

# **Developing a Comprehensive Research Plan for 'RuggedAtTick' Value Prediction**

## **Executive Summary**

This report outlines a comprehensive research plan to explore and evaluate advanced methodologies for predicting 'ruggedAtTick' values within the application, building upon and potentially enhancing the existing Bayesian indicator. The 'ruggedAtTick' metric, characterized by its sequential and potentially event-driven nature, is critical for app performance, user experience, and proactive system management. The current Bayesian indicator, while providing a statistical baseline, may possess limitations in capturing complex, non-linear dynamics or adapting to irregular data patterns.

The proposed plan is structured into distinct phases: initial data preparation and exploratory analysis to understand inherent time series properties, a detailed exploration of traditional statistical models, advanced machine learning regression techniques, and cutting-edge ensemble and hybrid forecasting paradigms. Emphasis is placed on meticulous feature engineering to transform raw, potentially irregular 'ruggedAtTick' data into suitable inputs for sophisticated models. Rigorous evaluation protocols, including walk-forward validation and a diverse set of performance metrics, are central to the plan. This systematic approach aims to identify methodologies that can significantly improve predictive accuracy, offer complementary insights into underlying app dynamics, and provide a more robust and adaptive framework for 'ruggedAtTick' prediction, ultimately supporting more informed decision-making and system optimization.

## **1. Introduction: Context and Objectives for 'RuggedAtTick' Prediction**

The accurate prediction of key operational metrics is paramount for maintaining the stability, performance, and user satisfaction of modern applications. Within this context, the 'ruggedAtTick' metric has been identified as a crucial data point, necessitating a deeper investigation into its predictive modeling. This section establishes the foundational understanding of 'ruggedAtTick', reviews the current predictive approach, and articulates the strategic objectives driving this research initiative.

### **1.1. Understanding the 'ruggedAtTick' Metric and its Business Importance**

The 'ruggedAtTick' metric represents a critical, sequential data point within the application's operational framework. The provided data format, `gameId`ruggedAtTick, strongly suggests that these are discrete observations tied to a specific `gameId`, implying a per-game or per-session context. The term "tick" typically refers to an instantaneous event or a discrete time step within a system, such as a game engine's internal clock or a financial market's trade events.<sup>1</sup> This indicates that 'ruggedAtTick' values are likely associated with specific, discrete events or measurements rather than continuous, uniformly sampled data.

A fundamental characteristic of such "tick" data, particularly in dynamic application environments, is its potentially discrete and event-driven nature. Unlike traditional time series that might exhibit measurements at fixed, consistent time intervals, 'ruggedAtTick' values could arrive unpredictably, akin to stock trades or sensor warnings.<sup>2</sup> This inherent irregularity, where intervals between consecutive observations can vary, significantly influences the selection and application of forecasting models. Classical time series models often assume regularly spaced data, which would necessitate substantial preprocessing, such as aggregation or interpolation, to conform to their input requirements. Conversely, specialized models designed to handle non-uniform time steps, or robust feature engineering that explicitly encodes temporal irregularities, may be required to preserve the granular information and true dynamics of the system.<sup>3</sup> The predictive accuracy of 'ruggedAtTick' values is crucial for various app functions, potentially influencing real-time resource allocation, optimizing user experience by anticipating performance bottlenecks, or enabling proactive anomaly detection. High-fidelity predictions can lead to more responsive and efficient system management, ultimately enhancing

overall app stability and user engagement.

## **1.2. Review of the Existing Bayesian Indicator Approach**

The current reliance on a "Bayesian indicator" for 'ruggedAtTick' prediction suggests a probabilistic framework that leverages prior knowledge and observed data to update beliefs about future values. This approach provides a strong statistical foundation, inherently quantifying uncertainty and allowing for the incorporation of domain expertise through prior distributions.

However, a purely Bayesian indicator, particularly if based on simpler statistical assumptions, may encounter limitations when confronted with complex or non-linear data dynamics. While Bayesian methods are robust for incorporating prior knowledge and quantifying uncertainty, simpler statistical models often assume linearity or specific distributional forms. Many time series contain underlying patterns and trends, such as a continuous increase, which result in statistical properties like mean or variance changing over time, leading to non-stationarity.<sup>4</sup> If 'ruggedAtTick' exhibits intricate non-linear trends, sudden shifts, or complex interactions that are not easily captured by predefined probabilistic structures, a simpler Bayesian formulation might struggle to adapt quickly or accurately. For instance, models like ARIMA, while powerful, have a "limited ability to model nonlinear relationships".<sup>5</sup> The request to "explore additional methodologies" suggests a perceived need to enhance or move beyond the existing capabilities, implying that the current Bayesian approach may not fully capture all the nuances of the data. The potential for non-linear or complex patterns in 'ruggedAtTick' data therefore necessitates the exploration of more flexible, data-driven machine learning models. These models are inherently better equipped to capture such relationships, thereby potentially improving predictive accuracy and adaptability beyond the existing linear or probabilistic assumptions.

## **1.3. Objectives of This Comprehensive Research Plan**

This research plan aims to systematically investigate and evaluate a diverse portfolio of advanced time series forecasting methodologies. The primary objective is to identify and implement methodologies that can either significantly enhance predictive

accuracy, offer complementary insights into the underlying system dynamics, or provide a more robust and adaptive framework for 'ruggedAtTick' prediction compared to the existing Bayesian indicator. The plan will encompass a thorough approach covering data preparation, advanced feature engineering, rigorous model selection, practical implementation strategies, and comprehensive evaluation protocols. This systematic exploration seeks to uncover the most effective and reliable methods for forecasting 'ruggedAtTick' values, thereby supporting more informed and proactive decision-making within the application environment.

## **2. 'RuggedAtTick' Data Characteristics and Preprocessing**

Effective predictive modeling begins with a profound understanding of the data's inherent characteristics and meticulous preparation. This section details the initial steps for handling the 'ruggedAtTick' data, focusing on its extraction, analysis of its time series properties, and outlining essential feature engineering techniques crucial for preparing the data for advanced modeling.

### **2.1. Initial Data Extraction and Structuring**

The raw data provided, `gameIdruggedAtTick020250717-73861dcdae764117319120250717-6cfae8adbd814e672172...`, clearly indicates a gameId followed immediately by a ruggedAtTick value. This structure implies a direct mapping between a game identifier and a specific metric reading. The ruggedAtTick values appear to be integers, while the gameId is a string identifier that likely contains a timestamp component (e.g., 20250717 suggests a date).

The first critical step involves parsing this raw string to accurately extract individual gameId and ruggedAtTick pairs. This requires a robust parsing script that can reliably separate the identifier from its corresponding numerical value. Once extracted, these pairs must be structured into a tabular format, such as a dataframe, ensuring correct data types for both components and maintaining chronological ordering if the gameId contains time information. This initial data handling step is fundamental, as all

subsequent analytical and modeling processes depend on the accurate transformation of this raw input into a usable, structured format. The following table illustrates a sample of the parsed data, demonstrating the successful transformation and structuring:

**Table 1: 'ruggedAtTick' Data Sample and Extracted Values**

Game ID (Partial)	Date Component (Extracted)	ruggedAtTick Value
020250717-7386...	20250717	319
020250717-6cfa...	20250717	217
020250717-b6bc...	20250717	323
020250717-5291...	20250717	784
020250717-6e03...	20250717	379

**2.2. Analysis of Time Series Properties: Regularity, Trends, Seasonality, and Outliers**

Before applying any predictive model, a thorough exploratory data analysis (EDA) is indispensable. This involves visualizing time plots and correlograms to identify key time series properties such as non-stationarity, underlying trends, seasonal patterns, and the presence of outliers.<sup>6</sup>

A crucial property for many traditional time series models, such as ARIMA, is stationarity.<sup>4</sup> Stationary series exhibit consistent statistical properties, like mean and variance, over time, whereas non-stationary series often reflect trends or seasonality, leading to changing properties and potentially unreliable predictions.<sup>6</sup> Techniques like differencing (calculating differences between consecutive observations) or logarithmic transformations can be employed to achieve stationarity, thereby facilitating accurate forecasts.<sup>4</sup> Statistical tests, including the Augmented Dickey-Fuller (ADF) and Kwiatkowski-Phillips-Schmidt-Shin (KPSS) tests, are critical tools to rigorously assess these properties and guide necessary data

transformations.<sup>6</sup>

Trends describe long-term patterns in the data where measurement points either consistently increase or decrease over extended periods.<sup>4</sup> Seasonality refers to recurring patterns that occur at fixed, predictable intervals, such as daily, hourly, or annually.<sup>4</sup> Cycles, while also fluctuations, are typically longer-term and less regular than seasonal changes, often linked to broader economic or system cycles.<sup>4</sup> The presence and nature of these patterns directly inform the choice of appropriate forecasting methods, such as specific exponential smoothing variants<sup>10</sup> or Seasonal ARIMA (SARIMA) models.<sup>5</sup>

A significant consideration for 'ruggedAtTick' data is its potential for irregularity. The "tick" nature strongly suggests that data points might represent discrete events that do not necessarily occur at fixed, regular time intervals.<sup>2</sup> Most classical time series models, including ARIMA and standard Exponential Smoothing, are designed for and perform optimally with regularly sampled data.<sup>3</sup> If 'ruggedAtTick' is indeed irregular, a direct application of these models without prior transformation would be inappropriate and could yield misleading results. This fundamental data characteristic necessitates a critical strategic decision: either transform the irregular data into a regular time series through aggregation (e.g., summing or averaging 'ruggedAtTick' values over fixed time windows) or interpolation, or alternatively, employ specialized models or feature engineering techniques inherently capable of handling non-uniform time steps. While aggregation or interpolation can regularize the data, interpolation must be approached cautiously, as it can introduce artificial patterns or smooth out important "rugged" events, potentially distorting the true dynamics.<sup>8</sup> The inherent irregularity of 'ruggedAtTick' data thus causes the need for sophisticated data transformation or specialized model selection. This choice, in turn, has significant implications for the complexity of the modeling pipeline and the potential fidelity of the predictions to the true underlying system behavior. Ignoring this irregularity could lead to spurious correlations or highly inaccurate forecasts, making this a critical area of focus in the research plan.

Finally, outliers, defined as irregular patterns in time series data that follow neither seasonality nor a trend, must be identified and appropriately handled.<sup>4</sup> Their presence can significantly skew model training and lead to inaccurate predictions, making outlier detection and robust treatment crucial for model stability and reliability.<sup>12</sup>

### **2.3. Essential Feature Engineering Strategies for 'ruggedAtTick'**

Feature engineering is a critical process for transforming raw time series data into meaningful inputs for predictive models, particularly for non-deep-learning models.<sup>5</sup> This process enhances the accuracy of forecasting by capturing complex patterns and relationships that are not immediately apparent in the raw data.<sup>13</sup> For 'ruggedAtTick' data, given its potential for irregularity and discrete events, feature engineering serves as a crucial bridge to enable the application of powerful machine learning algorithms.

The challenge extends beyond simply deciding which features to create (e.g., lags, rolling statistics, time-based components) to meticulously designing how to construct these features when the underlying time intervals are not constant. This means that for irregular 'ruggedAtTick' data, features must be explicitly designed to encode the temporal relationships and event characteristics. For instance, instead of just a simple lag\_1 feature, one might need time\_since\_last\_tick or ruggedAtTick\_per\_second\_in\_last\_minute. This transformation is a causal step: the quality and relevance of the engineered features directly determine the success or failure of the subsequent machine learning models in capturing the "rugged" patterns and making accurate predictions. Without this crucial step, the models will not be able to fully leverage their non-linear capabilities. The inherent nature of 'ruggedAtTick' as an irregular, event-driven time series therefore causes the indispensable role of advanced feature engineering. This engineering, in turn, causes the transformation of raw data into a format suitable for powerful regression models, thereby causing the potential for improved predictive accuracy and the ability to capture complex, non-linear dynamics.

### **2.3.1. Lag Features and Rolling Window Statistics**

Lag features incorporate previous 'ruggedAtTick' values as predictors, providing essential context and capturing short-term dependencies and cyclical patterns.<sup>5</sup> This is fundamental for models with autoregressive components. For example, using

ruggedAtTick values from  $t-1$ ,  $t-2$ , and so forth, to predict the value at time  $t$ .<sup>14</sup> The number of lags to include (

$p$  in ARIMA, or  $n_{in}$  in sliding window for Random Forest) is a critical parameter that

needs to be determined through analysis of autocorrelation plots.<sup>9</sup>

Rolling window statistics, such as moving averages, standard deviations, and trends, are computed over a sliding window of past observations. These statistics help to smooth out noise, highlight local trends, and capture the temporal dynamics and volatility of the data.<sup>5</sup> This enables models to adapt to changing patterns and detect anomalies, providing a smoothed representation of the metric's recent behavior.

### 2.3.2. Time-Based and Cyclical Features

Deriving features from the timestamp component within the gameld (e.g., day of the week, month, hour of day) can effectively model temporal effects and recurring patterns.<sup>5</sup> This is crucial for capturing seasonality that might influence 'ruggedAtTick' values. Trigonometric transformations (e.g., sine and cosine functions of time components) can encode periodic features more effectively, allowing models to learn smooth cyclical patterns without assuming linear relationships between time and the target variable.<sup>13</sup> Additionally, incorporating binary indicators for holidays or other special in-app events can capture external influences that might cause sudden shifts or anomalies in 'ruggedAtTick'.<sup>12</sup>

### 2.3.3. Handling Irregular Intervals and Missing Data

If 'ruggedAtTick' is inherently irregular, where data points do not arrive at fixed intervals, specific strategies are required. Standard approaches include filling missing values forward or backward, assuming values remain constant until a new entry appears, or employing interpolation techniques to estimate missing values based on surrounding data.<sup>2</sup> However, interpolation must be performed with caution to avoid introducing artificial data or smoothing out genuine "rugged" characteristics.<sup>8</sup>

Alternatively, and perhaps more effectively for truly irregular 'tick' data, features can be explicitly engineered to capture the *duration* between ticks, the *rate* or *density* of ticks within a given time period (e.g., ticks per minute), or the *change* in 'ruggedAtTick' value per unit of time.<sup>3</sup> These features directly convey the temporal irregularity to the model, allowing it to learn from the timing of events rather than forcing the data into a



regular grid.

Table 2: Key Feature Engineering Techniques for 'ruggedAtTick' Data

Feature Type	Description	Relevance to 'ruggedAtTick'	Relevant References
Lag Features	Previous values of 'ruggedAtTick' (e.g., ruggedAtTick_t-1, ruggedAtTick_t-2).	Captures short-term dependencies and autoregressive patterns crucial for predicting the next 'tick' value.	5
Rolling Window Statistics	Moving averages, standard deviations, min/max over a defined past window (e.g., last 5 ticks, last 60 seconds).	Smooths noise, highlights local trends, captures volatility, and helps detect anomalies in 'ruggedness'.	5
Time-Based Features	Components extracted from gameId timestamp (e.g., hour of day, day of week, month, minute of hour).	Models cyclical and seasonal patterns, accounting for predictable variations in app usage or event frequency.	5
Irregular Interval Features	Time elapsed since last tick, number of ticks in a fixed time window, rate of change per unit time.	Directly addresses the irregular nature of 'tick' data by encoding temporal gaps and event density.	3
Fourier Transforms	Decomposition of time series into frequency components.	Reveals underlying periodicities and cyclical behaviors not obvious in the time domain, useful for complex seasonality.	13

### 3. Exploration of Advanced Forecasting Methodologies

This section delves into various advanced forecasting methodologies, categorizing them into traditional statistical models, machine learning regression models, and advanced paradigms like ensembles and hybrid approaches. Each methodology is discussed in terms of its applicability, strengths, and limitations for 'ruggedAtTick' prediction, drawing heavily on existing research.

#### 3.1. Traditional Statistical Time Series Models

Traditional statistical models serve as foundational baselines and are particularly effective for capturing linear relationships within time series data.

##### 3.1.1. ARIMA/SARIMA Family: Strengths for Linear Dependencies and Stationarity

The ARIMA (AutoRegressive Integrated Moving Average) model is a class of statistical models widely used for analyzing time series data and predicting future values, especially when temporal autocorrelation is present.<sup>4</sup> ARIMA models are composed of three core components:

- **AR (Autoregressive):** The autoregressive component, denoted by the parameter  $p$  (also known as lag order), utilizes past values of the series to make future forecasts, thereby recognizing and incorporating existing temporal structures.<sup>4</sup>
- **I (Integrated):** The "integrated" component refers to the process of differencing, denoted by the parameter  $d$  (degree of differencing). This technique is applied to make non-stationary data stationary by removing underlying trends and patterns, which is a prerequisite for many statistical models.<sup>4</sup>
- **MA (Moving Average):** The moving average component, represented by the parameter  $q$  (order of moving average), accounts for the random errors or noise from previous periods, smoothing out fluctuations.<sup>4</sup>

The SARIMA (Seasonal ARIMA) model extends the ARIMA framework by explicitly accounting for seasonality, incorporating additional seasonal parameters ( $P$ ,  $D$ ,  $Q$ ,  $S$ )

to capture recurring patterns at fixed intervals.<sup>4</sup> SARIMAX further enhances this by including exogenous variables, allowing for the incorporation of external factors that might influence the time series.<sup>5</sup> These models are applicable for 'ruggedAtTick' if the data exhibits linear trends and seasonality, provided it can be made stationary through appropriate differencing.<sup>4</sup> They are effective in capturing clear temporal structures. However, a significant limitation of ARIMA models is their "limited ability to model nonlinear relationships"<sup>5</sup> and their struggles with complex non-linear patterns.<sup>5</sup> Crucially, these models typically assume regularly spaced data<sup>3</sup>; thus, irregular 'ruggedAtTick' data would necessitate aggregation or interpolation, potentially leading to a loss of granular information and a less accurate representation of the discrete events.

### 3.1.2. Exponential Smoothing Methods: Simple, Holt's, and Holt-Winters

Exponential smoothing methods are time series forecasting techniques that predict future values by using an exponentially weighted average of past observations, assigning greater weight to more recent data points.<sup>10</sup> These methods are generally recognized for providing accurate short-term forecasts.<sup>11</sup>

- **Simple Exponential Smoothing (SES):** Also known as single exponential smoothing, this is the most basic form and assumes that the time series has no underlying trend or seasonal pattern. It is suitable for stationary data.<sup>10</sup>
- **Holt's Linear Exponential Smoothing (Double ES):** This method is designed for time series data that displays a linear trend but lacks a seasonal pattern. It incorporates an additional smoothing parameter,  $\beta$ , specifically for the trend component.<sup>10</sup>
- **Holt-Winters' Exponential Smoothing (Triple ES):** This is the most comprehensive of the three, used for time series data that exhibits both a trend and a seasonal component. It utilizes three smoothing parameters: for the level, trend, and seasonality.<sup>10</sup>

These methods are useful for 'ruggedAtTick' if the metric displays clear trends and/or seasonality, as they are adaptive to changing trends.<sup>10</sup> However, they are primarily designed for univariate time series data<sup>11</sup> and tend to produce less reliable long-term forecasts.<sup>11</sup> Similar to ARIMA, their direct application assumes regular intervals. Therefore, irregular 'ruggedAtTick' data would require preprocessing, such as

aggregation or interpolation, before these models can be effectively applied.

The necessity to transform irregular 'ruggedAtTick' data into a regular, stationary series for these traditional models introduces a critical layer of abstraction. This transformation, often involving aggregation (e.g., counting ticks per minute) or interpolation, might inherently smooth out or distort the precise "rugged" events, thereby losing granular information and unique characteristics of the "tick" nature. This creates a trade-off: while these models offer simplicity and interpretability for linear patterns, achieving compatibility with irregular data may cause a loss of fidelity to the original data's dynamics. This could potentially lead to less accurate or less nuanced predictions compared to models that can inherently handle such data or are robust to complex feature engineering. The inherent limitations of traditional statistical models with irregular data therefore cause the necessity for extensive data transformation. This transformation, in turn, may cause a potential loss of granular information crucial for predicting specific "rugged" events, thereby pushing the research towards more flexible machine learning or specialized models.

### 3.2. Machine Learning Regression Models for Time Series

Machine learning regression models offer greater flexibility in capturing non-linear relationships and can effectively leverage complex features engineered from the time series data.

#### 3.2.1. Tree-Based Ensembles: Random Forest and Gradient Boosting (XGBoost, LightGBM)

Tree-based ensemble methods are powerful and widely used for predictive modeling.

- **Random Forest:** This is a popular and effective ensemble machine learning algorithm that constructs a multitude of decision trees during training.<sup>14</sup> It enhances robustness and reduces variance by employing bagging (bootstrap aggregation) and by selecting a random subset of input features at each split point.<sup>14</sup> For regression problems, predictions are derived by averaging the predictions from all individual trees within the ensemble.<sup>14</sup>
- **Gradient Boosting (XGBoost, LightGBM):** Gradient boosting is a machine

learning technique that builds a prediction model as an ensemble of weaker prediction models, typically decision trees, in a sequential manner.<sup>17</sup> Each new model is trained to correct the errors made by the previously built models, thereby iteratively refining the overall prediction. XGBoost trees generally grow depth-wise, while LightGBM trees grow leaf-wise, which often makes LightGBM faster and more efficient, especially for large datasets.<sup>19</sup>

These models are powerful for handling nonlinear relationships within time series data.<sup>5</sup> They can be effectively applied to time series forecasting by transforming the sequential data into a supervised learning problem using a sliding-window representation and lagged variables.<sup>14</sup> They excel when provided with feature-engineered time series, where lagged values, rolling statistics, and external variables serve as predictors.<sup>5</sup> Furthermore, they can be utilized to predict residuals from a simpler linear trend model, capturing the non-linear deviations.<sup>21</sup> A primary consideration, however, is that these models typically "requires manual feature engineering"<sup>5</sup> to convert time series data into a tabular format. They may also struggle with extrapolation beyond the range of the training data if underlying trends are not explicitly captured through differencing or specific trend features.<sup>21</sup>

### **3.2.2. Deep Learning Approaches: Long Short-Term Memory (LSTM) Networks**

Long Short-Term Memory (LSTM) networks are specialized recurrent neural networks (RNNs) designed to process sequential input data by looping over time steps and continuously updating an internal "memory" state.<sup>22</sup> This sophisticated memory system, facilitated by internal gating mechanisms, enables LSTMs to effectively capture long-term dependencies in sequential data, thereby mitigating the vanishing gradient problem commonly encountered in traditional RNNs.<sup>5</sup>

LSTMs are particularly well-suited for "complex time series with long-term dependencies and nonlinear relationships".<sup>5</sup> They possess the unique ability to process an entire sequence of data and associate historical memories with the current input, making them ideal for forecasting subsequent values of a time series using previous time steps as input.<sup>22</sup> The implied non-linear and potentially volatile nature of 'ruggedAtTick' makes machine learning models, especially LSTMs, highly effective candidates. Their effectiveness is contingent upon substantial data and computational resources for LSTMs, which causes a strategic decision point regarding complexity versus resource allocation. However, LSTMs are "computationally intensive

and requires large datasets" <sup>5</sup> for effective training. Their architectural complexity can also make them challenging to design and tune for optimal performance.

### 3.2.3. Other Relevant Models: Prophet and Support Vector Machines

Beyond tree-based ensembles and LSTMs, other machine learning models offer distinct advantages for time series forecasting:

- **Prophet:** Developed by Facebook, Prophet is an open-source tool specifically designed for forecasting time series that exhibit strong seasonality and may contain missing data.<sup>5</sup> It is particularly useful for business time series where trends and seasonality are prominent and often change over time.<sup>5</sup>
- **Support Vector Machines (SVM):** SVMs can be adapted for time series regression or classification tasks when combined with appropriate feature engineering.<sup>5</sup> They are effective in finding complex non-linear decision boundaries or regression functions. However, SVMs can be "sensitive to parameter tuning and feature scaling" <sup>5</sup>, requiring careful optimization for robust performance.

## 3.3. Advanced Forecasting Paradigms

These advanced paradigms combine the strengths of multiple models or address unique data characteristics like hidden states, offering enhanced robustness and predictive power.

### 3.3.1. Ensemble Learning: Bagging, Boosting, and Stacking for Enhanced Robustness

Ensemble learning involves combining the predictions of multiple individual models into a single, more robust forecast, thereby increasing overall prediction performance and stability.<sup>17</sup> This approach is particularly valuable for 'ruggedAtTick' prediction as it can significantly improve accuracy and robustness by mitigating the limitations of

single models, especially their susceptibility to noise and outliers.<sup>24</sup>

- **Bagging (Bootstrap Aggregation):** This technique reduces variance and helps prevent overfitting by training multiple "weak" models (e.g., decision trees) independently on different bootstrapped (randomly sampled with replacement) subsets of the training data.<sup>17</sup> The final prediction is then obtained by averaging the predictions from these individual models.<sup>17</sup> Random Forest is a prominent example of a bagging method.<sup>14</sup>
- **Boosting:** Boosting is an ensemble technique that constructs a strong predictive model by sequentially combining several weak models.<sup>17</sup> Each successive model is trained to correct the errors made by its predecessor, thereby iteratively reducing bias and improving overall accuracy. Instances that were misclassified by previous models are given higher weights, ensuring the subsequent models focus on more challenging cases.<sup>17</sup> XGBoost, LightGBM, and AdaBoost are well-known and highly effective boosting algorithms.<sup>17</sup>
- **Stacking (Stacked Generalization):** This is a more advanced ensemble method that leverages model diversity. It involves training multiple diverse "base models" (also known as level-0 models) on the same dataset. The predictions generated by these base models are then used as input features to train a "meta-model" (or level-1 model), which makes the final prediction.<sup>17</sup> Stacking can capture a wide range of patterns in the data by combining different types of models, often leading to improved performance over individual models.<sup>18</sup>

Ensemble methods are a powerful approach to combine different algorithms for a final forecast.<sup>25</sup> However, they can increase overall model complexity and may reduce interpretability compared to individual models<sup>17</sup>, requiring careful implementation and tuning of combining strategies or meta-models.

### 3.3.2. Hybrid Models: Combining Statistical and Machine Learning Strengths

Hybrid models combine different techniques or approaches to capture and forecast complex patterns and dynamics present in time series data.<sup>25</sup> They integrate the strengths of multiple methods to overcome the limitations of individual models, aiming to enhance overall forecasting performance.<sup>25</sup> This approach is highly promising if 'ruggedAtTick' is found to have both linear (e.g., underlying trends related to game progression) and nonlinear (e.g., sudden "rugged" events, player-induced volatility) components, which is common in real-world time series.<sup>26</sup> This approach can

compensate for the weaknesses of one model type with the strengths of another.<sup>26</sup>

Common hybrid types include combining ARIMA with Exponential Smoothing, or ARIMA with Neural Networks.<sup>25</sup> Another example is Exponential Smoothing combined with Recurrent Neural Networks.<sup>26</sup> A prevalent methodology involves fitting a statistical model (e.g., ARIMA) to capture the linear component of the time series. Subsequently, a machine learning model (e.g., a Neural Network or Gradient Boosting model) is trained on the residuals (the errors) from the statistical model to capture the remaining nonlinear relationships.<sup>26</sup> This methodology is plausible because real-world time series often comprise a combination of linear and nonlinear patterns, and hybridization provides a solution to the dilemma of linearity assumptions that often limit traditional approaches.<sup>26</sup> The primary limitation of hybrid models is their increased complexity in design, implementation, and maintenance compared to single-model approaches.

### 3.3.3. Hidden Markov Models (HMMs) for Event-Driven Sequences

Hidden Markov Models (HMMs) are statistical models employed to characterize systems undergoing changes in unobservable states over time.<sup>27</sup> They are particularly useful for modeling sequential data where the underlying process generating observations is obscured or unknown.<sup>27</sup> HMMs are composed of:

- **Hidden States:** These represent the underlying, unobservable structure of the system (e.g., different "game states," "server load regimes," or "system health phases" that influence 'ruggedAtTick' values).<sup>27</sup>
- **Observable States:** These are the values that can be directly measured (in this case, the 'ruggedAtTick' values themselves).<sup>27</sup>
- **Transition Probabilities:** These define the likelihood of moving from one hidden state to another.<sup>27</sup>
- **Emission Probabilities:** These describe the probability of observing a particular observable state given the current hidden state.<sup>27</sup>

HMMs are highly relevant if 'ruggedAtTick' values are generated by an underlying, unobservable process that transitions between different "regimes" or "states" within the app. For instance, if the "ruggedness" in 'ruggedAtTick' is not merely random noise but reflects a system transitioning between different operational modes (e.g., from a "stable" state to a "high load" state, or a "player engagement" state to a



"player disengagement" state), HMMs offer a unique and powerful capability. They can explicitly model these hidden states and their probabilistic influence on the observed 'ruggedAtTick' values. This moves beyond pure statistical correlation to inferring the underlying causal processes within the app. Furthermore, the ability of HMMs to handle sequences of events directly makes them particularly well-suited for the "tick"-based, potentially irregular nature of the data.<sup>3</sup> HMMs are excellent for handling uncertainty and temporal dependencies<sup>27</sup> and can predict future observations or categorize sequences based on the concealed process.<sup>27</sup> The potential existence of unobservable, underlying system states that cause the observed patterns in 'ruggedAtTick' makes HMMs a uniquely powerful candidate. They offer not just prediction but also inference about these hidden system dynamics, which can lead to a richer understanding and potentially more effective intervention strategies within the app. This provides a more mechanistic understanding than pure regression, directly impacting the ability to diagnose and respond to issues related to 'ruggedAtTick'.

However, HMMs require careful definition of the state space and observation space.<sup>28</sup> They can be computationally intensive and complex to parameterize and train, often utilizing algorithms like Baum-Welch or Viterbi.<sup>28</sup> The interpretability of the "hidden" states might also be challenging without deep domain expertise.

**Table 3: Comparison of Advanced Time Series Forecasting Methodologies**

Methodology	Strengths	Weaknesses/Limitations	Applicability to 'ruggedAtTick'	Relevant References
<b>ARIMA/SARIMA</b>	Strong for linear trends & seasonality; statistical rigor; interpretability.	Limited for non-linear relationships; assumes stationarity (requires differencing); typically assumes regular intervals.	Good for capturing linear trends and seasonality if data can be regularized and made stationary.	<sup>4</sup>
<b>Exponential Smoothing</b>	Simple & effective for short-term	Primarily univariate; less reliable for	Useful for clear trends/seasonality if data is	<sup>10</sup>

	forecasts; adaptive to changing trends; handles trends & seasonality (Holt-Winters).	long-term forecasts; assumes regular intervals.	regularized; adaptive for short-term changes.	
<b>Random Forest</b>	Handles non-linearity; robust to outliers; reduces variance (bagging); good with feature-engineered data.	Requires extensive manual feature engineering; may struggle with extrapolation beyond training range.	Powerful for non-linear 'rugged' patterns with strong feature engineering, especially for irregular data.	5
<b>Gradient Boosting (XGBoost, LightGBM)</b>	Powerful for non-linearity; high accuracy; reduces bias (boosting); efficient (LightGBM).	Requires extensive manual feature engineering; can be prone to overfitting without careful tuning.	Excellent for complex, non-linear 'rugged' patterns; can predict residuals; efficient for large datasets.	5
<b>LSTM Networks</b>	Captures long-term dependencies; models complex non-linear relationships in sequential data; "memory" system.	Computationally intensive; requires large datasets; complex to design and tune.	Ideal for highly complex 'rugged' patterns with long-term dependencies, especially if underlying processes are non-linear.	5
<b>Prophet</b>	Handles strong seasonality, holidays, missing data; user-friendly; robust to	Assumes additive relationships; less flexible for highly complex non-linear	Good for business-like time series with clear trends and multiple seasonalities,	5

	outages.	patterns.	robust to data gaps.	
<b>Hidden Markov Models (HMMs)</b>	Models systems with unobservable (hidden) states; handles uncertainty and temporal dependencies; suitable for event sequences.	Requires careful definition of state/observation space; computationally intensive for training; interpretability of hidden states.	Highly relevant if 'ruggedAtTick' is influenced by underlying, unobservable app/game states or event-driven processes.	27
<b>Ensemble Learning</b>	Combines predictions for enhanced accuracy & robustness; mitigates single model limitations; reduces overfitting/bias.	Increased model complexity; potentially reduced interpretability of the combined model.	Universal applicability to improve any base model's performance; useful for noisy or outlier-prone 'ruggedAtTick' data.	17
<b>Hybrid Models</b>	Combines strengths of statistical & ML models; captures both linear & non-linear components; compensates for weaknesses.	Increased complexity in design, implementation, and maintenance.	Highly promising if 'ruggedAtTick' has both linear trends and complex non-linear event-driven patterns.	25

## 4. Comprehensive Research Plan and Methodology

This section details a phased approach to implementing and evaluating the selected forecasting methodologies, ensuring a systematic and rigorous comparison against

the existing Bayesian indicator. The plan is designed to be iterative, allowing for refinement and adaptation based on empirical results.

#### **4.1. Phase 1: Data Preparation, Exploratory Analysis, and Baseline Establishment**

This initial phase is critical for establishing a robust foundation for all subsequent modeling efforts.

##### **4.1.1. Data Ingestion and Parsing**

The first activity involves developing and implementing robust scripts to accurately extract `gameId` and `ruggedAtTick` values from the raw data. This process will ensure that the extracted values are correctly typed (e.g., `ruggedAtTick` as numerical, `gameId` as string or parsed timestamp) and that the data is chronologically ordered. The parsing logic must account for the specific format of the `gameId` to derive precise timestamps, which are essential for time series analysis.

##### **4.1.2. Initial Exploratory Data Analysis (EDA)**

Following data ingestion, a thorough EDA will be conducted. This includes visualizing time plots of 'ruggedAtTick' values, potentially aggregated by `gameId` or across all games if a meaningful global trend exists. These visualizations will help identify visual trends, seasonal patterns, cyclical behaviors, and potential irregularities in the data.<sup>6</sup> Further analysis will involve examining statistical distributions, summary statistics (mean, variance, standard deviation), and autocorrelation function (ACF) plots to reveal dependencies and potential non-stationarity.<sup>6</sup>

##### **4.1.3. Stationarity Testing and Transformation**

A critical step in EDA is assessing data stationarity. Statistical tests such as the Augmented Dickey-Fuller (ADF) test and the Kwiatkowski-Phillips-Schmidt-Shin (KPSS) test will be applied to rigorously evaluate whether the time series properties (mean, variance, covariance) remain constant over time.<sup>6</sup> If the data is found to be non-stationary, appropriate transformations will be applied. Differencing, which calculates the differences between consecutive observations, is a primary technique to stabilize the mean and remove trends.<sup>4</sup> Logarithmic transformations can be used to stabilize variance.<sup>6</sup> The minimum required differencing to achieve a near-stationary series, where the ACF plot quickly reaches zero, will be determined.<sup>9</sup>

#### **4.1.4. Feature Engineering for Initial Models**

Based on the EDA, initial feature engineering will be performed. This includes creating lagged features of 'ruggedAtTick' values (e.g.,  $t-1$ ,  $t-2$ ) to capture short-term dependencies.<sup>5</sup> Rolling window statistics (e.g., moving averages, standard deviations over recent ticks) will also be generated to smooth noise and highlight local trends.<sup>5</sup> If the

gameId contains date/time information, time-based features such as hour of day, day of week, or month will be extracted to capture potential cyclical patterns.<sup>5</sup> For irregular data, initial features capturing the time intervals between ticks or the density of ticks within a fixed window will be explored to encode the temporal irregularity.<sup>3</sup>

#### **4.1.5. Establishing Baseline Performance**

The existing Bayesian indicator's performance will be thoroughly documented and quantified using appropriate time series evaluation metrics (e.g., RMSE, MAE, MAPE). This will serve as the crucial baseline against which all new methodologies will be compared. This step ensures that any proposed new model demonstrates a measurable improvement over the current approach.

### **4.2. Phase 2: Model Development and Iterative Refinement**

This phase involves the systematic development, training, and initial tuning of the selected forecasting methodologies.

#### **4.2.1. Selection of Candidate Models**

Based on the insights from Phase 1 EDA and the comparative analysis in Section 3, a subset of the most promising methodologies will be selected for in-depth development. This selection will prioritize models that align with the identified characteristics of 'ruggedAtTick' (e.g., non-linearity, irregularity, presence of hidden states). A multi-model approach, where various algorithms are applied for forecasting, is a good practice that can lead to improved accuracy.<sup>16</sup>

#### **4.2.2. Data Splitting and Walk-Forward Validation**

For time series data, traditional cross-validation methods are unsuitable as they can lead to data leakage by training on future information to predict the past.<sup>12</sup> Therefore, a rigorous walk-forward validation strategy will be employed.<sup>14</sup> The dataset will be split chronologically into training and test sets, ensuring that training data always precedes test data. The model will be trained on the initial training set, predict the next time step, and then the actual observation will be added to the training set for the next prediction. This iterative process will be repeated across the entire test period, providing a realistic assessment of out-of-sample performance.<sup>14</sup>

#### **4.2.3. Model Training and Hyperparameter Tuning**

Each selected model will be trained on the prepared and feature-engineered data. Hyperparameter tuning will be conducted systematically to optimize each model's performance. For instance, for ARIMA, the  $p$ ,  $d$ ,  $q$  orders will be determined using ACF/PACF plots and information criteria.<sup>9</sup> For tree-based models like XGBoost and

LightGBM, parameters such as

max\_depth, n\_estimators, and learning\_rate will be tuned.<sup>19</sup> For LSTMs, the number of hidden units, layers, and optimization algorithms will be explored.<sup>22</sup> Numerical optimization can be used to find optimal smoothing factors for exponential smoothing methods.<sup>11</sup> Early stopping rounds will be utilized during training of boosting models to prevent overfitting and improve performance.<sup>19</sup>

#### 4.2.4. Initial Performance Evaluation and Error Analysis

Following training, each model's performance will be evaluated using a comprehensive set of time series specific metrics. These include Mean Absolute Error (MAE), Root Mean Square Error (RMSE), Mean Absolute Percentage Error (MAPE), and Mean Absolute Scaled Error (MASE).<sup>24</sup> Error analysis will involve examining the residuals (the differences between actual and predicted values) to identify any systematic biases, patterns in prediction errors, or specific scenarios where models perform poorly. This analysis will inform subsequent refinement steps.

### 4.3. Phase 3: Advanced Model Integration and Comparative Analysis

This phase focuses on leveraging the strengths of multiple models through ensemble and hybrid approaches, and conducting a thorough comparative analysis.

#### 4.3.1. Implementation of Ensemble Methods

Ensemble learning techniques, including bagging, boosting, and stacking, will be implemented to combine the predictions from the best-performing individual models.<sup>17</sup>

- **Bagging:** Multiple instances of a base learner (e.g., decision trees) will be trained on bootstrapped samples of the data, and their predictions averaged.<sup>17</sup>
- **Boosting:** Sequential training of weak learners will be performed, with each

subsequent model focusing on correcting errors from previous ones (e.g., XGBoost, LightGBM).<sup>17</sup>

- **Stacking:** A meta-model will be trained on the predictions generated by several diverse base models, allowing the meta-model to learn the optimal way to combine their outputs.<sup>18</sup>

#### 4.3.2. Development of Hybrid Models

Hybrid models will be developed to combine the complementary strengths of different model types. A promising approach involves fitting a statistical model (e.g., ARIMA) to capture the linear components of 'ruggedAtTick' and then training a machine learning model (e.g., LSTM or Gradient Boosting) on the residuals of the statistical model to capture the remaining non-linear patterns.<sup>26</sup> This allows for a more comprehensive capture of complex time series dynamics that may contain both linear and non-linear elements.<sup>26</sup>

#### 4.3.3. Hidden Markov Model (HMM) Exploration

Given the discrete and potentially event-driven nature of 'ruggedAtTick', the application of Hidden Markov Models will be explored. This involves defining potential hidden states that could influence 'ruggedAtTick' values (e.g., different app operational modes, user engagement levels). The HMM would then model the transition probabilities between these hidden states and the emission probabilities of observing specific 'ruggedAtTick' values given a hidden state.<sup>27</sup> This approach offers a unique advantage by providing not just predictions but also inferences about the underlying system dynamics, which can be invaluable for understanding the causal factors behind 'ruggedness'.

#### 4.3.4. Comprehensive Comparative Analysis and Selection

All developed models (individual, ensemble, hybrid, and HMM) will undergo a



comprehensive comparative analysis against the established Bayesian indicator baseline and against each other. This will involve:

- **Performance Metrics:** Detailed comparison across all established metrics (RMSE, MAE, MAPE, MASE) to quantify predictive accuracy.
- **Robustness Analysis:** Evaluating model stability under varying data conditions, including periods of high volatility or sparse data.
- **Computational Efficiency:** Assessing the training and inference time for each model, a critical factor for real-time application deployment.
- **Interpretability:** Evaluating the degree to which each model's predictions and underlying mechanisms can be understood and explained, which is important for actionable insights.
- **Scalability:** Considering how well each model can scale with increasing data volume and complexity.

The selection of the final recommended methodology will be based on a multi-criteria assessment, balancing predictive performance, robustness, computational feasibility, and the ability to provide actionable insights into the 'ruggedAtTick' metric.

## Conclusions

The comprehensive research plan outlined herein provides a systematic framework for advancing the prediction of 'ruggedAtTick' values beyond the existing Bayesian indicator. The analysis of 'ruggedAtTick' data reveals its critical nature as a sequential, potentially irregular, and event-driven metric within the application. This inherent irregularity necessitates either careful data regularization or sophisticated feature engineering to transform the raw time series into a format suitable for advanced models. The potential for non-linear patterns in 'ruggedAtTick' dynamics further underscores the need to move beyond purely linear statistical models.

The exploration of advanced methodologies highlights a diverse set of powerful tools. Traditional statistical models like ARIMA and Exponential Smoothing offer strong baselines for linear patterns, but their direct application to irregular 'ruggedAtTick' data would require careful preprocessing that could potentially smooth out important granular information. Machine learning regression models, particularly tree-based ensembles (Random Forest, Gradient Boosting) and deep learning approaches (LSTMs), are well-suited to capture the complex, non-linear relationships implied by

the "rugged" nature of the data. Their effectiveness, however, is significantly contingent upon meticulous feature engineering, which serves as a crucial bridge between raw time series and model input.

Furthermore, advanced paradigms such as ensemble learning and hybrid models offer pathways to enhance robustness and leverage the complementary strengths of different model types. Ensemble methods can mitigate the limitations of single models, while hybrid approaches can effectively combine the interpretability of statistical models with the non-linear learning capabilities of machine learning models. Uniquely, Hidden Markov Models present a compelling opportunity to not only predict 'ruggedAtTick' values but also to infer the underlying, unobservable system states that may cause the observed patterns. This capability offers a deeper, more mechanistic understanding of app dynamics, moving beyond mere correlation to potential causal inference.

The proposed phased research plan, from detailed data preparation and feature engineering to rigorous walk-forward validation and multi-criteria comparative analysis, is designed to identify the most effective and robust methodology. The ultimate recommendation will balance predictive accuracy, computational feasibility, and the ability to provide actionable insights, ensuring that the chosen approach significantly enhances the app's capacity for proactive management and optimization based on 'ruggedAtTick' predictions.

## Works cited

1. Discrete-event simulation - Wikipedia, accessed July 17, 2025, [https://en.wikipedia.org/wiki/Discrete-event\\_simulation](https://en.wikipedia.org/wiki/Discrete-event_simulation)
2. Irregular time series - IBM, accessed July 17, 2025, <https://www.ibm.com/docs/en/informix-servers/12.10.0?topic=concepts-irregular-time-series>
3. Time series models with irregular time intervals : r/quant - Reddit, accessed July 17, 2025, [https://www.reddit.com/r/quant/comments/1ftpm3t/time\\_series\\_models\\_with\\_irregular\\_time\\_intervals/](https://www.reddit.com/r/quant/comments/1ftpm3t/time_series_models_with_irregular_time_intervals/)
4. ARIMA: A Model to Predict Time Series Data | Towards Data Science, accessed July 17, 2025, <https://towardsdatascience.com/arima-a-model-to-predict-time-series-data-a34c7638310b/>
5. Machine Learning Models for Time Series Analysis: A Comprehensive Guide - Medium, accessed July 17, 2025, <https://medium.com/@aleksej.gudkov/machine-learning-models-for-time-series-analysis-a-comprehensive-guide-e0e4d126debb>

6. An Introduction to Stationarity and Non-stationarity in Econometrics, accessed July 17, 2025,  
<https://www.econometricstutor.co.uk/time-series-analysis-stationarity-and-non-stationarity>
7. 8.1 Stationarity and differencing | Forecasting: Principles and Practice (2nd ed) - OTexts, accessed July 17, 2025, <https://otexts.com/fpp2/stationarity.html>
8. Time series - Wikipedia, accessed July 17, 2025,  
[https://en.wikipedia.org/wiki/Time\\_series](https://en.wikipedia.org/wiki/Time_series)
9. Time Series Analysis using ARIMA - Medium, accessed July 17, 2025,  
<https://medium.com/@danushidk507/time-series-analysis-using-arima-d246a82c43a6>
10. Exponential Smoothing: A Beginner's Guide to Getting Started ..., accessed July 17, 2025,  
<https://www.influxdata.com/blog/exponential-smoothing-beginners-guide/>
11. An Introduction to Exponential Smoothing for Time Series Forecasting in Python, accessed July 17, 2025,  
<https://www.simplilearn.com/exponential-smoothing-for-time-series-forecasting-in-python-article>
12. Feature Engineering for Time Series Forecasting | Train in Data, accessed July 17, 2025, <https://www.trainindata.com/p/feature-engineering-for-forecasting>
13. Feature engineering for time-series data - Statsig, accessed July 17, 2025,  
<https://www.statsig.com/perspectives/feature-engineering-timeseries>
14. Random Forest for Time Series Forecasting ..., accessed July 17, 2025,  
<https://machinelearningmastery.com/random-forest-for-time-series-forecasting/>
15. Random Forest for Time Series Forecasting for Data Science - Analytics Vidhya, accessed July 17, 2025,  
<https://www.analyticsvidhya.com/blog/2021/06/random-forest-for-time-series-forecasting/>
16. AI and Machine Learning for Networks: time series forecasting and regression - Codilime, accessed July 17, 2025,  
<https://codilime.com/blog/ai-ml-for-networks-time-series-forecasting-regression/>
17. Bagging, Boosting, and Stacking in Machine Learning | Baeldung on Computer Science, accessed July 17, 2025,  
<https://www.baeldung.com/cs/bagging-boosting-stacking-ml-ensemble-models>
18. Bagging, Boosting and Stacking: Ensemble Learning in ML Models - Analytics Vidhya, accessed July 17, 2025,  
<https://www.analyticsvidhya.com/blog/2023/01/ensemble-learning-methods-bagging-boosting-and-stacking/>
19. XGBoost & LGBM for Time Series Forecasting: How to – 365 Data ..., accessed July 17, 2025,  
<https://365datascience.com/tutorials/python-tutorials/xgboost-lgbm/>
20. medium.com, accessed July 17, 2025,  
[https://medium.com/@kylejones\\_47003/boosting-stacking-and-bagging-for-ensemble-models-for-time-series-analysis-with-python-d74ab9026782#:~:text=Bag](https://medium.com/@kylejones_47003/boosting-stacking-and-bagging-for-ensemble-models-for-time-series-analysis-with-python-d74ab9026782#:~:text=Bag)

[ging%20helps%20reducing%20overfitting%2C%20boosting.performance%20of%20time%20series%20models.](#)

21. [D] How does xgboost work with time series? : r/MachineLearning - Reddit, accessed July 17, 2025,  
[https://www.reddit.com/r/MachineLearning/comments/1aoo7gc/d\\_how\\_does\\_xgb\\_oost\\_work\\_with\\_time\\_series/](https://www.reddit.com/r/MachineLearning/comments/1aoo7gc/d_how_does_xgb_oost_work_with_time_series/)
22. Time Series Forecasting Using Deep Learning ... - MathWorks, accessed July 17, 2025,  
<https://www.mathworks.com/help/deeplearning/ug/time-series-forecasting-using-deep-learning.html>
23. Time Series Stock Prediction with LSTM | by Gabe Nosek - Medium, accessed July 17, 2025,  
<https://gabenosek.medium.com/time-series-stock-prediction-with-lstm-eb04f2224c22>
24. (PDF) Ensembles for Time Series Forecasting - ResearchGate, accessed July 17, 2025,  
[https://www.researchgate.net/publication/272833422\\_Ensembles\\_for\\_Time\\_Series\\_Forecasting](https://www.researchgate.net/publication/272833422_Ensembles_for_Time_Series_Forecasting)
25. Hybrid models in time series analysis - Kaggle, accessed July 17, 2025,  
<https://www.kaggle.com/discussions/general/501602>
26. huytjuh/Hybrid-Time-Series-Modeling - GitHub, accessed July 17, 2025,  
<https://github.com/huytjuh/Hybrid-Time-Series-Modeling>
27. Hidden Markov Model in Machine learning | by Palak Bhandari ..., accessed July 17, 2025,  
<https://medium.com/@palakvb02/hidden-markov-model-in-machine-learning-a744ada52166>
28. Hidden Markov Model in Machine learning - GeeksforGeeks, accessed July 17, 2025,  
<https://www.geeksforgeeks.org/machine-learning/hidden-markov-model-in-machine-learning/>