



## Team Details

- a. **Team name:** Zolo Hallucinators
- b. **Team leader name:** Aravind S
- c. **Problem Statement:** Signal Extraction from Market & News Data

## Brief about the idea

We have created an end-to-end solution that:

- Regularly ingest **market & news data** with an automated pipeline
- Maintain **historical data** in a structured database for reference
- Perform **transformations & feature generation** to create a **time-synced dataset**
- Run **ML models** to predict **tomorrow's price** & analyze **market sentiment**
- Provide **actionable signals**: short-term (tomorrow) & general trends
- Include **explainability** to trace **what influenced each prediction**
- Built a **backtest engine** to measure **strategy performance** & scoring

**Opportunities:** How different is it from any of the other existing ideas?  
USP of the proposed solution?

**Our Solution is has the following USPs, and the ability to break this and provide developers with an intermediate view of data:**

- **Unified pipeline:** Combines **market data + news sentiment** in one end-to-end workflow, unlike many tools that focus on only one data type
- **Time-synced dataset:** Historical + live data **aligned for ML models**, enabling more accurate predictions
- **Explainability built-in:** Can **trace exactly what influenced a prediction**, unlike black-box models in standard platforms
- **Backtesting engine:** Evaluate **strategy performance & rewards** systematically, not just raw predictions
- **Customizable features & models:** Flexibility to **add new indicators, news sources, or ML strategies**
- **Actionable signals:** Provides **both short-term (tomorrow) and general trend signals**, bridging the gap between analytics and trading decisions
- **Data-driven decision support:** Goes beyond dashboards—turns insights into **predictive guidance**

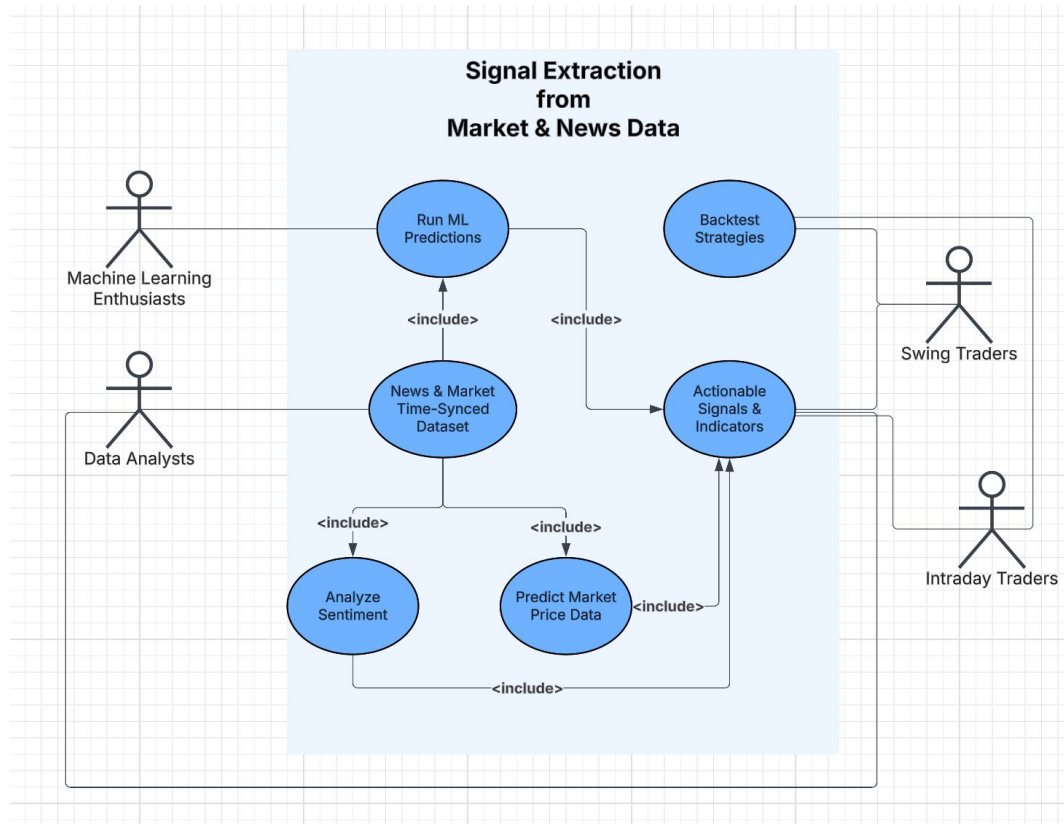
## Opportunities: How will it be able to solve the problem?

- **Unifies multiple data sources** (market + news) into a **single time-synced dataset** for consistent analysis
- Enables **ML-driven predictions** and sentiment scoring on a consolidated view
- Supports **strategy experimentation**: try different models to identify **winning trading strategies**
- Fills a **gap in the market**: currently, no dataset exists with this level of integration and flexibility for testing models
- **Upcoming Feature**: can plug in a **Hugging Face model link**, tune hyperparameters, and test directly on this dataset

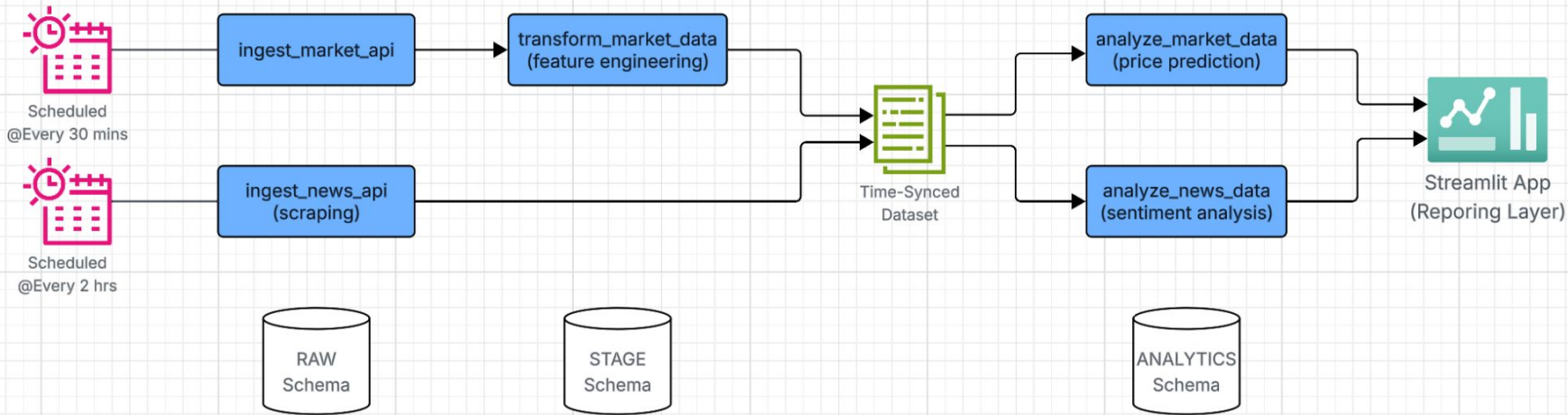
## List of features offered by the solution

- **Automated Data Ingestion** – market & news, scheduled regularly
- **Historical Data Storage** – all your past data in one place
- **Time-Synced Dataset** – aligns market & news for ML-ready analysis
- **Feature Engineering** – generate indicators, signals, sentiment scores
- **ML Predictions** – tomorrow's price & trend forecasting
- **Sentiment Analysis** – gauge market mood from news
- **Explainability** – see exactly **what drove each prediction**
- **Backtesting Engine** – test strategies, measure rewards & performance
- **Actionable Signals** – short-term & general trend guidance

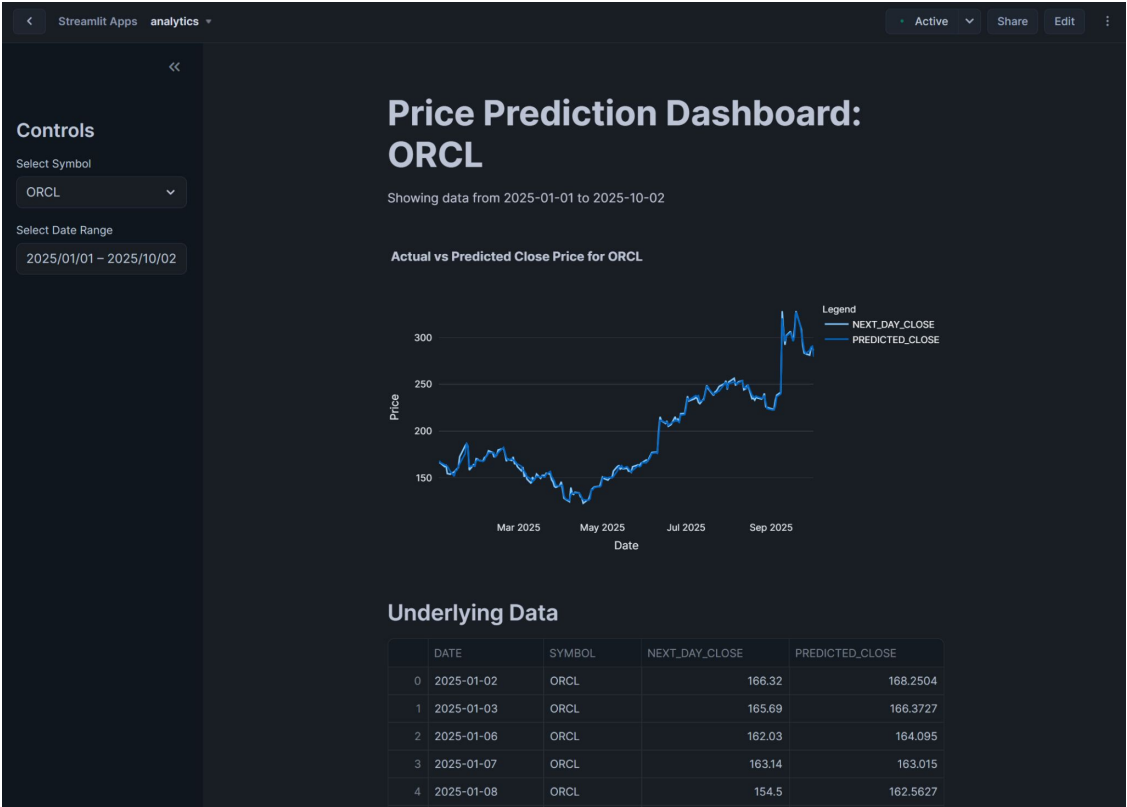
## Use-case Diagram



## Process flow diagram (inc. architecture)



Wireframes/Mock Diagrams of the proposed solution





## Technologies used in the solution

1. Ingestion & Transformation
  - Technologies: Python, APIs (Alpha Vantage, News API)
  - Components: Snowflake Tables, Medallion Arch., Snopipe, stages
2. Data Storage
  - Technologies: SQL, SQL Jobs
  - Storage: Snowflake Databases and Tables.
3. Machine Learning
  - Technologies: Python, ML Libraries, XGBoost
  - In-built functions: Snowpark ML, Cortex Search, Dataframes, Python UDFs
4. Explainability & Backtesting
  - Python, NLP Libraries, Statistics
5. Frontend/Dashboard/Reporting
  - Streamlit Application integrated in Snowflake

## Estimated implementation cost

Title	Comment	Importance	Cost
Snowflake Business Critical Edition	For the News Ingestion Pipeline, the unique number of sites/domains that we have to scrape is in the 200s. This is only possible within this edition to have more number of network policies.	High	4\$/credit ~200 credits/month. Totalling: 800\$/month
<a href="https://newsapi.org">newsapi.org</a> Business Edition	In the free tier, we can only get articles up to a month with 24 hr delay, which is not enough to create high quality datasets for our ML models that need at least 3-5 years of data.	High	449\$/month
<a href="https://alphavantage.co">alphavantage.co</a>	Free is at 25 API requests per day. Monthly plans for premium membership offers 75 API requests per minute	Low	50\$/month
-	-	<b>Total</b>	<b>1299\$/month</b>

Snapshots of the prototype

SIGNAL-EXTRACTION

docs

src

1\_ingestion

infra

local\_test

1\_ingest\_market\_api.ipynb

1\_ingest\_news\_api.ipynb

environment.yml

market\_config.json

news\_config.json

README.md

2\_transformation\_and\_feature\_engine...

1\_transform\_and\_feature\_engineeri...

1\_transform\_market\_data.ipynb

README.md

x1\_transform\_and\_feature\_engineeri...

3\_ml

1\_analyze\_news\_data.ipynb

1\_predict\_market\_data.ipynb

environment.yml

market\_config.json

README.md

4\_frontend

env

infra

utils

.gitignore

LICENSE

README.md

requirements.txt

Streamlit Apps analytics

Active Share Edit

Controls

Select Symbol

ORCL

Select Date Range

2025/01/01 – 2025/10/02

Price Prediction Dashboard: ORCL

Showing data from 2025-01-01 to 2025-10-02

Actual vs Predicted Close Price for ORCL

Price

300

250

200

150

Mar 2025

May 2025

Jul 2025

Sep 2025

Date

Legend

NEXT\_DAY\_CLOSE

PREDICTED\_CLOSE

Underlying Data

	DATE	SYMBOL	NEXT_DAY_CLOSE	PREDICTED_CLOSE
0	2025-01-02	ORCL	166.32	168.2504
1	2025-01-03	ORCL	165.69	166.3727
2	2025-01-06	ORCL	162.03	164.095
3	2025-01-07	ORCL	163.14	163.015
4	2025-01-08	ORCL	154.5	162.5627

9	2025-01-16	ORCL	161.03
...	...	...	...

Model Performance Metrics

RMSE  
2.855

MAE  
2.158

R<sup>2</sup>  
0.997

## Prototype Performance report/Benchmarking

- Major focus was on building the ingestion and transformation pipeline, creating the batch processed time-synced dataset.
- Model Used: XGBRegressor
- Currently the ML model is overfitting for our case, but with more research we can integrate multiple models and better features to get better scores.

9	2025-01-16	ORCL	161.03
<b>Model Performance Metrics</b>			
RMSE			
2.855			
MAE			
2.158			
R <sup>2</sup>			
0.997			

## Future Development

- **Clickable Explainability** – click any data point to see full **lineage of high-impact features**
- **Hugging Face Integration** – plug in your HF model link and test directly on our **time-synced dataset**
- **Strategy & Model Backtesting** – try different ML models & trading strategies to see **which performs best**
- **End-to-End Experimentation** – build, test, and compare models **right from the app frontend**

## Provide links to your:

### GitHub Public Repository:

<https://github.com/Zolo-Hallucinators/Signal-Extraction>

### Demo Video Link:

<https://drive.google.com/file/d/1iWiBB3lU3H3SMDZ82SmHFFqIPTmzbApU/view?usp=sharing>

### Final Product Link (draft):

[https://app.snowflake.com/us-east-1/lac70367/#/streamlit-apps/SIGNAL\\_EXTRACTION\\_DB.UTILS.AINREU5NXYDJBG2Y?ref=snowsight\\_shared](https://app.snowflake.com/us-east-1/lac70367/#/streamlit-apps/SIGNAL_EXTRACTION_DB.UTILS.AINREU5NXYDJBG2Y?ref=snowsight_shared)

**Acknowledgement:**

I would like to sincerely thank you for this opportunity. Working on this problem statement was highly engaging, and using Snowflake for the first time was an eye-opening experience. Coming from a Databricks background, I was impressed by the platform's capabilities, which are extensive and powerful. Overall, this project has been a great learning experience, offering valuable takeaways and insights that I will carry forward.

YOURSTORY

PRESENTS



# HELLO, GCC THE DEV PREMIER (LEAGUE);

CO-PRESENTED BY



INNOVATION PARTNER



Thank you