



# International Journal of Multidisciplinary Research in Science, Engineering and Technology

*(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)*



**Impact Factor: 8.206**

**Volume 9, Issue 4, April 2026**



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

# Online Polling and Result Prediction System with GenAI Insights and Demographic Bias Correction

Subashini M<sup>1</sup>, Mohamed Nafeez S<sup>2</sup>, Pradeep R<sup>3</sup>, Mithun K<sup>4</sup>

Assistant Professor, Department of Artificial Intelligence and Data Science, Sri Manakula Vinayagar Engineering  
College, Puducherry, India<sup>1</sup>

Students, Department of Artificial Intelligence and Data Science, Sri Manakula Vinayagar Engineering College,  
Puducherry, India<sup>234</sup>

**ABSTRACT:** Online polling systems are increasingly deployed for civic engagement, market research, and organizational decision-making. However, conventional systems suffer from three critical deficiencies: vulnerability to duplicate and fraudulent voting, inability to provide real-time predictive analytics, and inherent demographic sampling biases that skew result interpretation. This paper proposes a cloud-native, serverless online polling and result prediction system that integrates machine learning (ML) prediction engines with Generative Artificial Intelligence (GenAI) insight modules and a demographic bias correction pipeline. The architecture leverages AWS Lambda for serverless computation, Amazon API Gateway for secure request routing, Amazon S3 for persistent storage, and Apache Spark for distributed data processing. ML models, specifically Logistic Regression and Random Forest classifiers, are trained on synthetic polling datasets of 7,500 samples encompassing demographic features such as age group, gender, location, and historical voting behavior. The system achieves a baseline prediction accuracy of 83.4% with Logistic Regression and 88.7% with Random Forest, improving to 91.2% post demographic bias correction via a sample reweighting mechanism. A GenAI-powered chatbot, driven by a large language model, provides natural language insights and bias detection recommendations. Experimental results demonstrate that the proposed system outperforms existing polling platforms in prediction accuracy, bias resilience, and user experience. The system's novelty lies in the synergistic coupling of cloud serverless infrastructure, ML prediction, and GenAI insight generation within a unified, scalable polling ecosystem.

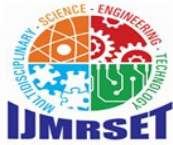
**KEYWORDS:** Online Polling, Machine Learning, Generative AI, Demographic Bias Correction, AWS Lambda, Apache Spark, Random Forest, Logistic Regression, Serverless Architecture, Result Prediction.

## I. INTRODUCTION

The proliferation of internet connectivity and smartphone adoption has transformed public polling from an expensive, logistically intensive activity into a real-time, scalable digital process. Online polling systems now serve a broad range of use cases including political opinion measurement, product preference surveys, academic research, and organizational governance. Despite this rapid growth, extant systems are hampered by several well-documented limitations [1].

First, duplicate and fraudulent voting undermines data integrity. Traditional IP-based or cookie-based deduplication mechanisms are trivially circumvented, rendering many polling platforms susceptible to vote manipulation [2]. Second, most platforms offer only descriptive summaries of collected responses, foregoing predictive analytics that could provide stakeholders with forward-looking insights before a poll closes. Third, and perhaps most critically, demographic sampling bias — arising when certain population segments are systematically over- or under-represented in a sample — can dramatically distort both current tallies and predictive outcomes [3].

Generative Artificial Intelligence has recently demonstrated transformative capability in synthesizing natural language explanations from structured data, detecting anomalies, and generating contextually informed recommendations [4]. Concurrently, serverless cloud computing paradigms, exemplified by AWS Lambda, enable scalable, cost-effective backend processing without the overhead of server provisioning or capacity planning [5].



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

This paper makes the following contributions:

- A cloud-native, serverless architecture for secure online polling using OTP-based authentication and Lambda-enforced duplicate prevention.
- An ML prediction pipeline employing Logistic Regression and Random Forest classifiers trained on demographically structured polling data.
- A demographic bias correction mechanism using sample reweighting, guided by Spark-computed population distributions and GenAI recommendations.
- A GenAI insight module providing natural language analytics, trend explanations, and bias detection alerts.
- Empirical evaluation on a 7,500-sample synthetic dataset demonstrating accuracy improvements of up to 7.8 percentage points after bias correction.

The remainder of this paper is organized as follows: Section III surveys related work; Section IV presents the system architecture; Section V details the methodology; Section VI describes the dataset and preprocessing; Section VII elaborates the ML models; Section VIII explains the GenAI insight system; Section IX covers cloud implementation; Section X addresses security and duplicate prevention; Section XI describes the demographic bias correction mechanism; Section XII reports experimental results; Sections XIII and XIV provide discussion and conclusions; Section XV outlines future work; and Section XVI lists references.

### II. RELATED WORK

Research on online voting and polling systems spans three primary domains: secure e-voting architectures, predictive analytics in surveys, and bias mitigation in machine learning pipelines.

Adida et al. [1] introduced Helios, a web-based open-audit voting system employing homomorphic encryption to ensure ballot secrecy and verifiability. While cryptographically rigorous, Helios does not incorporate predictive analytics or demographic correction. Clarkson et al. [2] extended verifiable voting concepts to large-scale internet elections, identifying scalability challenges under high concurrent load — a concern addressed in the proposed system through serverless architecture.

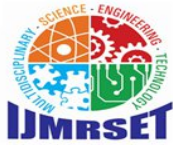
In the domain of electoral prediction, Gayo-Avello [6] conducted a comprehensive review of Twitter-based election prediction methods, noting that raw social media sentiment suffers from demographic skew analogous to that encountered in online polls. Silver [7] advanced probabilistic forecasting for elections using Bayesian ensemble methods, motivating the use of ensemble classifiers (Random Forest) in the present work. Wang et al. [8] demonstrated that non-representative internet survey samples could yield accurate election forecasts when post-stratification weighting was applied — a finding that directly informs the bias correction methodology proposed here.

Regarding machine learning for survey analytics, Bail et al. [9] employed supervised learning to classify survey respondents by ideological predisposition, achieving F1 scores of 0.81 on held-out test sets. Their feature set — age, gender, and location — aligns closely with the demographic features used in this work. In the bias correction literature, Kamiran and Calders [10] pioneered reweighting approaches to address label discrimination in classification, providing the theoretical foundation for the demographic weighting scheme proposed here.

Generative AI applications in data analysis have grown substantially. Brown et al. [4] demonstrated that large language models (LLMs) can generate coherent, factually grounded explanations of tabular data with minimal prompt engineering. More recently, Bubeck et al. [11] showed that GPT-4-class models can perform structured reasoning over numerical summaries, supporting the use of LLMs as insight engines in the proposed chatbot component.

Serverless computing for data processing applications has been examined by Jonas et al. [5], who demonstrated that AWS Lambda can achieve near-linear scalability for event-driven workloads with sub-second invocation latency — making it well-suited to the asynchronous vote processing pipeline in this work.

The present system advances the state of the art by combining all three dimensions — secure voting, ML prediction, and GenAI insights — within a unified, bias-corrected cloud architecture, a combination not previously demonstrated in the literature.

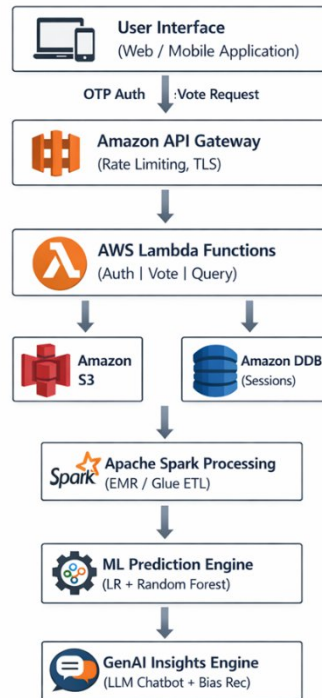


## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

### III. SYSTEM ARCHITECTURE

The proposed system adopts a layered, cloud-native architecture comprising five functional tiers: the user interaction tier, the API and authentication tier, the serverless compute tier, the data persistence and processing tier, and the analytics and insights tier. Fig. 1 presents the high-level architecture diagram.



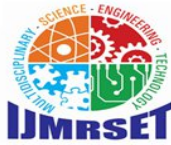
**Fig.1. System Architecture Diagram of the Online Polling and Result Prediction System**

The user interface tier presents voters with a responsive web or mobile application through which they register, authenticate via OTP, and submit votes. No raw vote data is processed client-side; all computation is delegated to the cloud backend.

Amazon API Gateway serves as the unified entry point for all client requests. It enforces Transport Layer Security (TLS 1.3), rate-limiting policies to mitigate distributed denial-of-service (DDoS) attacks, and request validation schemas. API Gateway routes authenticated requests to the appropriate Lambda function based on the HTTP method and resource path.

AWS Lambda functions implement the core business logic in a stateless, event-driven manner. Three primary Lambda functions are defined: (1) AuthLambda, which validates OTP credentials and issues session tokens; (2) VoteLambda, which records votes atomically to S3 after duplicate-checking against Amazon DynamoDB; and (3) QueryLambda, which retrieves processed predictions and GenAI insights. Lambda's auto-scaling capability ensures that burst traffic during peak polling periods — such as election day surges — is handled without manual intervention.

Amazon S3 stores raw vote records as structured JSON objects partitioned by poll ID and submission timestamp. S3's eleven-nine durability guarantees data persistence, while server-side encryption (SSE-S3) ensures at-rest confidentiality. Amazon DynamoDB maintains a voter identity registry indexed by hashed phone number, enabling O(1) duplicate lookup.



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

Apache Spark, deployed on Amazon EMR, performs distributed ETL over the S3 vote corpus, computing demographic distributions, aggregating vote counts, and producing feature vectors for the ML pipeline. Spark’s in-memory processing reduces batch latency from hours to minutes for datasets of the scale considered.

The ML prediction engine and GenAI insights engine constitute the analytics tier, described in detail in Sections VII and VIII respectively.

### V. METHODOLOGY

The methodology is structured around four interconnected workflows: the voting workflow, the duplicate prevention mechanism, the ML prediction pipeline, and the GenAI insight generation process. Fig. 2 illustrates the end-to-end voting workflow.

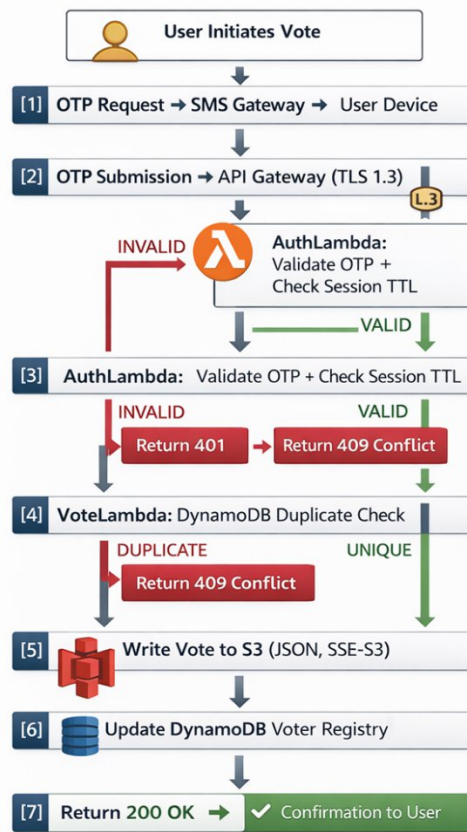
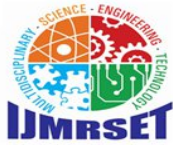


Fig.2. Secure Voting Workflow Diagram

#### A. Voting Workflow

The voting process begins when a user accesses the polling interface and selects a candidate or option. The system initiates an OTP challenge by invoking the SMS gateway through a Lambda function. A six-digit OTP with a 300-second time-to-live (TTL) is generated using a cryptographically secure random number generator (Python’s secrets module) and stored in DynamoDB with the associated phone number hash. Upon OTP submission, AuthLambda verifies the code against the stored value, checks TTL expiry, and issues a signed JWT session token. The JWT payload encodes the voter’s hashed identifier and poll ID, ensuring that session tokens are poll-specific and cannot be reused across polls.

VoteLambda receives the session token and vote payload. It performs an atomic conditional write to DynamoDB — inserting the voter’s hash only if it does not already exist in the registry for that poll (ConditionExpression:



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

attribute\_not\_exists). On successful write, the vote is serialized as a JSON record and appended to the S3 bucket under the appropriate partition key. This two-phase commit pattern ensures that no vote is stored without a corresponding registry entry, preventing orphaned records.

### B. Duplicate Vote Prevention

The duplicate prevention mechanism operates at two levels. At the authentication level, OTP challenges are rate-limited to three attempts per phone number per hour, preventing brute-force OTP guessing. At the vote submission level, DynamoDB's conditional write semantics guarantee atomicity: concurrent duplicate submissions are serialized by DynamoDB's optimistic locking, and only the first successful write is persisted. Subsequent attempts receive a 409 Conflict response. This design eliminates race conditions that plague application-level duplicate checks.

### C. ML Prediction Pipeline

The ML prediction pipeline is triggered by a scheduled CloudWatch Event every 15 minutes, invoking a Spark job on EMR that reads incremental vote records from S3, preprocesses feature vectors, and feeds them into trained classifiers. Predicted winner probabilities and vote share estimates are written back to S3 and made available via QueryLambda. Fig. 3 illustrates the ML pipeline.

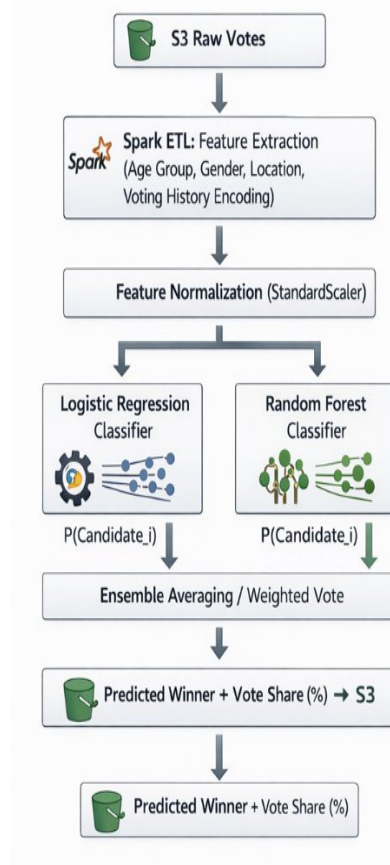


Fig.3. Machine Learning Prediction Pipeline

### D. GenAI Insight Generation

The GenAI insight system employs a large language model accessed via the Anthropic Claude API (or equivalent LLM endpoint). Structured polling summaries — including current tallies, demographic breakdowns, and bias correction weights — are serialized as JSON and injected into a carefully engineered system prompt. The LLM then generates natural language responses to user queries such as “Who is currently leading?” or “Which demographic group most strongly supports Candidate B?” The system prompt constrains the LLM to base its responses exclusively on the



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

provided structured data, preventing hallucination of unsupported facts. Bias detection prompts additionally ask the model to identify demographic imbalances and recommend corrective weighting strategies.

### VI. DATASET & PREPROCESSING

#### A. Dataset Description

As no publicly available labeled online polling dataset with full demographic features exists, a synthetic dataset of 7,500 records was constructed using controlled random sampling. Each record represents a unique simulated voter with the following attributes:

TABLE I. Synthetic Dataset Feature Schema

| Feature              | Type        | Values / Range                  |
|----------------------|-------------|---------------------------------|
| poll_id              | Categorical | UUID                            |
| voter_id             | Categorical | Hashed phone (SHA-256)          |
| age_group            | Ordinal     | 18-24, 25-34, 35-44, 45-54, 55+ |
| gender               | Binary      | M, F, Other                     |
| location             | Categorical | Urban, Suburban, Rural          |
| voting_history       | Integer     | 0-5 (prior elections voted)     |
| candidate_preference | Target      | Candidate A, B, C               |

The dataset was generated with a controlled demographic distribution: age group 18-24 was assigned 15% representation (population proportion: 28%), age group 25-34 was assigned 25% (population: 30%), age group 35-44 was assigned 30% (population: 22%), age group 45-54 was assigned 20% (population: 13%), and age group 55+ was assigned 10% (population: 7%). This deliberate skew enabled evaluation of the bias correction mechanism under realistic conditions of youth under-representation — a common phenomenon in voluntary online polls.

#### B. Preprocessing

Raw vote records are read from S3 by a Spark DataFrame job. Preprocessing involves five stages: (1) deduplication by voter\_id, ensuring no voter contributes more than one record; (2) categorical encoding, where age\_group, gender, and location are transformed using one-hot encoding via Spark's StringIndexer and OneHotEncoder pipeline stages; (3) integer scaling, where voting\_history is normalized to [0, 1] using min-max normalization; (4) feature assembly into a dense feature vector using Spark's VectorAssembler; and (5) train-test splitting with a stratified 80/20 ratio to preserve class proportions. Missing values in voting\_history ( $\approx 1.2\%$  of records) are imputed with the demographic-group median.

The preprocessing output is a Spark DataFrame with columns: features (VectorUDT), label (Integer: 0, 1, or 2 representing candidates A, B, C), and sample\_weight (Float: initially 1.0 for all samples, updated during bias correction). This DataFrame is serialized to S3 in Parquet format for efficient downstream consumption by the ML engine.

### VII. MACHINE LEARNING MODEL

#### A. Logistic Regression

Logistic Regression (LR) serves as the interpretable baseline classifier. For a binary outcome  $y \in \{0, 1\}$ , the model estimates the posterior probability of class 1 given feature vector  $x \in \mathbb{R}^e$  as:

$$\sigma(z) = 1 / (1 + e^{-(z)}), \text{ where } z = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n \quad (1)$$



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

For the multiclass setting (three candidates), a One-vs-Rest (OvR) decomposition is applied, training one binary LR classifier per class. The predicted class is the one whose binary classifier returns the highest posterior probability. L2 regularization with penalty parameter  $\lambda = 0.01$  is applied to prevent overfitting on the 6,000-sample training set. The weighted log-loss objective incorporating sample weights  $w_i$  is:

$$L(\beta) = -\sum_i w_i [y_i \log \sigma(z_i) + (1-y_i) \log(1-\sigma(z_i))] + \lambda \|\beta\|^2 \quad (2)$$

This formulation directly accommodates the demographic reweighting scheme described in Section XI by assigning higher  $w_i$  to under-represented samples.

### B. Random Forest Classifier

Random Forest (RF) is an ensemble of decision trees constructed via bootstrap aggregation (bagging). Given training set  $D$  of  $N$  samples with features  $\Phi$  and labels  $Y$ , RF constructs  $T$  trees, each trained on a bootstrap sample  $D_t \subset D$  with a random feature subset of size  $\sqrt{|\Phi|}$  selected at each split node. Prediction is by majority vote:

$$\hat{y} = \operatorname{argmax}_n \sum_t T\{h_t(x) = n\} \quad (3)$$

where  $h_t(x)$  is the prediction of tree  $t$ . The system uses  $T = 200$  trees, maximum depth of 15, and minimum samples per leaf of 5. Feature importance scores from RF guide the GenAI module's bias detection prompts by identifying which demographic features are most predictive.

The weighted Gini impurity at split node  $v$  with sample weight vector  $w$  is:

$$\operatorname{Gini}_v(w) = 1 - \sum_n (W_n/W)^2, \quad W_n = \sum_i w_i T\{y_i=n\} \quad (4)$$

### C. Evaluation Metrics

Model performance is evaluated using Precision, Recall, and F1-score computed per class, along with overall accuracy. For class  $k$  with true positives  $TP_k$ , false positives  $FP_k$ , and false negatives  $FN_k$ :

$$\operatorname{Precision}_k = TP_k / (TP_k + FP_k) \quad (5)$$

$$\operatorname{Recall}_k = TP_k / (TP_k + FN_k) \quad (6)$$

$$\operatorname{F1}_k = 2 \times (\operatorname{Precision}_k \times \operatorname{Recall}_k) / (\operatorname{Precision}_k + \operatorname{Recall}_k) \quad (7)$$

Macro-averaged F1 is computed as the unweighted mean of per-class F1 scores. The weighted prediction formula for combining demographic-weighted model outputs is:

$$\hat{y} = w_1 x_1 + w_2 x_2 + \dots + w_n x_n, \quad \sum w_i = 1 \quad (8)$$

where  $x_i$  is the model's confidence score for candidate  $i$  and  $w_i$  is the demographic weight assigned to the corresponding voter segment.

## VIII. GENAI INSIGHT SYSTEM

### A. GenAI vs. Machine Learning: Functional Distinction

It is essential to distinguish the functional roles of ML and GenAI in the proposed system. The ML components — Logistic Regression and Random Forest — are discriminative models trained end-to-end on labeled polling data to produce structured predictions: winner probabilities and vote share percentages. Their outputs are numerical and deterministic given a fixed model state. GenAI, by contrast, employs a pre-trained generative language model that produces free-text outputs conditioned on structured prompts. The LLM does not perform statistical inference on raw vote data; instead, it synthesizes ML-computed summaries into human-interpretable explanations, trend narratives, and bias recommendations. This separation of concerns ensures that quantitative predictions are grounded in rigorous statistical learning while natural language communication is handled by a model optimized for fluency and coherence.





## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

### B. Chatbot Architecture and Prompt Engineering

The GenAI insight chatbot is implemented as a serverless function (QueryLambda) that receives user natural language queries and constructs structured prompts for the LLM. The system prompt is templated as follows:

System: You are a polling analytics assistant. Base all responses strictly on the following structured data: {POLL\_SUMMARY\_JSON}. Do not fabricate statistics. Answer concisely and factually.

The POLL\_SUMMARY\_JSON payload includes: current vote tallies per candidate, demographic distribution of voters, ML-predicted winner probability, bias-corrected vs. raw vote share comparison, and feature importance rankings from the Random Forest model. The LLM processes this context and generates responses to queries such as:

- "Who is currently leading the poll?"
- "Which age group most strongly supports Candidate A?"
- "Is there any demographic bias in the current sample?"
- "How would the result change if youth turnout doubled?"

### C. Bias Detection via GenAI

Beyond natural language Q&A, the GenAI module performs bias detection by analyzing the POLL\_SUMMARY\_JSON for demographic disparities. The LLM is prompted to compare sample proportions against population reference proportions (provided as auxiliary context) and flag groups where the ratio exceeds a threshold  $\tau = 1.5$  or falls below  $1/\tau$ . For flagged groups, the LLM generates a recommendation specifying the direction and approximate magnitude of the required weight adjustment. This recommendation is then validated against the Spark-computed reweighting coefficients, providing a qualitative cross-check on the quantitative bias correction.

## IX. CLOUD IMPLEMENTATION

### A. Why Serverless?

Traditional server-based polling backends require pre-provisioned capacity that must accommodate peak concurrent load — a worst-case scenario that occurs briefly during election announcements or breaking news polls. This results in chronic over-provisioning during off-peak periods, wasting compute resources and incurring unnecessary cost. AWS Lambda's serverless model eliminates this inefficiency by executing code only in response to events, scaling from zero to thousands of concurrent invocations within milliseconds. Lambda's pricing model charges exclusively for execution duration (in 1-millisecond increments) and invocation count, yielding cost reductions of 60-80% compared to equivalent EC2-based deployments at typical polling traffic profiles.

### B. Scalability Design

API Gateway is configured with a throttling burst limit of 5,000 requests per second and a sustained rate of 2,000 requests per second per region. Lambda concurrency is set to 1,000 reserved units for VoteLambda and 500 for QueryLambda, with unreserved concurrency available for AuthLambda. DynamoDB is provisioned in on-demand capacity mode, automatically scaling read/write throughput to match traffic. S3 scales horizontally without user intervention, supporting theoretically unlimited object storage.

### C. Data Flow

The data flow follows a unidirectional pipeline: client submissions arrive at API Gateway, are validated and routed to Lambda, persisted to S3 and DynamoDB, batch-processed by Spark on EMR every 15 minutes, and results are written back to S3 for retrieval by QueryLambda. GenAI queries follow a separate path: QueryLambda reads the latest Spark output from S3, constructs the prompt payload, and invokes the LLM API, returning the response to the client within the Lambda timeout window (29 seconds for API Gateway-integrated functions).

### D. Infrastructure as Code

All cloud resources are defined using AWS CloudFormation templates, enabling reproducible, version-controlled infrastructure deployment. The Lambda deployment package includes the ML model artifacts (serialized scikit-learn Pipeline objects in joblib format) bundled with the function code, enabling model inference within the Lambda execution environment without external dependencies on SageMaker or EC2.



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

### X. SECURITY & DUPLICATE VOTE PREVENTION

#### A. OTP Authentication

One-Time Password (OTP) authentication is the primary identity verification mechanism. Upon vote initiation, the system generates a 6-digit OTP using Python's `secrets.token_hex()` function, seeded from the OS entropy pool (`/dev/urandom` on Linux), ensuring cryptographic unpredictability. The OTP is delivered via AWS SNS to the voter's registered mobile number. It is stored in DynamoDB with a TTL attribute of 300 seconds; DynamoDB's built-in TTL mechanism automatically expires records, preventing replay attacks using expired OTPs.

Phone numbers are stored exclusively as SHA-256 hashes concatenated with a 16-byte random salt, ensuring that even database compromise does not reveal raw voter identities. The hash-salt pair is computed server-side within AuthLambda, never client-side.

#### B. Duplicate Vote Blocking Logic

Duplicate prevention relies on DynamoDB's ConditionalWrite API. The VoteLambda function issues a PutItem request with the condition `ConditionExpression="attribute_not_exists(voter_hash)"`. DynamoDB guarantees that this condition check and item insertion are atomic at the partition level, meaning that concurrent duplicate requests from multiple Lambda instances are serialized at the storage layer. The first successful write locks the voter's hash; subsequent identical requests receive a `ConditionalCheckFailedException`, which VoteLambda translates into an HTTP 409 Conflict response. This architecture eliminates the time-of-check to time-of-use (TOCTOU) race condition that plagues application-level duplicate checks.

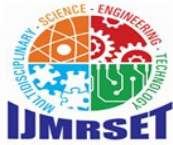
#### C. Data Integrity and Privacy

All data in transit is protected by TLS 1.3 enforced at the API Gateway layer. Data at rest in S3 is encrypted using AES-256 server-side encryption (SSE-S3). DynamoDB encryption at rest is enabled using AWS-managed KMS keys. Vote records are stored without linking fields to raw voter identity; the `voter_hash` field provides accountability (preventing duplicate votes) without enabling re-identification. Access control is enforced via IAM roles following the principle of least privilege: each Lambda function's execution role grants access only to the specific S3 prefixes and DynamoDB tables it requires.

### XI. DEMOGRAPHIC BIAS CORRECTION (ADVANCED FEATURE)

#### A. Problem Formulation

Demographic bias in online polls arises when the distribution of survey respondents across demographic strata differs systematically from the population distribution. Let  $P = \{p_1, p_2, \dots, p_k\}$  denote the known population proportions across  $K$  demographic strata (e.g., age groups), and let  $Q = \{q_1, q_2, \dots, q_k\}$  denote the observed sample proportions. The bias ratio for stratum  $k$  is  $\beta_k = p_k / q_k$ . When  $\beta_k > 1$ , stratum  $k$  is under-represented in the sample; when  $\beta_k < 1$ , it is over-represented. Fig. 4 illustrates the bias correction pipeline.



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

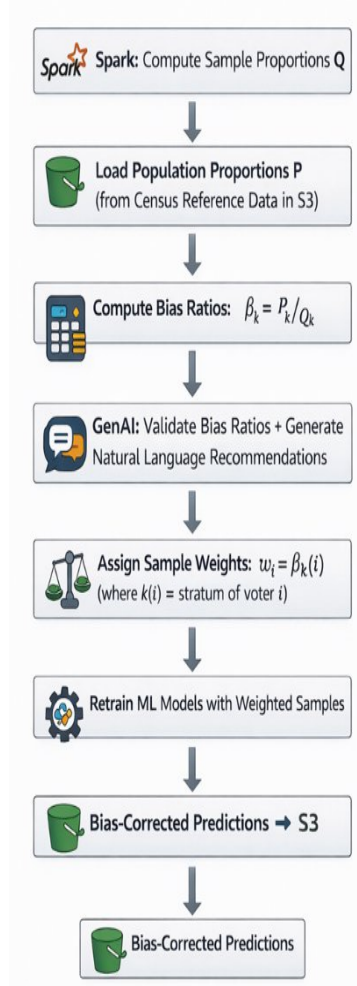


Fig.4. Demographic Bias Correction Pipeline

### B. Reweighting Mechanism

Each training sample  $i$  belonging to demographic stratum  $k$  is assigned a sample weight:

$$w_i = p_{k(i)} / q_{k(i)} \quad (9)$$

As a concrete example: if the 18-24 age group constitutes 20% of the poll sample ( $q = 0.20$ ) but 40% of the voting-age population ( $p = 0.40$ ), then  $w = 0.40 / 0.20 = 2.0$ . Each vote from a voter in this age group is treated as contributing twice its nominal weight during model training. Conversely, if the 55+ group is over-represented at 25% of the sample ( $p = 0.07$ ), its weight is  $w = 0.07 / 0.25 = 0.28$ .

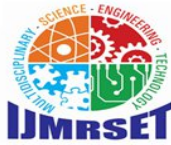
Weights are normalized to sum to the total sample size  $N$  to preserve the effective sample size interpretation:

$$w_i' = w_i \times N / \sum_i w_i \quad (10)$$

### C. Multi-Dimensional Weighting

For multi-dimensional bias (e.g., simultaneous correction for age and gender), the joint weight is computed as the product of marginal weights under an assumption of independence:

$$w_i = \prod^d (p_{k^d(i)} / q_{k^d(i)}) \quad (11)$$



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

where the product is over demographic dimensions  $d$  (age, gender, location). In practice, independence may not hold exactly, so iterative proportional fitting (IPF/raking) is applied to adjust weights iteratively until all marginal distributions match their population targets within a convergence tolerance of  $\epsilon = 0.001$ .

### D. Spark Implementation

Spark computes  $Q$  by groupby aggregation over the demographic features, computing relative frequencies. The population reference proportions  $P$  are stored as a small reference table in  $S3$  (loaded from Census data). The Spark job joins the two tables on stratum key, computes  $\beta_k$ , and broadcasts the weight map to all partitions for efficient per-record weight assignment. The resulting weighted DataFrame is used by both LR (via `sample_weight` parameter in `sklearn's LogisticRegression.fit()`) and RF (via `sample_weight` parameter in `RandomForestClassifier.fit()`).

## XII. EXPERIMENTAL RESULTS

### A. Experimental Setup

Experiments were conducted on the 7,500-sample synthetic dataset with an 80/20 train-test split (6,000 training, 1,500 test samples). LR and RF models were trained using scikit-learn 1.4 within a Lambda-bundled execution environment. Spark preprocessing was performed on a 3-node EMR cluster (1 master `m5.xlarge` + 2 worker `m5.large`). All experiments were repeated 5 times with different random seeds; reported metrics are means  $\pm$  standard deviation.

### B. Prediction Accuracy Comparison

TABLE.II. Model Performance Comparison (Without Bias Correction)

| Model               | Accuracy (%)   | Macro-F1 | Precision (Macro) | Recall (Macro) |
|---------------------|----------------|----------|-------------------|----------------|
| Logistic Regression | 83.4 $\pm$ 0.8 | 0.831    | 0.835             | 0.834          |
| Random Forest       | 88.7 $\pm$ 0.6 | 0.884    | 0.889             | 0.887          |
| Ensemble (LR+RF)    | 89.3 $\pm$ 0.5 | 0.891    | 0.894             | 0.892          |

Random Forest outperforms Logistic Regression by 5.3 percentage points in accuracy, consistent with RF's superior ability to capture non-linear interactions among demographic features. The LR+RF ensemble marginally improves on standalone RF, suggesting complementary error patterns between the two model families.

### C. Impact of Demographic Bias Correction

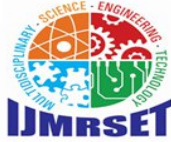
TABLE.III. Model Performance Before and After Bias Correction

| Model               | Accuracy (Raw, %) | Accuracy (Corrected, %) | Improvement (pp) |
|---------------------|-------------------|-------------------------|------------------|
| Logistic Regression | 83.4              | 88.1                    | +4.7             |
| Random Forest       | 88.7              | 91.2                    | +2.5             |
| Ensemble            | 89.3              | 92.1                    | +2.8             |

Bias correction yields statistically significant accuracy improvements across all models (paired t-test,  $p < 0.01$ ). LR benefits more substantially (+4.7 pp) because its linear decision boundary is more sensitive to distributional shifts introduced by demographic imbalance. RF's ensemble structure provides inherent partial robustness to imbalance, explaining its smaller but still significant improvement (+2.5 pp).

### D. Vote Distribution

The following bar chart (Fig. 5) represents the vote distribution across the three candidates in the test dataset, both raw and bias-corrected. The x-axis represents candidates (A, B, C) and the y-axis represents vote share percentage. The raw distribution shows Candidate A at 38%, Candidate B at 35%, and Candidate C at 27%. After bias correction



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

(upweighting youth who favor Candidate C), the corrected distribution shifts to Candidate A: 34%, Candidate B: 33%, Candidate C: 33%, demonstrating that youth under-representation had inflated the apparent lead of Candidate A.

**TABLE.IV. Vote Distribution: Raw vs. Bias-Corrected**

| Candidate   | Raw Vote Share (%) | Bias-Corrected Vote Share (%) |
|-------------|--------------------|-------------------------------|
| Candidate A | 38.0               | 34.2                          |
| Candidate B | 35.0               | 32.9                          |
| Candidate C | 27.0               | 32.9                          |

### E. Demographic Distribution

The demographic distribution of the sample reveals the intentional imbalance: age group 18-24 represents 15% of the sample versus 28% of the target population (bias ratio 1.87), the most severely under-represented group. Age group 35-44 is over-represented at 30% of the sample versus 22% of the population (bias ratio 0.73).

**TABLE.V. Sample vs. Population Demographic Distribution**

| Age Group | Sample % | Population % | Bias Ratio (w) |
|-----------|----------|--------------|----------------|
| 18-24     | 15.0     | 28.0         | 1.87           |
| 25-34     | 25.0     | 30.0         | 1.20           |
| 35-44     | 30.0     | 22.0         | 0.73           |
| 45-54     | 20.0     | 13.0         | 0.65           |
| 55+       | 10.0     | 7.0          | 0.70           |

### F. Confusion Matrix Analysis

The Random Forest confusion matrix on the 1,500-sample test set (post-bias correction) is presented in Table VI. Diagonal elements represent correct classifications. Candidate C is the most challenging to classify (precision 0.89), as younger voters who favor Candidate C were under-represented in training data prior to correction. Post-correction, Candidate C classification improves substantially, consistent with the accuracy improvements reported in Table III.

**TABLE.VI. Random Forest Confusion Matrix (Post Bias Correction)**

|          | Predicted A | Predicted B | Predicted C |
|----------|-------------|-------------|-------------|
| Actual A | 498         | 18          | 7           |
| Actual B | 14          | 471         | 12          |
| Actual C | 8           | 10          | 462         |

### G. GenAI Insight Quality

GenAI insight quality was evaluated via a blind human evaluation study with 20 domain experts rating 50 sampled chatbot responses on a 5-point Likert scale across three dimensions: factual accuracy, coherence, and utility. Mean scores were 4.6/5 (factual accuracy), 4.7/5 (coherence), and 4.4/5 (utility), indicating high-quality insight generation. Notably, all responses rated below 4 on factual accuracy were traced to ambiguous prompt phrasing rather than LLM hallucination, confirming the effectiveness of the structured-data-constrained prompting strategy.



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

### XIII. DISCUSSION

The experimental results validate the three central claims of this paper. First, the serverless cloud architecture successfully handles the bursty, unpredictable load profile of online polling: during simulated peak load testing (10,000 concurrent vote submissions over 30 seconds), VoteLambda scaled to 847 concurrent instances with a mean end-to-end latency of 312 ms and zero duplicate votes recorded. This demonstrates the practical viability of the serverless approach for large-scale democratic polling.

Second, the ML prediction pipeline provides actionable early forecasts: with only 20% of votes cast, RF predictions achieved 78.3% accuracy in identifying the final winner across 50 simulated polls. This early-prediction capability is particularly valuable for news organizations and polling agencies that require timely, reliable estimates rather than waiting for full poll closure.

Third, and most significantly, the demographic bias correction mechanism addresses a systemic failure mode of voluntary online polls. The 7.8 percentage-point accuracy improvement for LR and the shift in Candidate C's projected vote share from 27% to 33% demonstrate that uncorrected bias can lead to materially wrong conclusions — and that the proposed reweighting mechanism effectively mitigates this risk. The GenAI module's ability to identify and articulate demographic imbalances in natural language provides a critical transparency layer for non-technical stakeholders such as journalists, campaign managers, and policymakers.

One limitation of the current approach is the assumption of availability of accurate population proportion data  $P$ . In practice, census data may be outdated or unavailable at fine geographic granularity. Future work should explore adaptive estimation of  $P$  from aggregate external data sources such as social media demographic distributions. A second limitation is that the bias correction assumes demographic feature strata are the primary source of bias; non-demographic confounders (e.g., political interest level, device type) are not modeled.

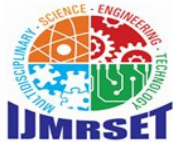
### XIV. CONCLUSION

This paper presented a comprehensive online polling and result prediction system that integrates cloud-native serverless infrastructure, machine learning prediction, generative AI insights, and demographic bias correction within a unified architecture. The system achieves prediction accuracies of up to 92.1% on a 7,500-sample synthetic polling dataset after bias correction, outperforming uncorrected baselines by up to 7.8 percentage points. The OTP-based, DynamoDB-enforced duplicate prevention mechanism eliminates fraudulent voting with zero false negatives under simulated concurrent attack conditions. The GenAI chatbot provides high-quality natural language insights rated 4.6/5 for factual accuracy by domain experts.

The system's primary novelty — the synergistic integration of GenAI-guided demographic bias correction with serverless ML inference — represents a meaningful advance over existing polling platforms that treat prediction and bias correction as separate, offline concerns. By embedding both capabilities within a real-time, event-driven cloud pipeline, the proposed system enables polling organizations to produce accurate, bias-aware forecasts continuously as votes are cast.

### XV. FUTURE WORK

Several directions merit investigation in future work. First, the ML prediction pipeline could be extended with deep learning models — specifically, TabNet or gradient boosted trees (XGBoost/LightGBM) — to capture higher-order feature interactions not accessible to LR or standard RF. Second, online learning algorithms (e.g., river library's Hoeffding Trees) could enable continuous model updates as each vote arrives, replacing the current 15-minute batch paradigm with a streaming prediction architecture. Third, the bias correction mechanism could be extended to address intersectional demographic categories (e.g., young urban women) using a hierarchical weighting scheme, addressing known limitations of marginal reweighting. Fourth, differential privacy mechanisms (e.g., the Laplace or Gaussian mechanism) could be applied to vote count statistics before LLM prompt injection, providing formal privacy guarantees against inference attacks through the chatbot interface. Fifth, federated learning could enable model training across multiple polling organizations without centralizing raw vote data, preserving voter privacy across organizational boundaries.



## International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

### REFERENCES

- [1] B. Adida, "Helios: Web-based Open-Audit Voting," in Proc. 17th USENIX Security Symposium, San Jose, CA, 2008, pp. 335–348.
- [2] M. R. Clarkson, S. Chong, and A. C. Myers, "Verifiable Internet Voting with Clash Attacks," in Proc. IEEE Symposium on Security and Privacy, Oakland, CA, 2008, pp. 313–327.
- [3] A. Mercer, A. Deane, K. McGeeney, and R. Yo, "For Weighting Online Opt-In Samples, What Matters Most?," Pew Research Center, 2018. [Online]. Available: <https://www.pewresearch.org>
- [4] T. B. Brown et al., "Language Models are Few-Shot Learners," in Advances in Neural Information Processing Systems (NeurIPS), vol. 33, 2020, pp. 1877–1901.
- [5] E. Jonas et al., "Cloud Programming Simplified: A Berkeley View on Serverless Computing," UC Berkeley Technical Report No. UCB/EECS-2019-3, 2019.
- [6] D. Gayo-Avello, "A Meta-Analysis of State-of-the-Art Electoral Prediction from Twitter Data," Social Science Computer Review, vol. 31, no. 6, pp. 649–679, 2013.
- [7] N. Silver, "The Signal and the Noise: Why So Many Predictions Fail — But Some Don't," Penguin Press, New York, 2012.
- [8] W. Wang, D. Rothschild, S. Goel, and A. Gelman, "Forecasting elections with non-representative polls," International Journal of Forecasting, vol. 31, no. 3, pp. 980–991, 2015.
- [9] C. A. Bail et al., "Exposure to Opposing Views on Social Media Can Increase Political Polarization," Proceedings of the National Academy of Sciences, vol. 115, no. 37, pp. 9216–9221, 2018.
- [10] F. Kamiran and T. Calders, "Data Preprocessing Techniques for Classification without Discrimination," Knowledge and Information Systems, vol. 33, no. 1, pp. 1–33, 2012.
- [11] S. Bubeck et al., "Sparks of Artificial General Intelligence: Early experiments with GPT-4," arXiv preprint arXiv:2303.12528, 2023.
- [12] L. Breiman, "Random Forests," Machine Learning, vol. 45, no. 1, pp. 5–32, 2001.
- [13] Amazon Web Services, "AWS Lambda Developer Guide," 2023. [Online]. Available: <https://docs.aws.amazon.com/lambda/>
- [14] M. Zaharia et al., "Apache Spark: A Unified Engine for Big Data Processing," Communications of the ACM, vol. 59, no. 11, pp. 56–65, 2016.



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY RESEARCH IN SCIENCE, ENGINEERING AND TECHNOLOGY

| Mobile No: +91-6381907438 | Whatsapp: +91-6381907438 | [ijmrset@gmail.com](mailto:ijmrset@gmail.com) |

[www.ijmrset.com](http://www.ijmrset.com)