# BIG DATA

## TOO BIG TO IGNORE
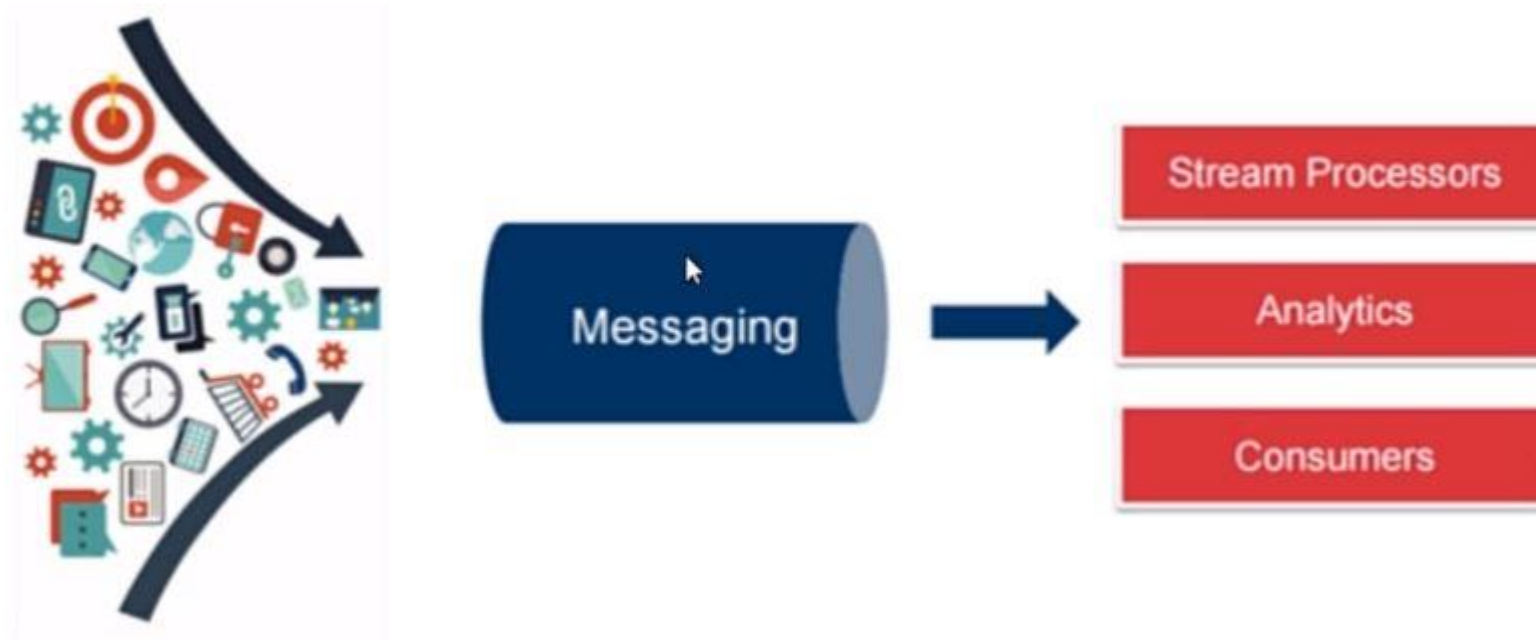
SÜMEYYE KAYNAK

# OUTLINE

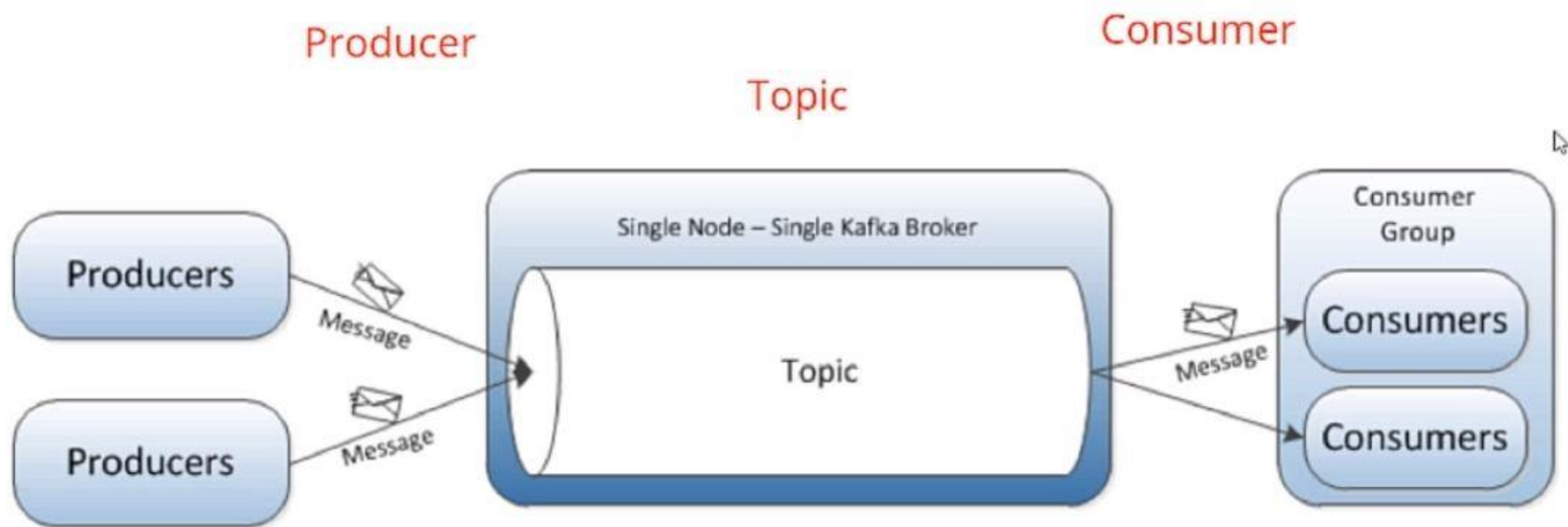Apache Kafka

Apache Zookeeper

Docker

# APACHE KAFKA

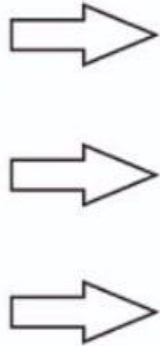- It is very important to collect and analyze big data quickly and without errors.
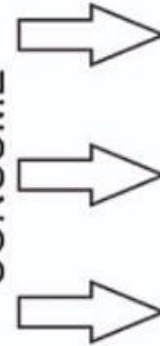
# APACHE KAFKA

# APACHE KAFKA

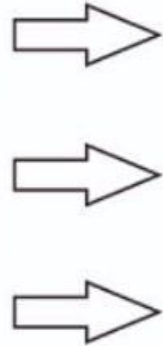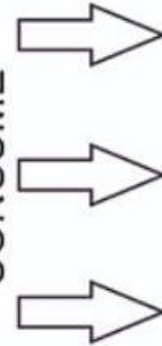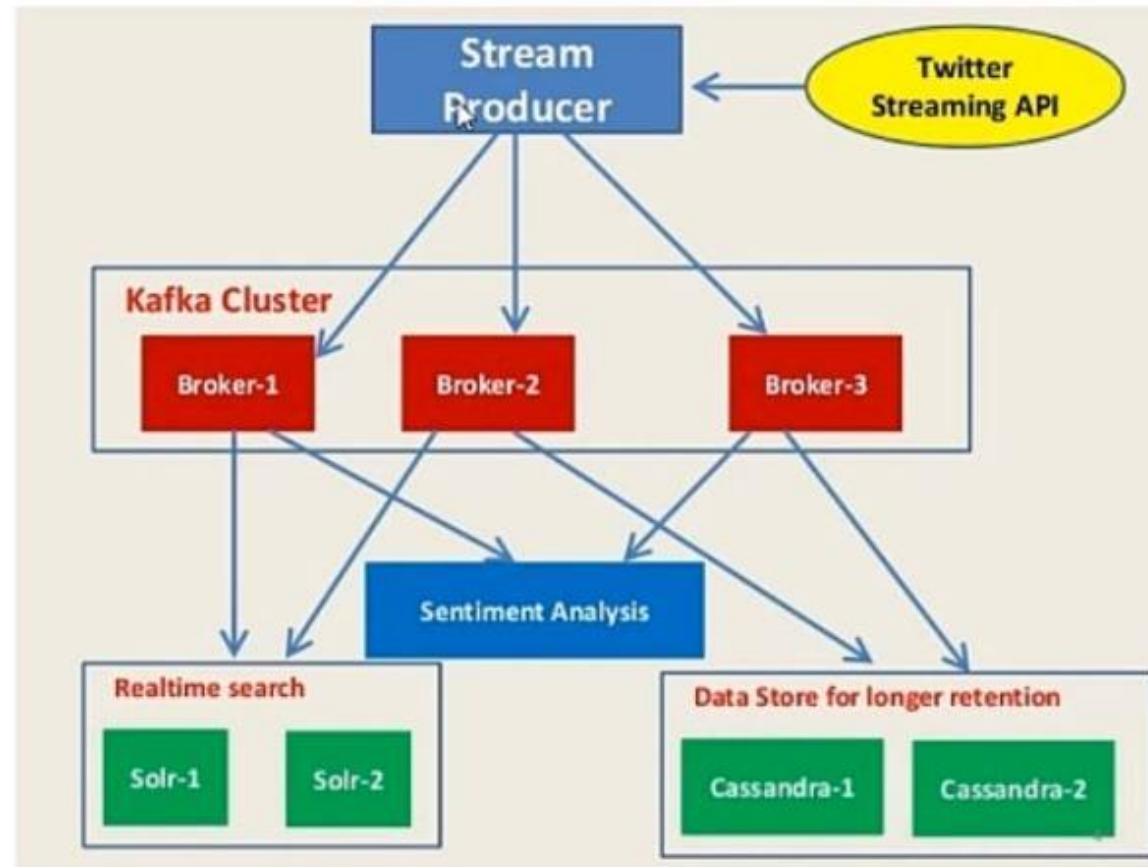# APACHE KAFKA

# APACHE KAFKA

Distributed data storage and replication are available in Kafka.

# APACHE KAFKA

- A file is divided into 3 parts.

- Partition 0 of the file is the leader in broker1.

- Replication factor value is 3.

- The leading part in each machine changes.

- Data loss is prevented.

Leader (red) and replicas (blue)

Client

write

Broker 1

Partition 0

replica write
replica write

Partition 1

Partition 2

Broker 2

Partition 0

Partition 1

Partition 2

Broker 3

Partition 0

Partition 1

Partition 2

# ZOOKEEPER



- Zookeeper coordinates resource management in distributed server architectures.
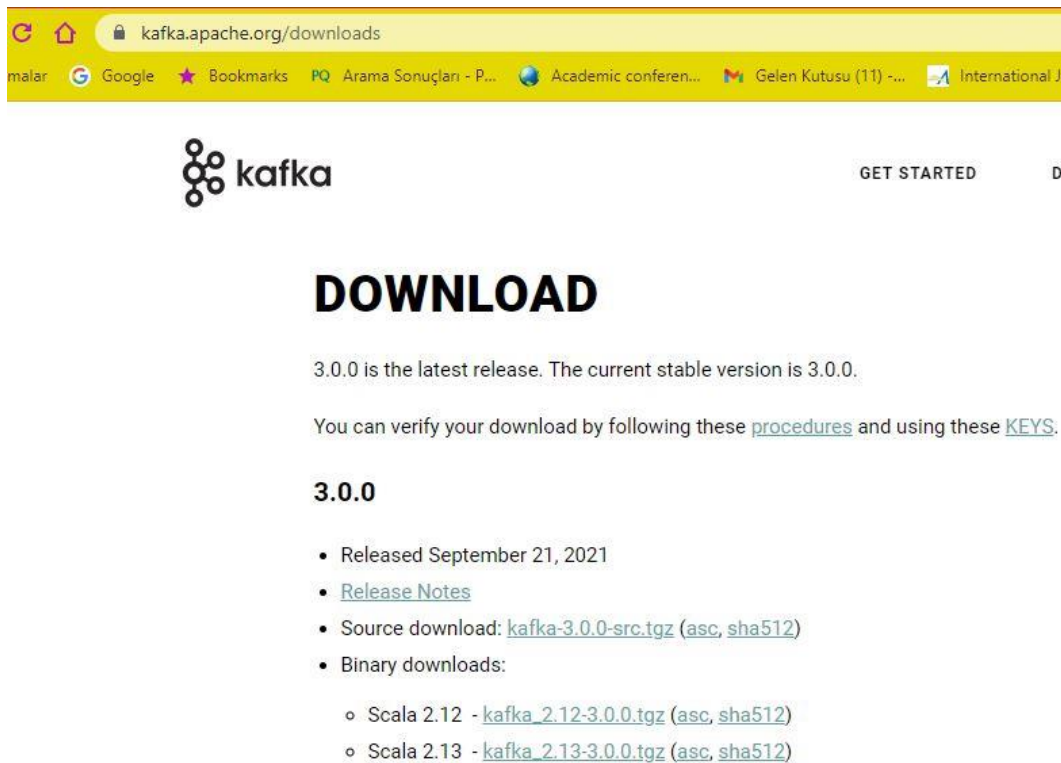- Zookeeper is generally used for configuration and keeps configuration files.

# DOWNLOAD KAFKA



- Download the appropriate version from the Kafka site.
- The downloaded zip file is extracted to the C: directory.
- Update the server.properties file in the config folder.

# DOWNLOAD ZOOKEEPER



- .Download the appropriate version from the Zookeeper site.

- The downloaded zip file is extracted to the C: directory.

- Update the zoo_sample.cfg's file name in the config folder and named zoo.cfg.

- Update the zoo.cfg

- Create folders named data.

# DOWNLOAD ZOOKEEPER

# JAVA_HOME IS NOT SET

# DOWNLOAD ZOOKEEPER

# RUNNING ZOOKEEPER

# RUNNING OF KAFKA



```
C:\kafka\bin\windows>kafka-server-start.bat C:\kafka\config\server.properties
[2021-11-22 20:39:40,214] INFO Registered kafka:type=kafka.Log4jController MBean
```

# DOCKER

- Docker is a technology that provides virtualization thanks to hundreds or even thousands of isolated and independent containers on the same operating system.

- It is an open source 'container' technology.

- Docker is a new technology that allows development teams to build, manage and secure applications anywhere.

# CONTAINER

- It is the name given to each of the processes that are run in isolation from each other in the Linux kernel by the Docker Daemon.

- A container is a special type of transaction that is isolated from other processes. Containers are assigned resources that no other process can access.

# CONTAINER-VIRTUAL MACHINE



**VM (Virtual Machine)**
- OS : **Tam işletim sistemi**
- İzolasyon : **Yüksek**
- Çalışır hale gelmesi : **Dakikalar**
- Versiyonlama : **Yok**
- Kolay paylaşılabilirlik : **Düşük**

**Docker**
- OS : **Küçültülmüş işletim sistemi imajı**
- İzolasyon : **Daha düşük**
- Çalışır hale gelmesi : **Saniyeler**
- Versiyonlama : **Yüksek**
- Kolay paylaşılabilirlik : **Yüksek**

# HYPERVISOR

- Hypervisor is a piece of code that allows multiple operating systems to run on the same hardware.

- This piece of code works directly on the hardware, allowing us to create multiple guest operating systems on our physical server.

# DOCKER IMAGE-DOCKER REGISTERY

- It is a structure that contains your application to run and the necessary operating system libraries running on your application's infrastructure.

- Images are kept in Docker Registries.

# DOCKER COMPOSE

- With Docker Compose, you can define more than one container in a single file, you can run the application by standing up all the requirements your application needs with a single command.

# KAFKA WITH DOCKER

- Apache Kafka is a distributed publish-subscribe messaging system that is designed to be fast, scalable, and durable.

- Kafka stores streams of records (messages) in topics. Each record consists of a key, a value, and a timestamp.

- Producers write data to topics and consumers read from topics.

# KAFKA WITH DOCKER

# TOPICS AND LOGS (KAFKA)

- A topic is a category to which records are published. It can have zero, one, or many consumers that subscribe to the data written to it.

- For each topic, the Kafka cluster maintains a partitioned log. Since Kafka is a distributed system, topics are partitioned and replicated across multiple nodes.



Anatomy of a Topic

# PRODUCERS AND CONSUMERS

- Producers publish data to the topics of their choice. It is responsible for choosing which record to assign to which partition within the topic.

- Consumer groups can subscribe to one or more topics. Each one of these groups can be configured with multiple consumers.

- Every message in a topic is delivered to one of the consumer instances inside the group subscribed to that topic. All messages with the same key arrive at the same consumer.

# MESSAGES

- Kafka is message-based, the main element that is processed is messages.

- Messages consist of a simple key-value pair. Both key and message contents can be anything that can be serialized.

- Examples of messages are page clicks, comments, and orders.

# PARTITION

- Kafka runs as multiple nodes on different physical/virtual machines for load distribution and fault tolerance.

- Topics are also physically divided into different partitions on these machines.

- Each partition can also be copied to other nodes as master/slave.

# DOCKER-COMPOSE.YML- RUNNING DOCKER-COMPOSE



```
docker-compose.yml - Notepad
File  Edit  Format  View  Help
version: "3"
services:
  zookeeper:
    image: 'bitnami/zookeeper:latest'
    ports:
      - '2181:2181'
    environment:
      - ALLOW_ANONYMOUS_LOGIN=yes
  kafka:
    image: 'bitnami/kafka:latest'
    ports:
      - '9092:9092'
    environment:
      - KAFKA_BROKER_ID=1
      - KAFKA_CFG_LISTENERS=PLAINTEXT://:9092
      - KAFKA_CFG_ADVERTISED_LISTENERS=PLAINTEXT://127.0.0.1:9092
      - KAFKA_CFG_ZOOKEEPER_CONNECT=zookeeper:2181
      - ALLOW_PLAINTEXT_LISTENER=yes
    depends_on:
      - zookeeper
```

```
Microsoft Windows [Version 10.0.18363.1556]
(c) 2019 Microsoft Corporation. All rights reserved.

C:\Users\Sumeyye>cd..

C:\Users>cd..

C:\>cd Users

C:\Users>cd Sumeyye

C:\Users\Sumeyye>cd Downloads

C:\Users\Sumeyye\Downloads>cd kafka-docker

C:\Users\Sumeyye\Downloads\kafka-docker>docker-compose up -d
[+] Running 2/2
 - Container kafka-docker-zookeeper-1   Started
 - Container kafka-docker-kafka-1       Started

C:\Users\Sumeyye\Downloads\kafka-docker>
```

```
C:\Users\Sumeyye\Downloads\kafka-docker>docker-compose stop
[+] Running 2/2
 - Container kafka-docker-kafka-1       Stopped
 - Container kafka-docker-zookeeper-1   Stopped

C:\Users\Sumeyye\Downloads\kafka-docker>
```

# VISUAL STUDIO

```csharp
namespace Consumer
{
    class Program
    {
        static async Task Main(string[] args)
        {
            var config = new ConsumerConfig
            {
                BootstrapServers = "localhost:9092",
                GroupId = "foo",
                AutoOffsetReset = AutoOffsetReset.Earliest
            };
            List<string> topics = new List<string>() { "testtopic" };
            using (var consumer = new ConsumerBuilder<Ignore, string>(config).Build())
            {
                consumer.Subscribe(topics);

                while (true)
                {
                    var consumeResult = consumer.Consume();
                    Console.WriteLine(consumeResult.Message.Value);
                    Console.WriteLine(consumeResult.Message.Timestamp.UtcDateTime);
                }

                consumer.Close();
            }
        }
    }
}
```

```csharp
namespace KafkaNetCore.Producer
{
    class Program
    {
        static async Task Main(string[] args)
        {
            var topicName = "testtopic";
            var kafkaUrl = "localhost";

            var config = new ProducerConfig() { BootstrapServers = "localhost:9092" };

            using (var producer = new ProducerBuilder<string, string>(config).Build())
            {
                while (true)
                {
                    Console.Write("Enter message: ");
                    var text = Console.ReadLine();

                    Message<string, string> message = new Message<string, string> { Value = text };
                    var deliveryResult = await producer.ProduceAsync(topicName, message);
                    Console.WriteLine($"Delivered to '{deliveryResult.TopicPartitionOffset}'");
                }
            }
        }
    }
}
```