

MLB Game Simulation Dashboard: A Lambda Architecture Approach

Aadarsha Gopala Reddy

Department of Computer Science & Engineering
Washington University in St. Louis
St. Louis, MO, USA
a.gopalareddy@wustl.edu

Eddy Sul

Department of Computer Science & Engineering
Washington University in St. Louis
St. Louis, MO, USA
eddysul@wustl.edu

Abstract—We present an interactive analytics dashboard for Major League Baseball (MLB) using *lambda architecture*. Our system combines batch processing of historical STATCAST data (2015–2025) with a streaming layer that enables realistic game simulation. The pipeline uses KAFKA for stream processing, SPARK for micro-batch operations, SNOWFLAKE for data warehousing, and AIRFLOW for orchestration. Users can explore historical games, analyze team matchups, and simulate past games pitch-by-pitch with configurable playback speeds.

Index Terms—baseball analytics, real-time streaming, lambda architecture, Apache Kafka, Spark Streaming, data visualization

I. INTRODUCTION

The introduction of Statcast in 2015 transformed statistical analysis across Major League Baseball (MLB). Before 2015, MLB analytics primarily relied on static, historical box-score data. Statcast facilitated a shift from 2D to 3D analytics, capturing granular pitching metrics such as velocity and spin rate, as well as exit velocity and launch angle for every batted ball. This abundance of new data has empowered teams to develop novel strategies for a competitive edge.

Most of this data is accessible via Baseball Savant, which provides advanced analytics and visualizations. Additionally, MLB Gameday allows viewers to watch live game simulations. However, existing tools often segregate these functions: MLB Gameday focuses on real-time streaming without deep historical context, while other platforms lack interactive simulation capabilities. To bridge this gap, we developed an interactive dashboard that integrates game simulation with statistical analysis backed by comprehensive historical data. Specifically, our system augments the simulation with aggregational statistics—such as a hitter’s performance against specific pitch types—and enables users to replay past games.

II. SYSTEM ARCHITECTURE

Our system follows lambda architecture (Fig. 1), combining two processing layers:

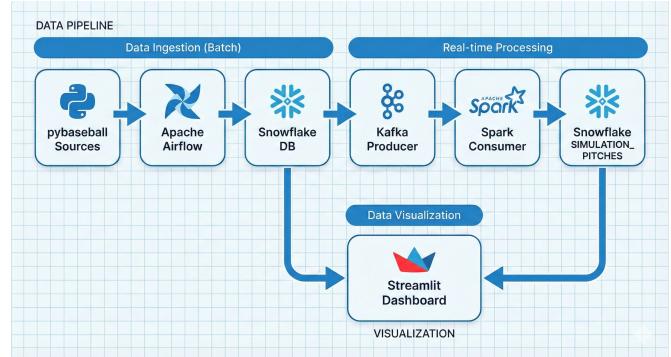


Fig. 1. Lambda architecture overview with batch and streaming layers.

- **Batch Layer:** Daily ingestion of historical data via AIRFLOW DAGs, stored in SNOWFLAKE using a star schema.
- **Streaming Layer:** KAFKA and SPARK Streaming process pitch events for game simulation updates.

III. DATA PIPELINE

A. Ingestion and Transformation

We use `pybaseball`, a Python wrapper for the MLB API, to fetch pitch-by-pitch logs. Each pitch contains over 118 fields; we filter these to retain only metrics relevant to our analysis. Although `pybaseball` does not support real-time streaming, it offers the most comprehensive daily-updated dataset dating back to 2015. We evaluated other APIs for real-time capabilities but found them either lacking in comprehensive metrics, incompatible with our stack, or prohibitively difficult to access.

The processing workflow:

1. Fetch daily data via `pybaseball`
2. Transform: remove deprecated fields, nulls, and out-of-scope metrics
3. Stage as Parquet files in SNOWFLAKE
4. Load into star schema tables
5. Execute data quality checks

B. Database Schema

We employ a *star schema* optimized for OLAP queries:

- **Fact table:** fact_pitches — one row per pitch
- **Dimensions:** player, team, pitch_type, game

C. Orchestration

AIRFLOW DAGs manage the pipeline (Fig. 2–3):



Fig. 2. Daily ingestion DAG.



Fig. 3. Historical backfill DAG (runs every 3 minutes).

IV. STREAMING LAYER

The streaming pipeline consists of three components:

1. **FASTAPI Producer:** Queries SNOWFLAKE and publishes pitch events to KAFKA
2. **KAFKA Broker:** Buffers events (Docker container with Zookeeper)
3. **SPARK Consumer:** Processes micro-batches and writes to Live_Pitches

Decoupling Strategy: Directly consuming the high-throughput stream in STREAMLIT can lead to dropped frames due to the framework's refresh latency. To mitigate this, we decoupled the visualization from the ingestion pipeline. The dashboard independently polls the Live_Pitches table, ensuring complete pitch sequence accuracy, albeit with marginally increased latency.

Trade-off: Data accuracy over minimal latency. Missing pitches would confuse viewers; slight delay is acceptable.

V. DASHBOARD FEATURES

A. Game Explorer

Select historical games and analyze pitcher statistics: velocity per pitch type, spin rate, and strikeout percentage. Includes strike zone visualization and velocity progression tracking.

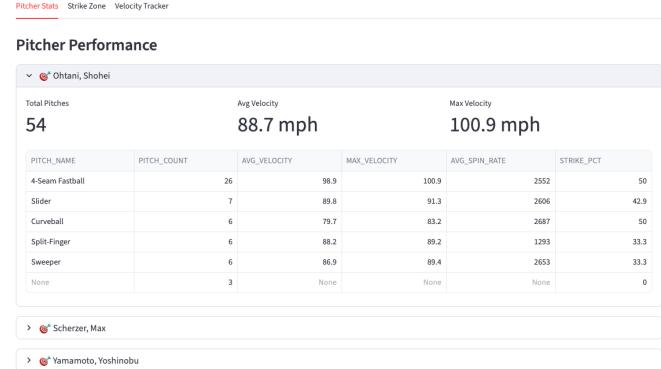


Fig. 4. Pitcher repertoire breakdown and key statistics.

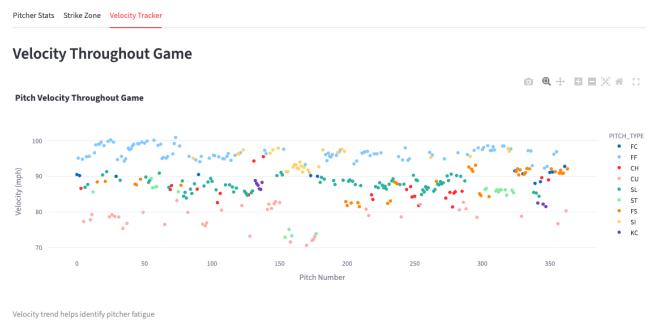


Fig. 5. Pitch velocity progression throughout the game.

B. Team Matchups

Compare head-to-head statistics between teams across a season (Fig. 6).

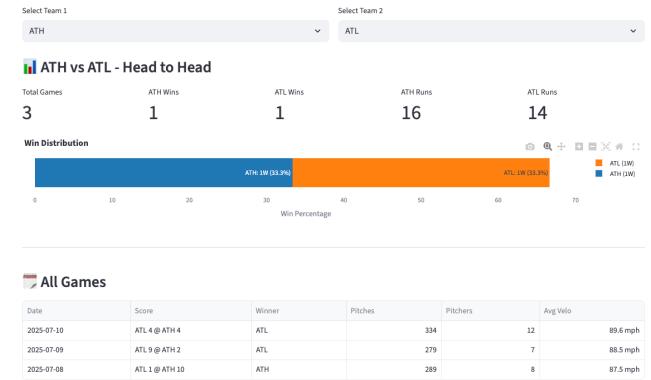


Fig. 6. Team matchup: pitcher vs. batter statistics.

C. Game Simulation

Watch historical games unfold pitch-by-pitch with playback controls (0.5–5 pitches/second). This feature allows users to replay any past game. Notably, while the replay speed is not a 1:1 real-time mirror, it provides a highly customizable visualization tool for tactical analysis.

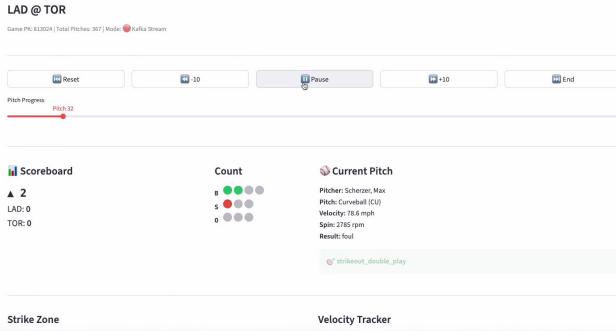


Fig. 7. Game simulation with game state and playback controls.

VI. CHALLENGES AND DESIGN DECISIONS

A. Schema Design

With 110+ fields per pitch, we invested significant effort selecting relevant columns and designing intuitive dimension tables.

B. Views vs. Materialized Tables

We chose SNOWFLAKE views for dashboard aggregations:

TABLE I VIEWS VS. MATERIALIZED TABLES TRADE-OFFS		
Factor	Views	Materialized
Auto-sync with data	✓	Requires refresh
Query performance	Slower	Faster
Staleness risk	None	Possible
Our choice	✓	—

For our workload (frequent writes, user-triggered reads), views provided acceptable latency with guaranteed consistency.

C. Dashboard Latency

STREAMLIT, depending on the system it is running on, imposes ~1 second minimum refresh. Our KAFKA producer can emit faster, but the dashboard cannot display sub-second updates. WebSocket-based alternatives (Dash, React) could improve this.

D. Historical Data Range

We focused on 2020–present for two reasons:

1. Data quality improved post-2019 with upgraded tracking hardware
2. Smaller backfill enabled faster development iteration

The pipeline supports 2015+ data with minimal modification.

E. Infrastructure Issues

We encountered significant environment-specific dependency conflicts when migrating AIRFLOW from a local development environment to the LinuxLab production server.

VII. FUTURE WORK

A. True Real-Time Integration

Current limitation: pybaseball provides data with ~1-day delay.

Required changes for live-season support:

1. WebSocket connection to MLB Gameday API
2. KAFKA producer emitting events as they arrive
3. Reconciliation logic: merge real-time events with next-day STATCAST data

Our lambda architecture already supports this—streaming layer handles live events; batch layer provides quality-assured corrections.

B. Data Corrections

MLB issues scoring changes and data corrections. Potential solutions:

- Change Data Capture (CDC) to detect source modifications
- Targeted updates to affected records
- Handle suspended games spanning multiple dates

C. Performance Optimizations

- **REDIS caching:** Lower latency for live updates
- **FLINK streaming:** Event-by-event processing (vs. SPARK micro-batches)
- **Materialized views:** For high-concurrency scenarios

D. Additional Data Sources

TABLE II
POTENTIAL DATA SOURCES AND APPLICATIONS

Data Type	Source/Library	Application
Weather	OpenWeatherMap API, Visual Crossing	Correlate temperature, humidity, and wind with pitch movement and home run distances
Stadium Info	MLB Park Factors, Chadwick Bureau	Normalize statistics across venues; account for altitude effects (e.g., Coors Field)
Historical Stats	Baseball Reference, Lahman Database	All-time head-to-head matchups; career splits analysis
Player Tracking	MLB Film Room API	Defensive positioning; sprint speed correlations

These integrations would enable advanced analytics such as weather-adjusted ERA, venue-normalized exit velocity, and predictive matchup modeling based on historical tendencies.

E. Planned Features

Batter splits against specific pitch types:

Ohtani vs. Righty Sliders: .333 Average

VIII. DEMONSTRATION

- **Video demo:** <https://www.youtube.com/watch?v=z01fHvzd-8w>
- **Source code:** <https://github.com/agopalareddy/CSE-5114-Project>

IX. CONCLUSION

We presented an MLB game simulation dashboard combining batch and streaming processing via lambda architecture. The system demonstrates effective integration of KAFKA, SPARK, SNOWFLAKE, and AIRFLOW for interactive baseball analytics, enabling users to replay historical games and explore detailed statistics. Future work includes true real-time integration during live seasons.

REFERENCES

- [1] MLB Advanced Media, “Statcast,” <https://www.mlb.com/glossary/statcast>.
- [2] Baseball Savant, “Statcast Search,” https://baseballsavant.mlb.com/statcast_search.
- [3] pybaseball, “Python package for baseball data,” <https://github.com/jldbc/pybaseball>.
- [4] Apache Kafka, <https://kafka.apache.org/>.
- [5] Apache Spark, <https://spark.apache.org/>.
- [6] Snowflake Inc., <https://www.snowflake.com/>.