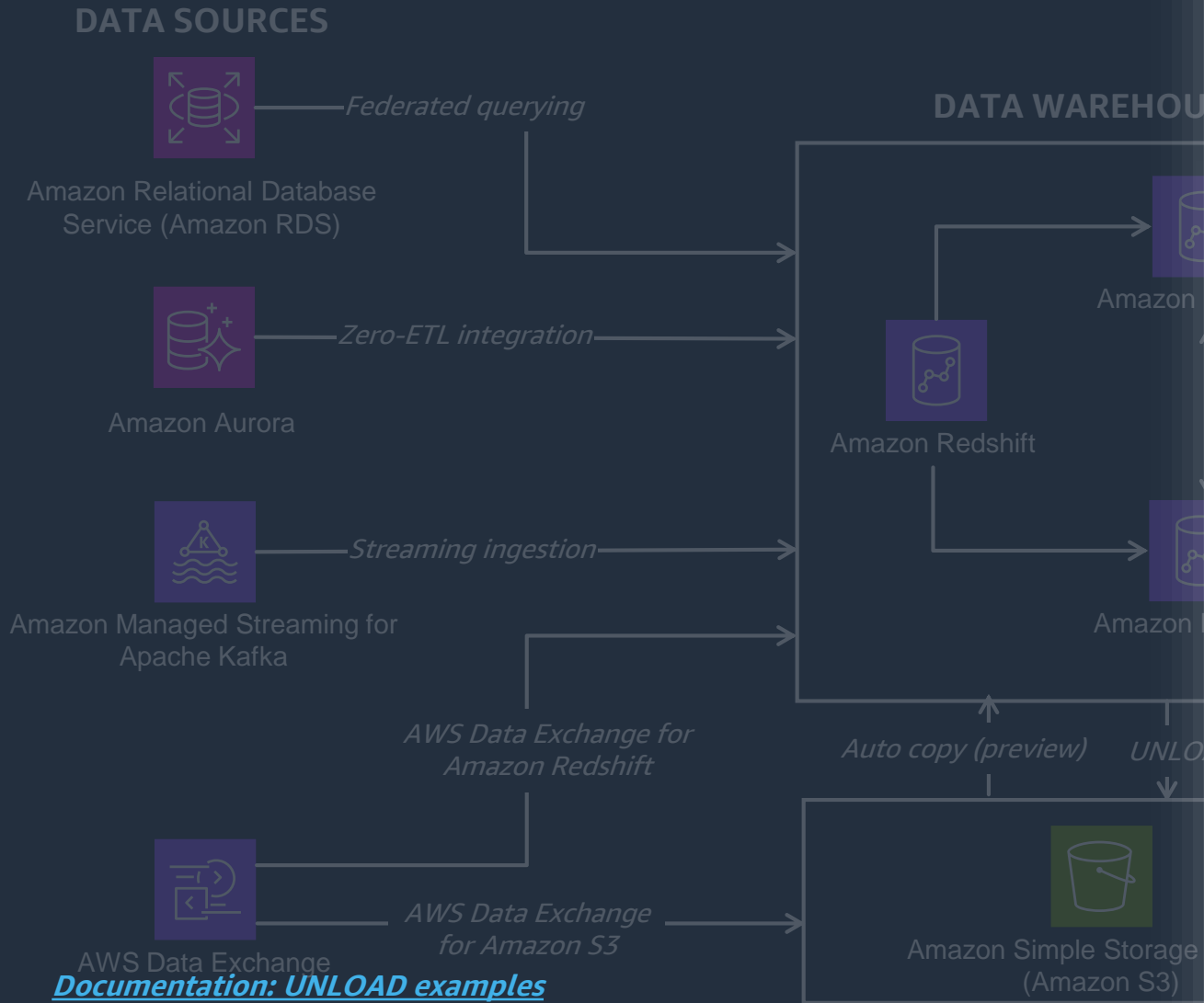
Using **AWS Data Exchange for Data APIs**, subscribers directly access data using their AWS IAM credentials and SDKs, no need for ETL pipelines.

With **AWS Data Exchange for AWS Lake Formation**, subscribers can share access to the data within their AWS account or across AWS organization.

# Zero ETL future enabled by service integrations

**DATA SOURCES**

**DATA WAREHOUSE**

Amazon Relational Database Service (Amazon RDS)

*Federated querying*

Amazon Aurora

*Zero-ETL integration*

Amazon Managed Streaming for Apache Kafka

*Streaming ingestion*

Amazon Redshift

Amazon Redshift

Amazon Redshift (Serverless)

*AWS Data Exchange for Amazon Redshift*

*UNLOAD*

AWS Data Exchange

*AWS Data Exchange for Amazon S3*

Amazon Simple Storage Service (Amazon S3)

**DATA LAKE**

# Unload data to your data lake

**DATA SOURCES**


Amazon Relational Database Service (Amazon RDS) — Federated querying


Amazon Aurora — Zero-ETL integration


Amazon Managed Streaming for Apache Kafka — Streaming ingestion

AWS Data Exchange for Amazon Redshift


AWS Data Exchange — AWS Data Exchange for Amazon S3

**DATA WAREHOUSE**

Amazon Redshift

Auto copy (preview)    UNLOAD

Amazon Simple Storage Service (Amazon S3)

### 1. UNLOAD file

```
unload ('select * from venue')
to 's3://mybucket/tickit/unload/venue_'
iam_role 'arn:aws:iam::0123456789012:role/MyRedshiftRole';
```

### 2. Delimited (default '|')

```
unload ('select * from venue')
to 's3://mybucket/tickit/venue/tab'
iam_role 'arn:aws:iam::0123456789012:role/MyRedshiftRole'
delimiter as '\t';
```

### 3. Create manifest

```
unload ('select * from venue')
to 's3://mybucket/tickit/venue_'
iam_role 'arn:aws:iam::0123456789012:role/MyRedshiftRole'
manifest;

{
  "entries": [
    {"url":"s3://mybucket/tickit/venue_0000_part_00"},
    {"url":"s3://mybucket/tickit/venue_0001_part_00"},
    {"url":"s3://mybucket/tickit/venue_0002_part_00"},
    {"url":"s3://mybucket/tickit/venue_0003_part_00"}
  ]
}
```

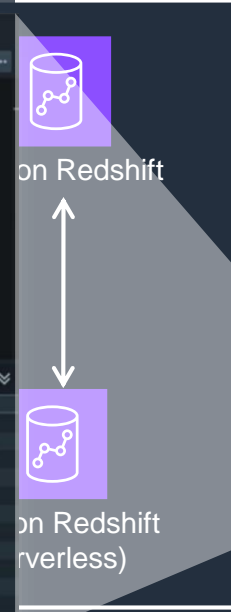# Zero ETL future enabled by service integrations



DATA SOURCES

Federated querying

DATA WAREHOUSE
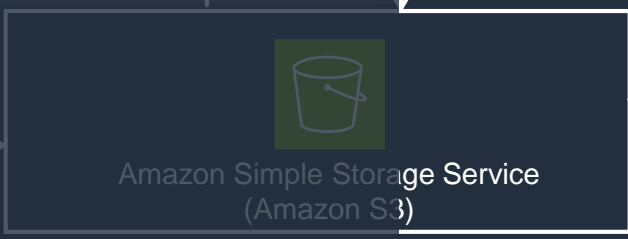
Amazon Redshift

Amazon Redshift
(Serverless)

Amazon Redshift

Auto copy (preview)    UNLOAD

AWS Data Exchange
for Amazon S3

AWS Glue
Data Catalog

crawlers    SQL/Spark

Amazon Simple Storage Service
(Amazon S3)

DATA LAKE

Amazon Athena

Zero ETL future enabled by service integrations

# Zero ETL future enabled by service integrations

**DATA SOURCES**

**CONSUMERS**

Amazon Relational Database Service (Amazon RDS)

*Federated querying*

Amazon Aurora

*Zero-ETL integration*

**DATA WAREHOUSE**

Amazon Redshift

Amazon Redshift

Amazon Managed Streaming for Apache Kafka

*Streaming ingestion*

Amazon Redshift (Serverless)

*AWS Data Exchange for Amazon Redshift*

*Spark native integration*

Amazon EMR

AWS Glue Data Catalog

Amazon QuickSight

*Auto copy (preview)*

*UNLOAD*

AWS Data Exchange

*AWS Data Exchange for Amazon S3*

Amazon Simple Storage Service (Amazon S3)

Amazon SageMaker

*crawlers*

*SQL/Spark*

Amazon Athena

**DATA LAKE**

# Connect all your data with native integrations



**DATA SOURCES**

**CONSUMERS**

**DATA WAREHOUSE**

Amazon Relational Database Service (Amazon RDS) — *Federated querying*

Amazon Aurora — *Zero-ETL integration*

Amazon Managed Streaming for Apache Kafka — *Streaming ingestion*

*AWS Data Exchange for Amazon Redshift*

AWS Data Exchange — *AWS Data Exchange for Amazon S3*

Amazon Redshift

*Data sharing* — Amazon Redshift

*Data sharing*

*Data sharing* — Amazon Redshift (Serverless)

*Redshift Data API* — Web Applications

*Spark native integration* — Amazon EMR

*Automatic mounting*

AWS Glue Data Catalog

Amazon QuickSight

Amazon SageMaker

Amazon Athena

*Auto copy (preview)*

*UNLOAD*

Amazon Simple Storage Service (Amazon S3)
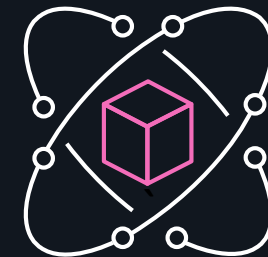
**DATA LAKE**

*crawlers*

*SQL/Spark*

# We are making data integration easier for you

**Zero ETL
future** enabled
by service integrations

**AWS Glue** for
value-add data
transformations and more

Services for
**connecting to 100s
of data sources**

# We are here to help

## Want to discover possibilities?

**aws** Well-architected Review

- ✓ Assess your workloads and learn how well your architecture aligns with cloud best practices and gain guidance for making improvements
- ✓ Use customer analytics lens to review your data workloads

Review and refine existing workloads

## Want to build a data vision and strategy?

**aws** Data Driven Everything

- ✓ Create an organizational vision for innovation with data to drive business outcomes
- ✓ Define the first pilot, learn, and build

Jump-start the data flywheel

## Want to modernize your data foundation?

**aws** Re:Imagine data

- ✓ Define the migrate & modernization strategy for a future data foundation
- ✓ Lower cost, increase capacity, and unlock business access

Migrate & Modernize data

# Thank you!

Sandipan Bhaumik

/in/sandipanbhaumik