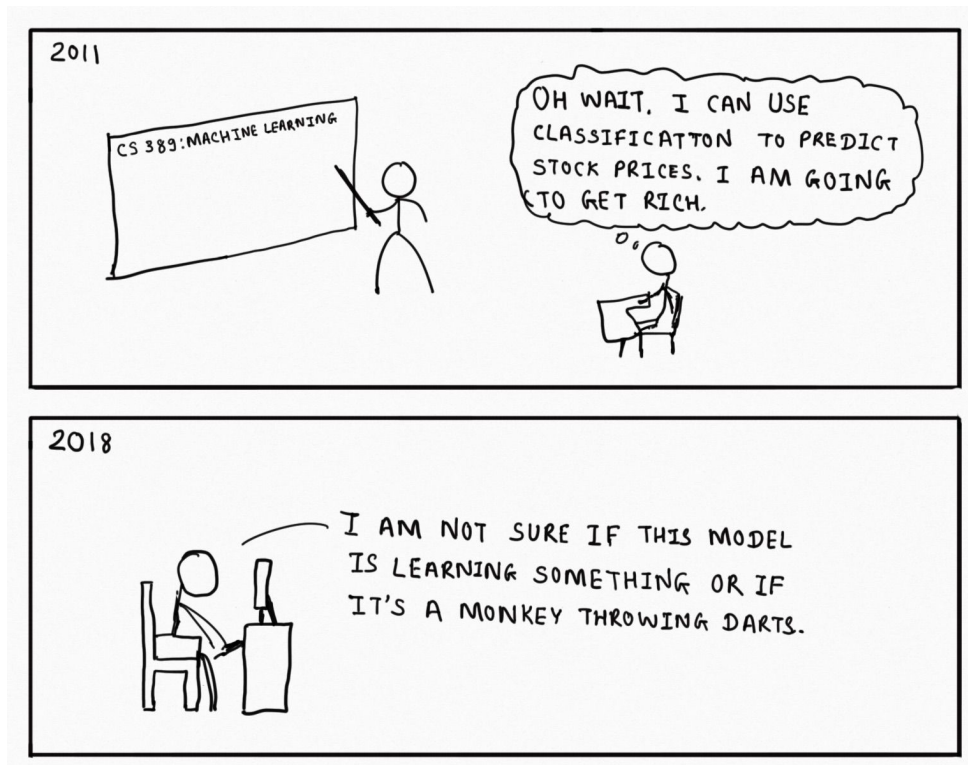


Workshop Report:

Stock Price Prediction System using FintHub Data-stream and Apache Spark

Anastasiya Merkulova, Raphael Waltenpül

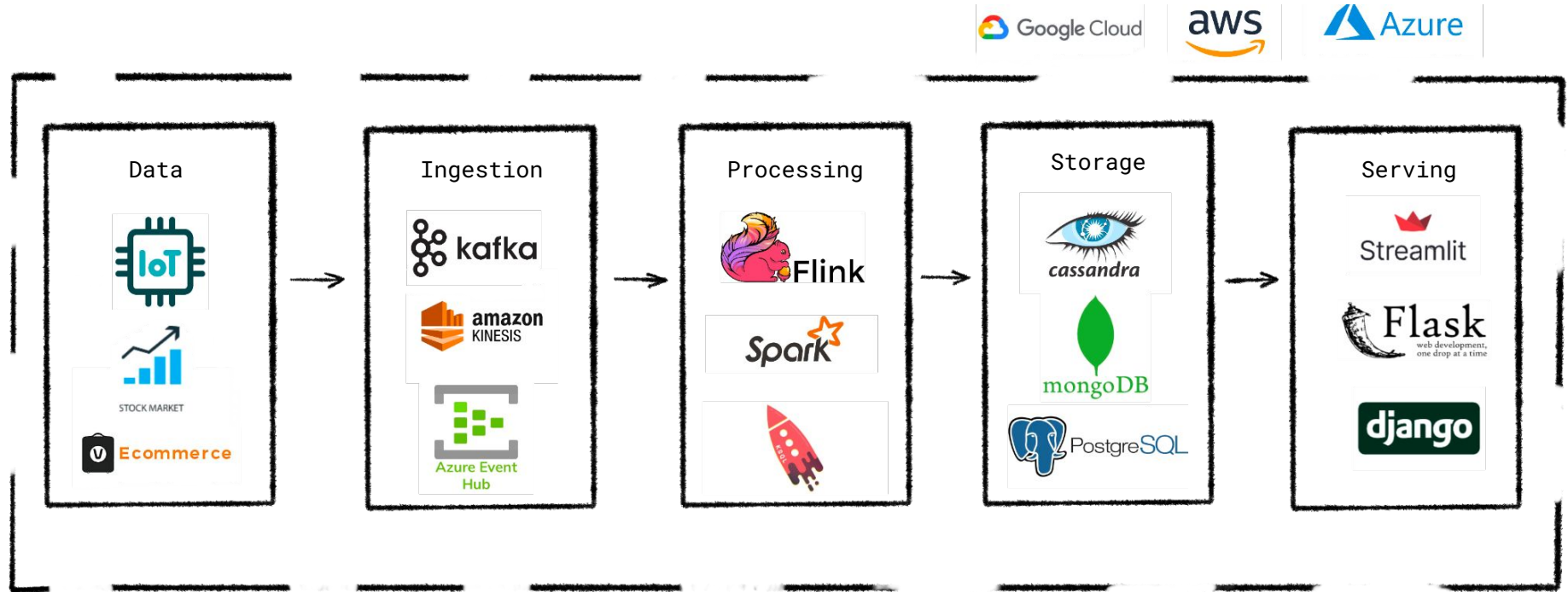


- **Goal:**

- Build a real-time distributed system to predict stocks price

- **Objectives:**

- Design a pipeline for real-time data processing and price prediction
- Identify and choose the optimal services, tools, and libraries
- Develop and integrate all components of the system
- Create a user-friendly interface to display data



Data streaming pipeline



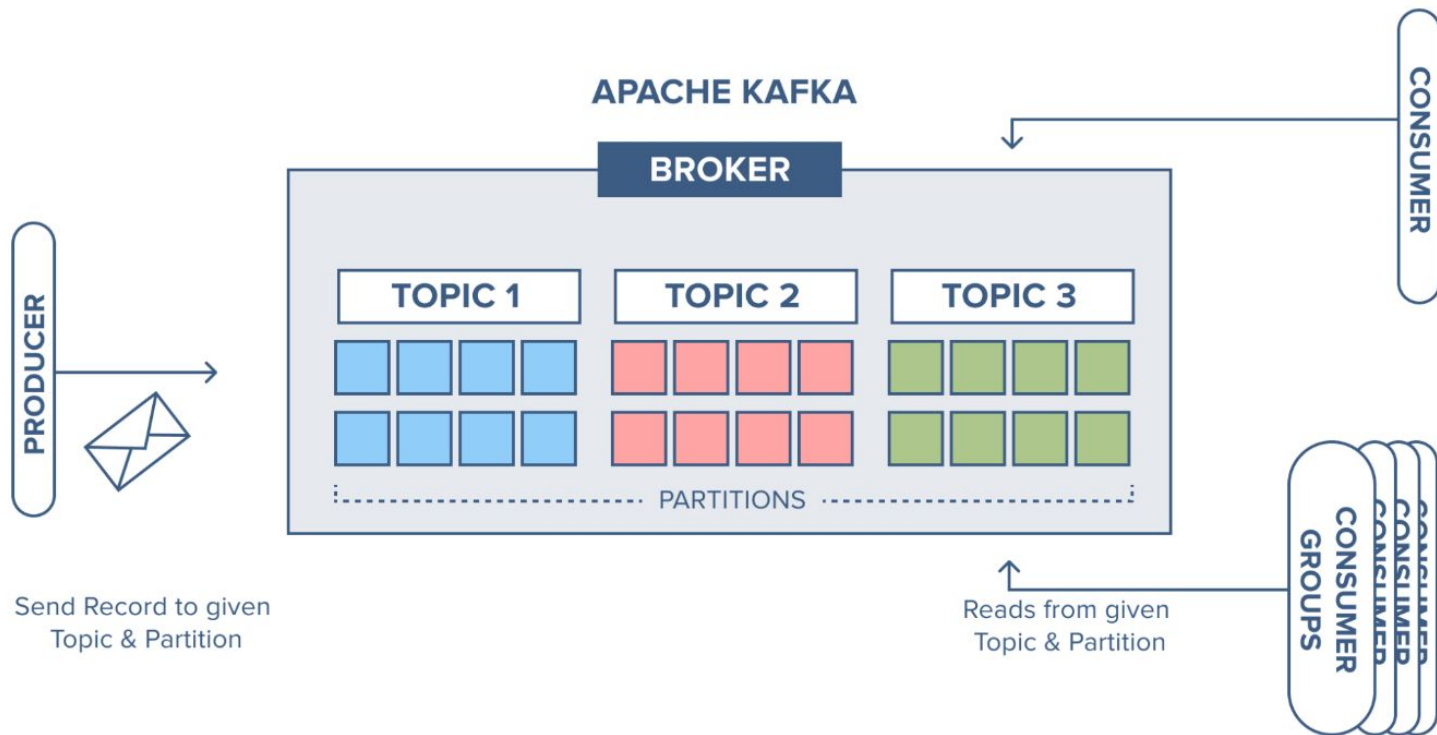
- **Why**

Finnhub?

- Data streaming via WebSockets, instant updates for stock prices
- Offers stocks, forex, and cryptocurrencies, covering global markets
- Free, detailed documentation, available for different languages



- Why Kafka?
 - Handles large volumes of data streams
 - Facilitates real-time data processing, enabling quick reaction to events
 - Data availability by replicating data across multiple brokers
 - Decouples producers and consumers, allowing asynchronous communication



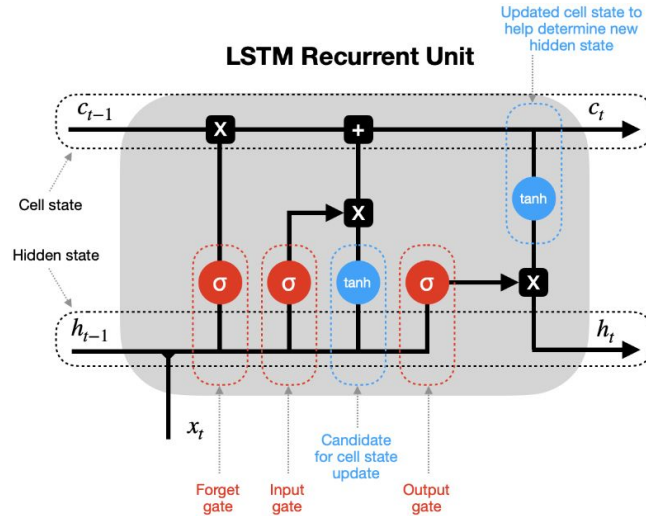


- **Why**

Apache

Hadoop?

- Distributed file system (HDFS) and processing framework (MapReduce) provide built-in fault tolerance
- Excels in processing vast amounts of data across distributed clusters
- Supports various models, including batch processing, interactive querying, and real-time processing



with  TensorFlow

- LSTM able to selectively forget information over time
- TensorFlow optimizes computation and allows easy experimentation with different model architectures and hyperparameters



- **Why**

MongoDB?

- Non-relational and document-oriented database
- High performance for read and write operations
- Horizontally scalable
- Support for real-time data processing and aggregation, well integrated with analytical frameworks



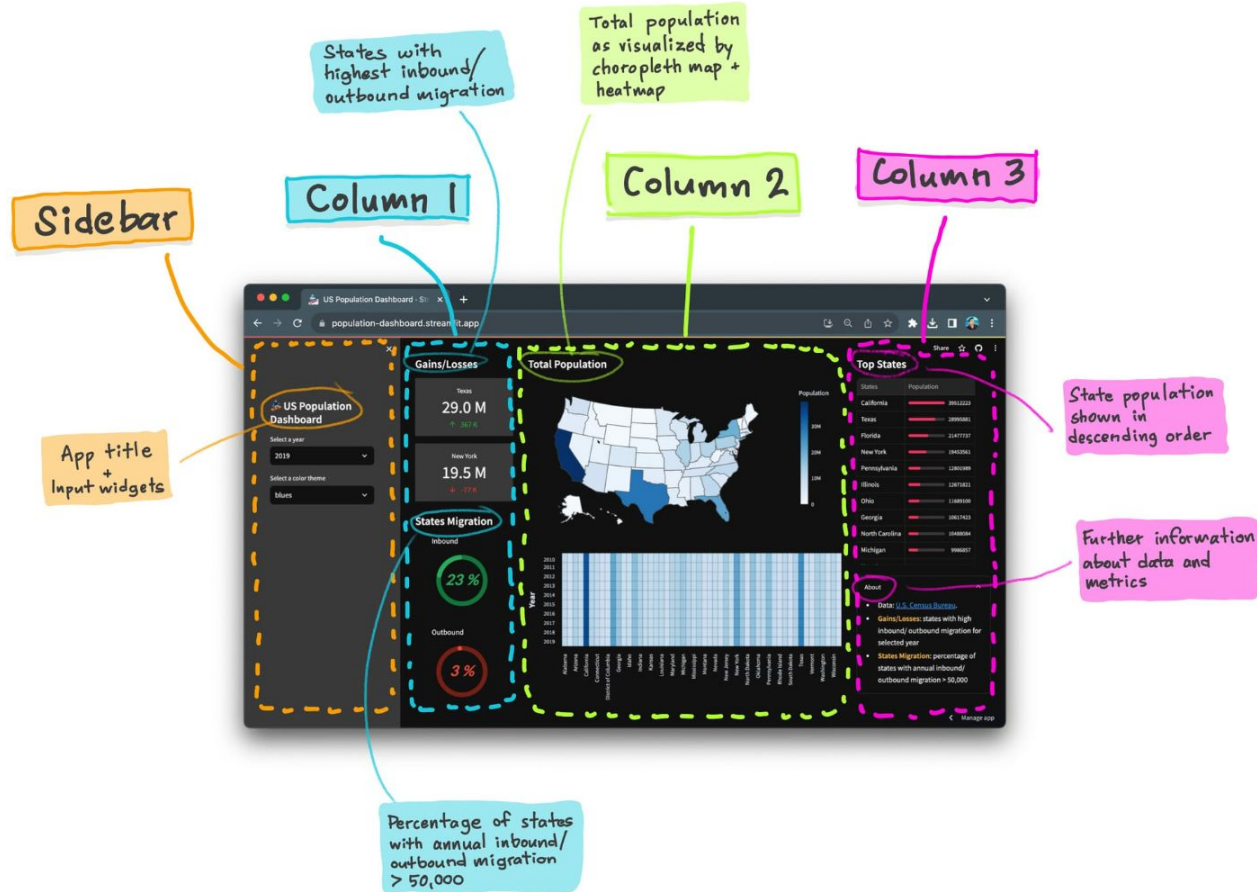
Streamlit

- **Why**

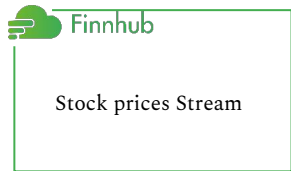
Streamlit?

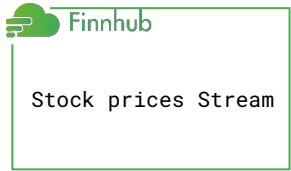
- Python-based framework, no need for JS or CSS
- Declarative syntax and automatic layout management
- Providing built-in support for popular data science libraries

VISUALIZATION: STREAMLIT



Implementation

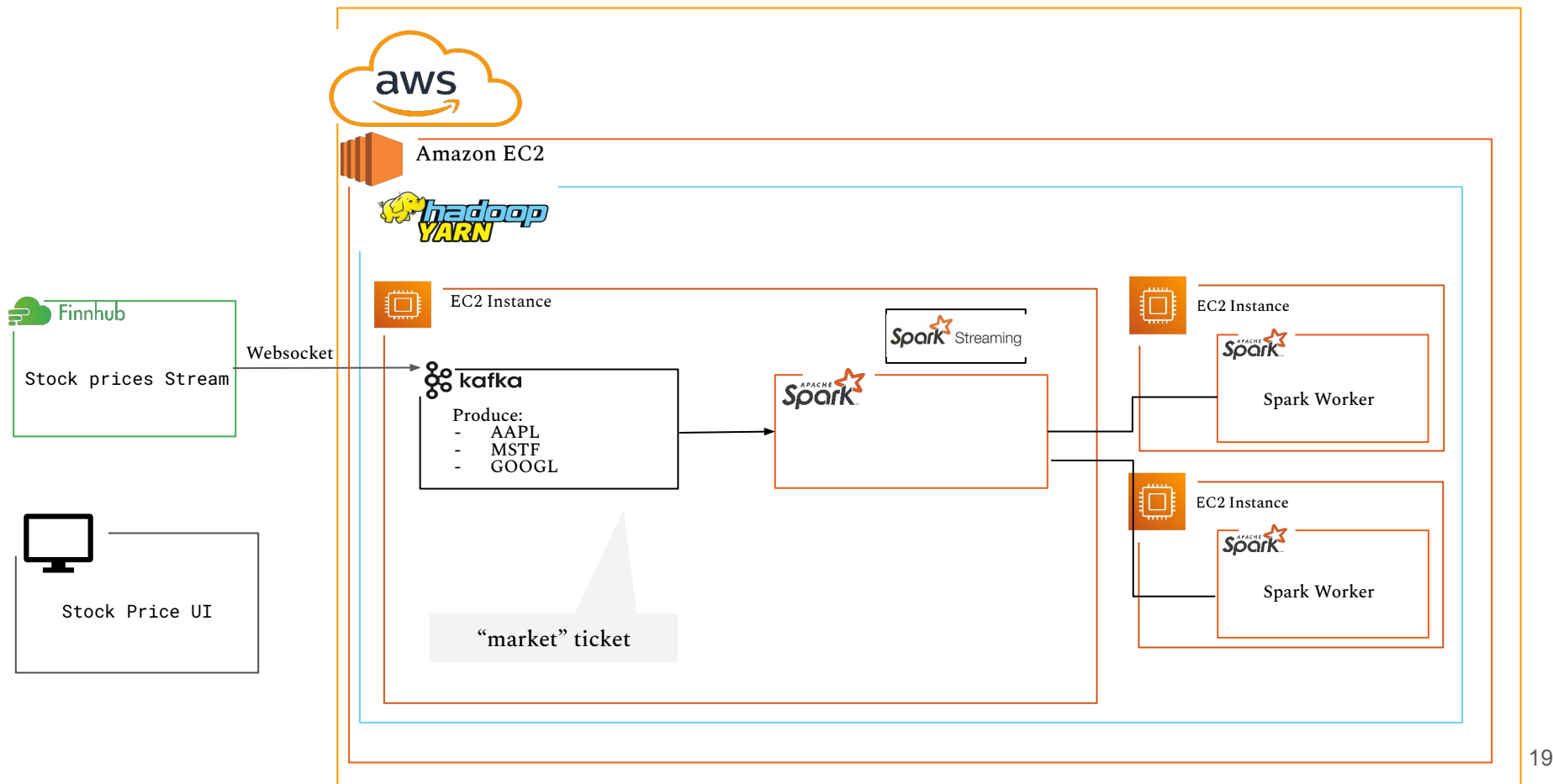


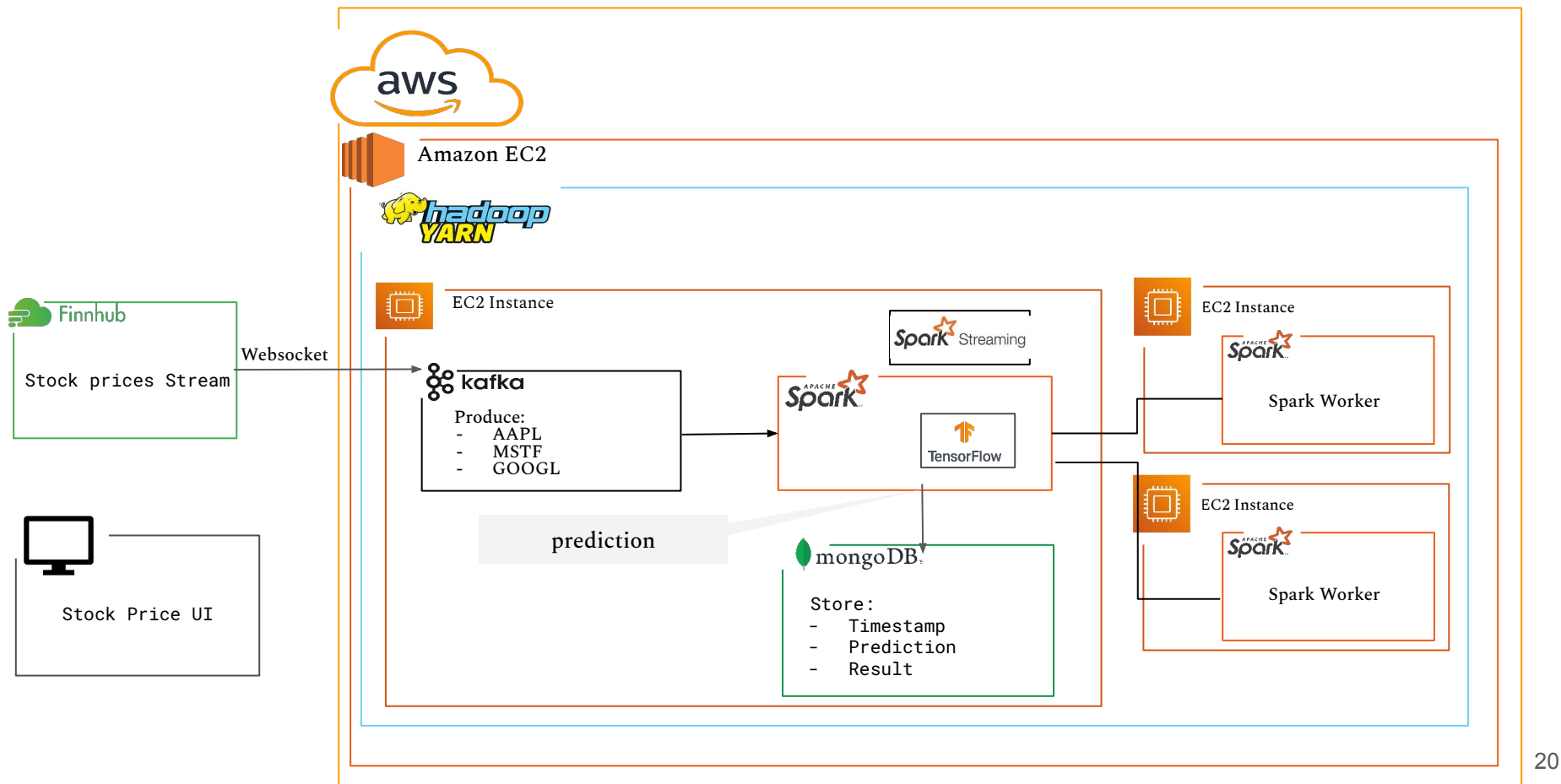


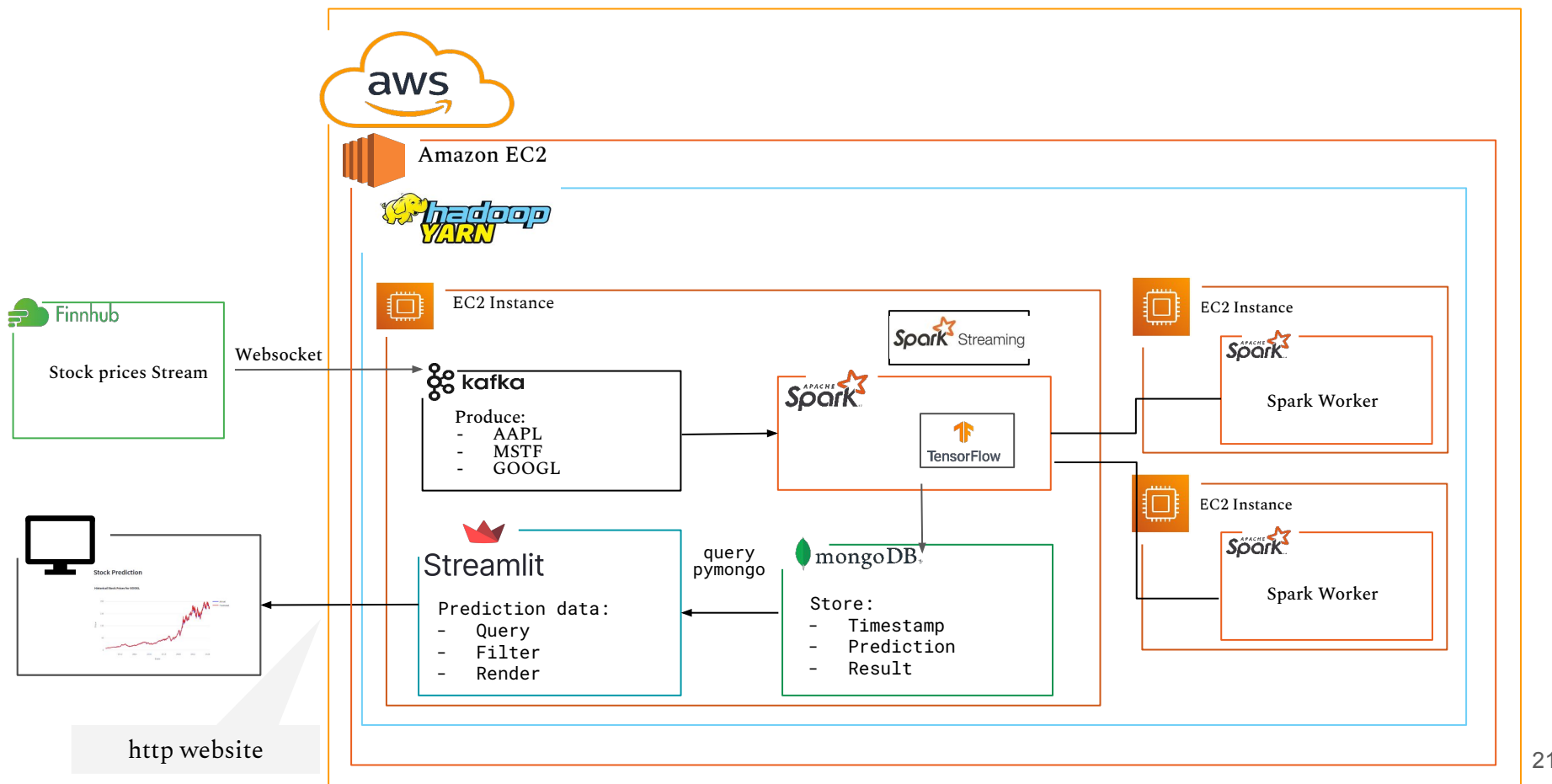












Demo

Results:

- We built a scalable distributed system for predicting stocks

Results:

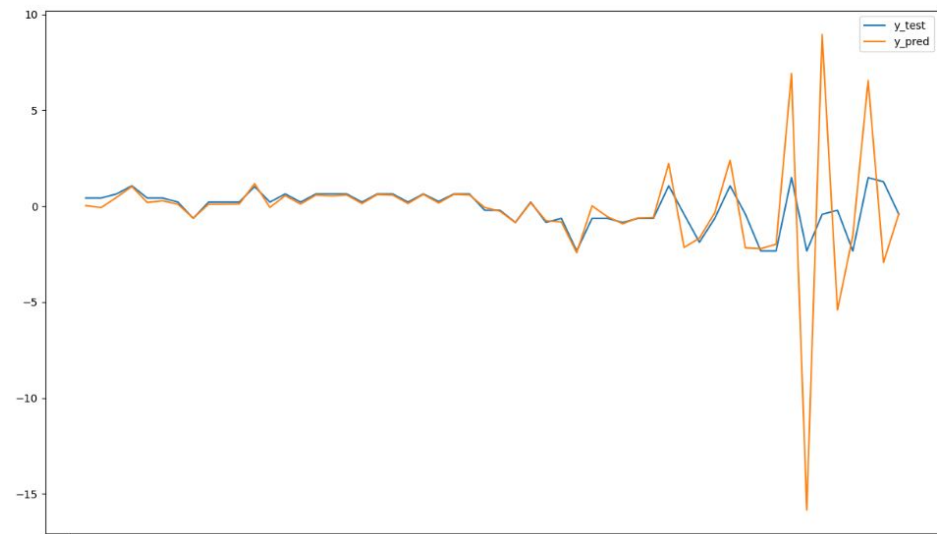
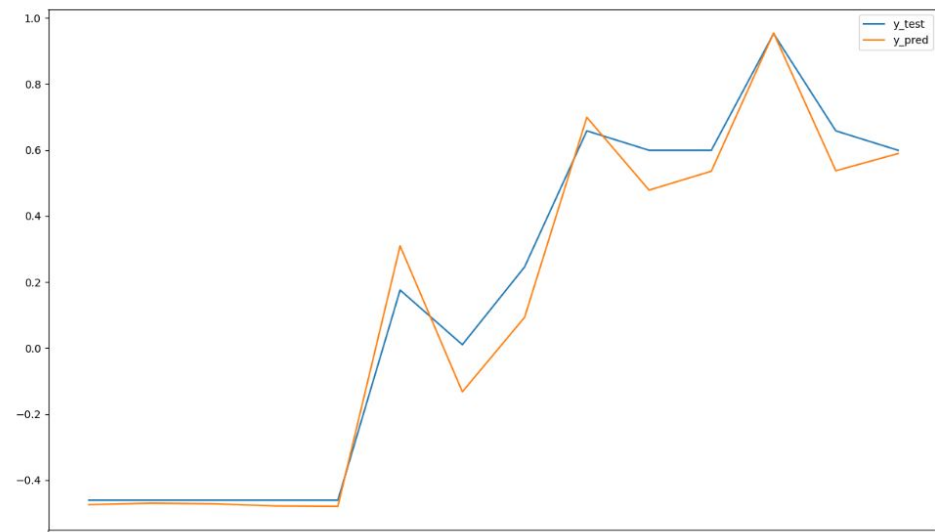
- We built a scalable distributed system for predicting stocks

Future work:

- Use Spark MLlib for distributed model training to improve scalability
- Containerize each part using Docker to ensure reproducibility
- Integrate CLIP or a similar model for real-time feature extraction on streaming data

Thanks for the attention!

Questions time!



Learnings:

- About RTFM
- About documentation.
- About Complexity.

