# A Strategic Blueprint for the Living Data Atlas

## Executive Summary: A Strategic Blueprint for Nigeria's Data Future

The "Living Data Atlas" project represents a highly promising and strategically aligned initiative aimed at building a centralized, open-source data repository for Nigeria. The project's mission is to address critical gaps in data transparency and accessibility by creating a "single source of truth" for key datasets spanning the economy, environment, health, and governance. This report provides an expert-level analysis of the project's phased roadmap, validating its technical architecture and identifying key strategic considerations and actionable recommendations.

The analysis indicates that the proposed project plan is both technically sound and directly responsive to Nigeria's documented need for a robust data infrastructure. The selection of a technology stack—comprising PostgreSQL, Python, SQLAlchemy, and Docker—is an excellent choice that judiciously balances a low barrier to entry for contributors with the enterprise-grade robustness required for long-term scalability. A critical challenge identified is not technical, but rather related to data sourcing, given the documented fragmentation and scarcity of official data, particularly in the governance and social sectors. Furthermore, the report underscores the importance of the project's open-source elements, noting that the choice of a visualization tool should prioritize the end-user's experience and the goal of self-service analytics over technical complexity. Finally, the report concludes that open-source documentation and community-building efforts are not a final phase but a foundational, ongoing effort essential for the project's long-term sustainability and success.

## The Strategic Imperative: Contextualizing the Need for Open Data in Nigeria

The development of the Living Data Atlas is not merely a technical undertaking; it is a direct response to a strategic imperative. The report highlights the foundational role of accessible data in fostering good governance, combating corruption, and enabling evidence-based policymaking. In Nigeria, weak data collection and poor management practices have historically hindered access to official data, which has eroded the nation's ability to enforce accountability and ensure good governance. This lack of reliable, timely data directly undermines transparency, particularly in public service delivery, where it can conceal the misuse of public funds and fuel inefficiencies. Without accessible data on government expenditures, performance metrics, and service outcomes, it becomes nearly impossible to monitor progress or hold public officials accountable.

A review of Nigeria's data openness across key sectors exposes significant disparities in transparency. While Nigeria ranks second in Africa for data openness with a score of 42%, it still lags significantly behind South Africa's 59%. Across sectors, data availability is fragmented: economic data is the most accessible at an average of 58.57%, while environmental and social data availability are alarmingly low, at 24% and 42%

respectively. The project's objective of creating a unified data repository is a direct solution to this national problem of data fragmentation. By standardizing and unifying disparate data points, the Living Data Atlas elevates itself from a simple technical task to a national strategic initiative. The economic implications are also profound. For example, traffic congestion in Lagos alone is estimated to cost the economy $1 billion annually in lost productivity, a tangible issue that can be addressed through bold, data-driven urban development strategies informed by robust data collection systems and analysis. The creation of a "single source of truth" makes the data required for accountability and development not only available, but also easily accessible and verifiable, directly attacking the issues of resource mismanagement and a lack of public trust.

# Architectural Blueprint: A Phased Breakdown of the MVP

## Phase 1: Foundations for a Robust Data System

The foundational phase of the project is designed to establish a robust and scalable architecture. This begins with the Extract, Transform, Load (ETL) process, which is the core of any data pipeline. ETL is the process of combining raw data from multiple sources into a central repository, or data warehouse, by first extracting it, then transforming it to a standardized format, and finally loading it into the target database. The project's choice of Python for this purpose is highly suitable, given its extensive ecosystem of libraries for data manipulation, API interaction, and database connectivity.

The selection of PostgreSQL as the primary data store is a prudent architectural decision. As an open-source relational database, it is renowned for its robustness and proven reliability under heavy workloads. PostgreSQL offers advanced features such as support for JSONB, full-text search, and partitioning, which make it well-suited for both transactional and analytical queries. The choice of SQLAlchemy for bridging Python applications and the database is also sound. SQLAlchemy functions as a SQL toolkit and an Object-Relational Mapping (ORM) system, allowing developers to interact with the database using familiar Python objects and methods instead of writing raw SQL queries. This abstraction provides database independence, meaning the codebase can be easily adapted to other SQL databases with minimal changes, and it significantly improves code readability and maintainability.

The project's inclusion of Docker for setting up the PostgreSQL environment is a highly strategic move that directly supports the "open-source friendly" goal. Docker is a containerization platform that packages an application and all its dependencies into a single, portable, and isolated unit. This approach effectively eliminates the common "it works on my machine" problem by ensuring a consistent environment across different developer machines and deployment platforms. By providing a Docker setup, the project proactively lowers the barrier to entry for potential contributors, who can get the database and its dependencies running with a single command, freeing them to focus on writing code rather than on complex environment configuration. This is a powerful, non-obvious step that fosters community growth, a key pillar for the long-term success of any open-source initiative.

## Phase 2: Ingestion Engine - Building the First Loader

This phase focuses on building the core ingestion pipeline, starting with a single data source to validate the architecture. REST APIs are the standard for web service communication, using HTTP requests (such as GET for reading data and POST for creating it) and typically returning data in a human- and machine-readable

JSON format. The project's use of the World Bank API for economic indicators serves as an excellent initial use case.

A key challenge with API data is that the responses often contain nested JSON objects which are not directly compatible with a flat relational database table. The pandas.json_normalize() function is a powerful tool for this purpose, designed to flatten these hierarchical JSON structures into a clean, tabular format that can be easily loaded into a database. The process for the World Bank loader would involve using a library like Python's

requests to fetch the data , followed by using

json_normalize to prepare the data for loading into Postgres. For instance, a simple API request to get GDP data returns a nested JSON object with indicator details and metadata, which must be normalized before storage.

| Technology | Rationale |
| --- | --- |
| **PostgreSQL** | Robust, open-source relational database with advanced features like JSONB support, well-suited for both transactional and analytical queries. |
| **Python** | Chosen for its ease of use and extensive ecosystem of libraries for data retrieval (requests), manipulation (Pandas), and database interaction (SQLAlchemy). |
| **SQLAlchemy** | Serves as a SQL toolkit and ORM, abstracting database-specific dialects and enhancing code readability and maintainability through a Pythonic domain language. |
| **Docker** | Simplifies environment creation and management by packaging applications and dependencies into isolated, portable containers, eliminating environment inconsistencies and lowering the barrier for new contributors. |

## Phase 3: Standardizing and Scaling Data Integration

The next major milestone is to scale the ingestion process by adding diverse data sources, specifically for weather and air quality. This phase presents the challenge of data heterogeneity, as different APIs, such as Open-Meteo and OpenAQ, have unique data structures and formats. OpenAQ, for example, provides data for specific pollutants and includes rich metadata about instruments, providers, and locations. The Open-Meteo API, by contrast, focuses on meteorological variables such as latitude, longitude, and elevation.

The core objective of this phase is not just to build these new loaders but to harmonize the data from these disparate sources. The proposed standardized schema (date, indicator, region, value, source, meta) is the

architectural linchpin that makes this possible. It ensures that data from all sources, regardless of its original format, can be consistently queried and compared within the single repository. The standardization process is a complex transformation step that must be carefully designed to avoid data loss and ensure data points are comparable. For example, while OpenAQ data is provided in physical units rather than air quality indices, a proper transformation is required to fit this into the unified schema. The long-term value of the Living Data Atlas lies in its ability to make collected data interoperable. This phase is where the project truly lives up to its name by creating a consistent framework for all future data sources, transforming disparate data points into a cohesive, queryable atlas.

## Phase 4: Expanding Horizons: Health and Governance

Building on the foundation laid in previous phases, the project aims to expand its scope to include health and governance data. This requires strategic sourcing from organizations such as the World Health Organization (WHO) and the United Nations (UN). The WHO's OData API provides a simple query interface for data and statistics on various indicators, while the UN Population Data Portal offers a RESTful API for retrieving specific indicators for different geographical areas and time periods.

While these global APIs provide a solid starting point, the research suggests a critical challenge in this phase: data scarcity and quality issues within Nigeria's specific context. The initial analysis reveals that governance and social data in Nigeria are often fragmented and difficult to access, with significant gaps in transparency. While the WHO and UN APIs offer broad, global datasets, the project's goal of creating a "single source of truth for Nigeria" may require supplementing this with local, potentially less structured, data. This phase introduces a notable risk of a fragmented or incomplete dataset for the target region. The project must be prepared to handle these data gaps and explicitly document sources, limitations, and collection methods in the

meta field of its unified schema. The project's documentation therefore becomes a critical feature for transparency and a key enabler of future data contributions and trust.

## Phase 5: The Exploration and Discovery Layer

The final component of the MVP is the exploration layer, which will make the collected data accessible and understandable to end-users. The project's success is dependent on its ability to make data "simple enough to digest," and the choice of a business intelligence (BI) tool is a strategic decision. The user's proposed options—Metabase, Superset, and Streamlit—each offer distinct trade-offs for this purpose.

Metabase is best suited for a non-technical audience, offering intuitive tools like an interactive query builder that allows users to create reports without writing any SQL. Its ease of setup and use makes it ideal for a startup or solo-developer environment, enabling team members to get answers from data on their own. However, it has limitations in advanced customization and lacks native Git version control, which some developers may find restrictive.

Superset, in contrast, is designed for highly technical teams with extensive data analysis needs. It offers a broad range of complex visualization options but requires significant technical expertise for deployment and management. Streamlit is positioned as a developer-centric tool for rapidly building and sharing data applications using Python. While it is excellent for prototyping and quick data apps, it may be less efficient for building large, complex dashboards compared to a dedicated BI tool.

The most valuable path forward is to select the tool that aligns most closely with the project's core mission of data democratization. The analysis indicates that Metabase is the most logical choice for the MVP, as it prioritizes self-service analytics and a low barrier to entry for non-technical users, thereby maximizing the project's potential impact on the widest possible audience. The following table provides a comparative breakdown of these tools.

| BI Tool | Ease of Use | Technical Expertise Required | Customization | Open-Source Readiness | Target Audience |
|---|---|---|---|---|---|
| **Metabase** | High | Low | Limited | Fully Open Source | Non-technical users, solo developers, startups seeking self-service analytics. |
| **Superset** | Low | High | Extensive | Fully Open Source | Highly technical teams with complex data analysis needs and resources for maintenance. |
| **Streamlit** | Medium | Medium | Medium | Fully Open Source | Developers seeking to rapidly prototype data apps using Python code. |

## Phase 6: Open-Source Readiness & Community Building

The final phase of the MVP roadmap centers on open-source readiness, a critical aspect that extends far beyond a simple checklist item. The project's use of GitHub is a foundational element of this strategy. GitHub is not merely a code repository; it is a collaborative platform that provides tools for issue tracking, pull requests, and code review, which streamline collaboration among a distributed community of developers. The pull request system, for example, allows maintainers to review and comment on proposed changes before they are merged, ensuring code quality.

Crucially, documentation is not an afterthought but a foundational product in its own right. A project's README.md and contribution guide are the first things a potential contributor sees, and they must be clear, concise, and easy to follow. Best practices for documentation include providing clear, working code examples, using screenshots to demonstrate features, and defining a common terminology to ensure the community speaks a common language. The research highlights that poor documentation is a primary barrier to entry for newcomers. The long-term viability of the Living Data Atlas is directly tied to its ability to attract and retain contributors. By treating documentation with the same rigor as the codebase itself, the project builds a platform for collaboration, not just a static repository. This re-frames the final phase from a task-oriented checklist to a strategic priority that will determine the project's ultimate success or failure.

# Conclusion and Next Steps: From Blueprint to Reality

The project plan for the Living Data Atlas is a robust, well-conceived roadmap that is strategically positioned to address a critical national need in Nigeria. The proposed technical architecture is sound, and the phased approach provides a clear path from a foundational MVP to a scalable and sustainable platform.

Based on this analysis, the following strategic recommendations are provided to ensure the project's long-term success:

- **Prioritize the Unified Schema:** The project's success hinges on its ability to create a truly "single source of truth." As new data sources are integrated, particularly those for governance and health, the unified schema must be flexible enough to handle data quality issues, scarcity, and disparate metadata from local sources. Investing time and resources in a robust metadata structure is crucial.
- **Embrace Documentation as a Product:** The open-source readiness phase should not be considered a final step. Documentation should be treated as an ongoing, first-class feature from day one. Implementing best practices, such as documenting new features as they are added and conducting regular reviews, will be essential for fostering a vibrant and engaged contributor community.
- **Select the Optimal BI Tool:** The choice of the exploration layer tool should be guided by the ultimate goal of democratizing data access. The analysis recommends Metabase for the MVP due to its low barrier to entry and focus on self-service analytics for non-technical users, which aligns perfectly with the project's mission of making data "simple enough to digest."

The Living Data Atlas holds immense potential to drive meaningful change by fostering transparency, accountability, and evidence-based decision-making. By meticulously executing this blueprint and treating its non-technical aspects with the same rigor as its technical ones, the team can lay the groundwork for a platform that will be invaluable for Nigeria's future.