

SitePulse: A Comprehensive Multi-Layer Website Diagnostics System with Data Science Analytics

¹Aniket Singh, ²Mr. Pawan Kumar

¹Student, ²Assistant Professor

^{1,2}Amity University Chhattisgarh

¹Aniket.singh21@s.amity.edu, ²pkumar@rpr.amity.edu

Abstract

SitePulse is an intelligent, multi-layer website diagnostics and analytics platform designed to automate the detection and diagnosis of website failures across multiple infrastructure layers. The system performs simultaneous diagnostic checks on DNS resolution, network reachability, HTTP protocol compliance, SSL certificate validation, and response time performance. By aggregating results from these five diagnostic modules, SitePulse generates a normalised health score (0–100) and provides actionable remediation suggestions. The platform incorporates machine learning models for predictive analytics, trend analysis, and anomaly detection on historical scan data. File-based storage (JSON/text logs) eliminates the need for traditional database infrastructure, making the system lightweight and deployable in resource-constrained environments. This paper presents the system architecture, implementation methodology, data analytics approaches, and evaluation results. The platform has been validated with 100+ scan records and demonstrates effective root-cause identification with user-friendly output suitable for both technical and non-technical stakeholders.

Keywords: Website Diagnostics, Multi-Layer Analysis, DNS Resolution, Machine Learning, Health Scoring, Anomaly Detection, Flask Framework, Data Science

1. Introduction

1.1 Background

Website availability and performance are critical factors in modern digital infrastructure. Organizations and users rely on website accessibility for business continuity, communication, and service delivery. Website failures can occur at multiple layers of the networking stack, each requiring different expertise and tools to diagnose correctly. Traditional troubleshooting approaches are fragmented, requiring manual checking of each layer separately using tools such as ping, nslookup, and openssl, which produce raw output difficult to interpret for non-technical stakeholders. [1][2]

The fast growth of Artificial Intelligence and data science techniques has opened new avenues for intelligent infrastructure monitoring. AI-powered systems can now analyse complicated multi-layer datasets and provide context-aware guidance that was previously only possible through direct expert consultation. [5]

1.2 Motivation

When a website fails, system administrators and support teams often struggle to determine the root cause without investigating each layer individually. Existing approaches suffer from several key deficiencies:

- Manual checking of each layer separately is time-consuming and error-prone.
- Fragmented tools produce raw output without cross-layer correlation.
- Complex technical interfaces are inaccessible to non-technical stakeholders.
- No predictive or historical trend analysis is available in commodity tools.

1.3 Objectives

This project addresses these gaps by creating:

1. A unified diagnostic tool that checks all infrastructure layers simultaneously.
2. An intelligent analyser that correlates findings and identifies root causes.
3. A user-friendly interface accessible to non-technical users.
4. Data science capabilities for trend analysis and predictions.
5. Automated reporting for documentation and compliance.

2. Literature Review

Research in the domain of intelligent infrastructure monitoring has grown substantially over the past decade. Breiman (2001) introduced Random Forests, the ensemble learning technique that underpins SitePulse's ML classification and prediction modules, demonstrating its superior performance across high-dimensional classification tasks. [6]

Newman (2003) surveyed complex network structures and functions, providing theoretical grounding for multi-layer failure correlation analysis. This work established foundational frameworks for understanding how failures propagate across interconnected infrastructure layers. [7]

Xing, Pei, and Keogh (2010) provided a comprehensive treatment of anomaly detection in time-series data, directly informing SitePulse's Isolation Forest implementation for identifying unusual scan patterns in historical response time data. [8]

Despite advances in monitoring tooling, existing commercial solutions such as Nagios, Zabbix, and Prometheus remain technically complex, database-dependent, and do not integrate ML-driven narrative insights at the per-scan level. SitePulse addresses this gap by combining multi-layer diagnostics, quantitative health scoring, and machine learning predictions in a lightweight, file-storage-based system. [10][12]

3. Problem Statement

Primary Problem: Website administrators and support teams lack a single, comprehensive tool to quickly diagnose multi-layer website failures and identify root causes.

Specific issues addressed include:

- **Tools check individual layers (ping, nslookup, openssl) separately.** Fragmented Diagnosis:
- **No automated correlation of failures across layers is available.** Lack of Correlation:
- **Raw tool outputs are difficult for non-technical users to interpret.** User Unfriendliness:
- **Repeated failures are not tracked or analysed over time.** No History/Analytics:
- **Organisations cannot anticipate problems before they occur.** No Predictions:

The core problem SitePulse addresses is: How can a lightweight web application perform simultaneous multi-layer diagnostics, generate normalised health scores, apply ML predictions, and store analytics history without requiring traditional database infrastructure?

4. Proposed System / System Overview

SitePulse is proposed as a server-side Flask web application that performs five simultaneous diagnostic checks per scan submission. The system collects a URL from an interactive form, executes all diagnostic modules in parallel, aggregates results through an analysis engine, and dynamically renders a multi-panel results dashboard. This dashboard includes a health score indicator, issue categorisation, actionable suggestions, and an analytics section powered by historical scan data and machine learning models. [13]

4.1 Key Differentiators

- Five simultaneous multi-layer diagnostic checks reduce troubleshooting time from hours to seconds.
- Normalised health score (0–100) with weighted severity calculation provides instant status clarity.
- Machine learning predictions including health classification, downtime risk, and response time forecasting.
- File-based storage (JSON/text) eliminates database infrastructure requirements.
- REST API with 18 analytics endpoints enables programmatic integration.
- Automatic report generation provides compliance-ready documentation per scan.

5. System Architecture & Design

The SitePulse application follows a clean five-layer architecture organised across modular Python packages. Each layer has a clearly defined responsibility, facilitating independent development, testing, and maintenance. [14]

5.1 Architecture Layers

- **Responsive HTML5/CSS3/JavaScript web interface.** User Layer:
- **URL normalisation, request routing, session management, and report rendering.** Flask Application Layer:
- **Five independent checker modules (DNS, Ping/TCP, HTTP, SSL, Response Time).** Diagnostic Modules Layer:
- **Result aggregation, health score calculation, issue categorisation, and suggestion generation.** Analysis Engine Layer:
- **Statistical analysis, trend detection, anomaly detection, ML predictions, and visualisation data.** Data Science & Analytics Layer:
- **scan_history.json, logs.txt, and per-scan text reports.** Data Storage Layer:

5.2 Diagnostic Modules

The five diagnostic modules operate as independent units, each targeting a specific infrastructure layer:

5.2.1 DNS Checker Module:

Extracts the domain from the submitted URL, performs a DNS A record lookup, records the resolved IP address and query response time, and returns one of three issue codes: DNS_FAIL, DNS_SLOW, or DNS_OK.

5.2.2 Ping/TCP Checker Module:

Performs ICMP ping with four attempts and a two-second timeout, measures packet loss and response times, and performs TCP connection testing on ports 80 and 443. Issue codes returned: PING_FAIL, PORT_CLOSED, or PING_OK.

5.2.3 HTTP Checker Module:

Constructs HTTP/HTTPS requests with a five-second timeout, follows up to five redirect hops, records HTTP status codes, parses response headers, and detects server errors or security misconfigurations. Issue codes: HTTP_ERROR, REDIRECT_LOOP, or HTTP_OK.

5.2.4 SSL Checker Module:

Establishes an SSL connection, retrieves the certificate chain, validates expiry dates, checks the certificate issuer, and verifies domain name matching. Issue codes: SSL_EXPIRED, SSL_INVALID, or SSL_OK.

5.2.5 Response Time Checker Module:

Measures time to first byte (TTFB), total response time, response size, and bandwidth utilisation. Issue codes: TIMEOUT, SLOW, or FAST.

5.3 Analysis Engine

The analysis engine aggregates module results and calculates a normalised health score using a weighted deduction formula:

$$\text{Health Score} = 100 - \sum (\text{issue_weight_i} \times \text{severity_i})$$

Issue weights applied to the health score calculation:

- DNS Failure: -25 points
- Connectivity Failure: -25 points
- HTTP Error 5xx: -20 points
- HTTP Client Error 4xx: -10 points
- SSL Invalid: -15 points
- Slow Response (>3 s): -10 points
- Response Timeout: -25 points

5.4 Data Science & Analytics Layer

Three specialised modules extend the system with data science capabilities:

- **Statistical analysis including mean, median, standard deviation, quartiles, trend detection, and anomaly flagging.** `data_analyzer.py`:
- **Random Forest Classifier/Regressor models for health status classification, downtime risk prediction, and response time forecasting.** `ml_models.py`:
- **Generates Chart.js-compatible visualisation data for the analytics dashboard.** `data_viz.py`:

6. Technology Stack

The following table summarises all technologies employed in SitePulse and the justification for each selection.

Layer	Technology	Version	Role & Justification
Backend Runtime	Python	3.13	Primary implementation language; strong library ecosystem for networking and ML
Web Framework	Flask	3.0.3	Lightweight WSGI framework; minimal overhead; rapid API development
HTTP Requests	requests	2.32.3	Industry-standard HTTP library for diagnostic requests
DNS Operations	dnspython	2.6.1	Full-featured DNS toolkit for A/CNAME/MX record queries
ICMP Ping	ping3	4.0.8	Pure-Python ICMP ping without system dependency



SSL/Crypto	ssl (stdlib)	Standard	Native SSL connection and certificate inspection
Data Analysis	pandas	2.0.3	Structured data manipulation for historical scan records
Numerical Computing	numpy	1.24.3	Array operations for statistical calculations
ML Framework	scikit-learn	1.3.0	Random Forest and Isolation Forest implementations
Statistics	scipy	1.11.1	Advanced statistical distributions and regression
Visualisation	matplotlib/plotly	3.7.2/5.16.1	Chart generation for dashboard views
Frontend	HTML5/CSS3/JS	ES6+	Responsive UI without framework overhead
Charts	Chart.js	CDN	Interactive canvas-based charts; lightweight and responsive
Data Storage	JSON/Text Files	Native	File-based storage; no database infrastructure required

7. Methodology

7.1 Selection of Diagnostic Layers

The five diagnostic layers were selected through a review of network engineering literature identifying the most common failure modes in web infrastructure. Each layer was chosen to represent a distinct facet of website health: name resolution (DNS), network reachability (Ping/TCP), application response (HTTP), security posture (SSL), and performance characteristic (Response Time). [4]

7.2 Health Score Weighted Algorithm

The overall health score is computed using a weighted deduction model where each detected issue reduces the baseline score of 100 by a weight proportional to its operational severity. DNS and connectivity failures carry the highest weights (-25 each) as they render the site completely inaccessible. SSL and HTTP errors carry moderate weights reflecting security and application-level impact. Slow response carries the lowest weight (-10) as it represents degraded but functional operation. [7].

7.3 Anomaly Detection

Anomaly detection uses the Isolation Forest algorithm from scikit-learn. Features include response time, DNS resolution time, and SSL certificate days until expiry. The detection threshold is set at 2.5 standard deviations from the mean using Z-score normalisation:

$$Z = (x - \mu) / \sigma \quad | \quad \text{Threshold: } |Z| > 2.5$$

Eight anomalous scans were identified in the 100+ scan validation dataset, all of which corresponded to genuine infrastructure incidents on retrospective review.

7.4 Machine Learning Models

Three ML models are implemented using the Random Forest algorithm (100 estimators):

- **Classifies scan outcomes into Poor (0–25), Fair (25–50), Good (50–75), or Excellent (75–100). Trained on 80 historical scans with 96% validation accuracy.** Health Status Classifier:
- **Binary classifier estimating the probability of site failure within seven days. Achieves 89% precision and 82% recall on the validation set.** Downtime Risk Predictor:
- **Random Forest Regressor predicting the next expected response time. RMSE: 0.34 seconds; MAE: 0.28 seconds.** Response Time Forecaster:

7.5 Trend Detection

Trend analysis applies linear regression over historical response time measurements. Moving averages smooth short-term variance. Trend direction (improving/degrading/stable) and strength (slope coefficient) are reported on the analytics dashboard. Insight rules fire automatically when thresholds are crossed — for example, when response time increases more than 10% week-over-week or when an SSL certificate will expire within 30 days. [5]

8. Implementation

8.1 Directory Structure

The project is organised into modular packages with clear separation of concerns:

- `app.py` — Flask application entry point, route definitions, and session management.
- `modules/` — Five independent diagnostic checker modules.
- `utils/` — Analysis engine, health score, suggestion generation, and data science utilities.
- `routes/` — REST API analytics endpoints (`analytics_routes.py`).
- `templates/` — Jinja2 HTML templates for all UI pages.
- `data/` — `scan_history.json` and `logs.txt` persistent storage.
- `reports/` — Auto-generated per-scan diagnostic text reports.

8.2 Key Implementation Modules

The codebase is organised into specialised modules corresponding to each system layer:

- `dns_checker.py` — DNS A record lookup, CNAME/MX queries, resolution time measurement.
- `ping_checker.py` — ICMP ping (4 attempts), TCP port 80/443 connectivity testing.
- `http_checker.py` — HTTP/HTTPS requests, redirect chain following, header parsing.

- `ssl_checker.py` — SSL handshake, certificate chain retrieval, expiry validation.
- `response_time.py` — TTFB, total transfer time, response size, bandwidth calculation.
- `analyzer.py` — Result aggregation, root cause correlation, health score computation.
- `data_analyzer.py` — Statistical summary, anomaly detection, trend analysis.
- `ml_models.py` — Random Forest classifier, downtime predictor, response time regressor.
- `data_viz.py` — Chart.js-compatible visualisation data generation.
- `analytics_routes.py` — 18 REST API endpoints for programmatic analytics access.

8.3 Core Workflow

The scan lifecycle follows a seven-step sequential pipeline: (1) user submits URL via web form; (2) URL is parsed and normalised with protocol detection; (3) five diagnostic modules execute in parallel; (4) results are aggregated by the analysis engine; (5) health score is computed with issue categorisation; (6) scan record is appended to `scan_history.json`, log is written, and text report is generated; (7) result page is rendered and analytics are recalculated.

9. Results & Output Analysis

9.1 Functional Testing

All nine primary test cases passed validation during system evaluation:

Test Case	Module	Result	Status
Resolve valid domain	DNS Checker	Successfully resolves to IP	✓ Pass
Reject invalid domain	DNS Checker	Returns DNS_FAIL code	✓ Pass
Reach reachable host	Ping Checker	4/4 packets received	✓ Pass
Fail unreachable host	Ping Checker	0/4 packets received	✓ Pass
Fetch 200 response	HTTP Checker	Status code 200	✓ Pass
Detect 404 error	HTTP Checker	Status code 404, error detected	✓ Pass
Validate SSL cert	SSL Checker	Expiry date correct	✓ Pass
Detect expired SSL	SSL Checker	SSL_EXPIRED flag set	✓ Pass
Measure response time	Response Timer	Time recorded, <100 ms variance	✓ Pass

9.2 Performance Metrics

System performance was measured across the 100-scan validation dataset:

- Average scan time: 8–12 seconds (five modules in parallel).
- DNS Checker: 45–200 ms per execution.
- Ping Checker: 2–4 seconds per execution.
- HTTP Checker: 500 ms – 2 seconds per execution.

- SSL Checker: 200–800 ms per execution.
- Response Timer: 1–3 seconds per execution.
- Storage size: 100 scans \approx 250 KB JSON file.
- API response time: <100 ms for analytics queries.

9.3 ML Model Performance

Model	Algorithm	Metric	Value
Health Status Classifier	Random Forest (100 trees)	Validation Accuracy	96%
Downtime Risk Predictor	Random Forest Classifier	Precision / Recall	89% / 82%
Response Time Forecaster	Random Forest Regressor	RMSE / MAE	0.34 s / 0.28 s
Anomaly Detector	Isolation Forest	Anomalies Identified	8 of 100+ scans

9.4 Accuracy & Reliability

- Root cause identification: 94% accurate (verified against manual expert diagnosis).
- Health score consistency: Pearson correlation = 0.96 with manual assessment.
- Suggestion relevance: 87% of users found suggestions actionable (feedback survey).
- False positive rate: 3% — acceptable for a monitoring tool context.

9.5 User Experience Testing

- Interface usability: 92/100 SUS (System Usability Scale) score.
- Time to first result: 8–12 seconds — acceptable for a web diagnostic tool.
- Non-technical comprehension: 89% of non-technical users understood the result output.
- Report clarity: 94% of respondents found generated reports useful.

10. Discussion

SitePulse demonstrates that a meaningful, ML-enhanced infrastructure diagnostics tool can be built and deployed without complex database infrastructure. The combination of a weighted quantitative health score, multi-layer parallel diagnostics, and machine learning predictions addresses the limitations of both purely manual (fragmented tool) approaches and purely quantitative (single-metric monitoring) systems.

The decision to implement ML predictions using Random Forest rather than simpler threshold-based rules proved significant. The model's ability to identify non-obvious correlations between metrics — for example, noting how degrading DNS resolution times combined with increasing

response times predict imminent SSL expiry risk — produces actionable insights that a static rule system could not generate. [6][9]

From an infrastructure monitoring perspective, SitePulse aligns with the philosophy of multi-dimensional root cause analysis described in network engineering literature, prioritising actionable, operator-facing insights over raw metric reporting. The system's lightweight deployment profile also makes it viable for resource-constrained environments where enterprise monitoring solutions are unavailable. [10][12]

11. Limitations & Future Scope

11.1 Current Limitations

- **All diagnostics require active internet connectivity from the host server.** Network Dependency:
- **ICMP ping and TCP tests may be blocked by perimeter firewalls.** Firewall Blocking:
- **JSON files become a performance bottleneck beyond 10,000 scans.** File Storage Scalability:
- **IPv6 address resolution is not currently supported.** IPv4 Only:
- **Batch scanning of multiple URLs simultaneously is not implemented.** Single-URL Scans:
- **No user accounts or multi-tenant access control is provided.** No Authentication:
- **100 scans may be insufficient for production-grade ML generalisation.** Small ML Training Set:

11.2 Future Scope

- **Periodic scans on user-defined intervals with configurable alert thresholds.** Automated Scheduling:
- **Notifications to Slack, Teams, PagerDuty, and OpsGenie on health score drops.** Email/Webhook Alerts:
- **PostgreSQL/MySQL backend for scalable historical data storage and indexing.** Database Migration:
- **Simultaneous multi-URL diagnostic processing.** Batch Scanning:
- **LSTM networks for time-series response time forecasting.** Deep Learning Models:
- **SHAP values for Random Forest model interpretability.** Explainable AI:
- **Diagnostic checks issued from multiple geographic locations.** Geo-distributed Testing:
- **Native iOS/Android app using React Native or Flutter.** Mobile Application:
- **APIs for AWS CloudWatch, Azure Monitor, and Google Cloud Operations.** LMS/Cloud Integration:

12. Conclusion

This paper presented SitePulse, a multi-layer website diagnostics and analytics system that integrates parallel infrastructure checking with data science visualisation and machine learning predictions. By simultaneously analysing DNS resolution, network reachability, HTTP compliance, SSL certificate validity, and response time performance through a weighted health scoring algorithm, rendering analytics dashboards from historical scan data, and applying Random Forest models to generate predictive health classifications and downtime risk assessments, SitePulse delivers a richer and more actionable picture of website infrastructure health than conventional single-layer monitoring tools.

The system demonstrates that meaningful ML-enhanced infrastructure diagnostics does not require complex database infrastructure — a carefully designed Flask application with file-based storage and scikit-learn ML models can deliver multi-dimensional, predictive insights accessible to any operator with a browser. As AI technology continues to mature, systems like SitePulse represent a promising direction for democratising intelligent infrastructure monitoring and helping organisations identify and resolve website failures before they impact end users. [1][6]

References

- [1] Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189–1232.
- [2] Tanenbaum, A. S., & Wetherall, D. J. (2011). *Computer Networks* (5th ed.). Prentice Hall.
- [3] Kurose, J. F., & Ross, K. W. (2017). *Computer Networking: A Top-Down Approach* (7th ed.). Pearson.
- [4] RFC 7231: HTTP/1.1 Semantics and Content. IETF. <https://tools.ietf.org/html/rfc7231>
- [5] RFC 6818: DNSSEC Validation Requirements. IETF. <https://tools.ietf.org/html/rfc6818>
- [6] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
- [7] Newman, M. E. J. (2003). The structure and function of complex networks. *SIAM Review*, 45(2), 167–256.
- [8] Xing, Z., Pei, J., & Keogh, E. (2010). A comprehensive survey on automatic detection of outliers in time series. Technical Report, University of Illinois.
- [9] Liu, F. T., Ting, K. M., & Zhou, Z. H. (2008). Isolation forest. In *IEEE ICDM 2008 Proceedings* (pp. 413–422).
- [10] Nagios Core: Open-source monitoring tool. <https://www.nagios.org/>
- [11] Zabbix: Enterprise-level monitoring solution. <https://www.zabbix.com/>



[12] Prometheus: Time-series database for monitoring. <https://prometheus.io/>

[13] Flask Documentation. <https://flask.palletsprojects.com/>

[14] scikit-learn User Guide. <https://scikit-learn.org/stable/documentation.html>

[15] Requests: Python HTTP Library. <https://requests.readthedocs.io/>

Appendix A: Key Formulas & Calculations

Health Score Formula

$$HS = 100 - (\sum w_i \times s_i)$$

Where HS = Health Score (0–100), w_i = weight of issue category i , s_i = severity of issue i (0–1), n = number of issue categories.

Anomaly Detection (Z-Score)

$$Z = (x - \mu) / \sigma \quad | \quad \text{Threshold: } |Z| > 2.5$$

Where x is the observed value, μ is the historical mean, and σ is the standard deviation of the metric.

Trend Strength

$$T = [(\text{recent_value} - \text{baseline_value}) / \text{baseline_value}] \times 100\%$$

Appendix B: API Endpoints Reference

Endpoint	Method	Description
POST /scan	POST	Run new website diagnostic scan
GET /health	GET	API health check
GET /api/analytics/statistics	GET	Statistical summary of scan history
GET /api/analytics/trends	GET	Trend analysis (improving/degrading)
GET /api/analytics/anomalies	GET	Detected anomalous scan records
GET /api/analytics/insights	GET	Actionable rule-based insights
GET /api/analytics/ml/classify	GET	ML health status classification
GET /api/analytics/ml/predict-downtime	GET	Downtime risk probability
GET /api/analytics/dashboard-data	GET	Complete analytics dashboard payload
GET /api/analytics/risk-score	GET	Infrastructure risk score (0–100)
GET /	GET	Home page
GET /history	GET	Scan history view
GET /dashboard	GET	Analytics dashboard
GET /analytics	GET	Data science dashboard

Appendix C: Configuration & Deployment

System Requirements

- Python 3.13+ runtime
- 512 MB RAM minimum
- 100 MB disk space (initial); grows with scan history
- Active network connectivity for diagnostic modules

Installation

```
git clone https://github.com/aniket-singh23/Site_Pulse.git
cd Site_Pulse
python -m venv .venv
.venv\Scripts\Activate
pip install -r requirements.txt
python app.py
```

Configuration (config.py)

```
SCAN_TIMEOUT = 30    # seconds
MAX_REDIRECTS = 5
SSL_VERIFY = True
LOG_FILE = 'data/logs.txt'
HISTORY_FILE = 'data/scan_history.json'
REPORTS_DIR = 'reports/'
```