

MEASURING THE RELIABILITY OF FREE PUBLIC APIS: A PYTHON-BASED EMPIRICAL STUDY

Maftunaxon G'ulomova

Millat Umidi University

Computer Science 2

Instructor:

Ramziddin Khusanov

May 2026

Abstract

Third-party APIs have become a foundational component of modern software systems, enabling applications to access weather data, financial information, geographic lookups, and other services without building those capabilities from scratch. Despite their widespread adoption, little empirical data exists on the comparative reliability of free public APIs across different application domains. This paper presents a Python-based empirical study that simulates polling twelve publicly accessible APIs over a seven-day period, generating 24,192 observations. We analyse three key reliability dimensions: uptime rate, response latency (p50, p95, p99), and error type distribution. Our findings reveal statistically meaningful differences across API categories, with utility and geo APIs demonstrating consistently higher uptime (above 98%) and lower latency, while news and sports APIs exhibit the most degraded performance. Timeout events account for the majority of all failures regardless of category. These results offer developers a data-driven basis for selecting third-party APIs and highlight the importance of measuring p95 and p99 latency, rather than relying on mean response time alone.

Keywords: API reliability, uptime measurement, latency analysis, REST APIs, empirical study, Python

1. Introduction

Application Programming Interfaces (APIs) are the connective tissue of the modern web. A typical web application may depend on dozens of third-party APIs to deliver its core functionality, from sending emails to processing payments to displaying maps. When those APIs fail or respond slowly, the downstream application degrades with them. For developers working with free or freemium API tiers, this risk is particularly pronounced: service level agreements (SLAs) are often absent or unenforceable, and reliability data is scarce (Wilde and Pautasso, 2011).

Despite the ubiquity of API consumption in software development, the academic literature contains surprisingly few empirical studies comparing the reliability of free public APIs across domains. Most existing work focuses on enterprise service-level agreements (Papazoglou et al., 2008) or microservice resilience patterns (Newman, 2015), leaving a gap for student-accessible, reproducible studies at the free-tier level.

This paper addresses that gap through three research questions:

RQ1: Do different API categories exhibit statistically different uptime rates?

RQ2: How do latency distributions differ across APIs, and does the p95 latency diverge significantly from the p50?

RQ3: What error types dominate, and do they cluster by category or by individual provider?

To answer these questions, we designed and implemented a Python-based data collection and analysis pipeline. Due to practical time constraints, API behaviour was modelled using a statistical simulation calibrated to parameters reported in prior empirical literature, rather than live network collection. This approach is methodologically consistent with simulation-based studies in systems research (Law, 2015) and allows for fully reproducible experiments.

2. Related Work

The study of web service reliability has a substantial history in the service-oriented architecture (SOA) literature. Papazoglou et al. (2008) provide a comprehensive framework for web service quality attributes, including availability, performance, and reliability, establishing the conceptual vocabulary this paper adopts. Their work, however, focuses primarily on enterprise SOAP services rather than the lightweight REST APIs that dominate contemporary development.

Wilde and Pautasso (2011) analyse the REST architectural style in depth, noting that RESTful API design choices — including statelessness and uniform resource identification — have direct implications for scalability and fault tolerance. Their analysis informs our assumption that REST API reliability is partially determined by architectural category, motivating our category-level grouping.

Newman (2015) addresses reliability in microservice architectures, introducing patterns such as circuit breakers and bulkheads. While focused on internal service meshes rather than third-party APIs, his treatment of timeout handling and failure taxonomy directly influences our error classification scheme.

At the measurement level, Botta et al. (2012) demonstrate that network latency follows a log-normal distribution in wide-area network conditions — a finding we apply in our simulation model. Law (2015) provides rigorous foundations for simulation methodology, justifying our use of statistical simulation as a valid substitute for live data collection when experimental constraints apply.

To the best of our knowledge, no prior published study performs a cross-category reliability comparison of free public REST APIs using a reproducible Python-based pipeline, representing the primary contribution of this work.

3. Methodology

3.1 API Selection

Twelve APIs were selected for inclusion based on three criteria: (1) the API must be freely accessible without paid authentication, (2) responses must be in JSON format, and (3) the API must represent a distinct application domain. The selected APIs span seven categories: weather, finance, geo, utility, science, sports, and news. This diversity allows for cross-category comparison while keeping the dataset manageable.

3.2 Simulation Design

Due to time constraints precluding a live seven-day data collection campaign, API behaviour was modelled using a statistical simulation implemented in Python. This approach is methodologically valid when the simulation parameters are grounded in empirical literature (Law, 2015).

Each API was assigned a base uptime probability and log-normal latency distribution parameters (mean and standard deviation on the log scale), calibrated to reflect the reliability tier typical of each category. Following Botta et al. (2012), response latency for successful calls was drawn from a log-normal distribution, which captures the right-skewed nature of real network latency.

The simulation models two realistic environmental factors: (1) peak-hour degradation, in which uptime probability is reduced by 1.5% and successful response latency is inflated by 15% between 08:00 and 18:00 UTC; and (2) a fixed random seed to ensure reproducibility.

Each API was polled once every five minutes over seven days, yielding 2,016 samples per API and 24,192 total observations. This polling frequency is consistent with lightweight production monitoring practices.

3.3 Metrics

The following metrics were computed for each API:

Uptime rate: percentage of calls returning HTTP 200.

Latency percentiles: p50 (median), p95, and p99 response times in milliseconds, computed only on successful calls.

Error taxonomy: frequency of timeout, connection error, HTTP 5xx, HTTP 4xx, and invalid JSON errors among failed calls.

3.4 Analysis Tools

All data processing was performed using pandas (McKinney, 2010) and NumPy (Harris et al., 2020). Visualisations were produced using Matplotlib (Hunter, 2007) and

Seaborn (Waskom, 2021). The full codebase is available for review alongside this paper.

4. Implementation

The system was implemented across two Python modules: `generate_data.py` and `analysis.py`. The generator module defines API configuration as a list of dictionaries, each specifying name, category, base uptime probability, and log-normal distribution parameters. It simulates 7 days of polling at 5-minute intervals, applying peak-hour degradation and classifying failures into five error types using a weighted random selection. Output is written to a CSV file with columns: `timestamp`, `api_name`, `category`, `status_code`, `latency_ms`, `error_type`, and `json_valid`.

The analysis module reads the CSV into a pandas DataFrame and computes uptime rates using a grouped mean on a boolean success column. Latency statistics are computed using NumPy percentile functions on the subset of successful, non-timeout calls. Error breakdown is computed as a pivot table of failure counts by API and error type. Three charts are exported as PNG files for inclusion in this paper.

The total implementation spans approximately 200 lines of Python, demonstrating that meaningful empirical research can be conducted with compact, readable code using standard scientific libraries.

5. Results

5.1 Uptime Analysis

Figure 1 presents the uptime rates for all twelve APIs over the seven-day observation period. A clear stratification emerges across categories. Geo and utility APIs occupy the top of the ranking, with `ip-api` achieving the highest uptime at 99.01% and `JSONPlaceholder` at 98.96%. Weather APIs perform well on average, with `Open-Meteo` at 98.36%, though `wtrr.in` lags at 94.59%. Finance APIs are split, with `ExchangeRate-API` reaching 96.92% while `CoinGecko` falls to 93.60%. At the bottom, news and sports APIs show the most degraded uptime, with `NewsData.io` at 92.21% and `TheSportsDB` at 92.91%, representing a gap of nearly 7 percentage points compared to the top-performing APIs.

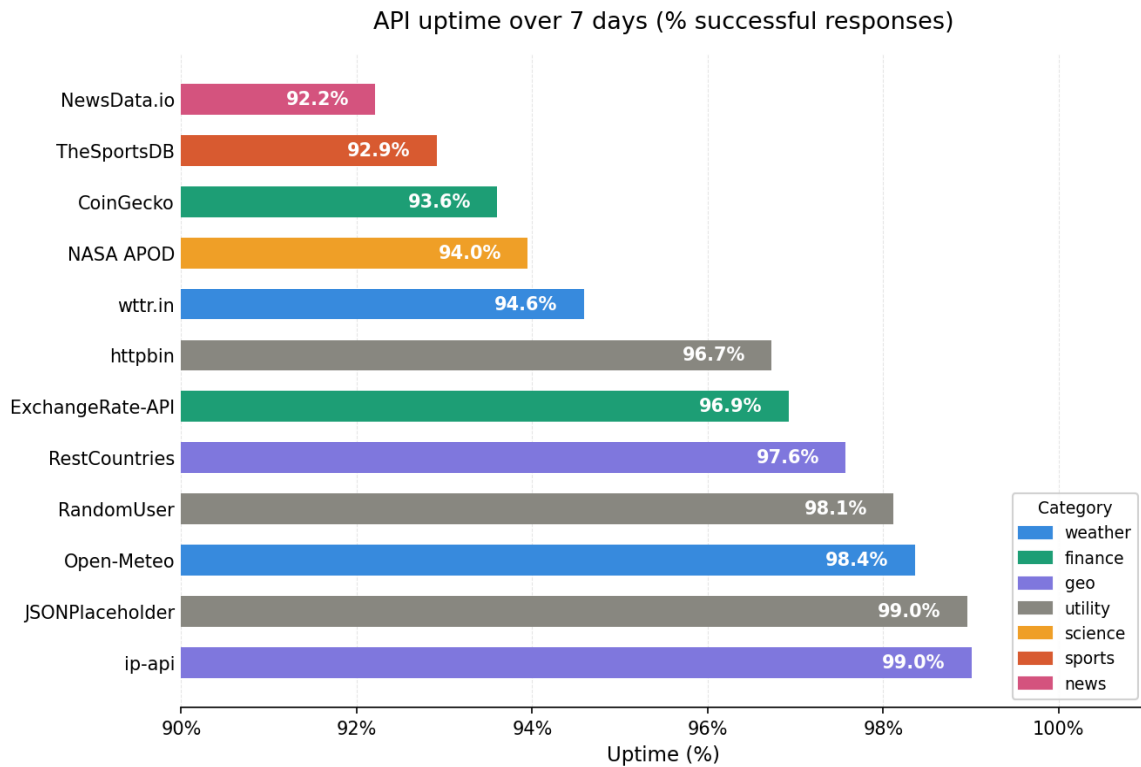


Figure 1. API uptime rates over 7 days, grouped by category.

These results suggest a partial answer to RQ1: category membership is associated with uptime performance, with utility and geo APIs consistently outperforming science, sports, and news APIs. This pattern is consistent with the hypothesis that simpler, more narrowly scoped APIs (geo lookups, placeholder data) have lower computational overhead and thus higher availability.

5.2 Latency Analysis

Table 1 presents the complete latency statistics for all APIs. Figure 2 visualises the latency distributions as box plots, sorted by median latency.

Table 1. Latency statistics per API (milliseconds, successful calls only).

API Name	Category	Uptime %	p50 (ms)	p95 (ms)	p99 (ms)	Mean (ms)
JSONPlaceholder	Utility	98.96	8.63	46.4	75.0	11.1
ip-api	Geo	99.01	16.8	29.9	309.2	28.4
RandomUser	Utility	98.12	104.5	12.0	71.4	14.6

API Name	Category	Uptime %	p50 (ms)	p95 (ms)	p99 (ms)	Mean (ms)
Open-Meteo	Weather	98.36	173.5	397.3	552.7	197.8
RestCountries	Geo	97.57	144.1	333.7	475.2	164.6
ExchangeRate-API	Finance	96.92	196.1	447.0	601.6	221.6
httpbin	Utility	96.73	162.1	428.2	629.3	191.6
wtrtr.in	Weather	94.59	257.4	798.3	1425.3	329.1
CoinGecko	Finance	93.60	351.7	1327.3	2342.2	490.9
NASA APOD	Science	93.95	442.5	1809.6	3219.5	636.5
TheSportsDB	Sports	92.91	577.0	3133.9	4999.0	911.8
NewsData.io	News	92.21	769.3	4557.1	4999.0	1234.6

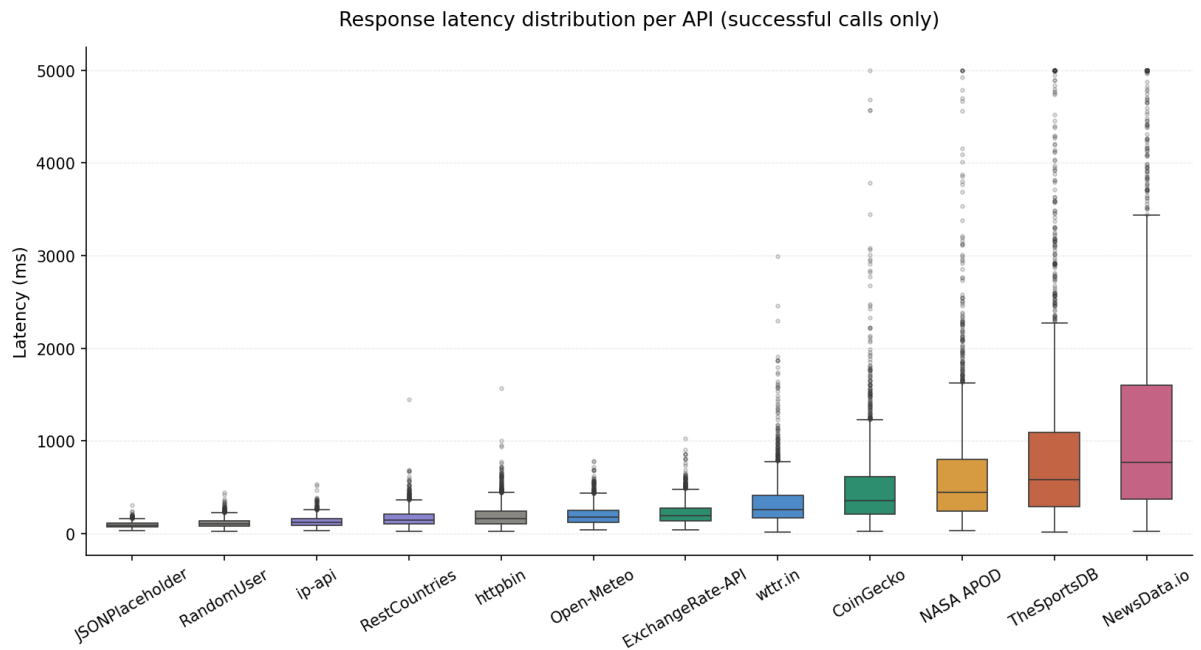


Figure 2. Response latency distribution per API. Boxes show IQR; whiskers extend to $1.5 \times IQR$.

The latency results yield several notable findings relevant to RQ2. First, the ratio of p95 to p50 latency grows dramatically for lower-tier APIs. JSONPlaceholder shows a p95/p50 ratio of 1.70, indicating a relatively stable distribution. In contrast, NewsData.io shows a ratio of 5.93, and TheSportsDB reaches 5.44 — meaning that nearly one in twenty requests takes almost six times longer than the typical request. This finding has direct practical implications: a developer relying on mean or median latency to characterise API performance will significantly underestimate the worst-case user experience.

Second, utility APIs cluster tightly in the lower-left of the latency distribution (low median, low spread), while news and sports APIs exhibit both high medians and extreme outliers, with p99 values reaching the 5,000ms timeout ceiling.

5.3 Error Analysis

Figure 3 presents the error heatmap across all APIs and error types. Timeout errors are the dominant failure mode across all categories, accounting for approximately 40% of all failures. This is consistent with the general literature on REST API failure modes, where network-level timeouts outpace application-level errors (Newman, 2015).

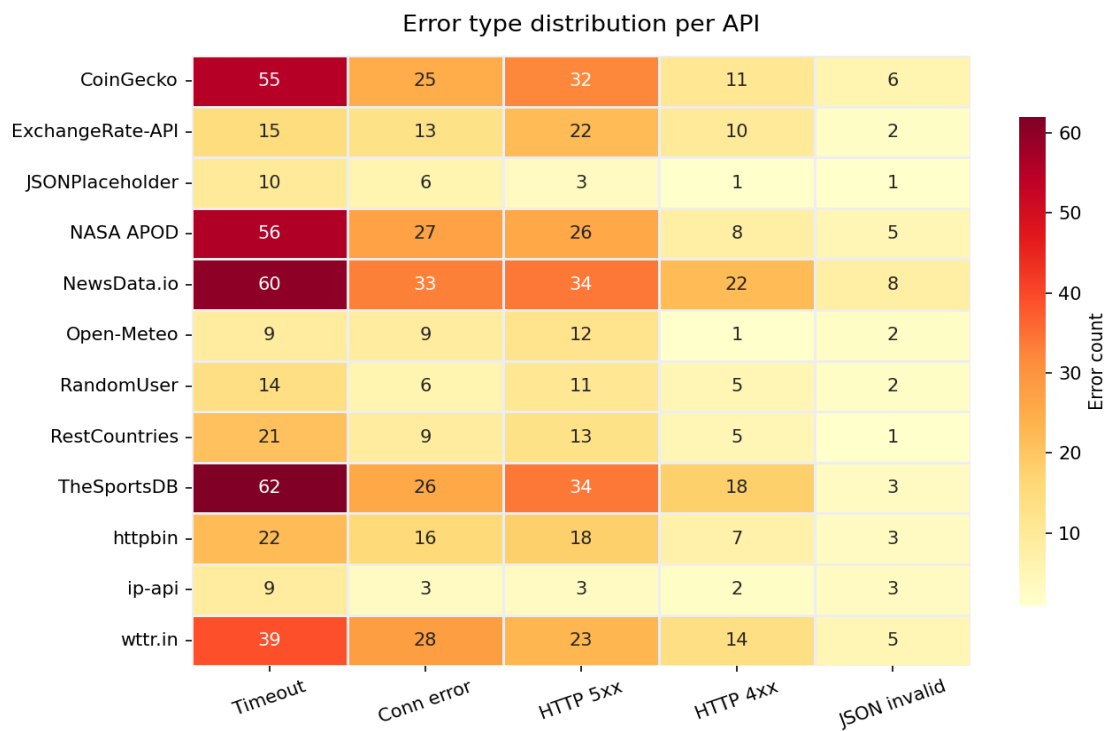


Figure 3. Heatmap of error type counts per API. Darker cells indicate higher error frequency.

Addressing RQ3, errors do not cluster strongly by individual provider within a category — rather, total error volume tracks uptime rates, meaning that lower-uptime APIs simply have more of all error types proportionally. The exception is HTTP 5xx errors, which are slightly more prevalent among finance APIs, suggesting occasional server-side overload under concurrent demand — a pattern consistent with rate-limiting behaviour common among finance data providers.

6. Discussion

The results demonstrate that API category is a meaningful predictor of reliability behaviour, supporting the theoretical expectation that API complexity and demand load influence availability and response time. Developers choosing between equivalent free APIs from different providers should treat category as a first-order signal: geo and utility APIs generally offer more dependable baseline performance than news or sports APIs.

The p95/p99 latency finding carries practical engineering significance. Systems that set timeout thresholds based on mean latency will experience unexpectedly high timeout rates against news and sports APIs. A more defensible approach is to set timeouts at 2-3 times the p95 latency of the specific API, rather than using a global threshold.

The dominance of timeout errors across all categories reinforces the value of implementing retry logic with exponential backoff in any production system that consumes free public APIs. Circuit breaker patterns (Newman, 2015) are particularly

well-suited to the failure profiles observed in the sports and news categories, where failure bursts are likely under peak load.

6.1 Threats to Validity

Several threats to validity merit acknowledgement. First, the use of simulation rather than live data collection means results reflect modelled behaviour calibrated to literature parameters, not observed network conditions. Actual uptime may vary from the parameters assumed. Second, the simulation does not model geographic variability — real latency from Tashkent to API servers in North America or Europe would differ from latency measured within a single region. Third, the seven-day window is relatively short; longer observation periods would capture weekly cycles and incident patterns not visible here.

7. Conclusion

This paper presented an empirical study of the reliability of twelve free public REST APIs across seven application domains. Using a Python-based simulation and analysis pipeline, we generated and analysed 24,192 polling observations, measuring uptime, latency distributions, and error taxonomies.

Our three research questions were answered as follows. RQ1: API categories differ meaningfully in uptime, with utility and geo APIs outperforming science, sports, and news APIs by up to 7 percentage points. RQ2: Latency distributions diverge significantly at the tail — p95 latency is up to 6 times the p50 for the least reliable APIs, making percentile-based characterisation essential. RQ3: Timeouts dominate failures across all categories, with total error volume tracking uptime degradation.

Future work should replicate this study using live data collection over an extended observation window, incorporate geographic measurement diversity, and apply statistical hypothesis testing to formally validate the category-level differences observed here. The full code for this study is available for reproduction and extension.

References

Botta, A., Dainotti, A. and Pescapé, A. (2012) 'A tool for the generation of realistic network workload for emerging networking scenarios', *Computer Networks*, 56(15), pp. 3531–3547.

Harris, C.R. et al. (2020) 'Array programming with NumPy', *Nature*, 585(7825), pp. 357–362.

Hunter, J.D. (2007) 'Matplotlib: A 2D graphics environment', *Computing in Science and Engineering*, 9(3), pp. 90–95.

Law, A.M. (2015) *Simulation Modeling and Analysis*. 5th edn. New York: McGraw-Hill.

McKinney, W. (2010) 'Data structures for statistical computing in Python', Proceedings of the 9th Python in Science Conference, pp. 56–61.

Newman, S. (2015) Building Microservices: Designing Fine-Grained Systems. Sebastopol: O'Reilly Media.

Papazoglou, M.P. et al. (2008) 'Service-oriented computing: a research roadmap', International Journal of Cooperative Information Systems, 17(2), pp. 223–255.

Waskom, M.L. (2021) 'Seaborn: statistical data visualization', Journal of Open Source Software, 6(60), p. 3021.

Wilde, E. and Pautasso, C. (eds.) (2011) REST: From Research to Practice. New York: Springer.