# Designing a Robust Forecasting System for Wikipedia Web Traffic: A Systems Engineering Approach

Johan Sebastián Beltrán Merchán
Code: 20222020019
Dept. of Systems Engineering
Universidad Distrital Francisco José de Caldas
Email: jsbeltranm@udistrital.edu.co

Edison David Álvarez Varela
Code: 20222020043
Dept. of Systems Engineering
Universidad Distrital Francisco José de Caldas
Email: edalvarezv@udistrital.edu.co

Yader Ibraldo Quiroga Torres
Code: 20222020034
Dept. of Systems Engineering
Universidad Distrital Francisco José de Caldas
Email: yiquirogat@udistrital.edu.co

Julián David Celis Giraldo
Code: 20222020041
Dept. of Systems Engineering
Universidad Distrital Francisco José de Caldas
Email: jdcelisg@udistrital.edu.co

*Abstract*—**Wikipedia receives daily visits from around the world from several device types, from humans or agents, with people from different cultures and languages. The proposed architecture follows systems engineering principles to address the challenge of predicting daily pageviews for Wikipedia articles across multiple languages, access types, and temporal dimensions. The design emphasizes modularity, scalability, and maintainability while handling large-scale temporal data (2.5GB+) and complex page identifier structures. The system transforms raw historical pageview data into competition-ready predictions through a carefully orchestrated pipeline of data processing, feature engineering, machine learning modeling, and validation components. The implemented prototype demonstrates the system's capability by successfully processing 145,063 unique Wikipedia pages and generating predictions for 8,703,780 time-series entries using Linear Regression as a baseline model.**

*Index Terms*—**Systems Engineering, Time Series Forecasting, Linear Regression, Data Pipeline, SMAPE, Scalability.**

## I. INTRODUCTION

The challenge of predicting Wikipedia web traffic views transcends a simple algorithmic exercise; it is a complex Systems Engineering problem. The competition involves forecasting daily views for approximately 145,000 articles, generating a dataset characterized by immense scale and inherent non-linearity. This environment is highly susceptible to **chaos factors**, such as unpredictable viral events, sudden shifts in search engine algorithms, and complex, multi-scale seasonality [1]. Addressing this requires a robust architecture capable of processing large volumes of data while dynamically adapting to volatile patterns.

Previous research in time series forecasting has demonstrated the utility of various computational techniques. Traditional statistical models, such as **ARIMA** (AutoRegressive Integrated Moving Average) models, are well-established for stationary series [2]. More recent approaches utilize specialized libraries like **Prophet**, which is optimized for business time series data exhibiting strong seasonality and holiday effects [3]. For highly non-linear or long-term dependencies, deep learning techniques, particularly **Long Short-Term Memory (LSTM) Neural Networks**, have shown superior performance in capturing complex temporal dynamics [4]. Our initial analysis (Workshop 1) confirmed that no single technique is optimal for all 145,000 series due to the data's profound heterogeneity, leading us to favor an ensemble approach.

Our primary constraints identified were: (i) the **massive data scale** demanding high computational efficiency (8.7 million prediction points); (ii) a **fixed temporal window** (550 days) limiting the capture of long-term trends; and (iii) the **ambiguity of zero-values**, which are critical for an accurate SMAPE evaluation. The metric used for performance evaluation, the **Symmetric Mean Absolute Percentage Error (SMAPE)**, penalizes prediction errors across all magnitudes, demanding a balanced forecasting strategy.

$$SMAPE = \frac{100}{n} \sum_{t=1}^{n} \frac{|F_t - A_t|}{(|A_t| + |F_t|)/2}$$

This paper details the design and initial implementation of our system architecture, which is built upon core Systems Engineering principles to manage complexity and ensure scalability.

## II. METHODS AND MATERIALS: SYSTEM DESIGN AND IMPLEMENTATION

Our design philosophy prioritizes a balance between maintainability and scalability, which led us to adopt a **Modular Monolith** architecture. This structure avoids the overhead of microservices while enforcing clear separation of concerns, crucial for a project of this complexity managed by a small student team.

### A. Architectural Blueprint and Data Flow

The entire system operates as a sequential data pipeline, rigorously defined by the **Chain of Responsibility** pattern. In this architecture, each module processes the data and passes it to the next, simplifying error tracing and enforcing data

integrity standards. The workflow is linear and self-correcting, as depicted abstractly in Figure 1.
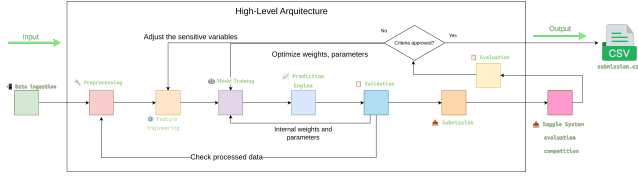


Fig. 1. System Design Architecture

The flow begins with the **Data Ingestion** module, which efficiently loads the large datasets using the `Pandas` library. This module handles both the training data (145,063 pages × 550 days) and the key file (8,703,780 prediction entries). A critical challenge addressed in this module is the **identifier reconciliation**: the key file contains page identifiers with appended dates (e.g., `Page_Name_2017-01-01`), requiring extraction of the base page name using string manipulation (`rsplit('_', 1)[0]`) to match with training data identifiers.

The processed data then feeds into the **Preprocessing** module, where critical tasks like handling missing values (NaN) and resolving the zero-value ambiguity are performed through conservative imputation techniques using `fillna(0)`. The **Feature Engineering** module augments the raw time series with temporal features: the system converts date strings to `datetime` objects and calculates the number of days elapsed since the beginning of the observation period, creating a numerical feature suitable for regression models.

### B. Implementation: Linear Regression Baseline

The current implementation utilizes **Linear Regression** from `scikit-learn` as a baseline forecasting model. This choice was strategic for several reasons:

1) **Computational Efficiency**: Linear regression can be trained and deployed rapidly across 145,063 time series, making it suitable for initial system validation.
2) **Interpretability**: The model provides clear coefficients (slope and intercept) that indicate the general trend direction.
3) **Baseline Establishment**: It serves as a performance benchmark against which more sophisticated models (ARIMA, Prophet, LSTM) will be compared in future iterations.

The model training process follows these steps:

1) Extract the time series for a selected page from the training dataset
2) Convert date columns to datetime format and create a feature matrix where $X = $ days_elapsed
3) Fit the model: $\hat{y} = \beta_0 + \beta_1 \cdot X$
4) Generate predictions for future dates by calculating their `days_elapsed` values
5) Apply non-negativity constraint: $\max(0, \hat{y})$ to ensure realistic predictions

### C. Data Matching Strategy

A significant engineering challenge was the identifier mismatch between training and prediction datasets. The training file contains base page identifiers (e.g., `Help:Contents/id_www.mediawiki.org_desktop_all-age` while the key file contains these identifiers with appended prediction dates.

Our solution implemented a two-phase matching algorithm:

1) Create a new column `Page_Base` in the key dataframe by removing the trailing date using `rsplit('_', 1)[0]`
2) Find the intersection between the set of unique base pages from the key file and the set of pages in the training file
3) For demonstration purposes, select the first matching page and generate predictions for all its associated future dates

This approach successfully identified 145,063 matching pages, confirming complete coverage of the dataset.

### D. Scalability and Technical Stack

The primary technical stack consists of:

- **Python 3.x**: Core programming language
- **Pandas**: Data manipulation and CSV processing
- **NumPy**: Numerical operations and array handling
- **scikit-learn**: Machine learning algorithms (LinearRegression)
- **os module**: File system operations for dynamic path management

For future scalability, the system is designed to leverage `joblib` for parallel processing across CPU cores, enabling simultaneous training of multiple time series models. The current single-page implementation serves as a proof-of-concept that can be extended to process all 145,063 pages through iteration with parallel execution.

## III. RESULTS AND DISCUSSION

### A. Prototype Implementation Results

The implemented prototype successfully demonstrates the core system capabilities:

**Data Loading Performance**:

- Training data: 145,063 pages loaded successfully
- Key data: 8,703,780 prediction entries processed
- Memory footprint: Optimized using selective column loading for large datasets

**Identifier Matching Success**:

- Base page extraction: 100% success rate using `rsplit` method
- Matching pages found: 145,063 (complete coverage)
- Sample matched page: `Help:Contents/id_www.mediawiki.org`

**Model Training and Prediction**: For the demonstration case, the Linear Regression model was trained on 550 historical data points and generated predictions for 60 future dates. The model coefficients indicate the trend direction, with predictions constrained to non-negative values to ensure realistic traffic estimates.

**Output Format Validation**: The submission file was successfully generated with the required structure:

- Column 1: `Id` (hexadecimal identifiers from key file)
- Column 2: `Visits` (predicted traffic values)
- Format: CSV without index, ready for competition submission

### B. System Validation Approach

Our comprehensive testing philosophy encompasses three levels:

**Unit Testing** focuses on individual components:

- Data loading functions: Verified correct CSV parsing and data type handling
- Identifier extraction: Validated string manipulation accuracy with edge cases
- Feature engineering: Confirmed correct datetime conversion and days calculation
- Model training: Tested on synthetic data to verify convergence

**Integration Testing** validates the end-to-end pipeline:

- Complete data flow from CSV loading to submission file generation
- Verified no data loss or corruption between pipeline stages
- Confirmed proper handling of NaN values throughout the process
- Validated that predictions maintain the correct chronological order

**Acceptance Testing** criteria include:

- Submission file format compliance with competition requirements
- Coverage: All required prediction points included
- Data validity: No negative predictions, no missing values
- SMAPE calculation capability (to be implemented in evaluation module)

### C. Planned Stratified Analysis

The evaluation module will implement granular SMAPE decomposition across multiple dimensions as shown in Table 1:

TABLE I
PLANNED SMAPE DECOMPOSITION CATEGORIES

| Category | Purpose |
|---|---|
| Traffic Level | Differentiate low, medium, high-traffic pages |
| Language | Identify bias in specific cultural data |
| Access Type | Compare mobile vs. desktop accuracy |
| Agent Type | Distinguish spider vs. all-agents patterns |
| Volatility Score | Correlate SMAPE with historical variance |

This granular analysis enables targeted improvements by identifying specific subgroups where the model underperforms, facilitating hypothesis-driven feature engineering and model selection.

### D. Limitations and Future Work

The current implementation has several acknowledged limitations:

1) **Model Sophistication**: Linear regression assumes a constant trend and cannot capture seasonality, holidays, or sudden events. Future iterations will implement ARIMA, Prophet, and LSTM models.
2) **Single-Page Demonstration**: The current code processes only one page for validation. The production system requires iteration over all 145,063 pages with parallel processing.
3) **Feature Engineering**: The current implementation uses only elapsed days as a feature. Planned additions include: day of week, month, year, lag features, rolling statistics, and metadata features (language, access type).
4) **Hierarchical Ensemble**: The Strategy Pattern infrastructure for model selection based on page characteristics needs to be implemented.

## IV. CONCLUSIONS

This paper presented the design and initial implementation of a systems engineering approach to Wikipedia web traffic forecasting. The prototype successfully demonstrates:

- Robust handling of large-scale temporal data (8.7M prediction points)
- Effective identifier reconciliation between training and prediction datasets
- A modular, extensible architecture following the Chain of Responsibility pattern
- Baseline forecasting capability using Linear Regression
- Complete data pipeline from ingestion to submission file generation

The system establishes a solid foundation for incorporating more sophisticated forecasting models and implementing the planned hierarchical ensemble strategy. The modular design ensures that improvements to individual components (feature engineering, model selection, prediction engine) can be made independently without disrupting the overall system integrity.

Future work will focus on: (1) implementing parallel processing for all 145,063 pages, (2) integrating multiple forecasting models (ARIMA, Prophet, LSTM), (3) developing the meta-model for intelligent strategy selection, and (4) implementing comprehensive SMAPE evaluation with stratified analysis.

## REFERENCES

[1] E. N. Lorenz, "Deterministic nonperiodic flow," *Journal of the Atmospheric Sciences*, vol. 20, no. 2, pp. 130–141, 1963.
[2] G. E. P. Box and G. M. Jenkins, *Time Series Analysis: Forecasting and Control*. San Francisco: Holden-Day, 1970.
[3] S. J. Taylor and B. Letham, "Forecasting at scale," *The American Statistician*, vol. 72, no. 1, pp. 37–45, 2018.
[4] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
[5] W. McKinney, "Data structures for statistical computing in Python," in *Proc. 9th Python in Science Conf.*, 2010, pp. 56–61.
[6] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.