

TrackMan API to PostgreSQL ETL Pipeline for Pomona-Pitzer Baseball

Capstone in Data Science

Z Skigen

November 21, 2025

1 Introduction

1.1 Executive Summary

This project creates an automated Extract–Transform–Load (ETL) pipeline that connects the Pomona-Pitzer Baseball TrackMan system directly to a PostgreSQL database. The goal is to move away from the current workflow, manual CSV downloads obtained through FileZilla and stored in Google Drive, toward a reproducible, cloud-based database that can be accessed and queried by multiple analysts simultaneously.

The pipeline authenticates with the TrackMan API, extracts raw JSON data, parses it into structured tables with Python and pandas, and loads it into PostgreSQL using SQLAlchemy. The centralized database is hosted on Railway, which provides a managed, persistent PostgreSQL instance accessible from anywhere. Analysts can run queries either locally (e.g., through VS Code) or through a Streamlit-based SQL playground that will be deployed on Streamlit Cloud.

To ensure the database remains continuously up to date, automated ingestion will be scheduled using PythonAnywhere or GitHub Actions, allowing the ETL script to run nightly without manual intervention.

Long-term maintainability is supported through a dedicated analyst hub at sagehenanalytics.sites.pomona.edu, which will provide:

- onboarding materials for new analysts,
- a SQL primer tailored to baseball data,

- a data dictionary describing every field and table,
- example queries and reproducible workflows, and
- documentation for sustaining and extending the ETL pipeline.

Together, the cloud infrastructure, centralized storage, and documentation ecosystem create a maintainable, MLB level infrastructure that supports automated reporting, advanced modeling, and long-term continuity for the program.

1.2 Motivation: Current Workflow and Limitations

TrackMan is an optical radar- and camera-based tracking system that measures pitch trajectories, velocities, spin rates, and batted-ball outcomes.

Pomona-Pitzer Baseball currently relies on manual processes to create TrackMan-based scouting reports, visualizations of baseball statistics. After each game, analysts download separate CSV files containing pitch-level data, many of which contain thousands of rows. These files must be cleaned, merged, and analyzed individually, with every analyst maintaining their own local copy of the data. This approach is slow, error-prone, and difficult to reproduce.

TrackMan's proprietary software provides visual summaries and leaderboards, but it does not grant analysts direct, flexible access to underlying pitch-level data. As a result, analysts rely on manual R/Python scripts to generate metrics like pitch usage, release velocity trends, or command consistency. These analyses are valuable but isolated; different students often compute similar metrics differently, leading to inconsistent results across reports.

The ETL pipeline addresses these issues by:

1. Centralizing all TrackMan data into a relational PostgreSQL database hosted on Railway.
2. Automating data extraction and cleaning directly from the TrackMan API.
3. Standardizing field names, data types, and relationships between games, players, and pitches.
4. Providing shared access through SQL, Python, and a Streamlit-based web interface.
5. Supporting long-term onboarding and sustainability through documentation at sagehenanalytics.sites.pomona.edu, including SQL examples and database-maintenance guides.

Once the pipeline is deployed, every analyst will query the same database. This eliminates redundant data wrangling and ensures reproducible, transparent analytics across seasons and analysts.

2 Data: Structure, Scope, and Organization

2.1 API Acquisition and Data Cleaning

The TrackMan Range API delivers data as deeply nested JSON rather than structured tables. To make these data usable for analysis, the ETL pipeline implements a complete acquisition, cleaning, and loading workflow:

1. **Authenticate with OAuth2** using a client-credentials flow.
2. **Fetch available sessions** in 30-day windows to remain within API rate limits (TrackMan 2025).
3. **Request /plays and /balls JSON objects** for each session.
4. **Flatten all nested JSON** using the `requests` library and `pandas.json_normalize()` (Goel 2020).
5. **Standardize fields**, including:
 - converting keys to `snake_case`,
 - parsing timestamps into UTC,
 - enforcing numeric types, and
 - removing duplicate or null rows.
6. **Load cleaned DataFrames into PostgreSQL** using SQLAlchemy (GeeksforGeeks 2021).

This transforms TrackMan's raw JSON into clean, queryable relational tables suitable for cloud storage, SQL analysis, and downstream modeling.

A temporary note: only a subset of 2024 sessions currently return complete `plays` and `balls` data. This appears to be a permissions or session-type filtering issue rather than a problem with the ETL code. I am actively working with TrackMan support to confirm the correct access scope.

2.2 PostgreSQL Database Creation

Using the official PostgreSQL documentation (PostgreSQL Global Development Group 2025), I created a dedicated database for the TrackMan pipeline. After cleaning, the ETL writes data into three core tables: `sessions`, `plays`, and `balls`. These tables reflect the natural hierarchy of baseball events:
`sessions` → plate appearances → pitches.

2.2.1 sessions

One row per TrackMan event (game, scrimmage, bullpen).

Includes:

- session ID
- date and time
- facility and field information
- team identifiers
- metadata anchoring all downstream tables

2.2.2 plays

One row per plate appearance from the `/plays` endpoint.

Includes:

- `playID` (TrackMan's identifier)
- batter and pitcher identifiers
- inning, outs, and count progression
- pitch sequence order
- tagger-behavior and contextual metadata

This table situates each pitch within its in-game context.

2.2.3 balls

One row per pitch from the `/balls` endpoint.

Includes detailed pitch-tracking metrics:

- release speed, height, side, and angles
- spin rate and spin axis
- induced vertical break, horizontal break, and movement
- pitch location and trajectory
- both `playId` and `sessionId` keys

This table captures all pitch-level physics and tracking information.

2.3 How the Tables Connect (Joins)

The relational schema is designed so analysts can reconstruct full pitch sequences and game contexts through standard SQL joins:

- `plays.sessionId = sessions.sessionId`
links plate appearances to their parent session.
- `plays.playID = balls.playId`
links each plate appearance to its individual pitches.
- `balls.sessionId` can also be used to join directly back to sessions if needed.

A typical combined query looks like:

```
SELECT
    p."taggerBehavior_pitchNo" AS pitch_no,
    to_timestamp(p."localDateTime", 'MM/DD/YYYY HH24:MI:SS') AS date,
    p."pitcher_name",
    b."pitch_release_relSpeed",
    b."pitch_movement_inducedVertBreak"
FROM plays p
JOIN balls b
    ON p."playID" = b."playId"
WHERE to_timestamp(p."localDateTime", 'MM/DD/YYYY HH24:MI:SS') >= '2024-01-01'
ORDER BY date, pitch_no;
```

This join structure forms the backbone of all downstream analyses. Additional tables—such as `players`, normalized `pitch_types`, or derived scouting-summary tables—can be added later without altering these core relationships.

2.4 Note on Limited 2024 Data

Only a small subset of 2024 sessions currently return complete pitch-level (`balls`) and plate-appearance (`plays`) data from the TrackMan API. This appears to be related to account permissions or session-type filtering, rather than any issue with the ETL pipeline itself. I am actively working with TrackMan support to verify the correct access scope.

Once the access configuration is resolved, the existing ETL infrastructure will ingest the full historical dataset without requiring architectural changes.

2.5 Cloud Architecture and Deployment Workflow

The pipeline is designed for lightweight, modular cloud deployment so that analysts can access the system without installing local databases or running ingestion scripts. Three platforms serve distinct roles in the architecture:

2.5.1 Railway (Managed PostgreSQL Hosting)

Railway hosts the centralized PostgreSQL database and provides:

- a persistent, fully managed PostgreSQL instance,
- built-in environment-variable management,
- automatic backups, and
- secure connections from VS Code, Python scripts, or web applications.

This ensures that all analysts query the same live database regardless of device or location.

2.5.2 Streamlit Cloud (Interactive SQL Playground)

Streamlit Cloud will host the team's web-based SQL playground. The interface will:

- allow analysts to run parameterized SQL queries,

- visualize pitch-level or player-level metrics,
- export results for scouting reports, and
- provide a low-barrier, MLB-style interface for interacting with TrackMan data.

This is especially useful for onboarding new analysts who may not yet be comfortable writing SQL directly.

2.5.3 PythonAnywhere (Optional Scheduled ETL Worker)

To keep the database continuously up to date, PythonAnywhere can execute the ETL script on a nightly schedule. The automated workflow will:

- authenticate with the TrackMan API,
- fetch new sessions, plays, and pitches,
- clean and validate the data,
- append new rows to the Railway PostgreSQL instance, and
- make updated data immediately available in Streamlit dashboards.

This removes the need for manual updates and supports a fully automated ingestion pipeline.

3 Analysis

The goal of this project is not to build statistical models or conduct performance analytics, but to create the infrastructure that enables those tasks. To demonstrate the practical value of the centralized PostgreSQL database, this section showcases examples of queries and workflows that were previously time-consuming or impossible under the manual CSV-based system.

Rather than presenting full analytic results, these examples illustrate how the new schema supports fast, reproducible, and consistent access to pitch-level data.

3.1 Example 1: Accessing All Pitches for a Game

A query that once required manually merging multiple CSVs can now be executed in a single step:

```
SELECT *
FROM balls
WHERE session_id = 'some_game_id';
```

3.2 Example 2: Linking Pitches to Plate Appearances

The relational structure allows analysts to connect events without manual joins:

```
SELECT p.pitch_number, p.pitch_type, pa.batter_id, pa.inning
FROM balls p
JOIN plays pa USING (session_id, pitch_number)
LIMIT 50;
```

3.3 Example 3: Checking Data Completeness

Because ingestion is automated, analysts can quickly verify which sessions have full play and pitch data:

```
SELECT session_id,
       COUNT(*) AS n_pitches
FROM balls
GROUP BY session_id;
```

3.4 Example 4: Preparing Data for Scouting Reports

Analysts can build consistent scouting templates directly from the database:

```
SELECT pitcher_id,
       AVG(release_speed) AS avg_velo,
       STDDEV(release_speed) AS velo_sd
FROM balls
GROUP BY pitcher_id;
```

These examples are not intended as full statistical analyses. Instead, they demonstrate that the new database functions as a reliable foundation for future analytics work, reducing time spent on data wrangling and ensuring that all analysts work from the same, consistent data.

4 Ethical Considerations

Building a live data pipeline from the TrackMan API introduces important ethical considerations related to data privacy, ownership, contractual compliance, access control, and responsible use. Because the pipeline automates acquisition and centralizes all pitch-level data, it is essential to ensure that athlete information is handled securely and used only for appropriate internal purposes.

4.1 Data Ownership and Privacy

TrackMan data belong to Pomona-Pitzer Baseball and its athletes. Although the ETL system streamlines ingestion, it does not change who controls the data or how the data may be shared. The database is designed only for internal use within the Pomona-Pitzer Baseball program.

To maintain privacy and comply with institutional expectations:

- API credentials are stored securely through environment variables, never in source code.
- Database access is restricted to verified analysts through password-protected accounts.
- Personally identifiable information (PII), such as student names or ID numbers, is never included in public exports or external reports; all analysis is conducted using anonymized player identifiers.
- No raw or processed TrackMan data are publicly distributed or indexed on the open internet.

These safeguards ensure that the pipeline enhances accessibility without compromising privacy or exposing information.

4.2 Contractual and Compliance Considerations

TrackMan requires users of its Data API to agree to a formal Terms of Service (ToS). As of now, I am awaiting a copy of the ToS from our TrackMan representative. Once received, I will review the document to determine:

- whether data may be shared internally across analysts,
- what forms of storage or distribution are prohibited,
- whether the database must remain private,
- any restrictions on derived metrics or visualization,
- any expectations regarding secure data handling or retention.

Until the ToS is reviewed, the system is intentionally designed to restrict access to authorized internal users only. This conservative default minimizes the risk of accidental violation of TrackMan's contractual requirements.

4.3 Fairness and Responsible Use

Performance data can influence evaluation and playing-time decisions, which makes responsible interpretation essential. Raw pitch metrics do not capture context such as injury history, mechanical adjustments, or coaching strategy, and could be misused if treated as objective rankings.

To promote fair use:

- The database supports development-focused analysis, not public comparison, ranking, or punitive evaluation.
- Summaries and dashboards are contextualized with coaching input.
- Metrics are interpreted as tools for improvement rather than indicators of athlete value.
- Analysts are encouraged to consider uncertainty, variability, and sampling issues in their conclusions.

This approach reduces the risk of overinterpreting noisy performance data or unintentionally creating inequitable evaluation environments.

4.4 Security and Access Control

Centralizing data introduces technical risks that must be mitigated. To ensure secure operation:

- Access follows least-privilege principles: each analyst receives only the permissions necessary for their role.
- Backups are encrypted, and the database is never left open to public traffic.
- API usage logs and ingestion events are monitored for anomalies.
- Credentials are rotated when personnel change.

These safeguards protect both athletes and the institution from data breaches, unauthorized access, or unintentional exposure.

4.5 Cloud Platform Considerations

The infrastructure supporting this project relies on modern cloud services, including Railway for database hosting, Streamlit Cloud for interactive SQL exploration, and potentially PythonAnywhere for scheduled ingestion or lightweight API tasks. Each of these services offers different strengths, Railway provides a simple and secure managed PostgreSQL instance, Streamlit Cloud enables accessible web-based dashboards for analysts, and PythonAnywhere supports automated Python execution without requiring local resources.

However, using third-party cloud platforms introduces additional ethical responsibilities. To prevent accidental public exposure of athlete data:

- all Railway database instances are configured to require password-authenticated connections,
- the Streamlit Cloud application will authenticate users before granting access to any query functionality, and
- no TrackMan-derived data will be written to publicly accessible URLs, buckets, or web endpoints.

If PythonAnywhere or similar services are used for scheduled ingestion, credentials will remain stored in environment variables, not in source code, and logs will not contain sensitive player information. These platform decisions are made for operational simplicity, but each is configured to preserve the privacy and confidentiality of athlete performance data and to prevent unauthorized external access.

4.6 Potential Risks With Automation

The pipeline is designed for automated updates, which introduces new ethical and operational risks:

- If TrackMan mislabels sessions (e.g., bullpens vs. games), automated ingestion may import incorrect or unintended data.
- Errors in the API or incomplete sessions could propagate into the database without being manually reviewed.
- Analysts could rely on real-time metrics without understanding their limitations or uncertainty.

To address these risks:

- Automated updates will include logging, flags for new sessions, and procedures for roll-back.
- Analysts will receive documentation explaining how ingestion works, what assumptions are made, and how to detect anomalies.
- Any automated reporting tools will emphasize transparency and uncertainty to avoid overconfidence in raw metrics.

5 Next Steps and Automation

The core ETL pipeline has been implemented and tested locally, but several major components remain before the system can operate as a complete, production-ready analytics platform. The following next steps focus on resolving current data-access limitations, building analyst-facing tools, and establishing long-term maintainability for the program.

5.1 1. Resolve the Limited 2024 Data Returned by the TrackMan API

The most immediate priority is diagnosing why the TrackMan API is returning only a small subset of 2024 sessions with complete pitch-level data. I am currently:

- reviewing the TrackMan Data API Terms of Service once it is provided,
- confirming the access scope associated with the Pomona-Pitzer TrackMan account,
- contacting our TrackMan representative to clarify available endpoints and permissions, and
- testing alternative approaches to querying sessions or session types.

Once access details are confirmed and expanded, the existing ETL code will be able to ingest the full historical dataset without requiring structural changes.

5.2 2. Build a Front-End Interface for Analysts

A lightweight Streamlit Cloud application will serve as the primary interface for team analysts who may not want to write SQL or Python. Planned features include:

- an authenticated login page,
- a SQL query sandbox connected to the central Railway database,
- built-in visualizations of pitch-level metrics,
- the ability to export data for scouting reports, and
- tools for monitoring data completeness and ingestion logs.

This web app will mirror internal MLB databases and significantly lower the barrier to using TrackMan data effectively.

5.3 3. Create a Persistent Analyst Hub (sagehenanalytics.sites.pomona.edu)

To ensure long-term sustainability and institutional memory, I will build a documentation hub hosted at:

sagehenanalytics.sites.pomona.edu

This site will contain:

- detailed documentation on how to access, query, and maintain the PostgreSQL database,
- a simple introduction to SQL for baseball analytics,
- a data dictionary for all tables produced by the ETL pipeline,

- setup instructions for new analysts,
- version-controlled documentation for future modifications, and
- guidance on data privacy, permissions, and operational best practices.

This hub will serve as the permanent home for Sagehen Baseball's analytics workflows.

5.4 4. Implement Automated Data Refreshes

A major goal of the system is to remove the need for manual data downloads entirely. Scheduled ingestion will be implemented using GitHub Actions or PythonAnywhere to automatically run the ETL script on a nightly basis (Sahu 2024). This automated process will:

- authenticate with the TrackMan API,
- detect and ingest new sessions,
- append updated play and pitch data,
- log ingestion events for debugging, and
- notify analysts if anomalies arise.

Automation is a crucial step: once implemented, the database will remain up to date without requiring analyst intervention.

5.5 5. Finalize Cloud Deployment Workflow

The last stage of development involves fully transitioning from local testing to secure cloud deployment:

- deploying the ingestion workflow to GitHub Actions or PythonAnywhere for scheduled execution (Joe Nelson, Steve Chavez 2025),
- hosting the SQL playground and visualization tools on Streamlit Cloud (Streamlit, Inc. 2025),
- ensuring secure connectivity between Streamlit, PythonAnywhere, and Railway,
- rotating credentials and storing them via environment variables, and

- performing multi-user testing to confirm reliability, latency, and usability.

Together, these steps will produce a robust, maintainable, and accessible cloud-based analytics ecosystem for the Pomona-Pitzer Baseball program (Plotly Technologies Inc. 2025).

6 Limitations

1. TrackMan API documentation changes periodically, requiring manual field validation.
2. Large JSON files can exceed local memory, necessitating batch ingestion.
3. Historical data from earlier seasons may use inconsistent field names or units.
4. Operational data (e.g., bullpen sessions, scrimmages) may need manual tagging to distinguish from official games.

Despite these limitations, the project's modular structure will allow straightforward debugging and schema expansion as the database scales.

6.1 Incomplete 2024 TrackMan Data (Under Investigation)

During development, I discovered that the TrackMan API was only returning a small subset of 2024 sessions with complete play- and ball-level data. Most of the pitch-level data currently available through the `/plays` and `/balls` endpoints falls between January 25–28, 2024, even though the API reports thousands of session identifiers for the year.

This does not appear to be an error in the ETL pipeline itself. Instead, it is likely related to the scope of the Pomona–Pitzer TrackMan account, the API's session-type filtering behavior, or access permissions defined in the TrackMan Data API Terms of Service. I am actively working to diagnose the cause by reviewing the ToS, communicating with TrackMan support, and testing alternative session-querying strategies.

As this investigation continues, the current database reflects the subset of 2024 data that the API makes available with the present access configuration. Once access details are confirmed and expanded, the pipeline will be able to ingest the full historical dataset without requiring structural changes.

7 Conclusion

The TrackMan ETL pipeline will modernize Pomona-Pitzer Baseball's data operations. Instead of dozens of fragmented CSVs, analysts will have a unified database accessible through SQL, Python, or a web app.

Once fully implemented, the system will:

1. Eliminate redundant manual cleaning.
2. Provide live, reproducible access to pitch-level data.
3. Serve as a foundation for advanced models and automated scouting reports.

This infrastructure will not only support the current analytics team but also create a sustainable framework for future seasons and future analysts.

References

- GeeksforGeeks. 2021. “Python SQLAlchemy Introduction.” <https://www.geeksforgeeks.org/python/sqlalchemy-introduction/>.
- Goel, Ankit. 2020. “How to Parse JSON Data with Python Pandas?” *Medium: TDS Archive*. <https://medium.com/data-science/how-to-parse-json-data-with-python-pandas-f84fb0b1025>.
- Joe Nelson, Steve Chavez. 2025. *PostgREST Documentation*. <https://postgrest.org/>.
- Plotly Technologies Inc. 2025. *Dash User Guide and Documentation*. <https://dash.plotly.com/tutorial>.
- PostgreSQL Global Development Group. 2025. *PostgreSQL: CREATE DATABASE — SQL Command*. <https://www.postgresql.org/docs/current/sql-createdatabase.html>.
- Sahu, Satyam. 2024. “Automate Data Ingestion: Connecting REST APIs to Your Data Warehouse with Python.” Medium. <https://medium.com/towards-data-engineering/automate-data-ingestion-connecting-rest-apis-to-your-data-warehouse-with-python-eb889fb668f7>.
- Streamlit, Inc. 2025. *Streamlit Documentation*. <https://docs.streamlit.io/>.
- TrackMan. 2025. *TrackMan Range API Documentation*. <https://docs.trackmanrange.com/>.