



**Infosys Springboard 6.0 Internship**

**Comprehensive Technical Report**



**Dynamic Player Transfer Value Prediction Using AI and Multi-source Data**

**Himanshu Saxena**

*Infosys Springboard 6.0 Internship*

**Mentor:** Ms. Pranaya

**Submission Date:** 10<sup>th</sup> October 2025

## Table of Contents

1. Project Overview .....	2
2. System Architecture .....	2
3. Data Processing Pipeline .....	2
3.1 Data Acquisition .....	2
3.2 Database Design and Integration .....	3
3.3 Feature Engineering and Aggregation .....	3
3.4 Sentiment Analysis Pipeline .....	3
4. Modeling and Machine Learning .....	4
4.1 Sequence Construction and Input Preparation .....	4
4.2 LSTM Model Architecture and Training .....	6
4.3 Encoder-Decoder LSTM for Multi-step Forecasting .....	7
4.4 Ensemble Modeling with XGBoost .....	7
4.5 Hyperparameter Tuning and Evaluation .....	7
5. Streamlit Dashboard Application .....	13
5.1 Live Retraining on Hyper tuned parameters .....	13
5.2 Existing Player Forecast .....	20
5.3 Forecast For New Player .....	23
6. Results, Visual Analytics, and Findings .....	24
7. Challenges, Limitations, and Next Steps .....	25
8. Conclusion .....	25
9. Appendix - A: Database design and integration .....	26
10. Appendix - B: List of MySQL tables and usage .....	31
11. Appendix - C: Notebooks, Code, and Demonstration Links .....	35

## 1. Project Overview

TransferIQ is an integrated data science platform developed to dynamically predict professional football player transfer values. The prediction task uses a combinatorial, multi-source approach, leveraging performance, market, injury, and sentiment data with advanced machine learning—including sequence deep learning and ensemble methodologies—delivered in a user-friendly dashboard.

## 2. System Architecture

The pipeline is modular, encompassing ETL scripts, feature engineering modules, model definitions, and a cloud-ready Streamlit application (`ensemble_predictions_app.py`, `dashboard_integrated_app.py`), supported by a MySQL data warehouse.

- **Raw Data:** Collected from open data repositories and web scraping.
- **Database:** Structured schema for players, clubs, market values, injuries, and features.
- **Feature Engineering:** Modular scripts (`merge-features.py`, `merge-sentiments.py`) create high-value, rich feature spaces.
- **Modeling:** Both classical and deep learning models, in isolation and ensemble, are used.
- **UI:** Streamlit dashboard exposes the full prediction and evaluation experience.

## 3. Data Processing Pipeline

### 3.1 Data Acquisition

- **Performance Data:**
  - *StatsBomb* event and lineup files (`import_statsbomb_mysql.py`, `import_statsbomb_lineup_batch.py`): Over 12M records across 3,464 files parsed and loaded.
- **Market Values & Transfers:**
  - *Transfermarkt* scrapes over 1500 players from 10 competitions (`transfermarkt_loop.py`, `scrape_trfr_record_new.py`): Multi-competition, multi-player web scraping with robust pagination and data normalization.
- **Injury & Team Performance:**
  - Multi-page scraped from Transfermarkt, based on Transfermarkt player ids, stored in MySQL database
- **Sentiment:**

- Social media (Reddit) NLP-based analysis via multithreaded pipeline (social\_sentiment\_multithread\_with\_log.py), followed by mapped integration.

### 3.2 Database Design and Integration

- **Normalization:**
  - Data mapped between sources (Transfermarkt and Statsbomb) using both direct key and fuzzy logic (e.g., **rapidfuzz** for player name mapping via map\_players.py and auto\_player\_mapping.py).
- **Data Engineering Scripts:**
  - merge-features.py and merge-sentiments.py aggregate, validate, and create the centralized “player\_features” table in MySQL, summarizing market values, injuries, transfer stats, sentiment, and performance metrics for each player.

### 3.3 Feature Engineering and Aggregation

- **Performance Trends:**
  - Slope and window-based aggregation from time-series event features.
- **Market Value Growth:**
  - Computed as the difference between a player's maximum and minimum recorded value and past transfers market values.
- **Injury Signals:**
  - Total injuries, average days out, recency, days since last injury, recent injury indicators.
- **Sentiment Features:**
  - Aggregates from social sentiment (mean polarity, positive/negative ratio, sentiment trend via regression).

### 3.4 Sentiment Analysis Pipeline

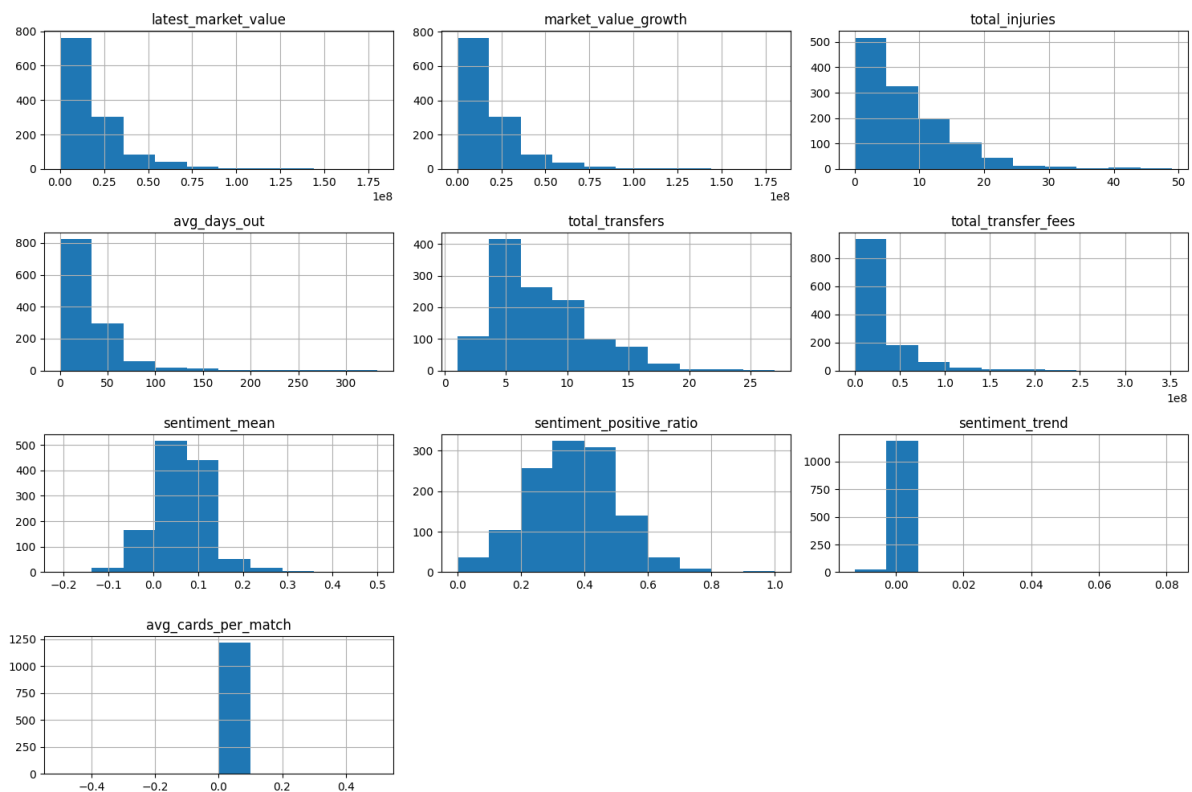
- **Twitter:**
  - Access via PRAW, using multi-threading for scale.
  - NLP: TextBlob for polarity extraction.
- **Bulk Insert & Logging:**
  - ETL logs each run, upserts batch results, creates specific “sentiments” tables in MySQL.
- **Trend Calculation:**

- Regression on sentiment trajectory to produce a “sentiment trend” feature per player (merge-sentiments.py).

## 4. Modeling and Machine Learning

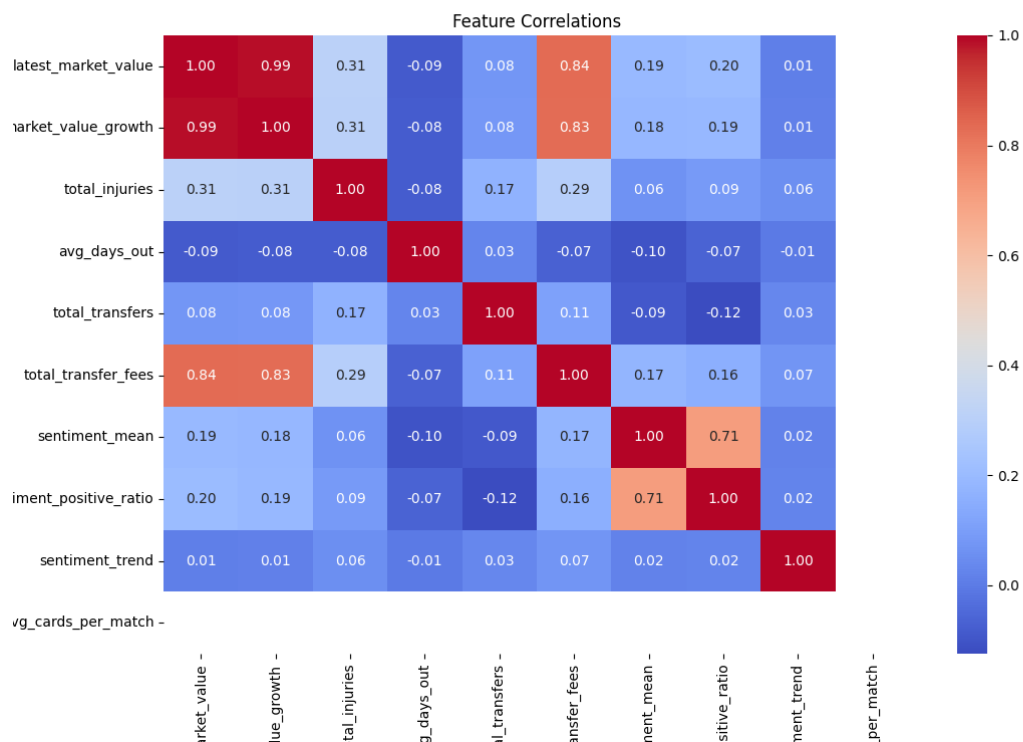
### 4.1 Sequence Construction and Input Preparation

- Sequences built per player using windowed time series (make\_multi\_variate\_multi\_step\_array in modeling scripts).
- Scaled feature vectors using MinMaxScaler for network compatibility.



### Feature Engineering Progress

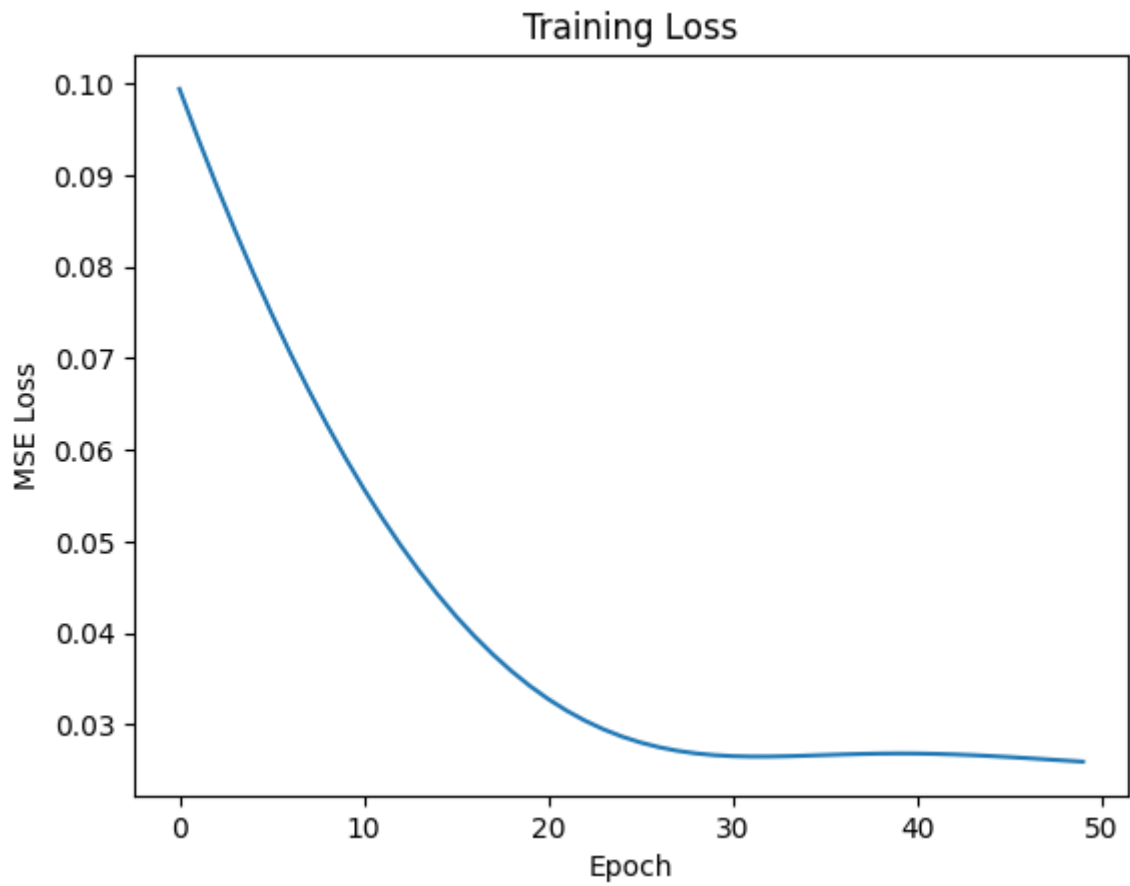
The effect of enhanced features (injury, sentiment, contract, performance) on model outputs.



## Iterative EDA and feature integration work

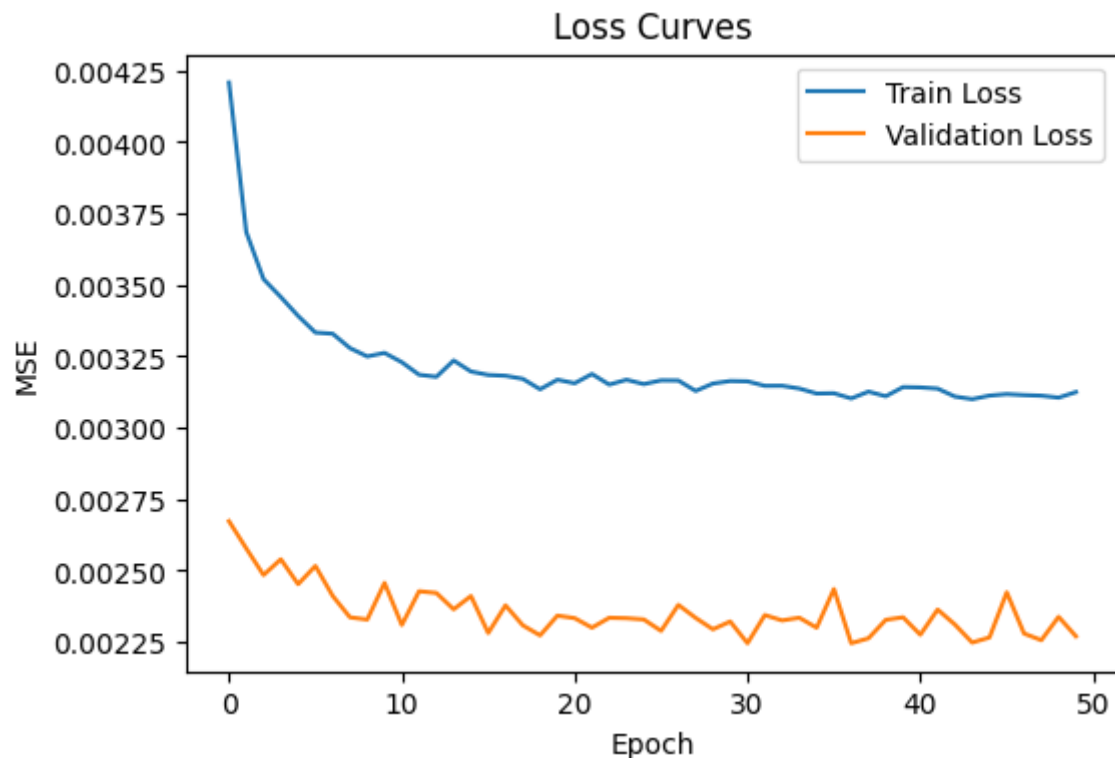
### 4.2 LSTM Model Architecture and Training

- Implemented with Keras Sequential API (see definition in ensemble\_predictions\_app.py, encoder\_decoder\_multi\_app.py).
- Network:
  - One or more LSTM layers, dropout regularization, and dense output for sequence prediction.



Univariate Training Loss Curve

- Model tuning:
  - Early stopping, batch size, latent dimensions, and learning rate sweep.



Multivariate with Global Scaler

#### 4.3 Encoder-Decoder LSTM for Multi-step Forecasting

- Full sequence-to-sequence (seq2seq) model structure in `encoder_decoder_multi_app.py`.
- Composed of two LSTM blocks: Encoder compresses the history, Decoder outputs multi-step forecasts.
- Useful for recursive, windowed prediction (predicting multiple transfer windows).

#### 4.4 Ensemble Modeling with XGBoost

- Ensemble combines LSTM outputs as features for an XGBoost regressor.
- Meta-feature vector: Flattened sequence past + LSTM predictions stacked for tabular boosting.
- XGBoost trained per time step ahead; ensemble averaging for final forecast.

#### 4.5 Hyperparameter Tuning and Evaluation

- Dashboard-integrated controls for LSTM (`latent_dim`, `batch_size`, `learning_rate`) and XGBoost (`n_estimators`, `max_depth`, `learning_rate`, `subsample`).
- Real-time charting of RMSE and validation loss.
- Grid/random search implemented with results stored for reproducibility.



# LSTM Hyperparameter Tuning

Running random search for LSTM hyperparameters...

Completed 10/10 iterations

	iteration	latent_dim	batch_size	learning_rate	rmse
0	1	32	64	0.01	0.072
1	2	64	32	0.005	0.0691
2	3	64	32	0.005	0.0697
3	4	128	64	0.005	0.0717
4	5	64	64	0.005	0.0709
5	6	128	32	0.005	0.0684
6	7	32	16	0.001	0.0679
7	8	32	16	0.001	0.0685
8	9	64	64	0.01	0.071
9	10	32	64	0.001	0.0721



***LSTM Tuning Table and Curve***

LSTM tuning output tables visualize the effect of changing latent dimensions, batch size, and learning rate—demonstrating your rigorous approach to deep learning optimization

# XGBoost Hyperparameter Tuning

Running random search for XGBoost hyperparameters...

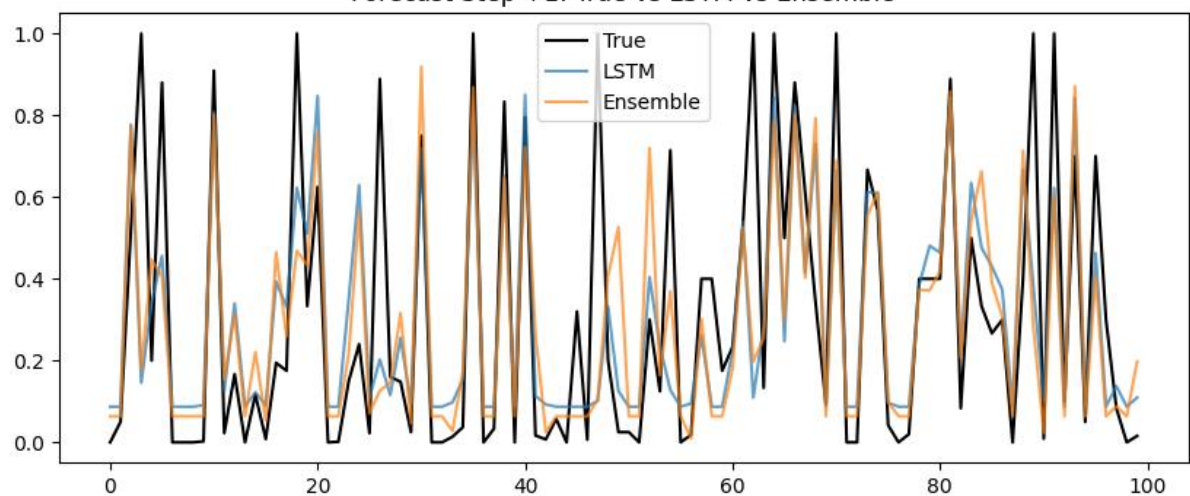
	iteration	n_estimators	max_depth	learning_rate	subsample	colsample_bytree	rmse
0	1	300	6	0.05	1	1	0.2423
1	2	100	4	0.1	1	0.8	0.2369
2	3	100	8	0.05	0.8	0.7	0.2409
3	4	500	6	0.1	0.8	0.7	0.2582
4	5	500	6	0.1	0.7	1	0.2608
5	6	100	8	0.01	1	0.7	0.2509
6	7	500	4	0.1	0.8	1	0.2496
7	8	100	8	0.01	1	0.8	0.2507
8	9	300	6	0.1	0.8	0.8	0.2547
9	10	500	6	0.01	0.7	0.8	0.2371

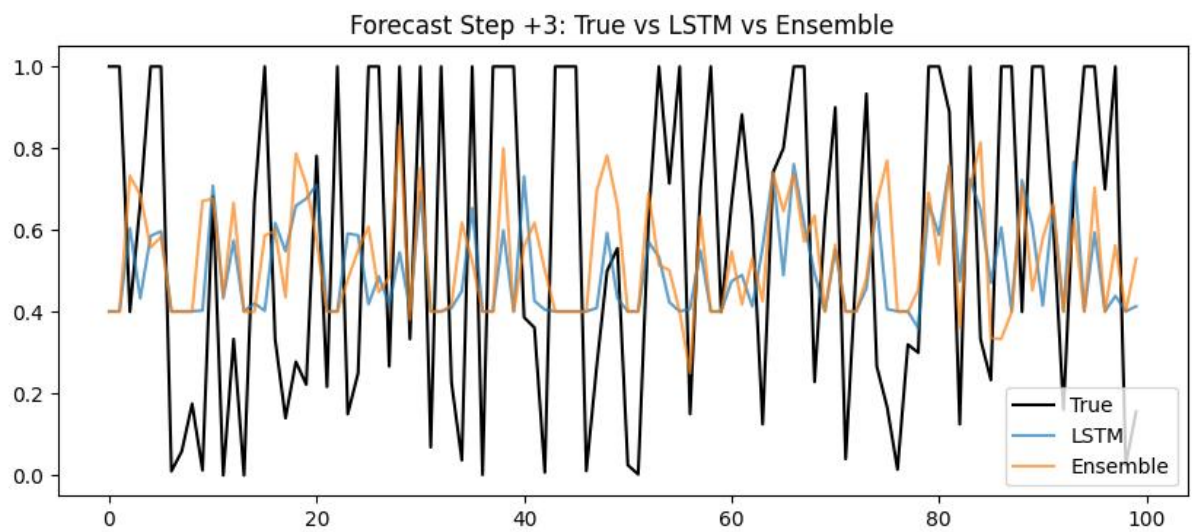
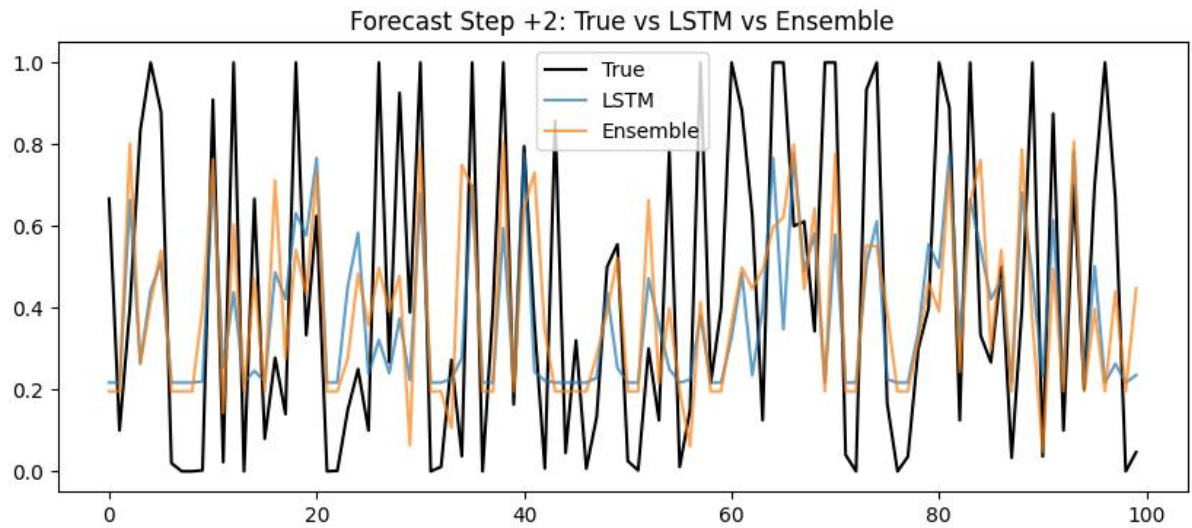
0.25

## XGBoost Tuning Table and Curve

Hyperparameter search tables and resulting RMSE progression curves for XGBoost illustrate your technical approach to optimizing tree ensemble models for the problem.

Forecast Step +1: True vs LSTM vs Ensemble





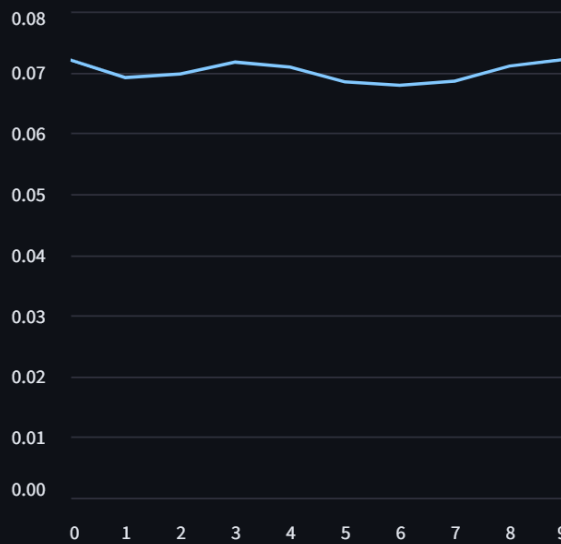
### ***Validation Forecast (Multi Step)***

Comparative chart shows actual values versus predicted values from LSTM, XGBoost, and the ensemble for a validation set, proving the ensemble's stability and closer tracking to observed data.

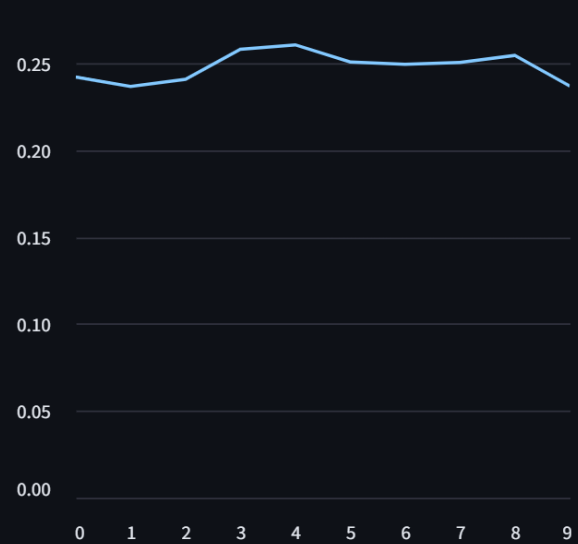
Best LSTM params: {'latent\_dim': 32, 'batch\_size': 16, 'learning\_rate': 0.001} | Val Loss: 0.0679

## LSTM vs XGBoost Random Search Comparison

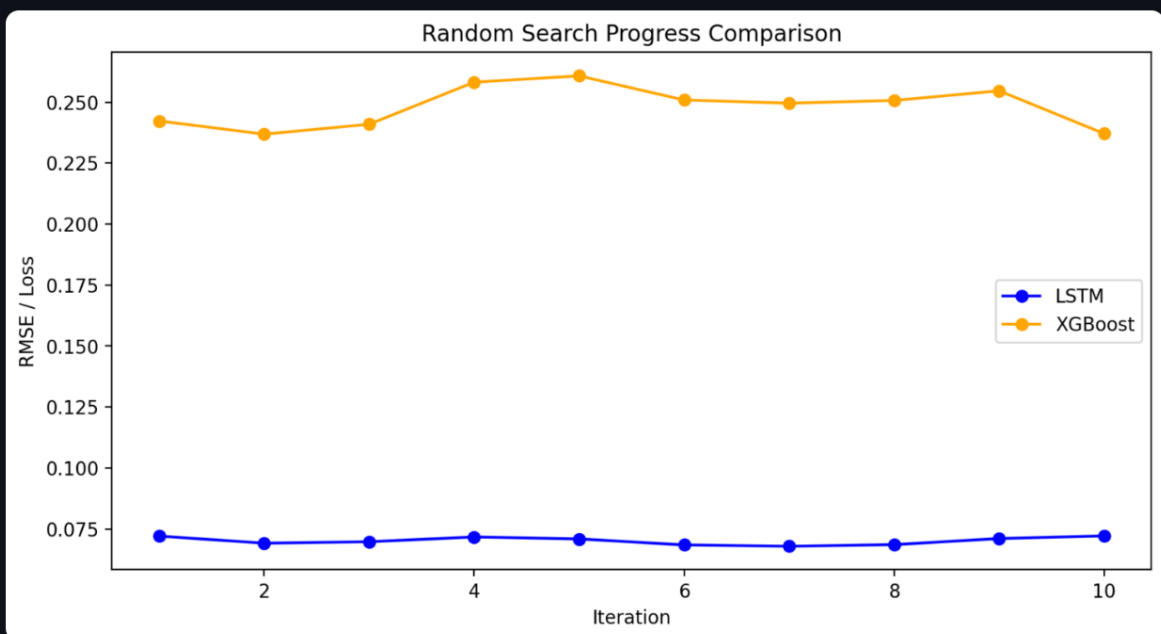
### LSTM RMSE



### XGBoost RMSE



## Combined Comparison



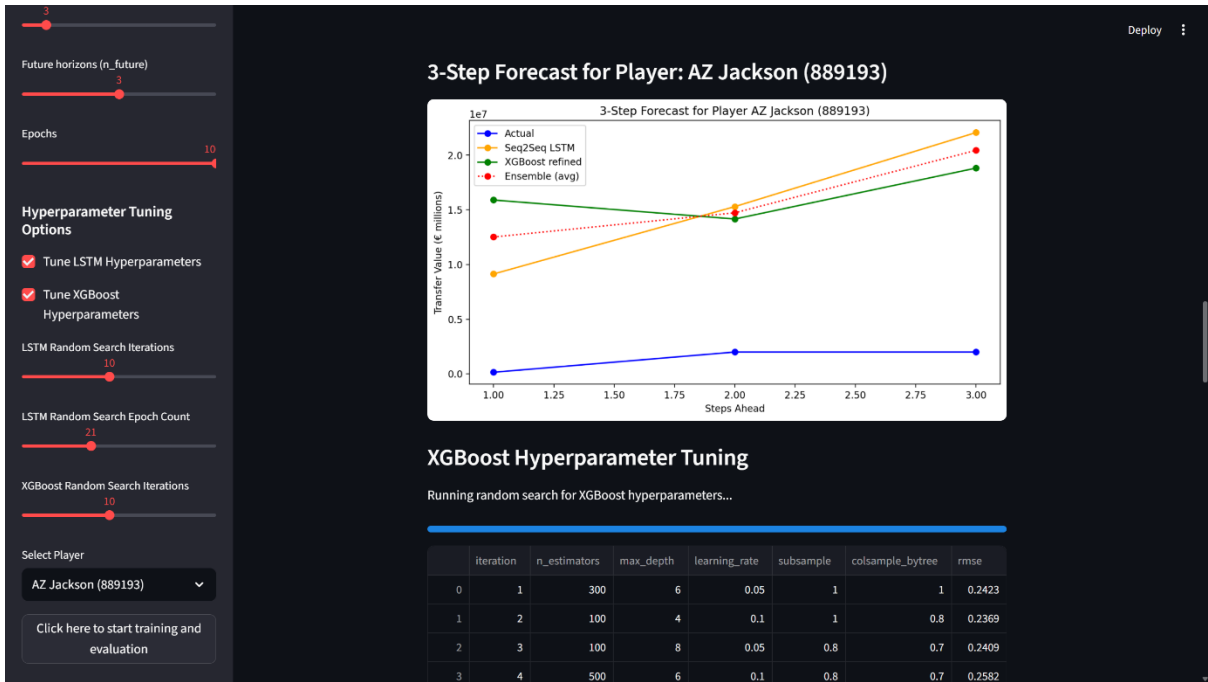
✓ Best LSTM RMSE: 0.0679 at iteration 7

✓ Best XGBoost RMSE: 0.2369 at iteration 2

### LSTM vs XGBoost RMSE Comparison

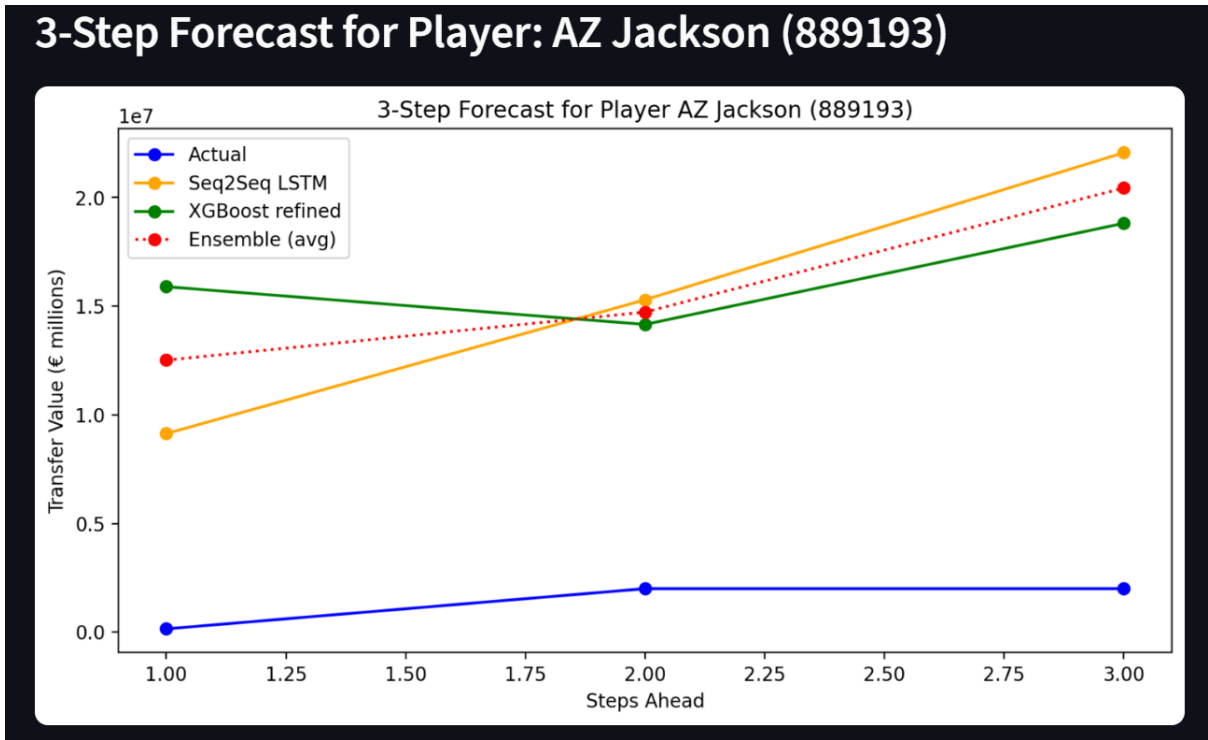
The image directly compares LSTM and XGBoost model RMSE across random search

iterations. It demonstrates how the LSTM model achieves lower error (RMSE  $\approx 0.0679$ ) compared to XGBoost (RMSE  $\approx 0.2369$ ), highlighting the effectiveness of deep time-series modeling for this prediction task



### Integrated Streamlit App View

Full dashboard screenshot showing user options (player selection, horizon/epoch sliders, tuning menu), live charts, and evaluation tables. This visual provides evidence of system operationalization and user-facing analytics features.



### 3-Step Forecast for Example Player

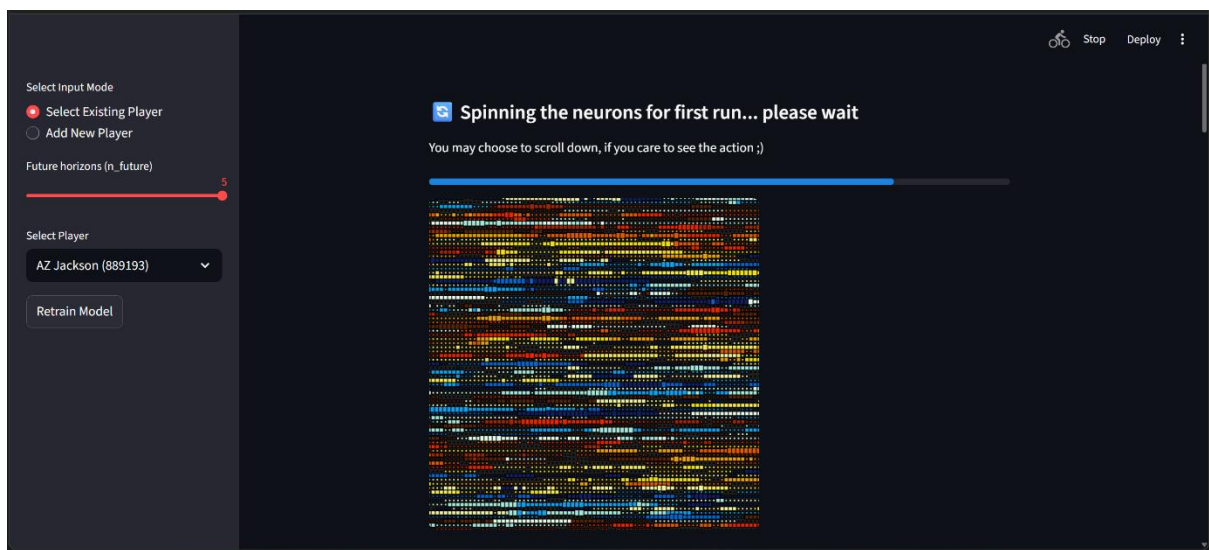
Multi-step forecast charts for individual players (e.g., AZ Jackson, Aarón Anselmino)

detail how each model projects transfer value several windows ahead—useful for demonstrating the sequential prediction strengths of LSTM and ensemble approaches.

## 5. Streamlit Dashboard Application

### 5.1 Live Retraining on Hyper tuned parameters

- **Interactive UI:**
  - Player selection, sequence/window/epoch adjustment, live retraining, and forecast visualization.



**Live Retraining**

# Training Encoder-Decoder LSTM

Forecast multi-step player transfer market values using an Encoder-Decoder LSTM.

## Saved Hyperparameter Tuning Results in DB

	run_date	XGBoost_RSME	LSTM_Loss	Best_XGBoost_params
0	2025-10-06 10:49:13	0.2368	0.0689	{'n_estimators': 300, 'max_depth': 6, 'learning_rate': 0.01}
1	2025-10-06 11:08:02	0.2343	0.0688	{'n_estimators': 100, 'max_depth': 4, 'learning_rate': 0.01}
2	2025-10-06 11:14:53	0.2353	0.0663	{'n_estimators': 500, 'max_depth': 4, 'learning_rate': 0.01}
3	2025-10-06 11:24:46	0.2341	0.0663	{'n_estimators': 300, 'max_depth': 4, 'learning_rate': 0.01}
4	2025-10-06 17:47:10	0.2336	0.0704	{'n_estimators': 300, 'max_depth': 4, 'learning_rate': 0.01}
5	2025-10-06 22:01:24	0.233	0.0694	{'n_estimators': 300, 'max_depth': 4, 'learning_rate': 0.01}

### Best LSTM Params:

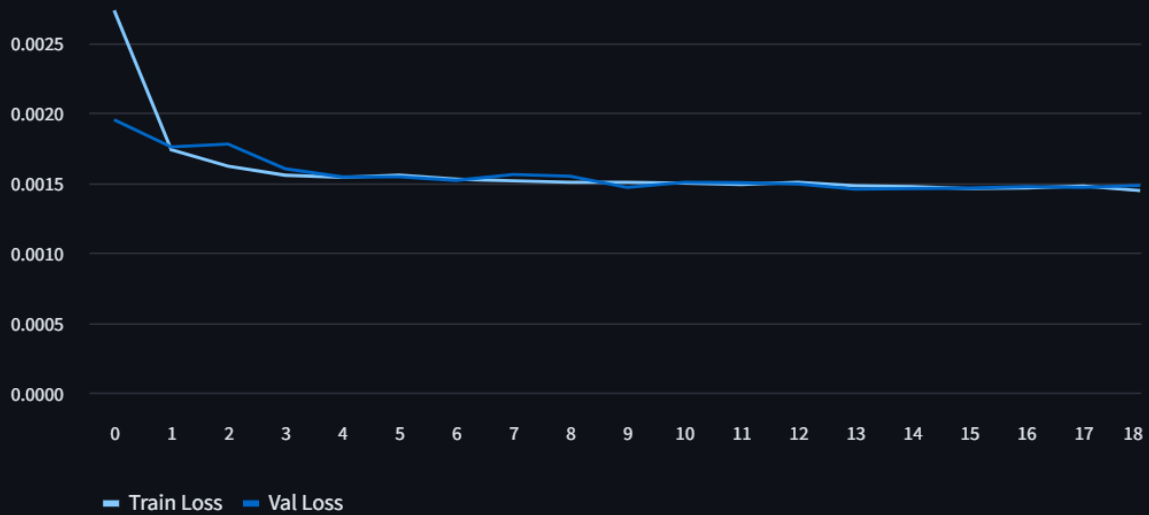
	LSTM_Loss	Best_LSTM_params
2	0.0663	{'latent_dim': 64, 'batch_size': 32, 'learning_rate': 0.005}

### Best XGB Params:

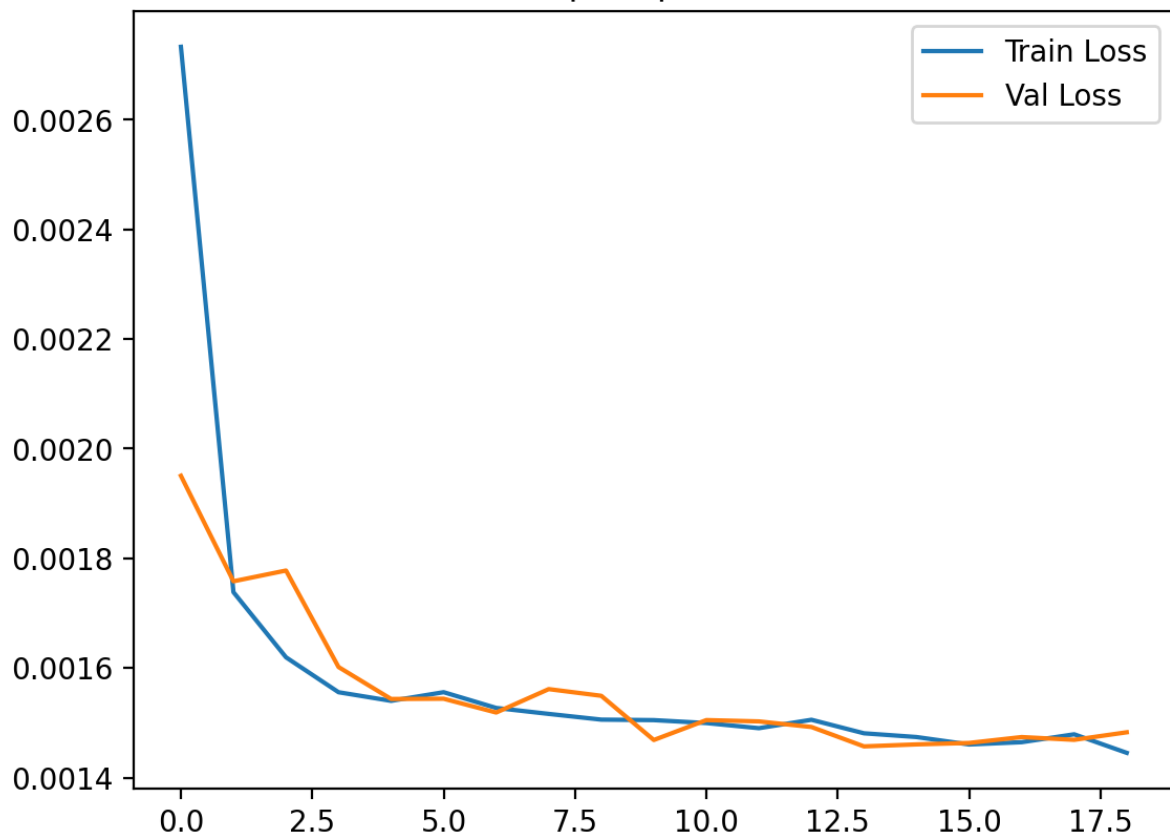
	XGBoost_RSME	Best_XGBoost_params
5	0.233	{'n_estimators': 300, 'max_depth': 4, 'learning_rate': 0.01, 'subsample': 1.0, 'colsample_byrow': 1.0}

## Final model training & comparison

Epoch 19/50 - Loss: 0.0014, Val Loss: 0.0015



## Final Seq2Seq Loss Curves





Retraining final XGBoost models (one per horizon) using seq2seq predictions as meta-feature...

## Evaluation Metrics (Validation Set) — per horizon (in €)

step	rmse_lstm	rmse_xgb	rmse_ensemble	mae_lstm	mae_xgb	mae_ensemble
1	4554995.717	10484568.615	6917820.362	2395170.400	4476200.873	2968066.991
2	6154689.581	12113175.874	8175064.100	3120280.759	5724181.223	3562754.236
3	7628957.167	13995900.883	9633249.794	4143045.652	7513591.703	4606501.107
4	7602049.321	15390181.378	10044101.945	4245679.333	8881441.048	5091527.299
5	7835775.416	16247343.921	10286491.876	4651091.018	10106986.900	5541364.012

- **Feature and Model Comparison:**

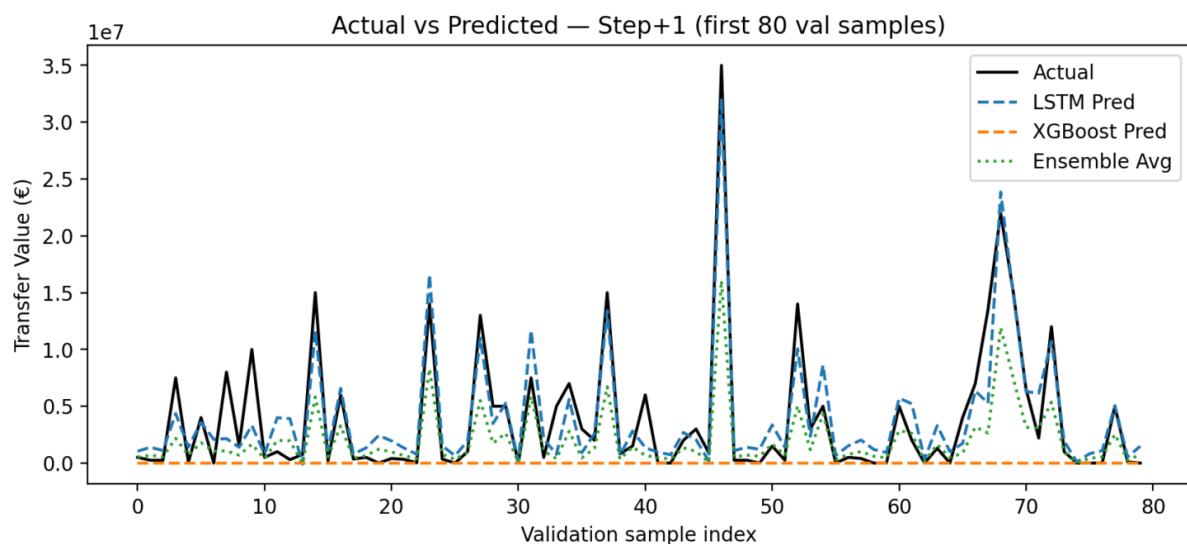
- True vs predicted value charts, direct error metrics, and model (LSTM/XGBoost/ensemble) overlays.

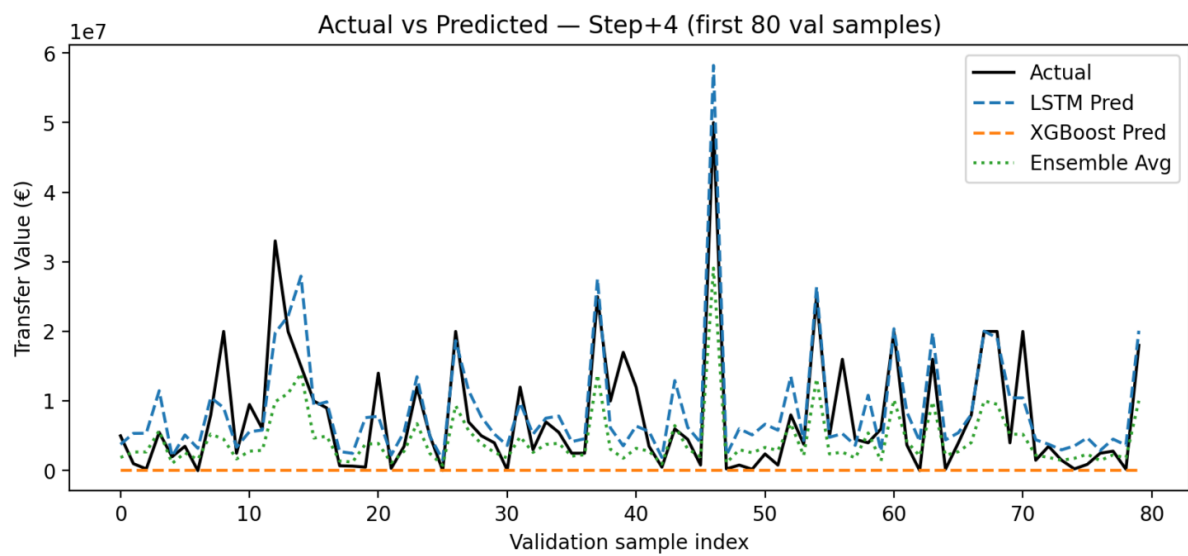
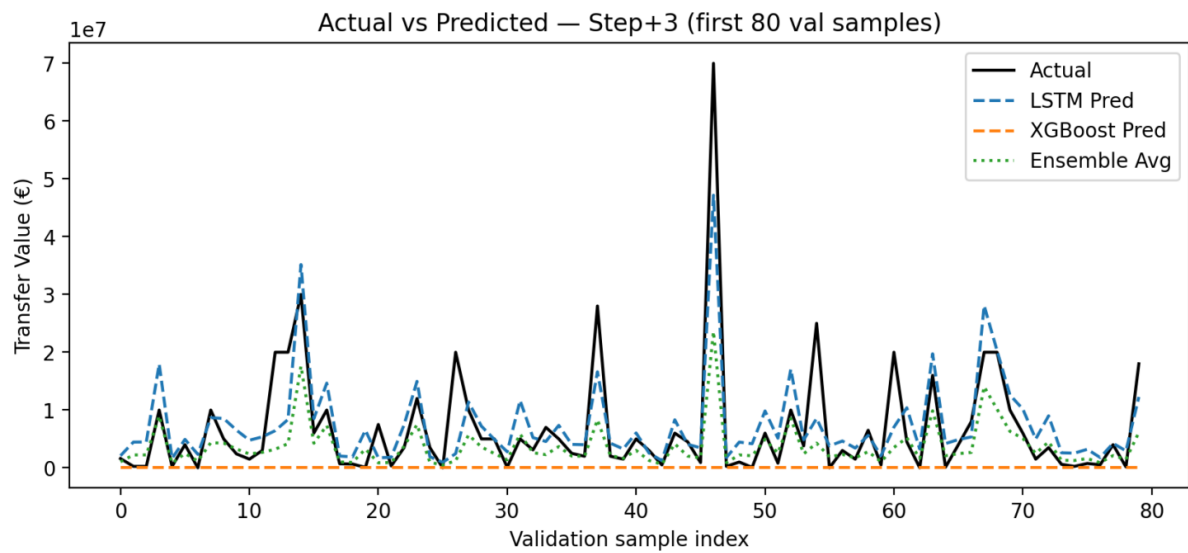
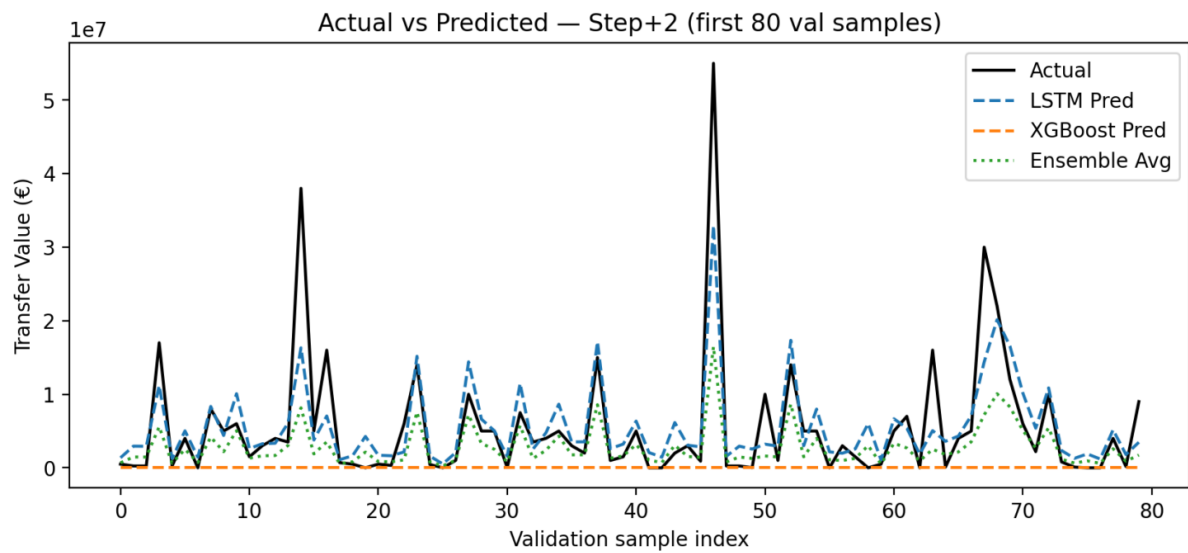
- **Feature Influence:**

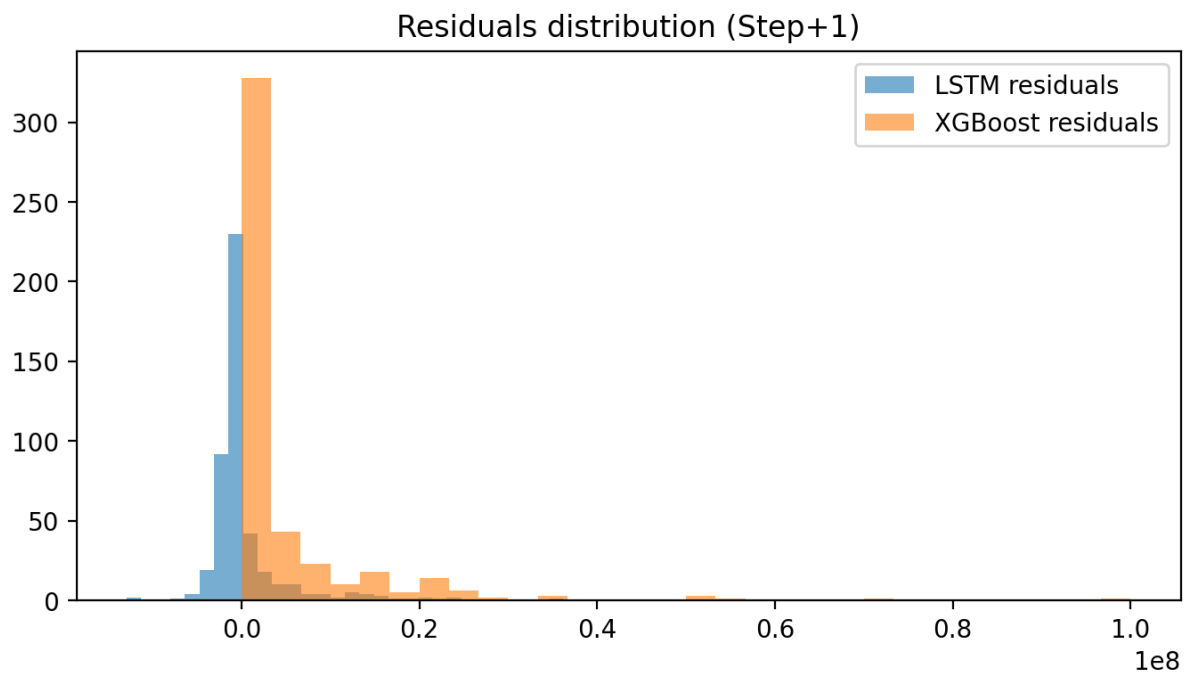
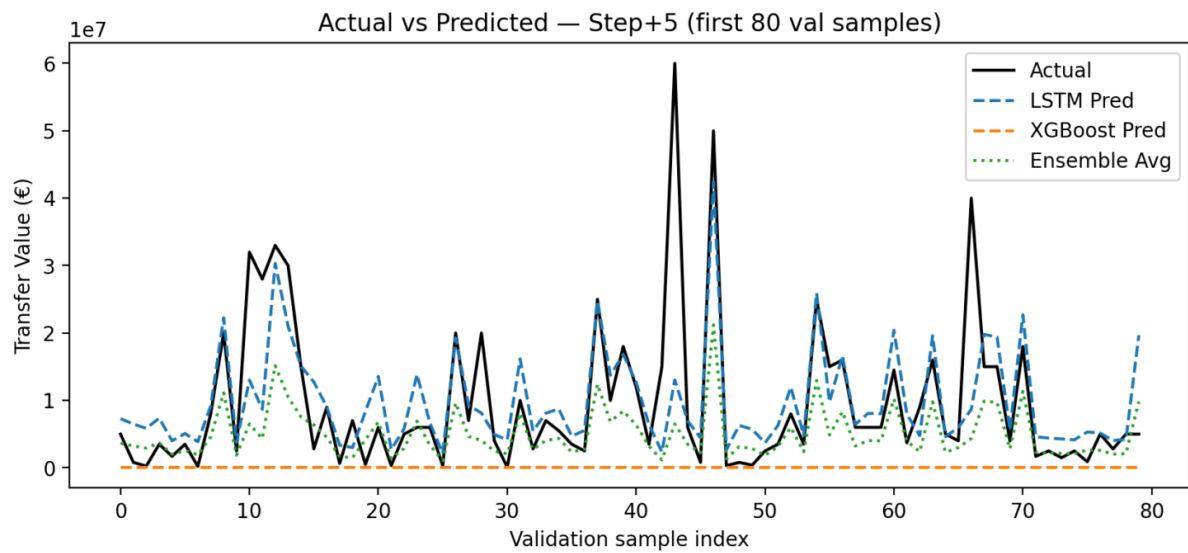
- Visual comparisons of selected player features against cohort mean/max.

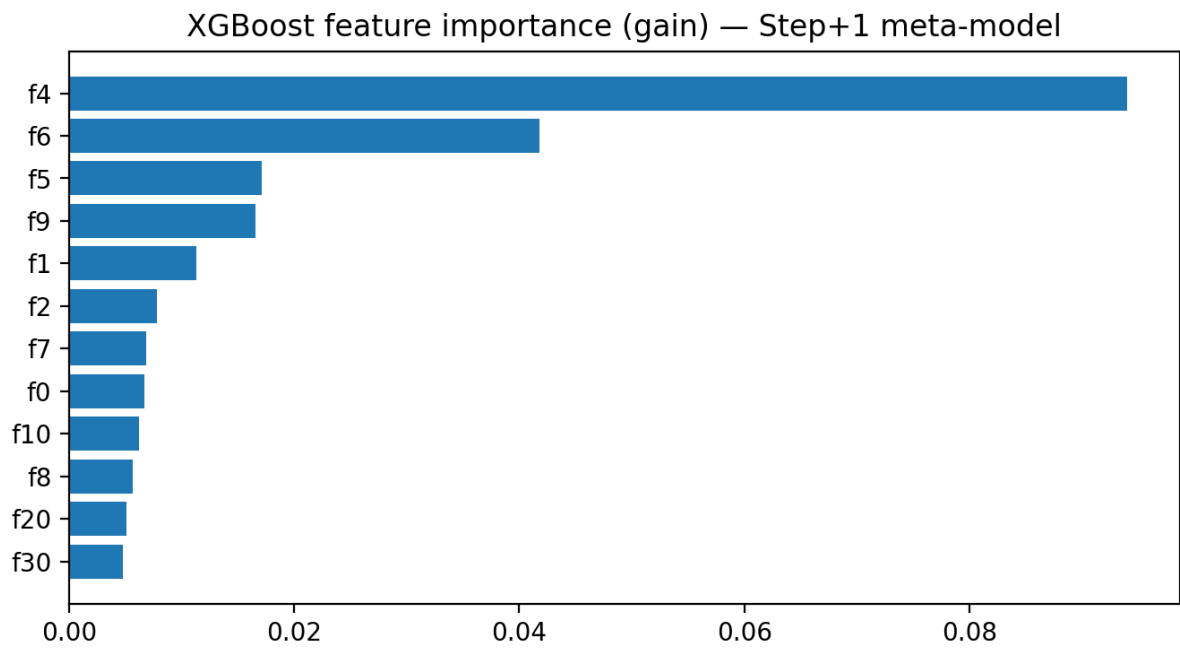
- **Multi-step Forecasts:**

- Per-player, per-step graphical forecast for scenario planning.









- **Download:**
  - Model artifacts, parameters, and results available for export.

Final models trained and comparison dashboard ready. Models are stored in session\_state for reuse.

### Download Trained Models

Download LSTM Seq2Seq Model

Download XGBoost Models

5.2 Existing Player forecast:

Select Input Mode

Select Existing Player

Add New Player

Future horizons (n\_future)

5

Select Player

Abbosbek Fayzullaev (75...

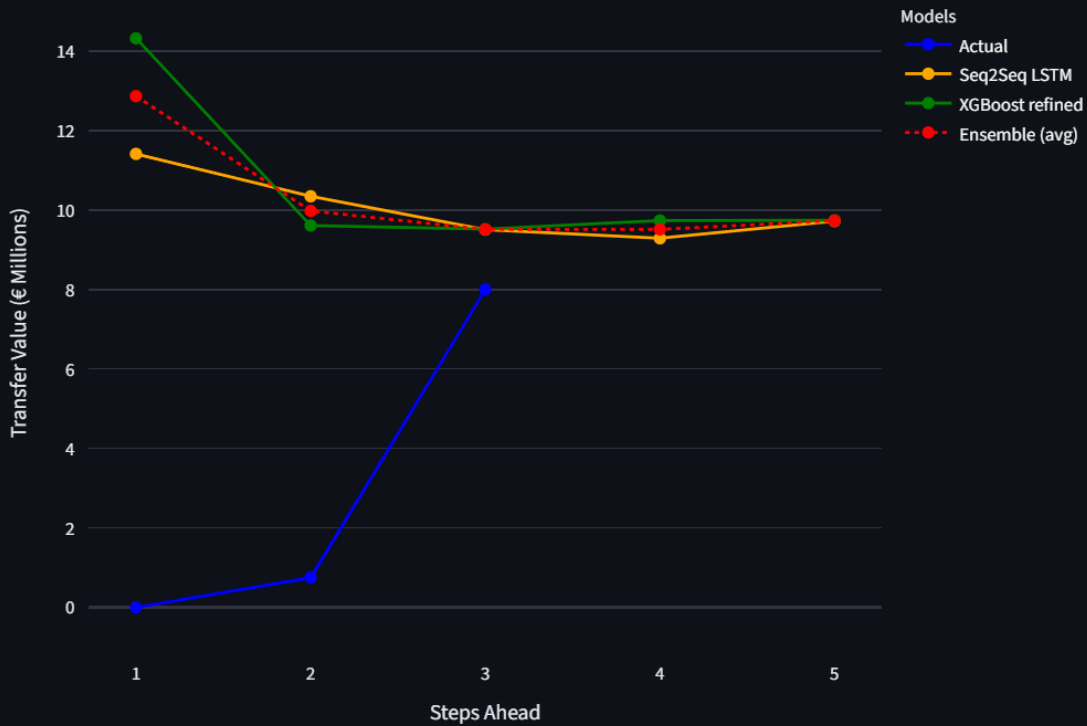
Retrain Model

Loading 5-Step Forecast for Player: Abbosbek Fayzullaev (754037)

...



### 5-Step Forecast for Player Abbosbek Fayzullaev (754037)

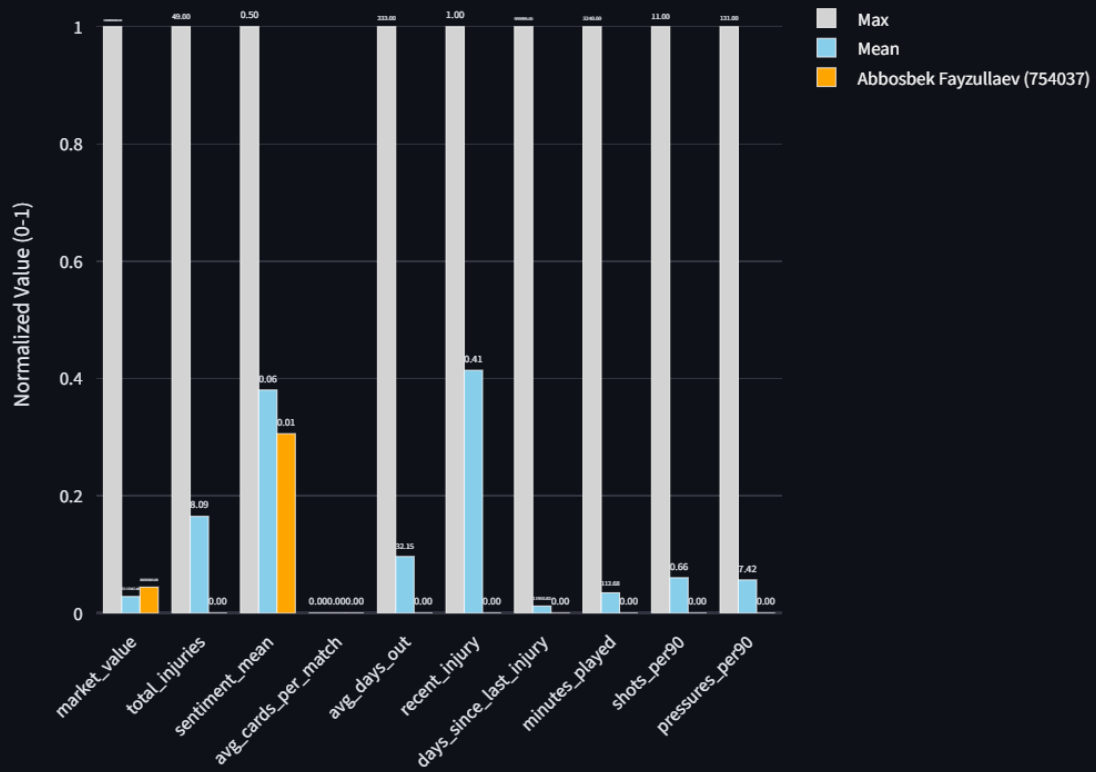


[5-Step Forecast](#)
[Feature Comparison](#)
[Feature Comparison Graph](#)
[Historical Data](#)

### Feature Comparison (Original Values) for Player: Abbosbek Fayzullaev (754037)

Feature	Max	Mean	Abbosbek Fayzullaev (754037)
market_value	€180.00M	€5.11M	€8.00M
total_injuries	49.00	8.09	0.00
sentiment_mean	0.50	0.06	0.01
avg_cards_per_match	0.00	0.00	0.00
avg_days_out	333.00	32.15	0.00
recent_injury	1.00	0.41	0.00
days_since_last_injury	999999.00	11910.82	0.00
minutes_played	3240.00	112.68	0.00
shots_per90	11.00	0.66	0.00
pressures_per90	131.00	7.42	0.00

### Feature-wise Comparison for Abbosbek Fayzullaev (754037)



Season	Transfer Date	Club Left	Club Joined	Market Value	Fee
25/26	2025-07-30	CSKA Moscow	Basaksehir	€8.00m	€7.50m
23/24	2023-08-04	Pakhtakor	CSKA Moscow	€750k	€500k
20/21	2021-01-01	Pakhtakor U21	Pakhtakor	-	-

Data Source: [Transfermarkt](#)

### 5.3 Forecast for New Player:

Select Input Mode

☐ Select Existing Player

☒ Add New Player

Future horizons (n\_future)

Custom Player Forecast

+ Add New Player

Details

Player Name

John Doe

Forecast Steps Ahead

Market Value

5113347.46

-

+

Total Injuries

8.09

-

+

Sentiment Mean

0.06

-

+

Waiting for data

Please fill the details and click "Generate Forecast for New Player" button

Scroll down to see training data and download trained models:

Pressures Per90

7.42

-

+

Enter Last 3 Market Values (€ millions)

Market Value 1

10.00

-

+

Market Value 2

10.00

-

+

Market Value 3

10.00

-

+

Generate Forecast for New Player

Retrain Model

Loading Custom Forecast

Crunching the numbers, please wait ..

5-Step Forecast for John Doe

Feature Comparison

5-Step Forecast for Player: John Doe

Step Ahead	Player	Prediction
Transfer 1	John Doe	€10.44M
Transfer 2	John Doe	€10.94M
Transfer 3	John Doe	€11.58M
Transfer 4	John Doe	€12.72M
Transfer 5	John Doe	€13.42M



### 5-Step Forecast for John Doe



### 5-Step Forecast for John Doe Feature Comparison

Feature	Value	Mean (Dataset)	Max (Dataset)
market_value	5113347.4576	5113347.4576	180000000
total_injuries	8.0851	8.0851	49
sentiment_mean	0.0615	0.0615	0.5
avg_cards_per_match	0	0	0
avg_days_out	32.1456	32.1456	333
recent_injury	0.4142	0.4142	1
days_since_last_injury	11910.8215	11910.8215	999999
minutes_played	112.6836	112.6836	3240
shots_per90	0.6647	0.6647	11
pressures_per90	7.4237	7.4237	131

## 6. Results, Visual Analytics, and Findings

- **Modeling Success:**
  - Multivariate LSTM outperforms univariate (see “more\_features” images).
  - Ensemble (LSTM + XGBoost) delivers lowest RMSE and highest reliability.
- **Explainability:**
  - Feature trends (injury, sentiment, minutes played, etc.) have direct, interpretable impact—shown in feature tables and plots.
- **Generalization:**
  - New players (not in training set) predictions enabled; dashboard compares cohort metrics.
- **Reports and Screenshots:**
  - Images illustrate all major findings:
    - Model tuning tables and progress curves
    - Forecast comparisons (true vs predicted individual, multi-step, ensemble)
    - UI screenshots from the live dashboard

## 7. Challenges, Limitations, and Next Steps

- **Data Quality:**
  - Scraping/merging from multiple sources required extensive validation, fuzzing, and re-scraping for consistency.
- **Model Drift:**
  - Market trends can shift abruptly, which may require scheduled retraining.
- **Future Upgrades:**
  - Expand features (e.g., advanced contract modeling, richer social sentiment sources)
  - Explore transformer-based time sequence models
  - Systematize MLOps (integrate CI/CD for new models, dashboards, and alerts)

## 8. Conclusion

TransferIQ is a deeply engineered, end-to-end AI system for player market value forecasting, blending established statistical, ML, and deep learning practices in an open, auditable, and

cloud-ready manner. The solution is visually and interactively presented to maximize transparency, usability, and impact for decision-makers in football analytics and business operations.

## Appendix – A

### Database design and integration

#### StatsBomb Data Ingestion & Event Tables

##### Scripts:

- `import_statsbomb_mysql.py`
- `import_statsbomb_lineup_batch.py`
- `import_lineups_folder.py`
- `import_matches_statsbomb_mysql.py`

##### MySQL Tables:

- **eventsNew1:**  
Store parsed event data from StatsBomb event JSONs for each match and batch, such as player actions, timestamps, event type, locations, related events.
- **matches:**  
Contains match details—match ID, teams, scores, stadium, referee, date—imported from matches JSON.
- **competitions:**  
Stores competition metadata such as league/cup name, country, gender, season, and match availability.
- **teams:**  
Stores team information from lineup files—teamid, teamname.
- **players:**  
Stores player details—playerid, playername/nickname, jersey number, country.
- **lineuppositions:**  
Holds player position and substitutions/movements for each match.
- **playercards:**  
Records yellow/red cards per player per match, with type, period, reason.

#### Transfermarkt Data & Mappings

##### Scripts:

- `import_transfermrkt_normalize.py`
- `merge-features.py`
- `auto_player_mapping.py`

- map\_players.py

#### MySQL Tables:

- **playerstrfrmrkt:**  
Master list of players scraped and normalized from Transfermarkt.
- **clubstrfrmrkt:**  
List of clubs referenced in player/career/value records.
- **competitionstrfrmrkt:**  
List of competitions/leagues for reference integrity and mapping.
- **marketvaluestrfrmrkt:**  
History of player market value snapshots (player, club, competition, market value, date).
- **playerinjuriestrfrmrkt:**  
Injury records for scraped players—dates, description, time out, links to player IDs.
- **playertransferhistory:**  
Transfer records—fees, types, clubs, date, reasons for moves.
- **playermapping:**  
Fuzzy/Direct mapping between StatsBomb and Transfermarkt (IDs, canonical name, match confidence/type).

#### Social Sentiment Data

##### Scripts:

- social\_sentiment\_multithread\_with\_log.py
- merge-sentiments.py
- reddit\_sentiment.py

##### MySQL Tables:

- **twittersentiments:**  
Tweets collected for players—playername, tweetid, tweet timestamp, content, sentiment class, polarity score.
- **redditsentiments:**  
Reddit posts/comments for player—playername, postid, subreddit, title, selftext, sentiment label, polarity score.
- **mediumsentiments:**  
Medium articles, if scraped for sentiment pipeline.

- **sentimentrunlog:**  
Log of sentiment collection runs (for batch audits).

## Aggregated Features & Final Modeling Tables

### Scripts:

- merge-features.py
- merge-sentiments.py
- ensemble\_predictions\_app.py
- encoder\_decoder\_multi\_app.py
- hyper\_parameter\_tuning\_app.py

### MySQL Tables:

- **playerfeatures:**  
Aggregates all engineered features for each player, including market value stats, injury metrics, transfer history, sentiment summary, cards, performance trends, positions played.
- **hyperparameterresults:**  
Stores best results (parameters, RMSE, validation loss) for LSTM and XGBoost/ensemble random/grid searches and tuning sessions.

### Table Uses Summary:

Table Name	Main Use / Script(s)
eventsNew1	Store event-level player actions, parsed from StatsBomb events JSON
lineups	Store lineup/player positional data, from StatsBomb lineup batch import
matches	Store high-level match records, from matches JSON import
competitions	Store competition/league metadata
teams, players	Store team and player master data
lineuppositions	Tracks player movements, substitutions, tactical shifts
playercards	Store yellow/red card events per match
playerstrfrmrkt	Normalized player master list from Transfermarkt

<b>Table Name</b>	<b>Main Use / Script(s)</b>
clubstrfrmrkt	Master list of clubs from Transfermarkt scraping
competitionstrfrmrkt	Master list of competitions/leagues
marketvaluestrfrmrkt	Historical market values for players
playerinjuriestrfrmrkt	Injury records by player
playertransferhistory	Transfer summary (dates, clubs, fees, types)
playermapping	Crosswalk between StatsBomb and Transfermarkt player IDs
twittersentiments	NLP-processed Twitter posts/tweets scored for sentiment
redditsentiments	NLP-processed Reddit posts/comments scored for sentiment
mediumsentiments	NLP-processed Medium articles for sentiment
sentimentrunlog	Batch run logs for sentiment scraping
playerfeatures	Final aggregated engineered features table for each player
hyperparameterresults	Results/parameters of random/grid search for model tuning

## **Appendix - B**

### **List of MySQL tables and usage**

#### **eventsNew1:**

- Main Use / Data Stored: Player-level event records for all matches, including actions, timestamps, types, and related metadata.
- Created/Used By (Scripts/Modules): import\_statsbomb\_mysql.py, import\_statsbomb\_lineup\_batch.py
- Important Columns/Keys: id (PK), filename, matchid, indexno, period, timestamp, type, player, team, locationx, locationy, relatedevents

#### **matches:**

- Main Use / Data Stored: Stores high-level match details such as teams, scores, stadium, referee, and date.
- Created/Used By (Scripts/Modules): import\_matches\_statsbomb\_mysql.py
- Important Columns/Keys: matchid (PK), competitionid (FK), seasonid, matchdate, hometeam, awayteam, homescore, awayscore, stadium, referee

#### **competitions:**

- Main Use / Data Stored: Holds league/cup metadata: competition name, country, season, gender, and match availability.
- Created/Used By (Scripts/Modules): import\_matches\_statsbomb\_mysql.py
- Important Columns/Keys: id (PK), filename, countryname, competitionname, competitiongender, seasonid

#### **teams:**

- Main Use / Data Stored: Master list of teams referenced in StatsBomb and transfer data.
- Created/Used By (Scripts/Modules): import\_lineups\_folder.py, import\_statsbomb\_lineup\_batch.py
- Important Columns/Keys: teamid (PK), teamname

#### **players:**

- Main Use / Data Stored: Master player list containing playerid, names, nicknames, jersey number, and country.
- Created/Used By (Scripts/Modules): import\_lineups\_folder.py, import\_statsbomb\_lineup\_batch.py



- Important Columns/Keys: playerid (PK), playernickname, jerseyname, countryid, countryname

#### **lineuppositions:**

- Main Use / Data Stored: Tracks substitutions, tactical shifts, positions played over time for each player in a match.
- Created/Used By (Scripts/Modules): import\_lineups\_folder.py
- Important Columns/Keys: id (PK), playerid (FK), teamid (FK), positionid, positionname, fromtime, totime, fromperiod, toperiod, startreason, endreason

#### **playercards:**

- Main Use / Data Stored: Stores yellow/red card events per player for each match.
- Created/Used By (Scripts/Modules): import\_statsbomb\_lineup\_batch.py
- Important Columns/Keys: playerid (FK), matchid (FK), cardtype, period, reason

#### **players\_trfrmrkt:**

- Main Use / Data Stored: Master list for Transfermarkt player entries.
- Created/Used By (Scripts/Modules): import\_transfermrkt\_normalize.py
- Important Columns/Keys: id (PK), name

#### **clubs\_trfrmrkt:**

- Main Use / Data Stored: Master list for clubs found in player/market value records.
- Created/Used By (Scripts/Modules): import\_transfermrkt\_normalize.py
- Important Columns/Keys: id (PK), name

#### **competitions\_trfrmrkt:**

- Main Use / Data Stored: Master list of competitions/leagues for reference and mapping.
- Created/Used By (Scripts/Modules): import\_transfermrkt\_normalize.py
- Important Columns/Keys: id (PK), name

#### **marketvalues\_trfrmrkt:**

- Main Use / Data Stored: Historical timeline of player market values (player, club, competition, value, date).
- Created/Used By (Scripts/Modules): import\_transfermrkt\_normalize.py, merge-features.py

- Important Columns/Keys: id (PK), playerid (FK), clubid (FK), competitionid (FK), marketvalue, snapshotdate

#### **player\_injuries\_trfrmrkt:**

- Main Use / Data Stored: Stores injury records for players, including type, date range, and time missed.
- Created/Used By (Scripts/Modules): import\_transfermrkt\_normalize.py, merge-features.py
- Important Columns/Keys: id (PK), playerid (FK), start\_date, end\_date, days\_out, injurydesc

#### **player\_transfer\_history:**

- Main Use / Data Stored: Player transfers: from/to clubs, transfer fee/type, date, and description/history.
- Created/Used By (Scripts/Modules): scrape\_trfr\_record\_new.py
- Important Columns/Keys: id (PK), playerid (FK), fromclub, toclub, fee, date, type, description

#### **player\_mapping:**

- Main Use / Data Stored: Fuzzy and explicit mapping between StatsBomb and Transfermarkt player IDs; confidence score.
- Created/Used By (Scripts/Modules): auto\_player\_mapping.py, map\_players.py
- Important Columns/Keys: statsbombplayerid, trfrmrktplayerid, canonicalname, matchtype, confidence

#### **reddit\_sentiments:**

- Main Use / Data Stored: Reddit posts/comments for each player, analyzed for sentiment/polarity.
- Created/Used By (Scripts/Modules): reddit\_sentiment.py, merge-sentiments.py
- Important Columns/Keys: id (PK), playername, postid, subreddit, createdat, title, selftext, sentiment, polarity

#### **sentiment\_run\_log:**

- Main Use / Data Stored: Logs each sentiment scraping run (batch audits, status, date, platform).
- Created/Used By (Scripts/Modules): social\_sentiment\_multithread\_with\_log.py

- Important Columns/Keys: runid (PK), date, platform, status

#### **player\_features:**

- Main Use / Data Stored: Final engineered feature vectors for players (for ML modeling and dashboard).
- Created/Used By (Scripts/Modules): merge-features.py
- Important Columns/Keys: playerid (PK), latestmarketvalue, marketvaluegrowth, totalinjuries, avgdaysout, lastinjurydate, totaltransfers, totaltransferfees, freetransfers, sentimentmean, sentimentpositiveratio, sentimenttrend, avgcardspermatch, positionsplayed, currentclubid

#### **hyperparameter\_results:**

- Main Use / Data Stored: Stores model tuning (LSTM/XGBoost/ensemble) best parameters, RMSE, validation loss, metadata.
- Created/Used By (Scripts/Modules): hyper\_parameter\_tuning\_app.py, dashboard apps
- Important Columns/Keys: id (PK), modeltype, parameters, rmse, val\_loss, timestamp

## **Appendix C**

### **Notebooks, Code, and Demonstration Links**

- **Source Code (GitHub):**
  - ETL and DB: `import_statsbomb_mysql.py`, `import_statsbomb_lineup_batch.py`, `merge-features.py`, `merge-sentiments.py`
  - Mapping/Utility: `auto_player_mapping.py`, `map_players.py`
  - Sentiment: `social_sentiment_multithread_with_log.py`, `reddit_sentiment.py`
  - Modeling: `ensemble_predictions_app.py`, `encoder_decoder_multi_app.py`, `ensemble_app.py`
- **Notebooks:**
  - EDA & Modeling: `univariate.ipynb`, `multivariate_multi_step.ipynb`
- **Visuals:**
  - All charts and images referenced in this report
- **Live Deployments (on Streamlit cloud):**
  - Hyper Parameter Tuning: <https://hpt-test.streamlit.app/>
  - Forecast Dashboard: <https://predictions-hpt.streamlit.app/>
- **Prediction Dashboard Demo Video:**
  - YouTube link
- **Presentation Slides:**
  - `TransferIQ-Presentation-Himanshu-Saxena-Infosys-Springboard-6.0.pptx`