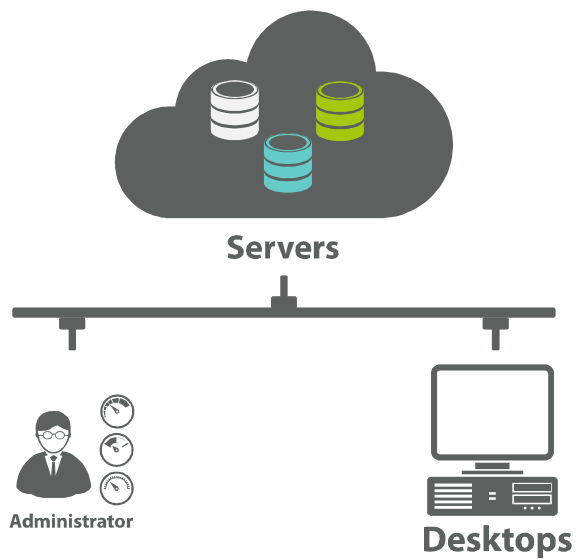




# 클라우드 기반 데이터레이크 및 분석

Data Lake를 잘 활용하기 위한 방안은 무엇인가?





# Index

- 01. Apache Spark 소개
- 02. AWS Glue 소개
- 03. Amazon RDS 소개
- 04. Amazon MSK 소개

# Apache Spark 소개



# Apache Spark



- <https://spark.apache.org/>
- Unified engine for large-scale data analytics.

## Apache Spark (Cont.)



- 빅데이터 워크로드에 쓰이는 오픈 소스 고속 통합 분석 엔진
- 2009년 UC Berkeley에서 개발
- 데이터 처리 분야에서 가장 규모가 큰 오픈 소스 프로젝트
- Netflix, Yahoo, eBay와 같은 인터넷 대기업들이 대규모로 Spark를 사용하고, 8000개가 넘는 클러스터에서 PiB 규모의 데이터를 처리
- 현재 250개 이상의 조직에서 1000명 이상이 Contributor로 활동
- 인메모리 기반의 데이터 처리로 Hadoop MapReduce 대비 100배(디스크 사용시 10배) 빠름
- DAG 스케줄러, 쿼리 최적화 도구, 물리적 실행 엔진을 사용하여 고성능을 제공
- 다양한 데이터 저장소와 생태계가 잘 구축되어 있음.
- 배치 및 실시간 처리 뿐만 아니라 머신 러닝 빌드 애플리케이션 지원(MLlib), 그래프 철(GraphX) 지원

## Apache Spark (Cont.)



- Spark Core API(일반 실행)
  - Spark Core는 Spark 플랫폼의 기본 일반 실행 엔진
- Spark SQL + DataFrame(구조화된 데이터)
  - Spark SQL은 구조적 데이터 처리를 위한 Spark 모듈
- Stream(Stream 분석)
  - Spark의 사용 편의성과 내고장성을 그대로 활용하면서도 Stream 데이터와 과거 데이터에 강력한 인터랙티브 분석 애플리케이션을 지원
- MLlib(머신러닝)
  - 확장 가능한 머신 러닝 라이브러리로, 고급 알고리즘과 빠른 속도 제공
- GraphX(그래프 계산)
  - Spark를 기반으로 한 그래프 계산 엔진으로, 사용자가 대규모의 구조화된 그래프 데이터를 상호작용 방식으로 구축, 변환하고 추론할 수 있도록 지원

## Apache Spark (Cont.)



- Spark 장점

- 속도

- 여러 개의 병렬 작업에 걸쳐 데이터를 메모리에 캐시하여 빠른 실행 속도를 자랑,
    - Hadoop MapReduce 대비 최고 100배 빠르고 디스크에서 처리시 10배 빠름.

- 실시간 스트림 처리

- 실시간 스트리밍을 처리하기도 하고, 다른 프레임워크와 통합 가능

- 통합 엔진(여러 워크로드 지원)

- SQL 쿼리, 스트리밍 데이터, 머신러닝과 그래프 처리 지원 및 높은 수준의 라이브러리 패키지 제공을 통해 생산성 향상 및 복잡한 워크플로 구현

- 사용 편리성 증가

- Java, Scala, Python, R 여러가지 프로그래밍 언어를 지원, 데이터 변환을 위한 100개 이상의 연산자 컬렉션과 반구조화된 데이터 조작에 흔히 사용

# AWS Glue 소개





## AWS Glue



- 분석 사용자가 여러 소스의 데이터를 쉽게 검색, 준비, 이동, 통합할 수 있도록 하는 서버리스 데이터 통합 서비스
- 작성, 작업 실행, 비즈니스 워크플로 구현을 위한 추가 생산성 및 데이터 운영 도구 제공
- 70개 이상의 다양한 데이터 소스 연결 지원
- 중앙 집중식 데이터 카탈로그에서 데이터 관리
- ETL 파이프라인을 시각적으로 생성, 실행, 모니터링 가능
- DataLake에 데이터를 로드하거나 Athena, EMR, Redshift Spectrum을 사용하여 카탈로그화된 데이터를 즉시 검색하고 쿼리 가능
- ETL, ELT, Streaming과 같은 모든 워크로드를 하나의 서비스에서 유연하게 지원

# AWS Glue (Cont.)



- AWS Glue 용어

- AWS Glue Data Catalog

- AWS Glue의 영구적 메타데이터 스토어.
    - 테이블 정의, 작업 정의 및 기타 관리 정보를 포함하여 AWS Glue 환경을 관리.
    - AWS 계정의 Region당 하나

- Classifier

- 데이터 스키마를 결정.
    - CSV, JSON, AVRO, XML 등과 같은 일반 파일 형식에 대한 분류자 뿐만 아니라 JDBC 연결을 사용한 일반 관계형 데이터베이스 관리 시스템을 위한 분류자를 제공

- Connection

- 특정 데이터 스토어에 연결하는 데 필요한 속성을 포함하는 Data Catalog 객체

- Crawler

- 데이터 스토어(소스 또는 대상)에 연결하는 프로그램
    - 분류자의 우선 순위 지정 목록을 통해 데이터의 스키마를 결정한 다음 AWS Glue Data Catalog에 메타데이터 테이블을 생성

## AWS Glue (Cont.)



- Event-driven ETL

- AWS Glue를 사용하면 새 데이터가 도착하는 대로 추출, 변환, 적재 작업을 실행할 수 있다.
- 예를 들어 S3에서 새 데이터를 사용할 수 있게 되는 즉시 실행할 ETL 작업을 시작하도록 AWS Glue를 구성할 수 있다.

- Data Catalog

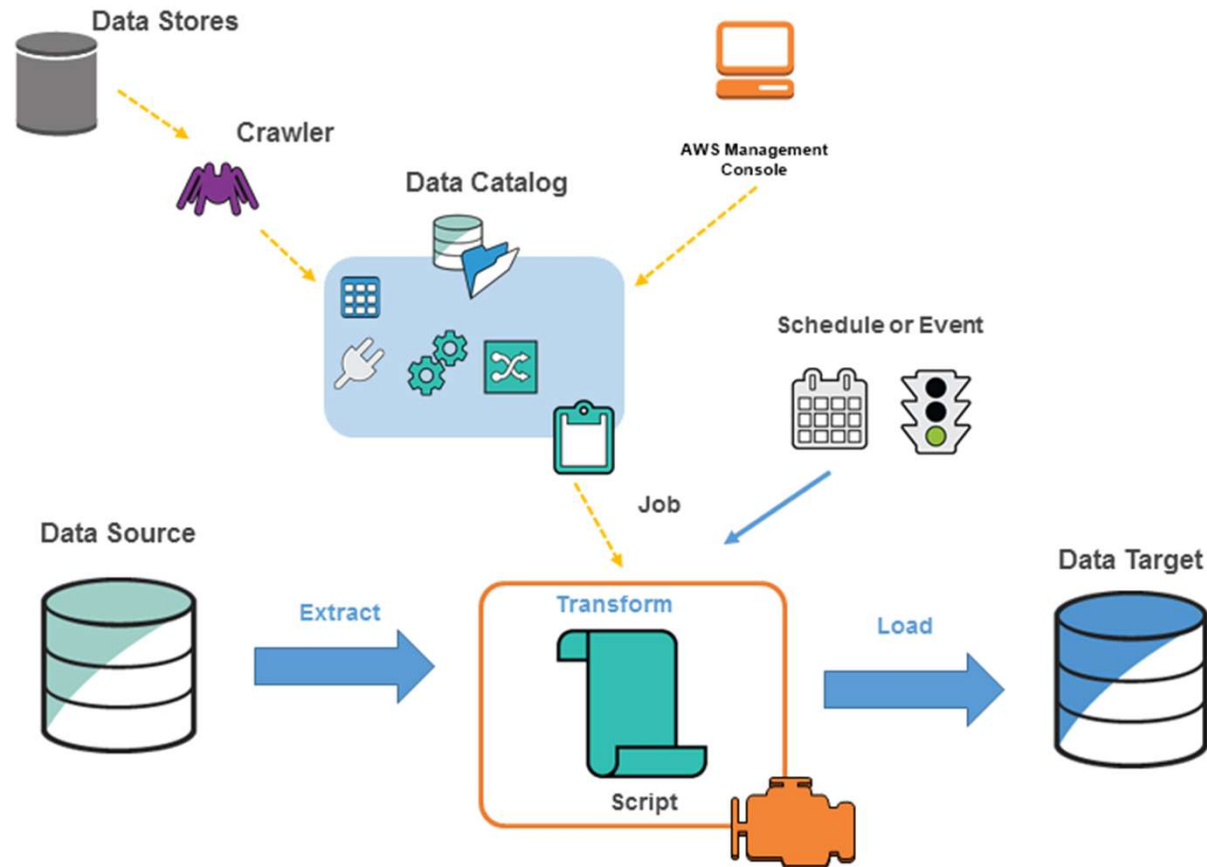
- 데이터 카탈로그를 사용하면 데이터를 이동하지 않고도 여러 AWS 데이터 세트 전체에서 신속하게 데이터를 검색할 수 있다.
- 일단 데이터가 카탈로그에 저장되면 Amazon Athena, Amazon EMR, Amazon Redshift Spectrum에서 즉시 검색 및 쿼리에 데이터를 사용할 수 있다.

## AWS Glue (Cont.)



- No-code ETL jobs
  - AWS Glue Studio를 사용하면 AWS Glue ETL 작업을 시각적으로 간편하게 생성, 실행 및 모니터링할 수 있다.
  - Drag & Drop 방식의 편집기를 사용하여 데이터를 이동 및 변환하는 ETL 작업을 구축할 수 있으며, AWS Glue가 자동으로 코드를 생성한다.
- Data Preparation
  - AWS Glue DataBrew를 사용하면 Amazon S3, Amazon Redshift, AWS Lake Formation, Amazon Aurora 및 Amazon RDS를 비롯한 DataLake, Data Warehouse 및 Database에서 직접 데이터를 탐색하고 데이터로 실험할 수 있다.
  - DataBrew의 사전 구축된 250여 개의 변환 중에서 선택하여 이상 항목 필터링, 형식 표준화, 잘못된 값 수정 등의 데이터 준비 작업을 자동화할 수 있다.

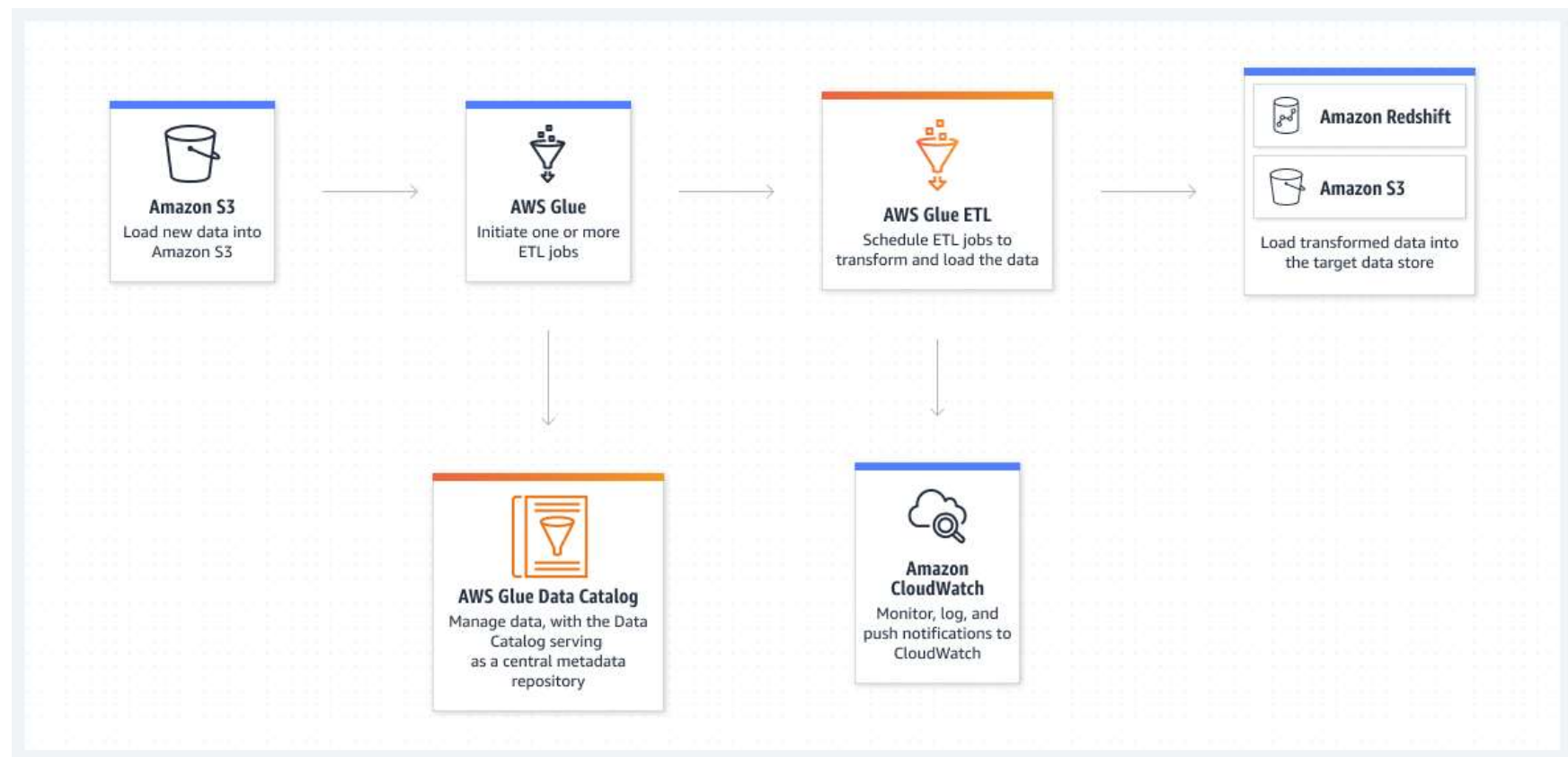
## AWS Glue (Cont.)



<https://docs.aws.amazon.com/glue/latest/dg/components-key-concepts.html>

## AWS Glue (Cont.)

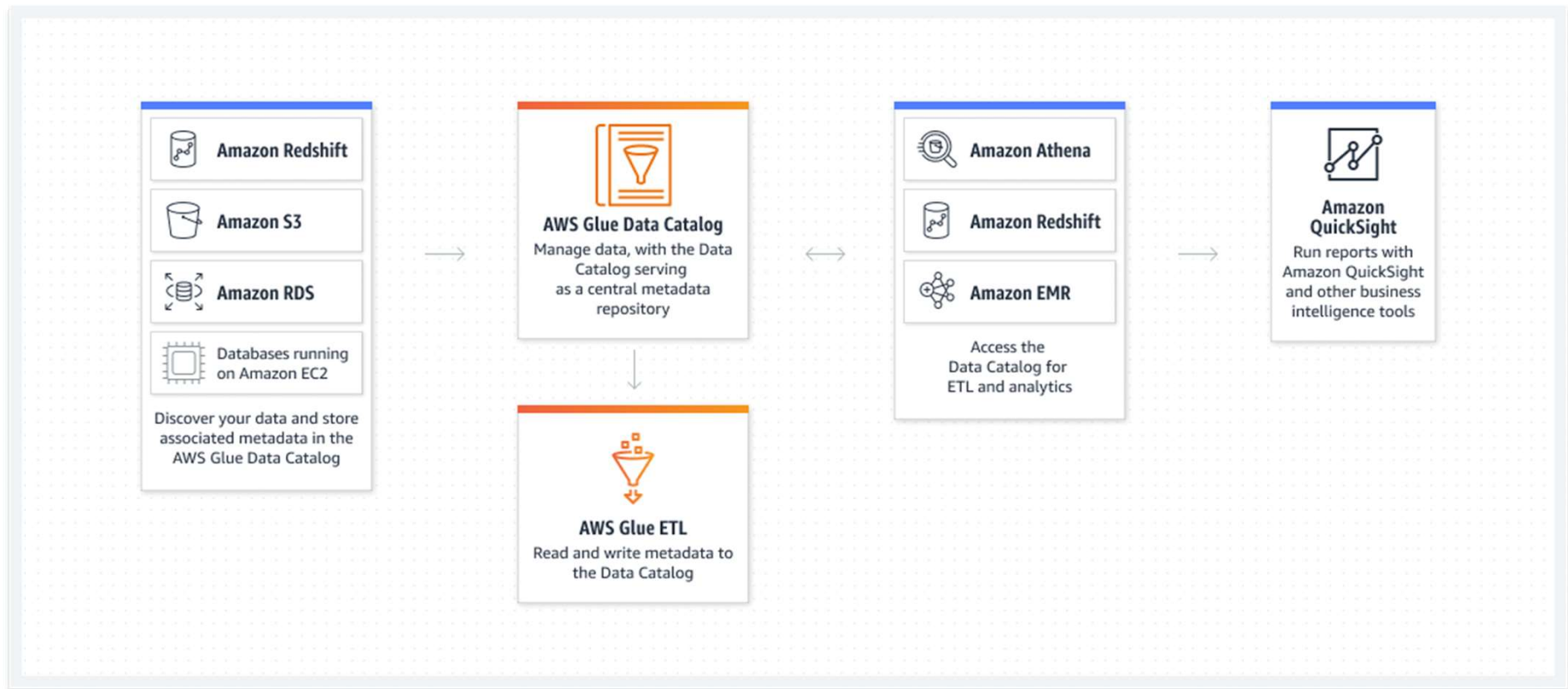
- How it works – Event-driven ETL



## AWS Glue (Cont.)

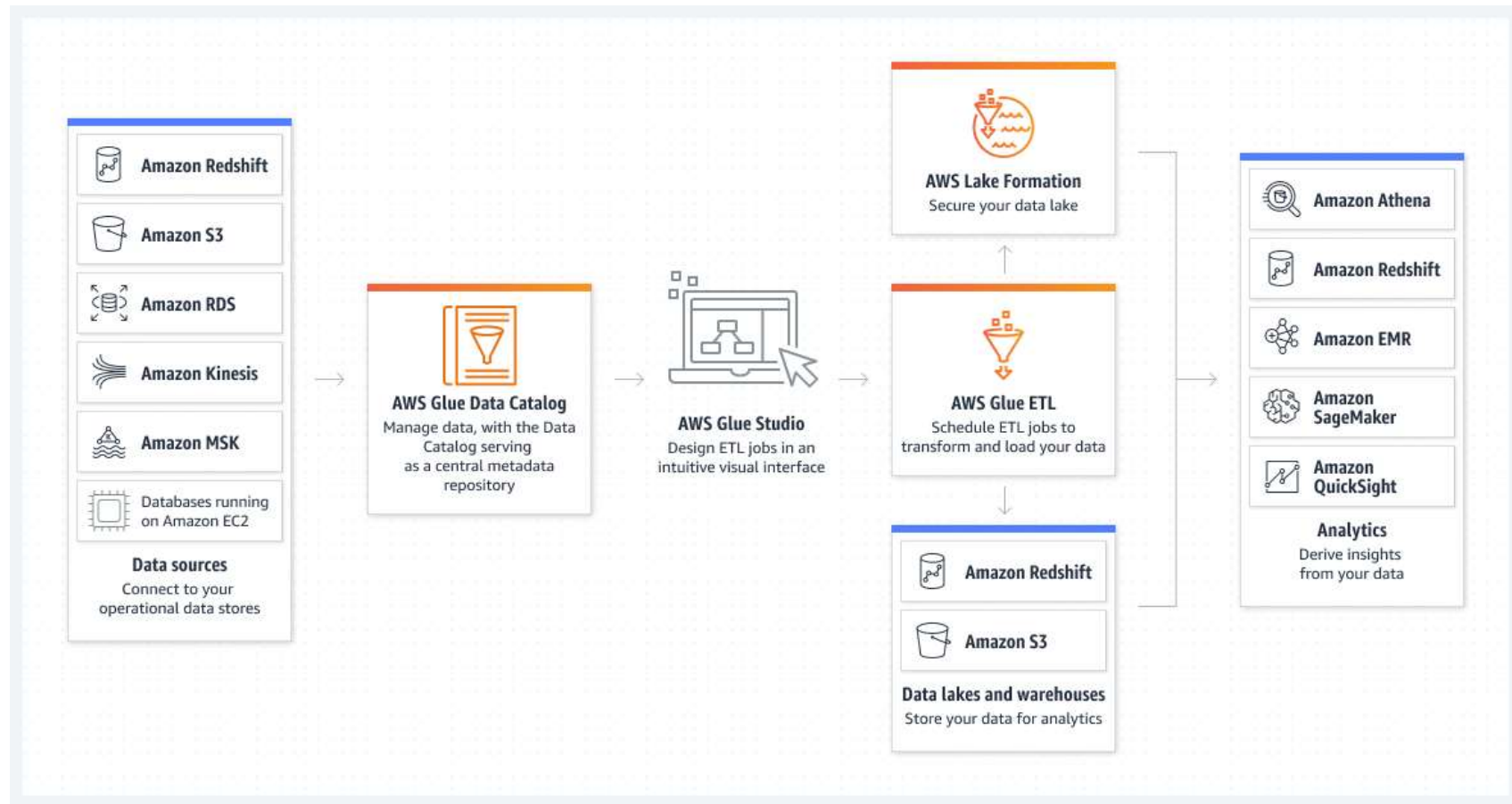


- How it works – AWS Glue Data Catalog



## AWS Glue (Cont.)

- How it works – No-code ETL jobs





## AWS Glue (Cont.)

- How it works – Data preparation



# Amazon Athena 소개

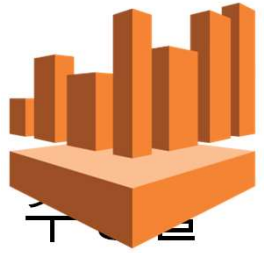


## Amazon Athena



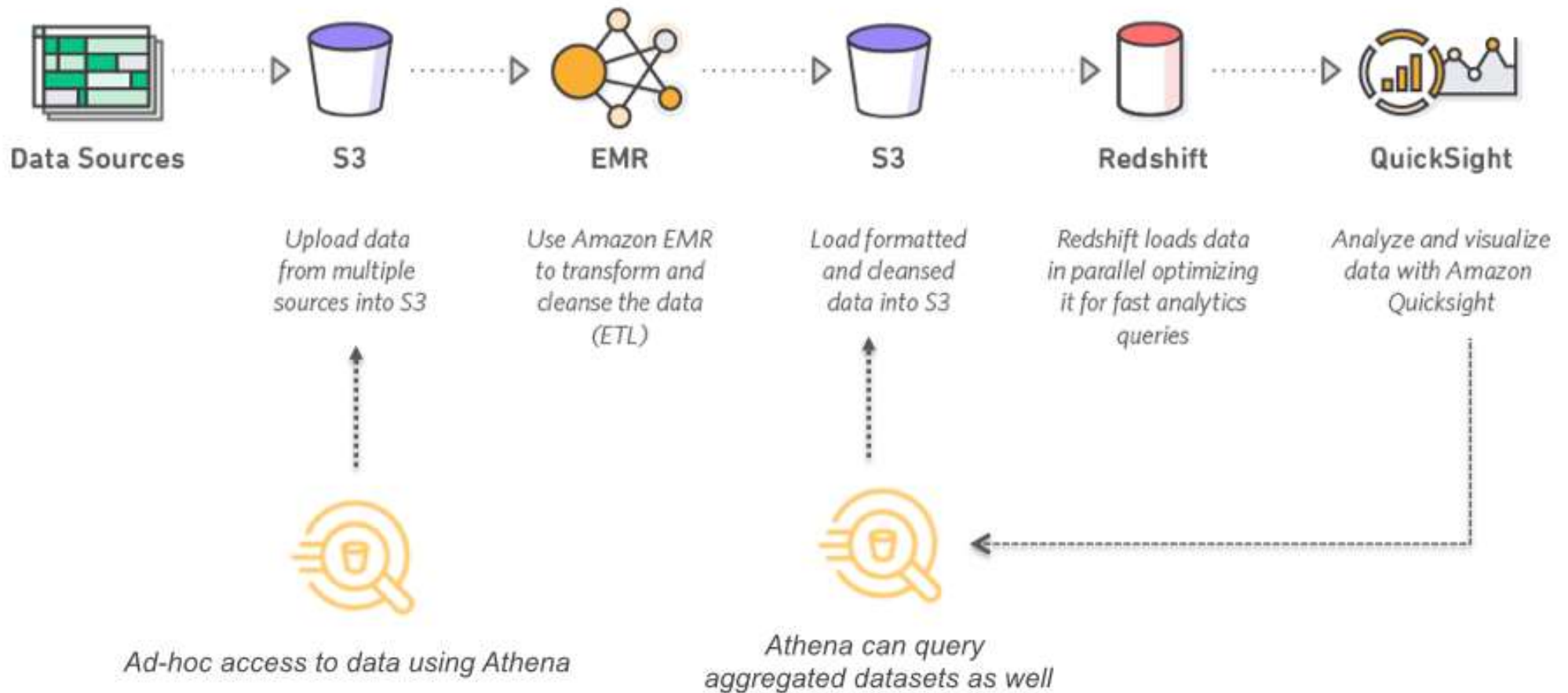
- 표준 SQL를 사용하여 Amazon S3에 있는 데이터를 직접 간편하게 분석할 수 있는 대화형 쿼리 서비스
- AWS Management Console에서 몇 가지 작업을 수행하면 Athena에서 S3에 저장된 데이터를 저장하고 표준 SQL을 사용하여 Adhoc Query를 실행하여 몇 초 안에 결과를 얻을 수 있음.
- Athena는 Serverless 서비스이므로 설정하거나 관리할 인프라가 없음.
- 비용은 실행한 쿼리에 대해서만 과금됨
- Athena는 자동으로 확장되어 쿼리를 병렬로 실행하여 대규모 데이터 집합과 복잡한 쿼리에서도 빠르게 결과를 얻을 수 있음.
- 일반적으로 비정형, 반정형 및 정형 데이터를 분석하는 데 도움(ex. CSV, JSON, Parquet, ORC)
- 다양한 데이터 시각화 도구와 연결을 지원.

## Amazon Athena - Features



- Athena는 PiB 규모의 데이터에 대해 표준 SQL문에 기반한 질의를 수행할 수 있다.
- S3를 스토리지로 사용하기 때문에 99.999999999%에 달하는 S3의 내구성이 그대로 데이터에 적용
- 데이터 소스에 대응하는 테이블 메타 정보만 생성하면 바로 쿼리를 수행할 수 있으며, 쿼리 수행 속도 또한 매우 빠르다.
- S3에서 스캔하는 데이터 1TB당 5달러로 매우 저렴한 가격(매번 쿼리를 수행할 때 스캔하는 데이터의 양에 따라 과금되며, 미리 서버를 준비할 필요가 없어 고정 비용이 발생하지 않음)
- Presto, Hive 크게 두 가지의 오픈 소스 기술이 적용되어 있음.
  - Presto
    - In-Memory 분석 쿼리 엔진으로 ANSI-SQL 호환
  - Hive
    - DDL 관련 기능을 처리하는 것을 담당, 복잡한 데이터 타입, 여러 포맷, 데이터 파티셔닝, 테이블 생성 등

# Amazon Athena(Cont.)



# Amazon RDS 소개



# Amazon RDS



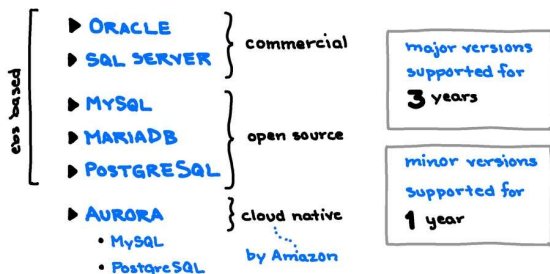
- AWS에서 관계형 데이터베이스를 더 쉽게 설치, 운영 및 확장할 수 있는 웹서비스
- 산업 표준 관계형 데이터베이스를 위한 경제적이고 크기 조절이 가능한 용량을 제공하고 공통 데이터베이스 관리 작업을 관리





# Amazo

## ENGINES



## PRICING

**COSTS:** (Oregon)

**DB INSTANCE HOURS** - per engine/instance type billed hourly, round up

**GP2** ... 11.54/GB/Mo

**STORAGE**

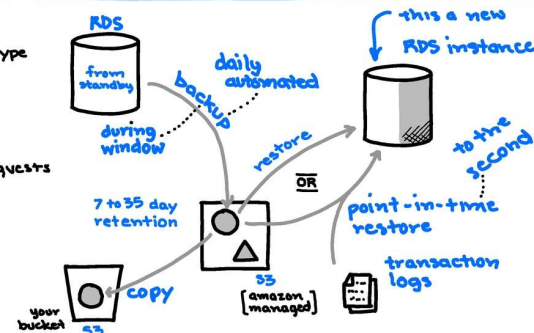
**IO1** ... 12.54/GB/Mo + 104/10/Mo

**Mag** ... 104/GB/Mo + 104/10/Mo Requests

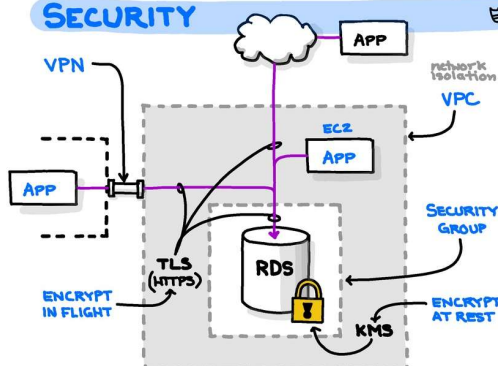
**BACKUP STORAGE** 100% of DB size free 9.54/GB/Mo after

**DATA TRANSFER** varies by destination

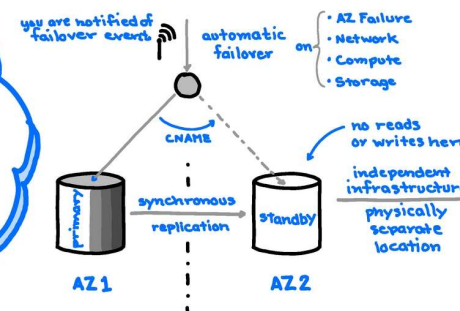
## BACKUPS



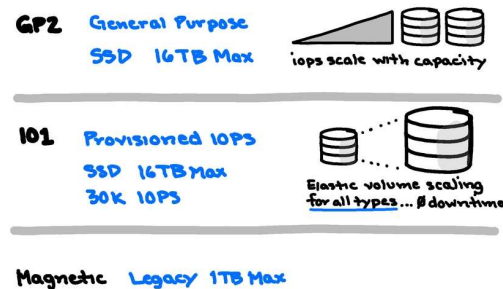
## SECURITY



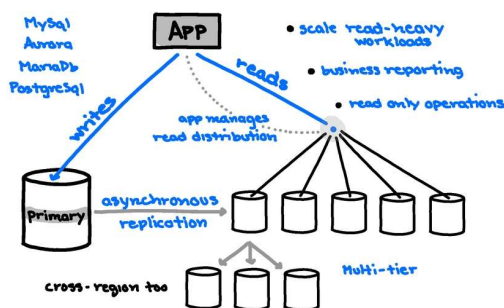
## MULTI-AZ



## STORAGE



## READ REPLICAS



## INSTANCES

SMALL WORKLOADS		
<b>T2</b> burst	1 vCPU 1GB RAM	8 vCPU 32 GB RAM
CPU INTENSIVE WORKLOADS		
<b>M3/M4</b> general purpose	2 vCPU 8GB RAM	64 vCPU 256 GB RAM
QUERY INTENSIVE WORKLOADS		
<b>R3/R4/X1(c)</b> mem optimized	2 vCPU 16 GB RAM	128 vCPU 3904 GB RAM



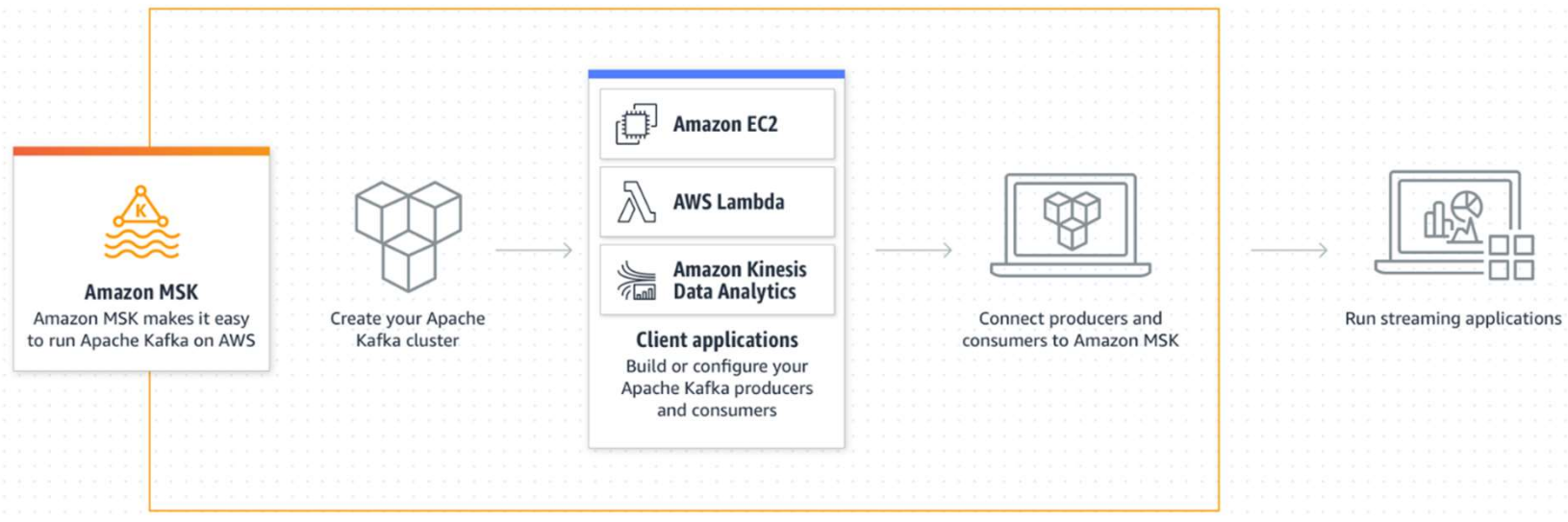
# Amazon MSK 소개



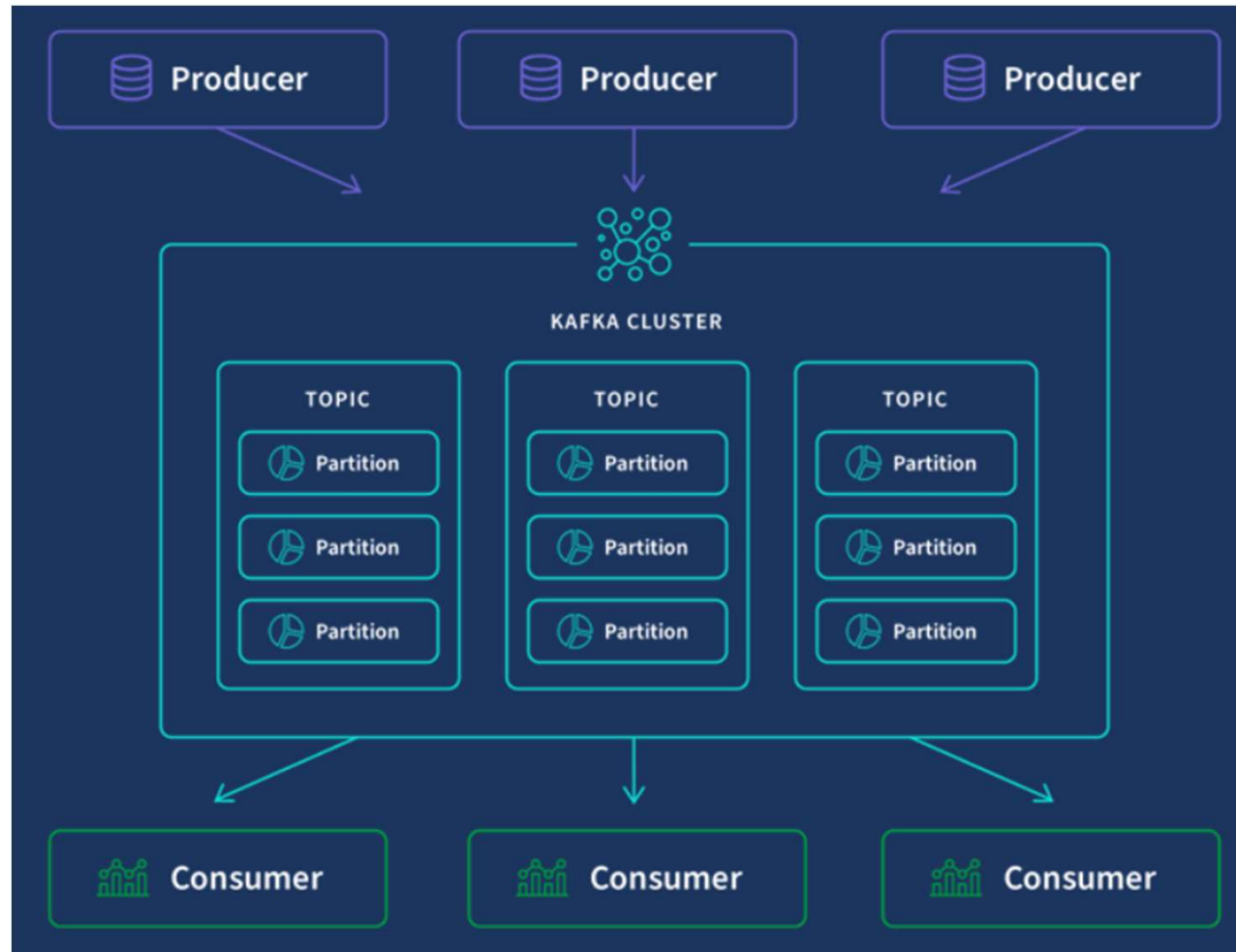
# Amazon MSK



- Kafka 인프라와 운영을 관리하는 AWS 스트리밍 데이터 서비스
- Kafka 운영 관련 전문 지식이 없는 개발자 및 Devops 관리자도 손쉽게 AWS에서 Apache Kafka 애플리케이션과 Kafka Connect 커넥터를 실행할 수 있도록 지원
- Amazon MSK는 Apache Kafka 클러스터를 운영, 유지 관리, 크기 조정하고, 즉시 사용 가능한 엔터프라이즈급 보안 기능을 제공하며, 스트리밍 데이터 애플리케이션 개발 속도를 높여 주는 내장 AWS 통합 기능 제공



## Amazon MSK (Cont.)



<https://www.qlik.com/us/streaming-data/apache-kafka>

## Amazon MSK (Cont.)



- Topic
  - 메시지를 전송하기 위해 사용되는 파티션의 집합
- Partition
  - 메시지를 병렬적으로 처리하기 위한 분산 저장 단위
  - 파티션 내부는 Queue로 되어 있어 순서를 보장하지만, 파티션 간 순서는 보장하지 않음.
- Record
  - 파티션에 들어가는 Byte 배열
  - Key/Value/Timestamp로 구성
- Offset
  - 파티션의 각 레코드를 식별할 수 있는 유일한 값
- Replica
  - 레코드의 복제본으로 이벤트 유실 방지

## Amazon MSK (Cont.)



- Broker
  - Kafka Client와 데이터를 주고받기 위해 사용하는 주체
  - 데이터를 분산 저장하여 장애가 발생하더라도 안전하게 사용할 수 있도록 도와주는 요소
- Producer
  - 데이터를 Kafka로 전송하는 주체
- Consumer
  - 데이터를 Kafka로부터 소비하는 주체
- ZooKeeper
  - Cluster의 설정 정보 관리, 동기화 등 Cluster 서버들이 공유하는 데이터 관리