

Task → Install Kafka on AWS EC2 ✓

ML System

Kafka

+ Distributed
+ fault-tolerant
+ high-throughput ✓

Streaming platform

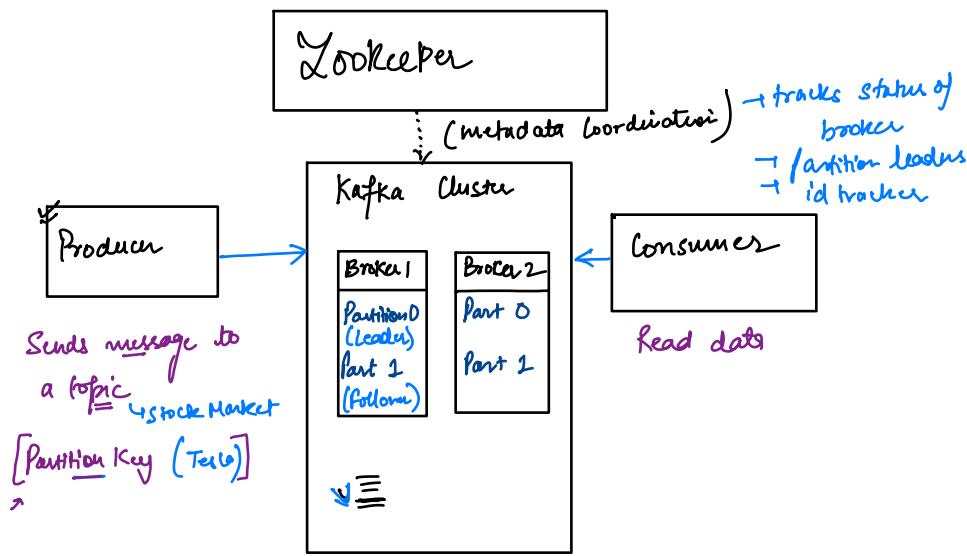
→ Real-time data feeds

→ Not a simple message-Queue

+ distributed
+ Optimised for horizontal scaling
+ durable
+ low-latency ✓

- ① Persistent Storage
- ② Scalable → Add brokers
- ③ Fault tolerant →
- ④ Ordering guarantee





Machine Learning System design for Real time Stock market Analysis

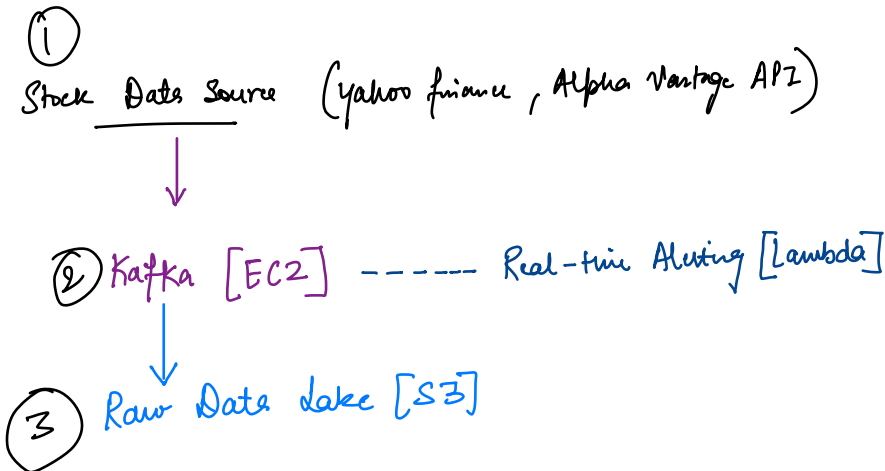
Functional Req.

- ① Real-time data ingestion
↳ ingest prices, volume, metadata
every 1-sec
- ② Stream processing → Clean, aggregate & transform data in real time.
- ③ Batch & real-time storage →
Store raw & processed data
in query-optimized format.

- ④ Real-time Alerts → Detect anomaly
(volume spikes, price drops)
- ⑤ Batch Analytics → SQL queries on historical data
- ⑥ ML predictions → Predict short-term stock trends
- ⑦ Dashboards → Visualize trends

New functⁿ

- ① low-latency → < 5 sec [End to end]
data ingestⁿ → predictions
- ② Scalability → 1000+ stocks with 1 sec updates
- ③ fault tolerance
- ④ Cost efficiency. [S3, EC2]
VM



Load Kafka on EC2

- ① Launch an EC2 instance → Amazon Linux 2
- ② SSH client details
- ③ Make a connection using key-value pair
 - local terminal
 - SSH command

1) `sudo yum install java`

2) `wget https://archive.apache.org/dist/kafka/3.3.1/kafka_2.12-3.3.1.tgz` (download Kafka)

3) `tar -xzf kafka_2.12-3.3.1.tgz` (unarchive)
4) `cd kafka_2.12-3.3.1/` (go inside the folder)

5) Kafka uses Zookeeper for distributed coordination. Start Zookeeper in the background:
`bin/zookeeper-server-start.sh config/zookeeper.properties`

6) open new terminal

Increase memory → `export KAFKA_HEAP_OPTS="-Xmx256M -Xms128M"`

7) Start the Kafka server:
`kafka-server-start.sh config/server.properties`

✓ Raw S3 file (csv / parquet / json)

3

haha

[`s3://stock-data/raw/date=2023-10-10/hour=12/haha.csv`]
partitions = date hour

✓ AWS Glue ETL ✓

④

Sample

Job1 Convert JSON to parquet, sanitize data
(handle null value)

Job2 feature engineering (moving Avg, RSI, MACD)



⑤ S3://stock-data/processed/date/hour/ —

--	--	--	--

⑥ SQL for analytics

AWS Athena

```
-- Create a table referencing S3 data
CREATE EXTERNAL TABLE my_table (
  id INT,
  name STRING,
  age INT
) ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION 's3://my-bucket/data/';
```

⑦ optional → Alerts → AWS Lambda

```
def lambda_handler(event, context):
    for record in event['records']:
        data = json.loads(record['value'])
        if data['volume'] > 1000000:
            sns.publish(TopicArn='arn:aws:sns:...', Message=f"Volume spike: {data['symbol']}")
```

SNS alert
↳ mailbox

⑧ Feature engineering

→ Moving Avg

`df['5min_MA'] = df['close'].rolling(window=300).mean() # 300 seconds = 5 mins`

⑨ Modelling

Model	Pros	Cons	Best Used
XGBOOST	Fast, Tabular data ✓	Manual feature Engineering	Intraday predictions
Prophet	Handles seasonality, easy to use	Struggle Volatile stocks	Long term trends
LSTM	Handle temporal patterns	Needs large data, slow training	High-frequency trading ✓

⑩ Minute level predictions

Deploy → Sagemaker endpoint

→ Predictions to Kafka topic

↳ predictions for real-time dashboards.

```
model = Sequential([
    LSTM(128, input_shape=(60, 5), # 60 time steps, 5 features
    Dropout(0.3),
    Dense(1)
])
model.compile(loss='mae', optimizer='adam')
```

Problem → Prediction ~ 1 sec

LSTM can take time

→ [Prophet ~ for alerts]

