

Overview Big Data on



Which tools should I use?

Let's solve the first question that might come to your mind: what's the right tools for building that pipeline? And the answer I found while building mine was:

“ There is not a right tool or architecture, it will always depend on your needs! “

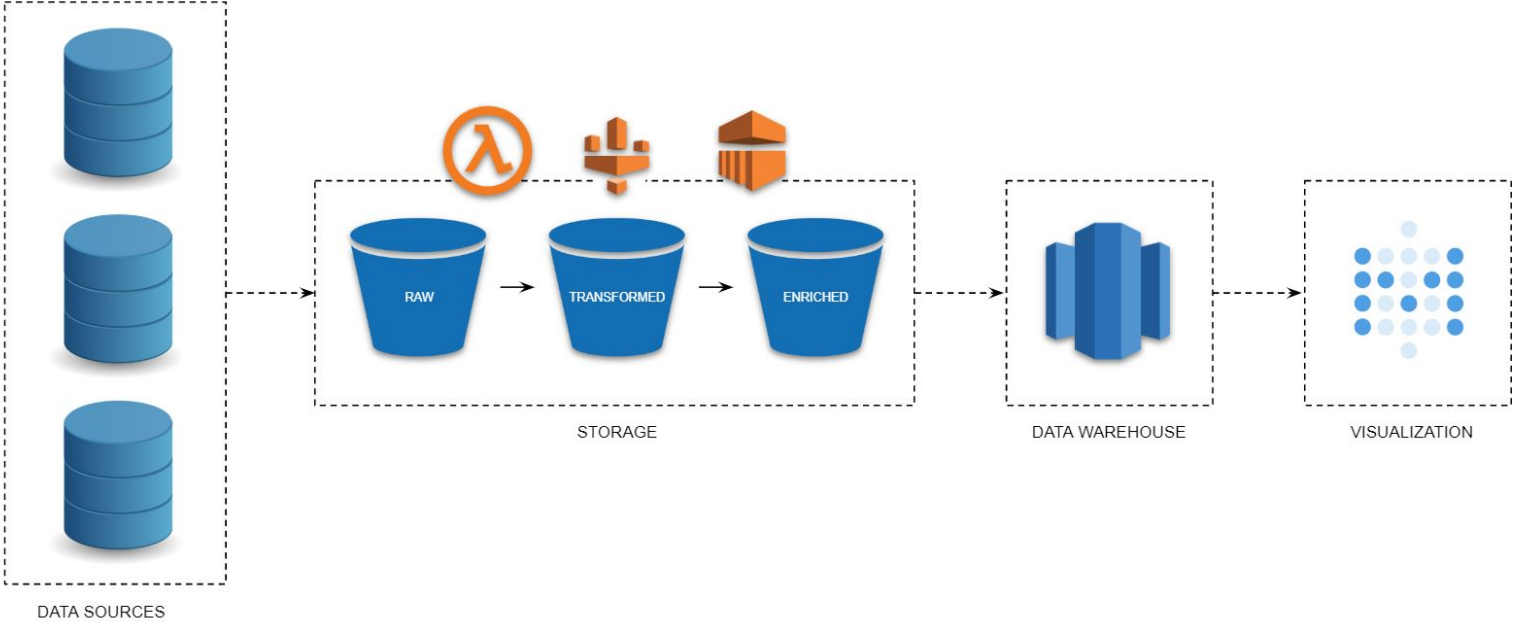
The challenge

By the time I got into the company, there was a big problem: the data was too isolated. Analyzing data was too slow and difficult that people could not find the motivation to do it.

And the challenge was: centralize that data and promote **data democratization** on the company in order to **empower** people! A big challenge, right?

The pipeline

The proposed pipeline architecture to fulfill those needs is presented on the image bellow, with a little bit of improvements that we will be discussing.

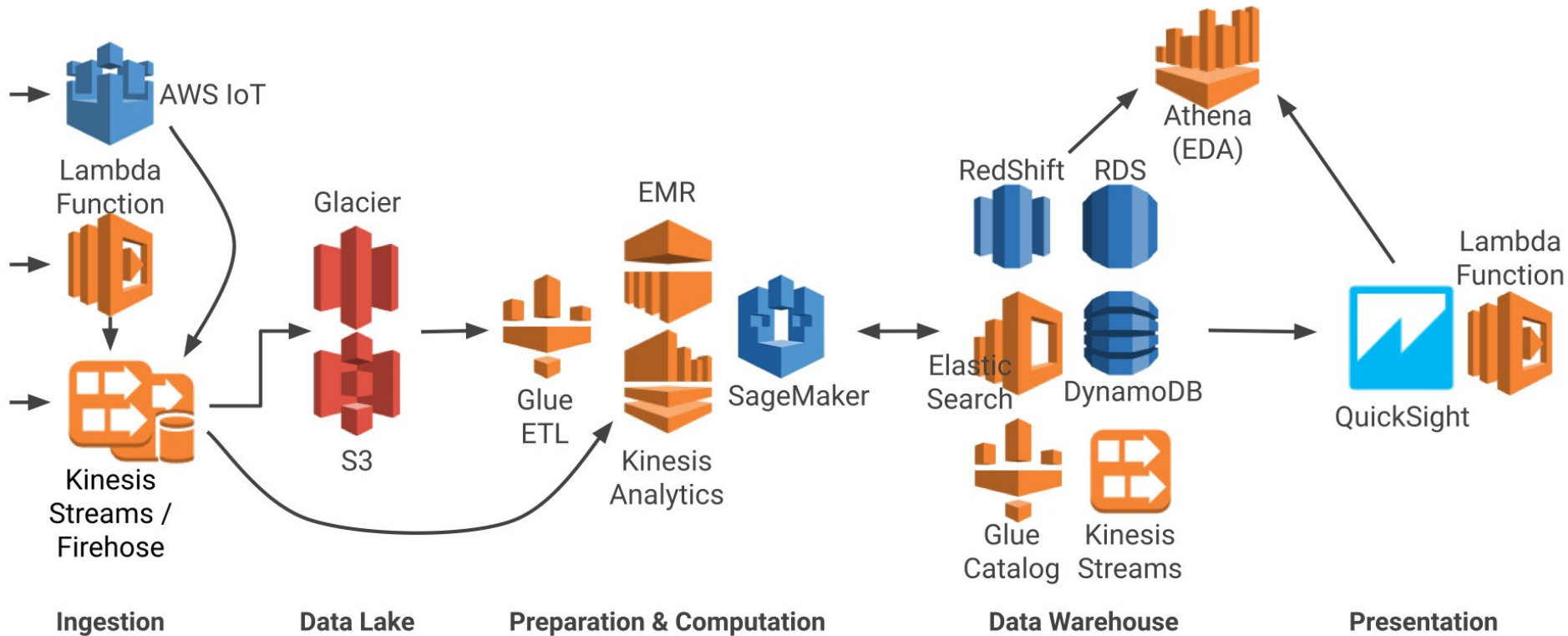


Category	Use cases	AWS service
Analytics	Interactive analytics Big data processing Data warehousing Real-time analytics Operational analytics Dashboards and visualizations Visual data preparation	<u>Amazon Athena</u> <u>Amazon EMR</u> <u>Amazon Redshift</u> <u>Amazon Kinesis</u> <u>Amazon Elasticsearch Service</u> <u>Amazon Quicksight</u> <u>AWS Glue DataBrew</u>

Category	Use cases	AWS service
Data movement	Real-time data movement	<u>AWS Glue</u> <u>Amazon Managed Streaming for Apache Kafka (MSK)</u> <u>Amazon Kinesis Data Streams</u> <u>Amazon Kinesis Data Firehose</u> <u>Amazon Kinesis Video Streams</u>

Category	Use cases	AWS service
Data lake	Object storage	<u>Amazon S3</u> <u>AWS Lake Formation</u>
	Backup and archive	<u>Amazon S3 Glacier</u> <u>AWS Backup</u>
	Data catalog	<u>AWS Glue</u> <u>AWS Lake Formation</u>
	Third-party data	<u>AWS Data Exchange</u>

Category	Use cases	AWS service
Predictive analytics and machine learning	Frameworks and interfaces Platform services	<u>AWS Deep Learning AMIs</u> <u>Amazon SageMaker</u>



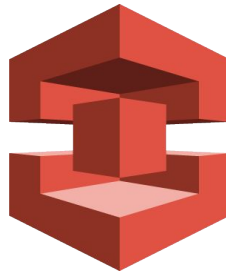
Data Movement



**AWS Direct
Connect**



**AWS
Database
Migration
Service**

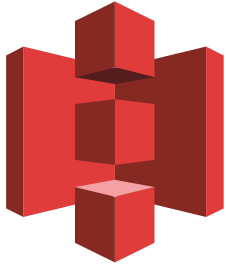


**AWS
Import/Export
& Snowball**



**AWS
Storage
Gateway**

Storage and Databases



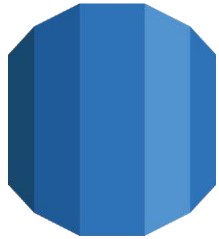
Amazon S3

- Store unlimited number of objects
- Designed for 99.999999999% durability
- As Data Lake with integration with other AWS services (Amazon Kinesis, Amazon Redshift, Amazon EMR, etc.)
- Low cost with tiered-storage (Standard, IA, Amazon Glacier) via life-cycle policy
- Secure – SSL, client/server-side encryption at rest



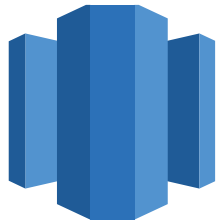
**Amazon
DynamoDB**

- Fully Managed NoSQL Database
- Fast consistent performance (single-digit millisecond latency at any scale)
- Highly scalable - automatic scaling of throughput capacity
- Highly available and durability
- Store unlimited number of data



**Amazon
Aurora**

- Fully Managed Relational Database Service
- MySQL and PostgreSQL compatible relational database with up to 5x better performance running on the same hardware
- Security, availability, and reliability of commercial databases at 1/10th the cost
- Designed to offer greater than 99.99% availability.
- Automatically grows storage as needed, from 10GB up to 64TB
- Achieve up to 500,000 reads and 100,000 writes per second



**Amazon
Redshift**

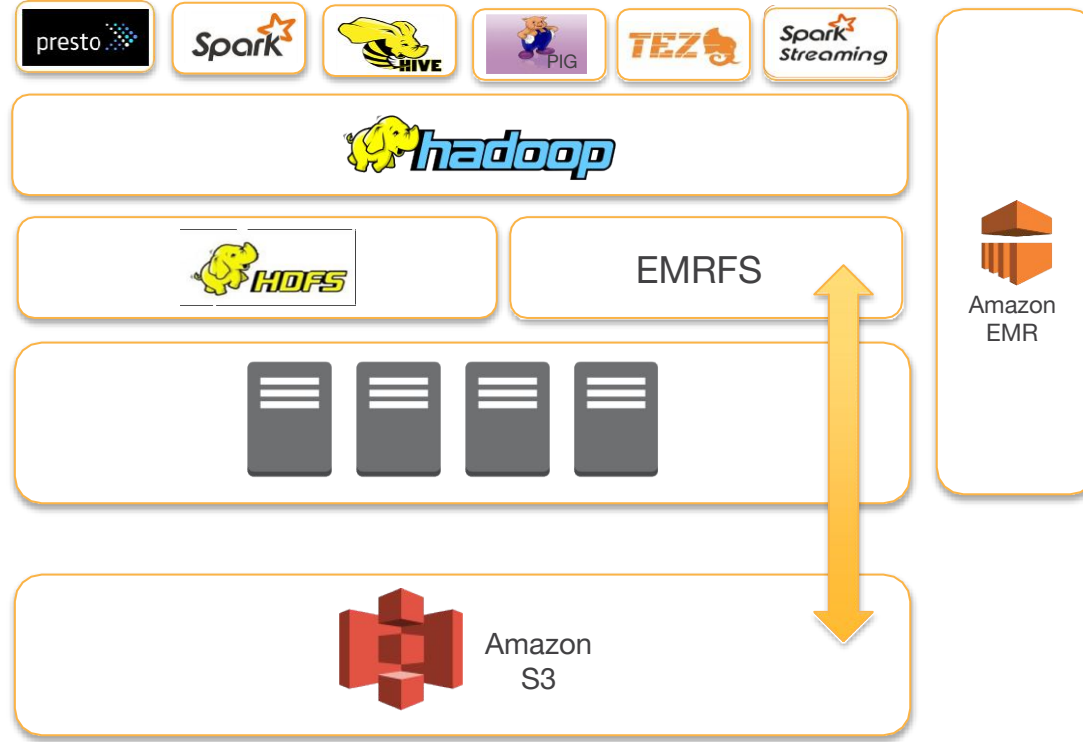
- Fully managed petabyte-scale relational, MPP, data warehousing
- Built-in end-to-end security, including SSL connections and cluster encryption
- Fault-tolerant - automatically recovers from disk and node failures
- Data automatically backed up to Amazon S3
- \$1,000/TB/Year; start at \$0.25/hour. Provision in minutes; scale from 160 GB to 2 PB of compressed data with just a few clicks

Analytic Frameworks



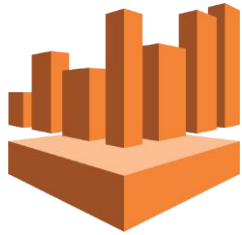
- Managed Hadoop framework
- Apache Hadoop, Hive, Spark, Zeppelin, Presto, HBase, Phoenix, Tez, Flink, etc.
- Auto Scaling clusters with support for on-demand and spot pricing
- Support for end-to-end encryption, IAM/VPC, S3 client-side encryption with customer managed keys and AWS KMS
- Integrates with Amazon S3, Amazon DynamoDB, Amazon Kinesis and Amazon Redshift

Amazon EMR





- Fully managed, reliable, and scalable Elasticsearch service
- Support for ELK
- Integration options with other AWS services (CloudWatch Logs, Amazon DynamoDB, Amazon S3, Amazon Kinesis)
- Use Case: log analytics, full text search, application monitoring, and more.



**Amazon
Athena**

- Serverless query service for querying data in S3 using standard SQL with no infrastructure to manage
- Support for multiple data formats include text, CSV, TSV, JSON, Avro, ORC, Parquet
- Pay per query only when you're running queries based on data scanned. If you compress your data, you pay less and your queries run faster

Familiar Technologies Under the Covers



Used for SQL Queries

In-memory distributed query engine

ANSI-SQL compatible with

extensions

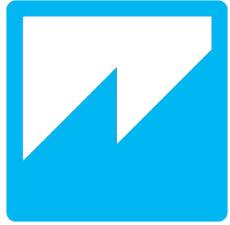


Used for DDL functionality

Complex data types

Multitude of formats

Supports data partitioning



**Amazon
QuickSight**

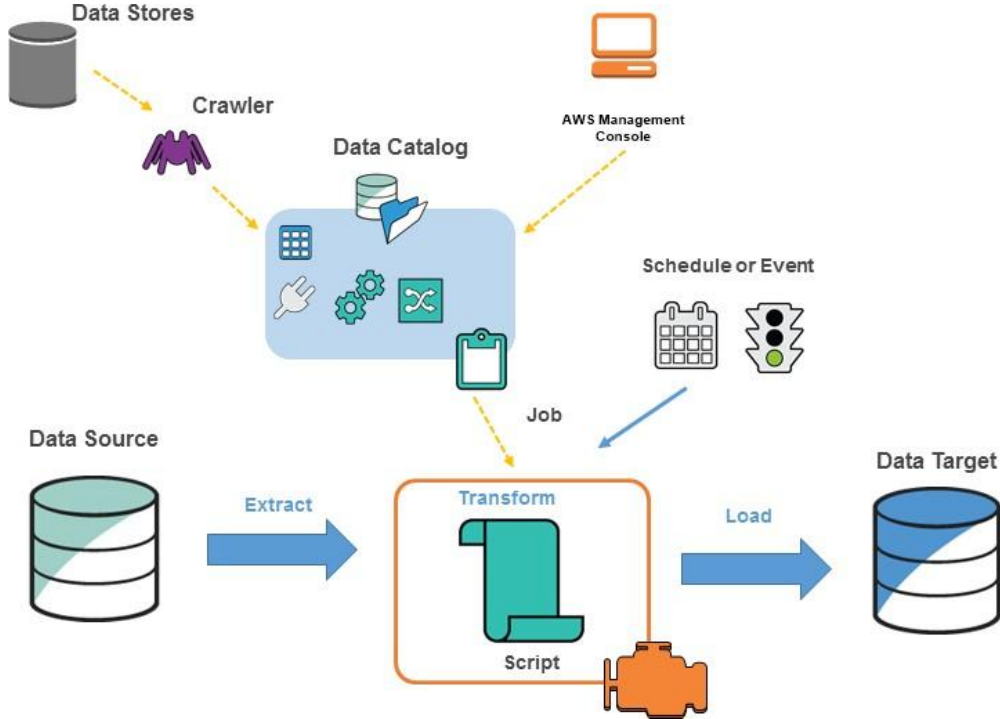
- Fast and cloud-powered Business Analytics
- Easy to use, no infrastructure to manage
- Quick calculations with SPICE
- 1/10th the cost of legacy BI software
- Accessed from any browser or mobile device



AWS Glue

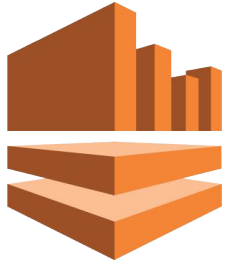
- Fully managed ETL (extract, transform, load) service
- Integrated data catalog, automatic schema discovery, ETL code generation, flexible job scheduler
- Integrated across a wide range of AWS services (Amazon RDS, Database running on Amazon EC2, Amazon Athena, etc.)

How AWS Glue Works



1. Build your data catalog
2. Generate and Edit Transformations
3. Schedule and Run Your Jobs

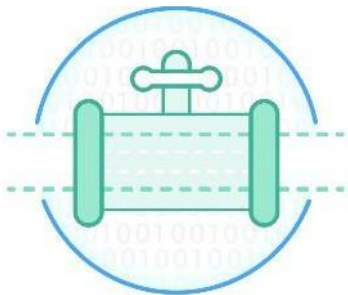
Real-time Analytics



**Amazon
Kinesis**

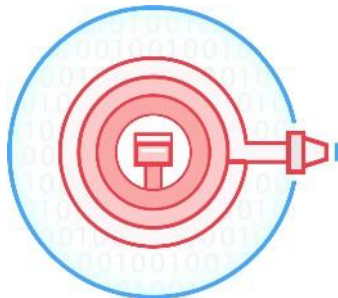
- Fully managed streaming application
- Scalable – handle any amount of streaming data
- Ingest, buffer and process data in real-time
- React quickly – derive insight in seconds

Amazon Kinesis



Amazon Kinesis Streams

Build your own custom applications that process or analyze streaming data



Amazon Kinesis Firehose

Easily load massive volumes of streaming data into Amazon S3, Amazon Redshift, and Amazon Elasticsearch

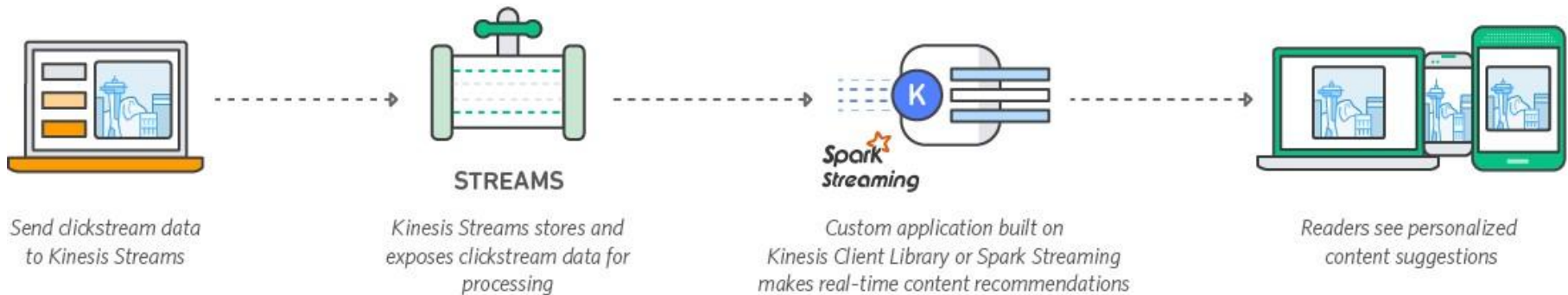


Amazon Kinesis Analytics

Easily analyze data streams using standard SQL queries

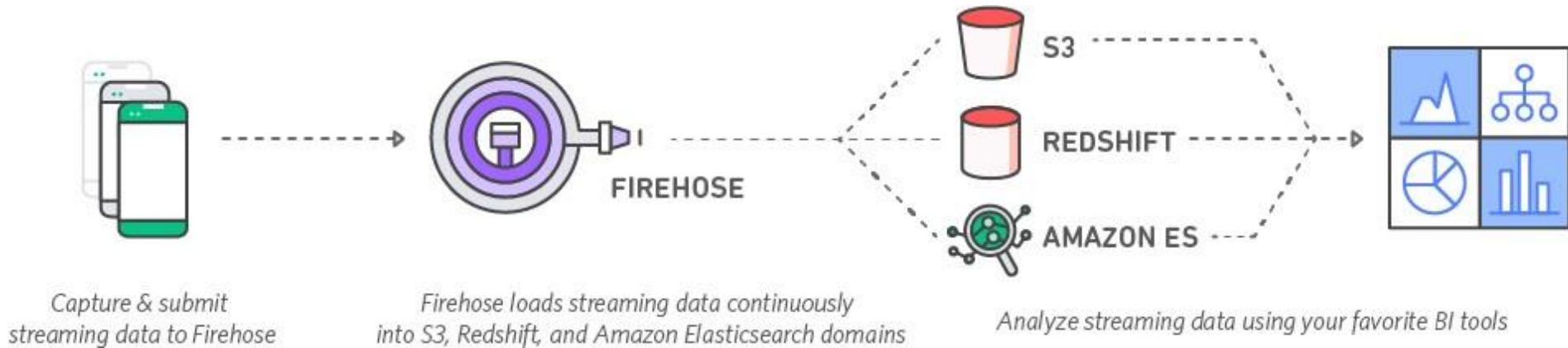
Amazon Kinesis Streams

- Reliably **ingest and durably store** streaming data at low cost
- Build **custom real-time applications** to process streaming data



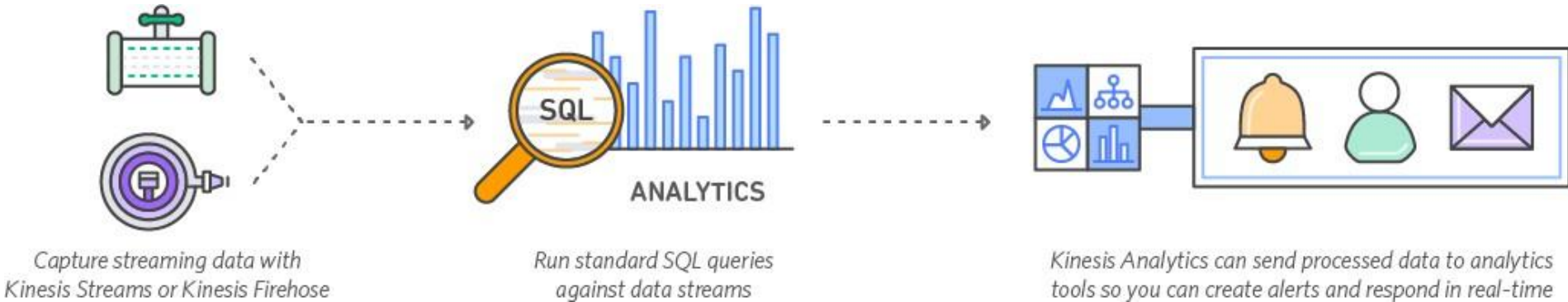
Amazon Kinesis Firehose

Reliably **ingest** and **deliver** **batched, compressed, and encrypted data** to **S3, Amazon Redshift, and Amazon Elasticsearch Service**



Amazon Kinesis Analytics

Interact with streaming data in real time using SQL



AWS Market Place for Big Data Solution



Hundreds of big data products are immediately available through the AWS marketplace

Database and Data Enablement



Advanced Analytics



Business Intelligence

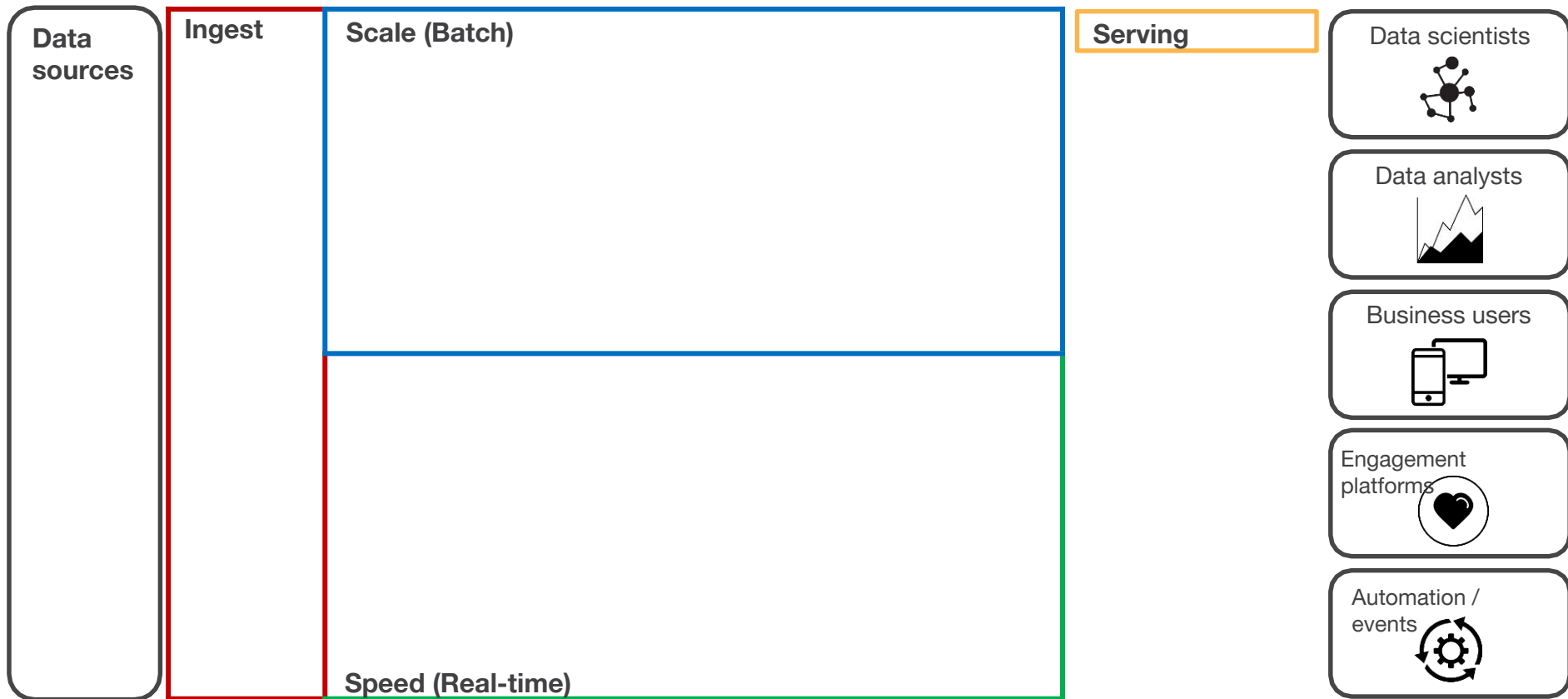


Fully Integrated | 1-click deployment |
Pay-as-you-go pricing

Modern Data Analytics Architecture on AWS

Modern data architecture

Insights to enhance business applications, new digital services



Modern data architecture

Insights to enhance business applications, new digital services

Data sources

Transactions



ERP



Web logs / cookies



Connected devices



Social media



Ingest

Scale (Batch)

Serving

Data scientists



Data analysts



Business users



Engagement platforms



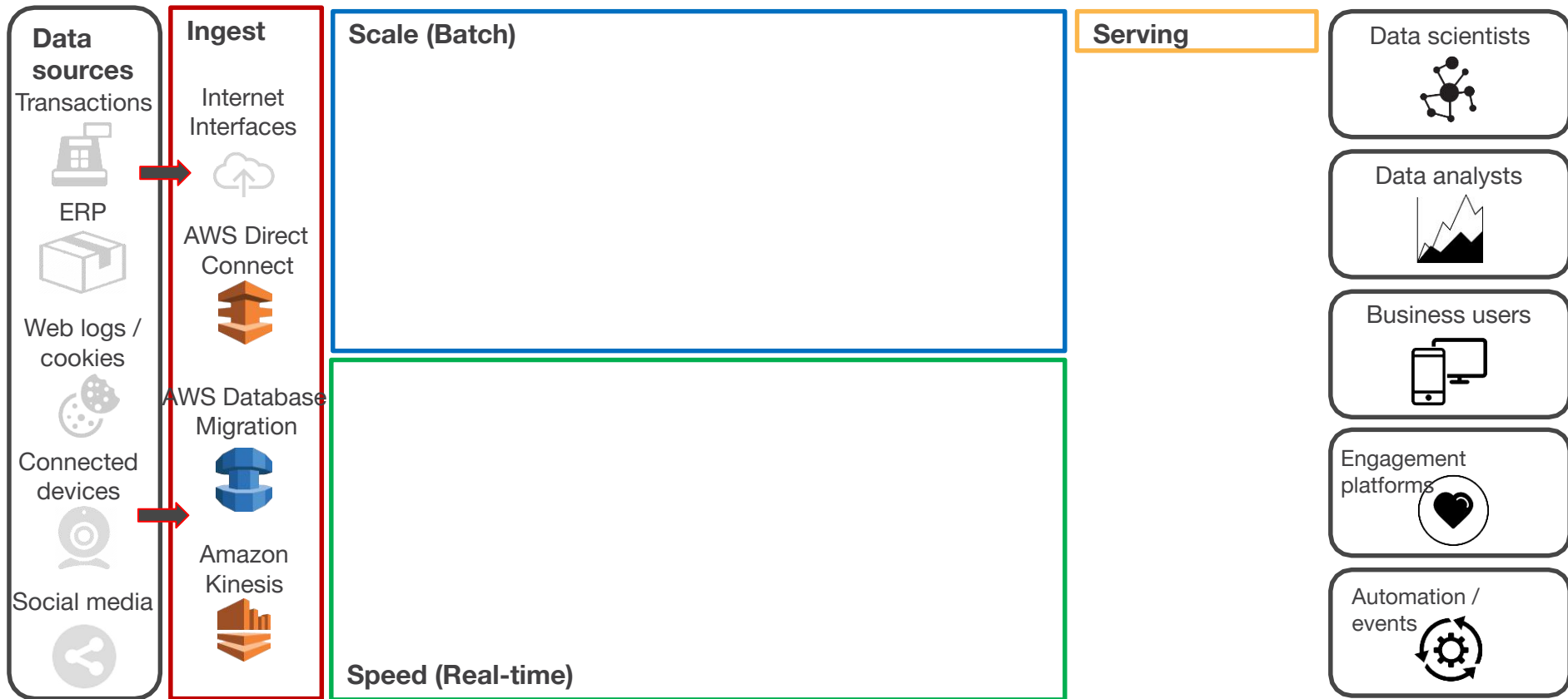
Automation / events



Speed (Real-time)

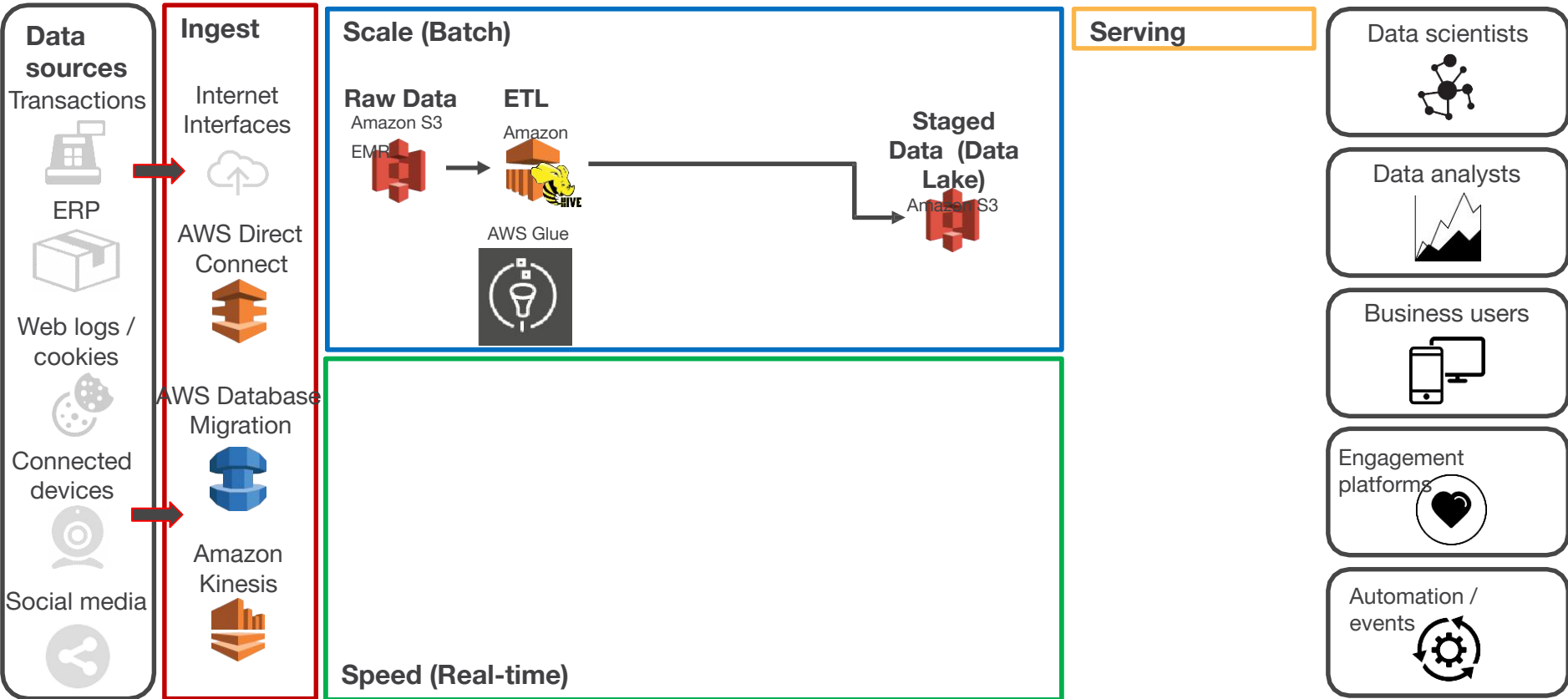
Modern data architecture

Insights to enhance business applications, new digital services



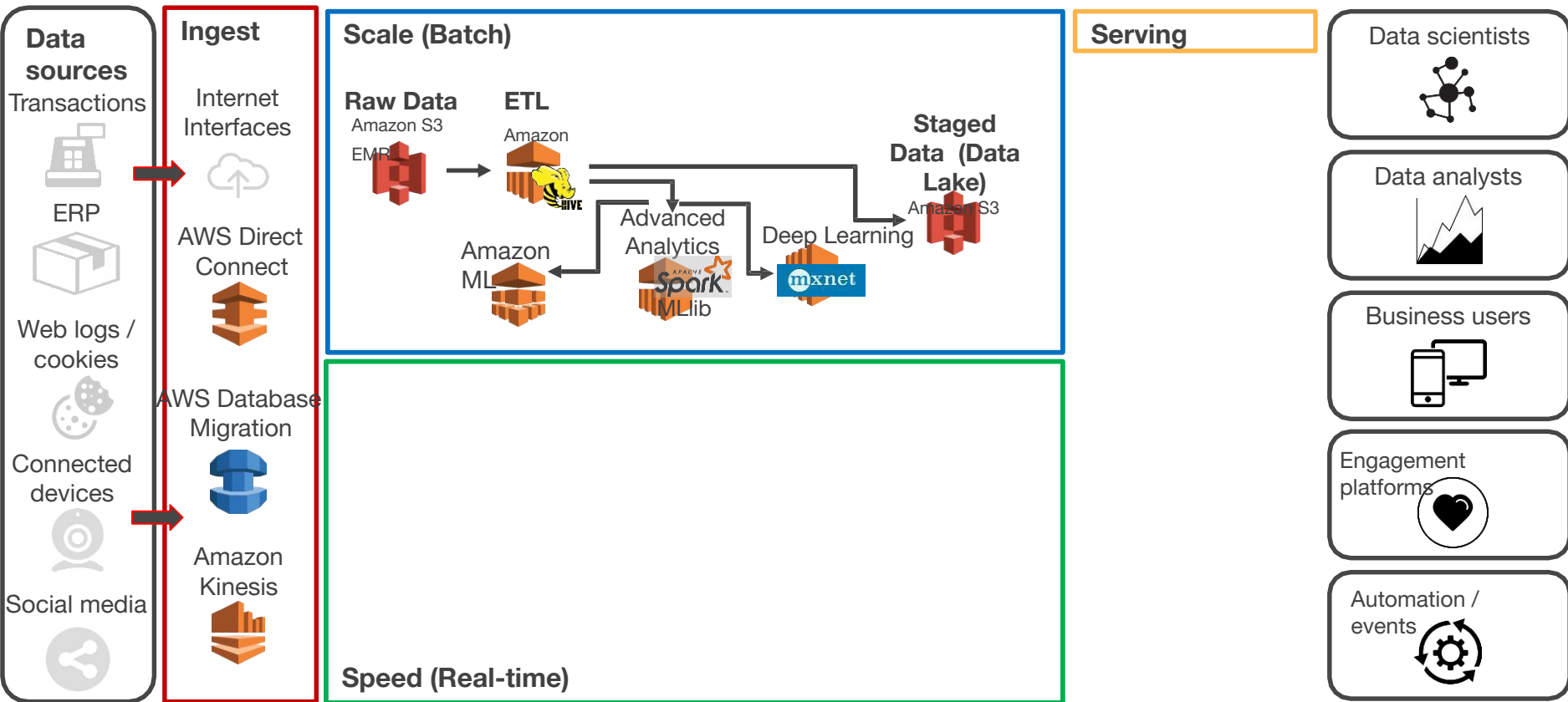
Modern data architecture

Insights to enhance business applications, new digital services



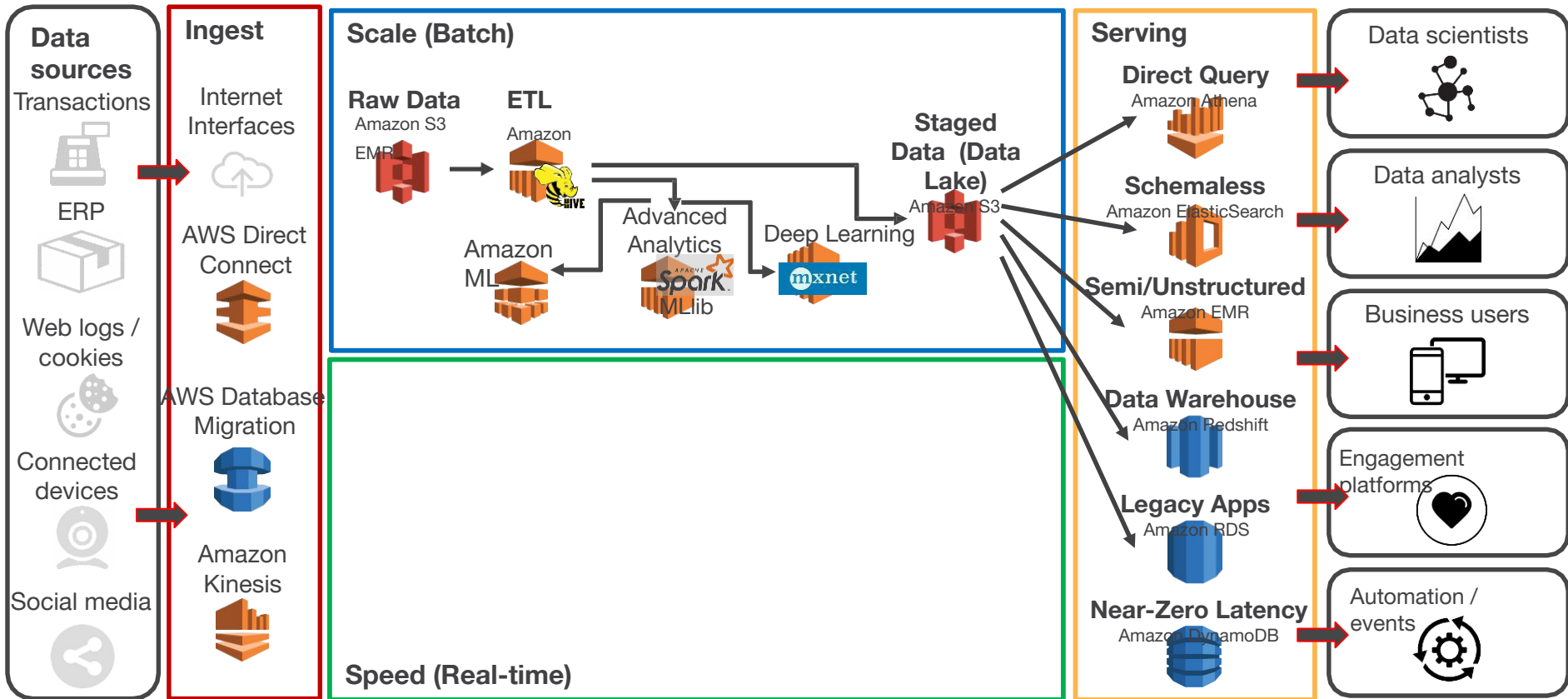
Modern data architecture

Insights to enhance business applications, new digital services



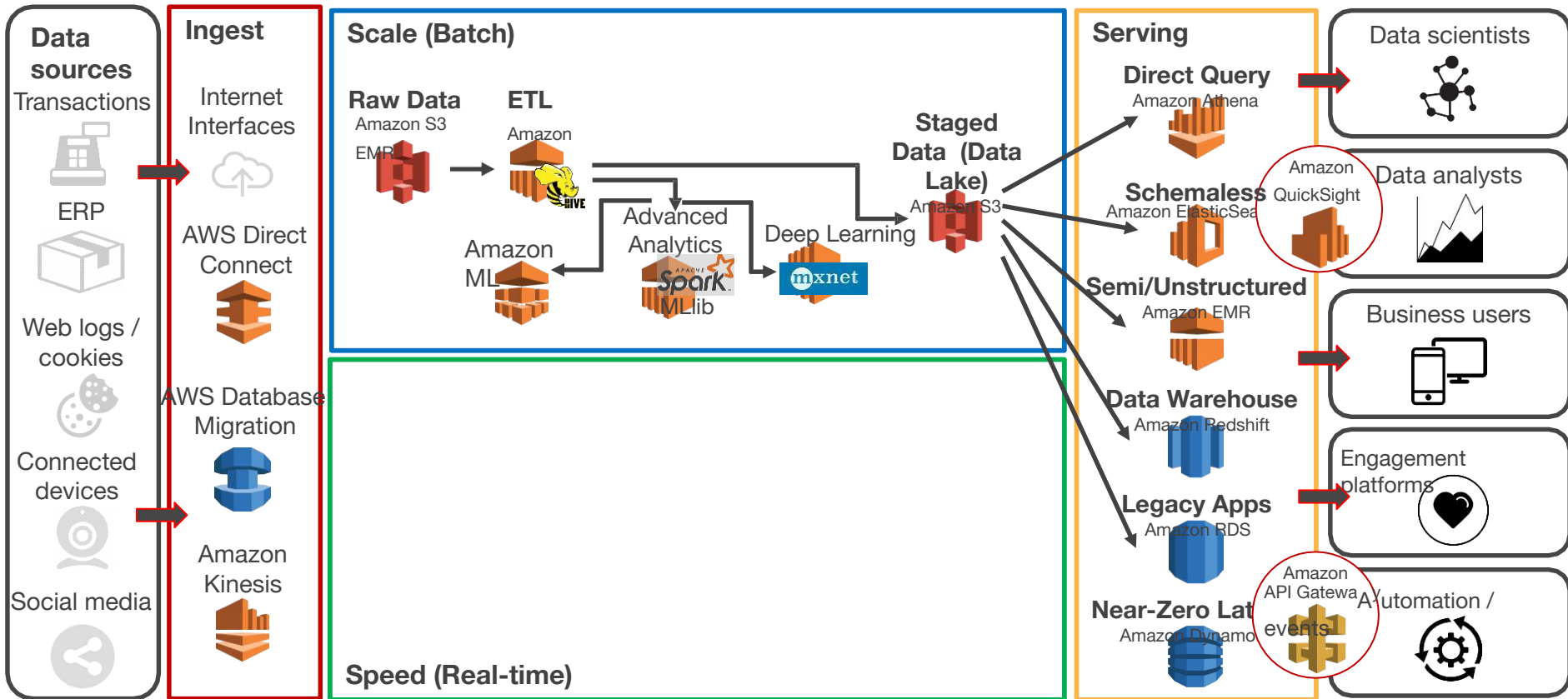
Modern data architecture

Insights to enhance business applications, new digital services



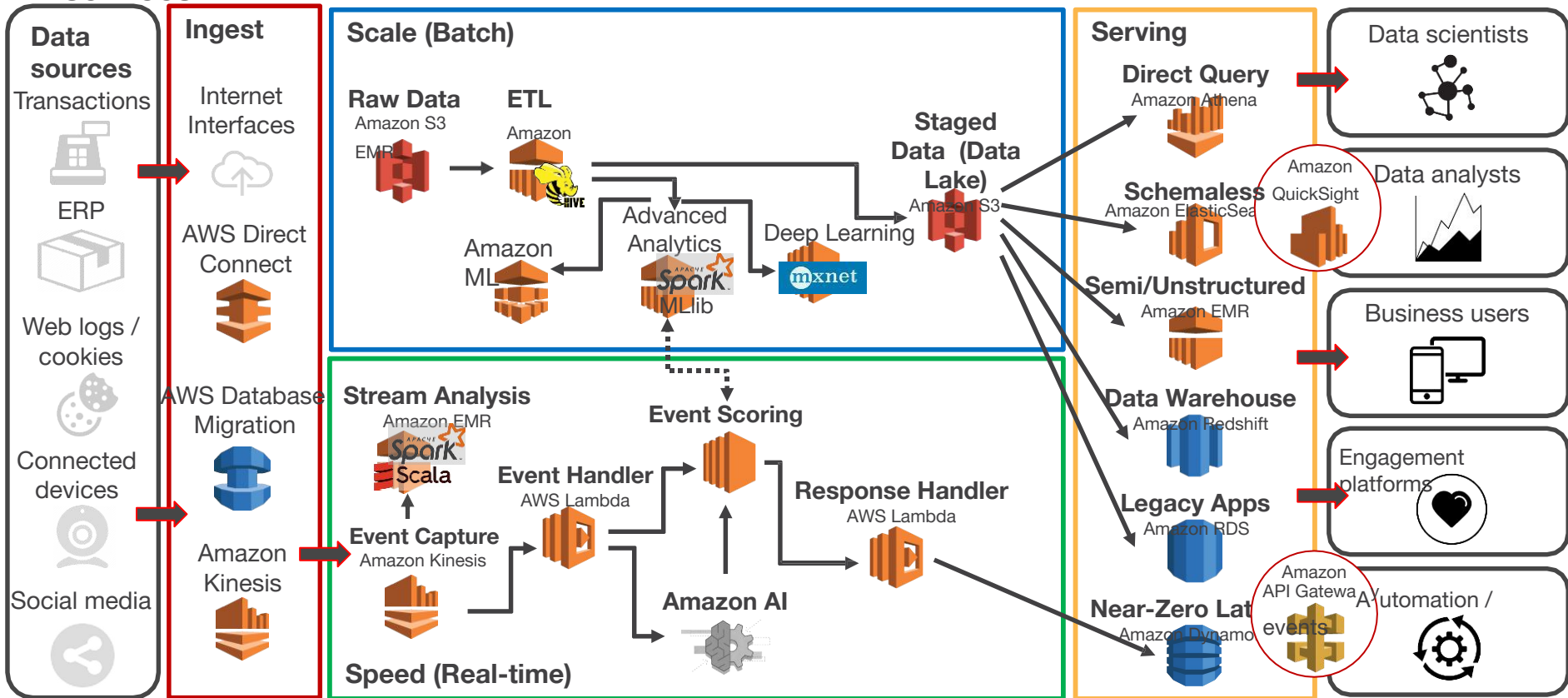
Modern data architecture

Insights to enhance business applications, new digital services



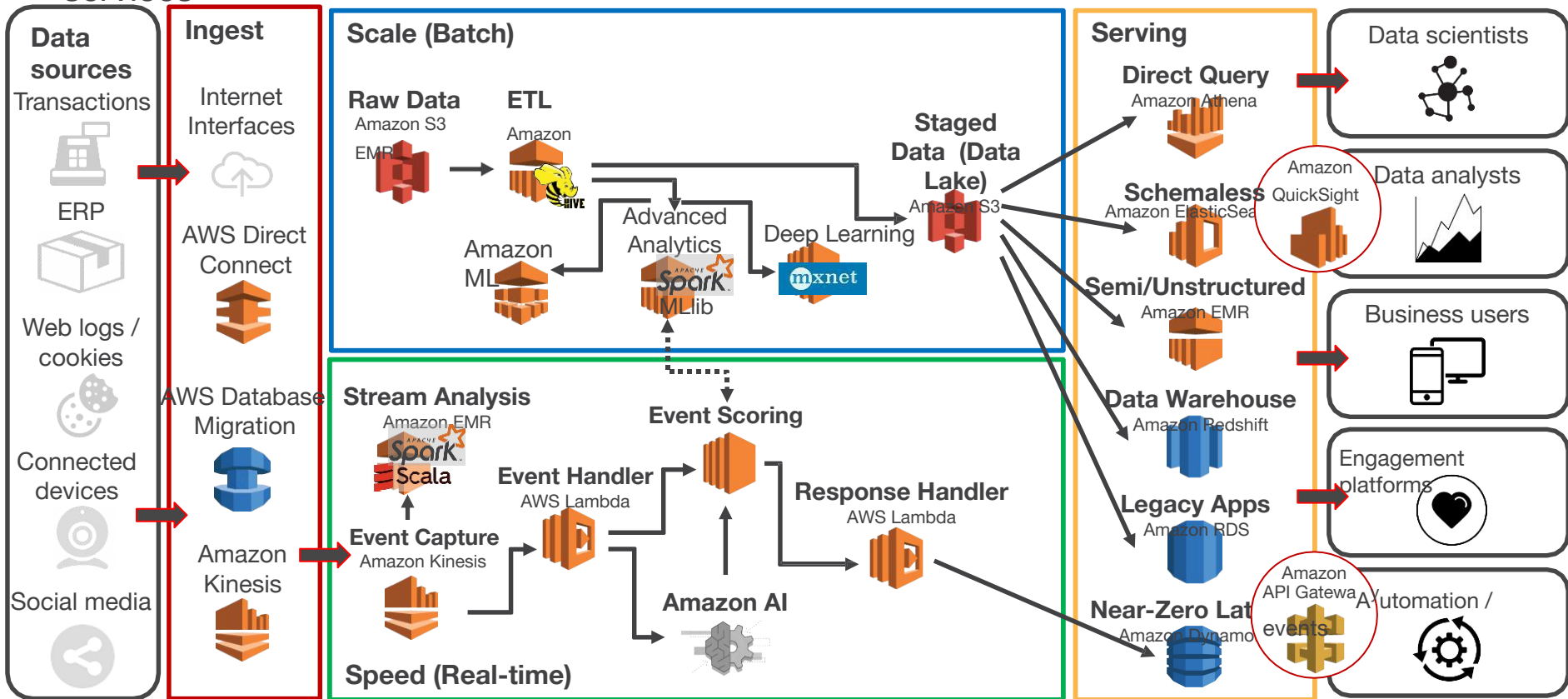
Modern data architecture

Insights to enhance business applications, new digital services



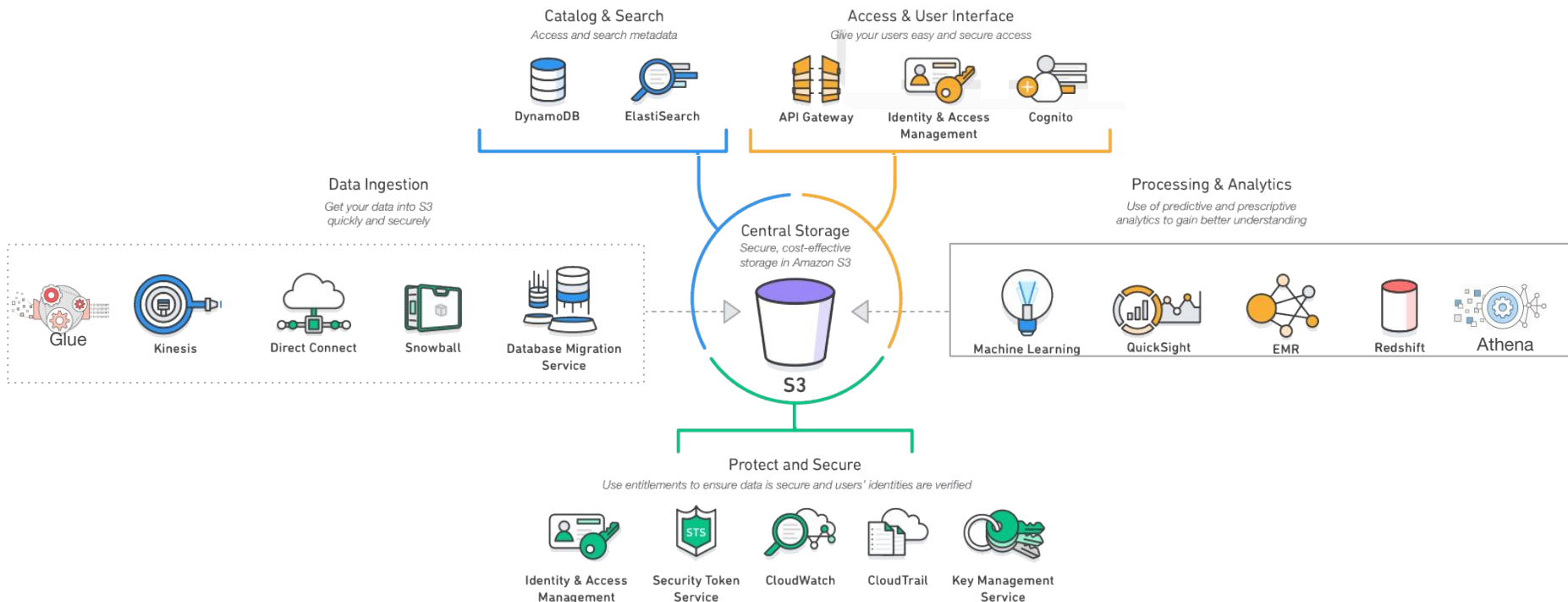
Modern data architecture

Insights to enhance business applications, new digital services

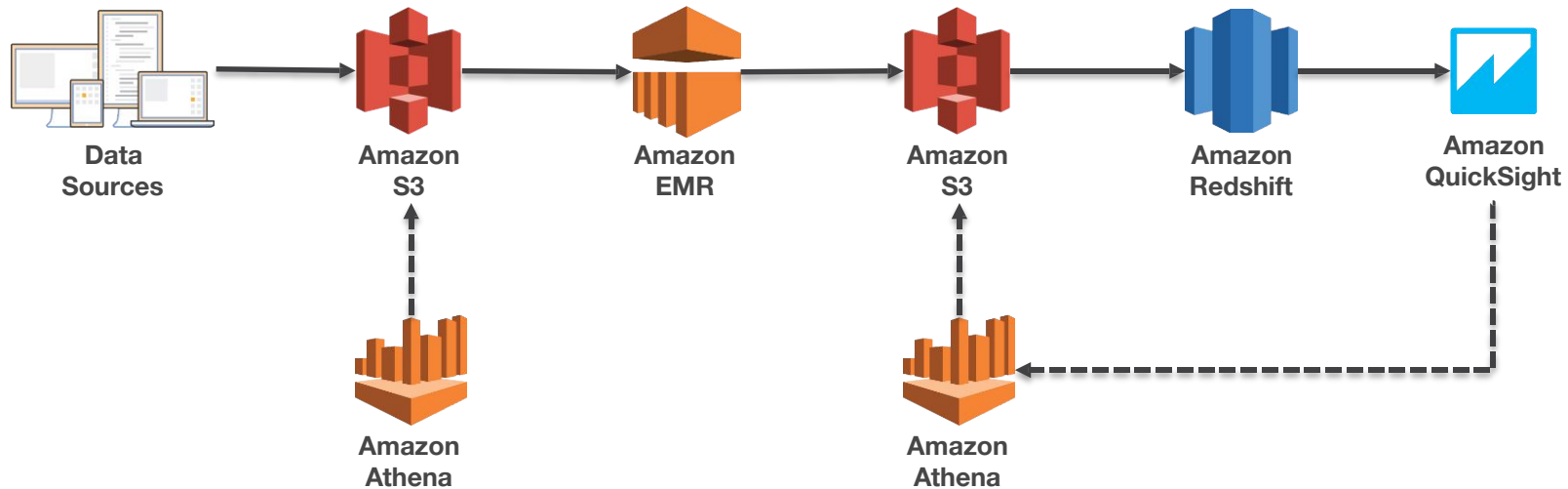


Reference Architecture

Sample Reference Architecture: Data Lake



Enterprise Data Warehouse



NORDSTROM

