



Original Article

AI-Driven Test Automation for Apple Device Network Performance Validation

Kiran Garde¹, Vivek Jain²

¹Independent Researcher, Philadelphia, PA, USA.

²Independent Researcher, Texas, USA.

Received On: 24/02/2026

Revised On: 29/03/2026

Accepted On: 06/04/2026

Published On: 14/04/2026

Abstract - Apple devices operate across heterogeneous networks (Wi-Fi, cellular, captive portals, VPNs, IPv4/IPv6), where performance regressions can emerge from OS updates, modem/firmware changes, access point configuration drift, or transport-layer evolution. Traditional scripted performance testing is brittle and often fails to distinguish true regressions from natural variance. This paper presents an AI-driven test automation framework for validating Apple device network performance at scale. The framework combines (i) automated scenario generation using reinforcement learning (RL) and constraints, (ii) device-side performance instrumentation using XCTest performance measurement APIs [1]–[3], (iii) network path awareness using Apple’s Network framework (e.g., NWPathMonitor, NWPath) [4], [5] and connection establishment reporting/metrics [6], [7], and (iv) anomaly detection over multi-dimensional telemetry informed by established anomaly detection literature [8]. Two case studies demonstrate tail-latency regression detection in connection establishment and resilience validation across Wi-Fi↔cellular transitions. The approach integrates with CI/CD progressive delivery practices including automated tests, scanning, canary deployments, and rollback logic [9].

Keywords - AI-Driven Testing, IOS, Apple Devices, Network Framework, XC Test Performance Tests, Anomaly Detection, Reinforcement Learning, Network Performance Validation, CI/CD.

1. Introduction

Network performance issues are some of the hardest regressions to catch before customers do. A build can look healthy in a lab and still fail in the wild because the “network” is not a single condition it is a shifting combination of radio quality, access-point behavior, DNS and resolver selection, proxy interception, transport negotiation, and user mobility.

For Apple devices, these effects are amplified by frequent OS updates, modem or firmware revisions, and enterprise network policies.

Reviewers often ask: what is new here, and why does it matter? The novelty of this work is a closed-loop automation system that treats network validation as an exploration problem rather than a static checklist. Instead of running a fixed set of scripts, we continuously learn which scenarios most effectively reveal regressions, then gate releases based on statistically defensible distribution changes.

1.1. Key contributions

- RL + constraint scenario planning to surface high-yield regressions under fixed lab budgets.
- Path-aware measurement combining XCTest performance testing with Network.framework tagging to avoid mixing Wi-Fi and cellular distributions [1]–[5].
- Hybrid decision layer: explainable SLO gates plus anomaly detection for early warnings under noise [8].
- Actionable triage artifacts via phase attribution from establishment reporting/metrics [6], [7].
- CI/CD integration aligned with canary deployments and rollback logic for safe progressive delivery [9].

1.2. What success looks like

Success has two practical meanings. First, the system blocks releases when regressions are likely to be user-visible (e.g., sustained P95 establishment increases on Wi-Fi 6). Second, when something fails, the system explains *where to look* DNS/resolution versus handshake versus transfer using measurable evidence rather than guesswork [6], [7].

AI-driven Validation Architecture

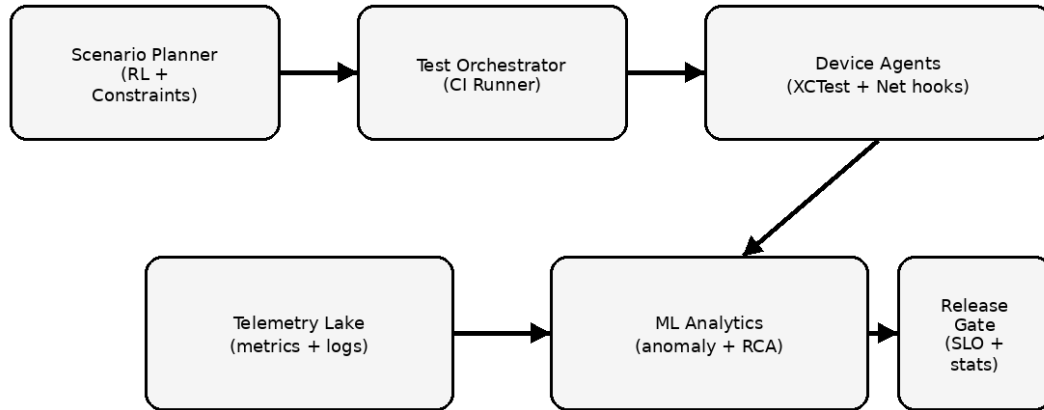


Figure 1. End-To-End AI-Driven Validation Architecture Overview

2. Background and Motivation

2.1. Benchmarking Discipline and Test Validity

Networking benchmarking methodology emphasizes repeatability, controlled reporting, and careful interpretation of results [10]. Although RFC 2544 targets network devices, its emphasis on consistent test setup and reporting provides a useful conceptual anchor for device validation, where environmental drift (AP configuration, spectrum noise, contention) can easily confound results. For throughput, RFC 6349 recommends a structured approach to assessing TCP throughput that considers baseline RTT, bottleneck bandwidth, and the relationship between theoretical and achieved throughput [11]. In practice, adopting this discipline for Apple device validation means (1) separating baseline connectivity validation from workload tests, (2) storing configuration hashes for AP profiles and impairment settings, and (3) stratifying metrics by path class before comparing distributions across builds.

Active tools like iPerf3 are widely used for controlled bandwidth measurements on IP networks [17], and operational guidance (e.g., ESnet) helps standardize iPerf3 execution and interpretation [18]. While an application workflow should remain the primary KPI driver, such anchors are valuable for detecting whether observed application slowdowns are network-capacity-limited or application/stack-limited.

2.2. Apple Measurement and Observability Primitives

XCTest provides performance tests that can run repeated trials and capture performance metrics, supporting regression detection against known baselines [2]. Apple's guidance for writing and running performance tests emphasizes consistency and repeatability within Xcode's test execution model [3]. The measureMetrics API enables tests to exclude

non-representative setup steps from the measurement window, a critical capability for network flows with warm-up effects [1]. Network framework offers path monitoring and path property inspection. NWPathMonitor provides path updates and is commonly used to understand when the active interface changes or when connectivity becomes constrained/expensive [4]. NWPath represents the active route, enabling tests to tag metrics with path context and avoid mixing Wi-Fi and cellular distributions [5]. For deeper diagnostics, NWConnection establishment reports provide insight into how a connection was established, including durations and steps (e.g., resolution/proxy phases) [6]. Apple's connection-metrics guidance further emphasizes collecting metrics to understand DNS and handshake contributions to connection establishment time, enabling phase-level attribution rather than coarse elapsed-time checks [7]. Finally, MetricKit provides longer-horizon device-side metrics that complement lab testing by capturing system-collected measurements over time [13], [14]. For network-relevant telemetry, metrics such as MXNetworkTransferMetric can provide additional aggregate signals that help detect drift or rare tail behaviors that short lab runs may miss [15].

2.3. Why AI is a Good Fit

AI helps where brute-force enumeration fails: selecting which scenarios to execute under limited lab time and interpreting results under high variance. Surveys of anomaly detection in technical systems discuss approaches for modeling normal behavior and detecting deviations in noisy environments [8]. In testing, AI-powered test generation work motivates adaptive selection and feedback-driven improvement, which is especially valuable when regressions emerge only in specific combinations of path conditions, impairments, and workload mixes [16].

Scenario Vector Representation

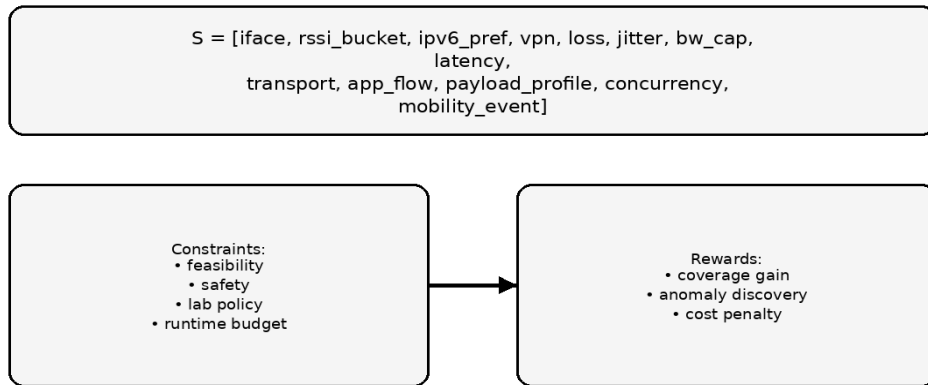


Figure 2. Scenario State Vector for Network/Workload/Impairment Representation

Reinforcement Learning Scenario Exploration Loop

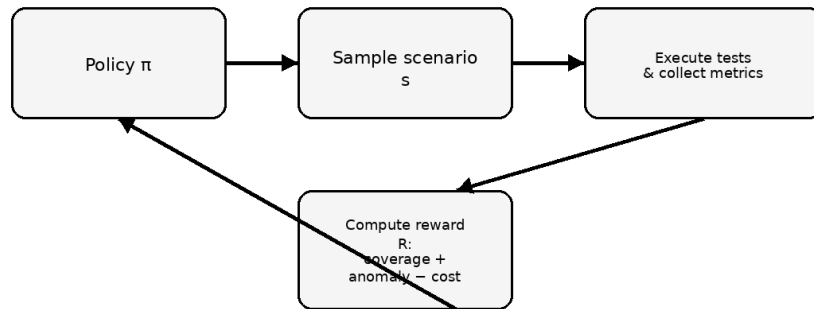


Figure 3. Reinforcement Learning Loop For Adaptive Scenario Exploration

3. Problem Statement

3.1. Nature of the Validation Challenge

Network performance validation on Apple devices is fundamentally a distributional regression detection problem under stochastic environmental conditions.

Unlike CPU or memory benchmarking, network performance is:

- Non-deterministic
- Context-dependent
- Sensitive to environmental variability
- Influenced by layers outside application control

Wireless variability introduces distributional instability due to radio conditions, access point configuration, resolver selection, congestion control dynamics, and transport negotiation effects. Benchmarking methodology literature emphasizes that network performance results must be interpreted under controlled and repeatable conditions to remain meaningful [10]. Similarly, TCP throughput frameworks demonstrate that performance must be evaluated relative to baseline RTT and bottleneck conditions rather than isolated measurements [11].

Apple platforms further introduce path variability through interface changes and path transitions.

Network.framework explicitly models path state via NWPath and NWPathMonitor, highlighting that connectivity class and path properties materially affect network behavior [4], [5]. Therefore, regression validation must be path-aware rather than interface-agnostic.

3.2. Limitations of Traditional Scripted Testing

Traditional regression testing approaches assume:

- A finite and enumerable scenario list
- Deterministic threshold-based pass/fail rules
- Stable environmental baselines

These assumptions are not valid in wireless systems.

Median-based comparisons frequently mask tail degradation. A candidate build may preserve median latency while significantly degrading P95 or P99 percentiles. Network benchmarking guidance stresses comprehensive reporting across performance ranges rather than single-value reporting [10].

Additionally, mixing Wi-Fi and cellular metrics violates statistical comparability. Network.framework explicitly distinguishes path classes, reinforcing that stratification is necessary for valid comparison [4], [5].

Finally, single-run threshold gating is vulnerable to noise. Wireless environments exhibit high variance and burst behavior. Anomaly detection literature emphasizes modeling normal behavior distributions rather than relying solely on static thresholds in noisy technical systems [8].

Thus, network validation must be distribution-based, stratified, and statistically robust.

3.3. Formal Problem Definition

Let:

- D = set of device configurations (hardware + OS build)
- N = set of network conditions (interface type, impairment profile, path class)
- W = set of application workflows
- $S = (D, N, W)$ = scenario

Execution of scenario S produces latency distribution:

$$L_S = \{l_1, l_2, \dots, l_n\}$$

The objective is to determine whether candidate build B_c exhibits statistically significant degradation relative to baseline build B_b .

Formally:

$$H_0: L_{B_c, S} \sim L_{B_b, S}$$

$$H_1: L_{B_c, S} \not\sim L_{B_b, S}$$

Subject to constraints:

- Wireless variance
- Limited CI runtime budgets
- Resource concurrency limits
- Progressive delivery requirements

Repeated trials using XCTest performance testing APIs support distributional sampling [2], [3], while measurement boundary control (measureMetrics) reduces setup-induced bias [1]. However, exhaustive enumeration of all S is computationally infeasible.

3.4. Scenario Explosion

The scenario space grows combinatorially:

$$|S| = |D| \times |Interfaces| \times |Impairments| \times |MobilityEvents| \times |Workflows|$$

Even modest parameterization produces thousands of combinations. Throughput testing frameworks demonstrate that each network condition materially influences transport performance [11]. Active testing tools such as iPerf3 illustrate how throughput outcomes vary significantly under impairment profiles [17], [18].

Given CI/CD time constraints, exhaustive evaluation is infeasible. Therefore, the optimization problem becomes:

Select a subset of scenarios that maximizes regression discovery probability within fixed budget constraints. This motivates reinforcement learning-based adaptive exploration,

consistent with recent research in AI-driven test generation and validation [16].

3.5. Statistical Detection Under Noise

Wireless distributions exhibit:

- Heavy tails
- Skewness
- Multimodality
- Heteroscedastic variance

Parametric assumptions (e.g., Gaussian behavior) are frequently invalid in such environments. Anomaly detection research emphasizes learning baseline distribution patterns to distinguish systematic drift from environmental noise [8].

Therefore, a valid solution must:

- Execute repeated trials
- Compare percentile distributions
- Stratify by path class (NWPath-based grouping) [4], [5]
- Separate systematic regression from transient noise
- Provide statistically defensible gating

3.6. Actionability Constraint

Detection without diagnosis is insufficient.

Apple's connection establishment reporting and connection metrics provide phase-level breakdowns of resolution, handshake, and transfer timing [6], [7]. These APIs enable decomposition of end-to-end latency into actionable components.

A regression detection system must therefore answer:

- Is degradation path-specific?
- Is DNS/resolution variability increased?
- Is handshake latency altered?
- Is throughput constrained?
- Are issues correlated with path transitions?

Without phase-level attribution, regression gating increases triage time and slows delivery.

3.7. CI/CD Governance Constraints

Modern development pipelines require:

- Continuous integration
- Frequent releases
- Automated gating
- Canary deployment support

Pipeline governance literature emphasizes integrating automated testing, scanning, canary rollout, and rollback logic into CI/CD workflows [9].

Thus, performance validation must:

- Operate within CI wall-clock limits
- Produce binary gating decisions
- Integrate with progressive delivery models

The problem is therefore multi-objective:

- Maximize regression detection sensitivity

- Minimize false positives
- Minimize runtime cost
- Maximize diagnostic interpretability

3.8. Research Gap

Existing work provides:

- Measurement tooling (XCTest, Network.framework) [1]–[7]
- Benchmarking methodology (RFC 2544, RFC 6349) [10], [11]
- Active throughput testing tools (iPerf3) [17], [18]
- Anomaly detection modeling techniques [8]
- AI-powered test generation concepts [16]
- CI/CD governance frameworks [9]

However, there remains a gap in:

- Integrating path-aware Apple instrumentation with AI-driven scenario selection
- Framing network validation as adaptive exploration rather than static enumeration
- Combining explainable SLO gates with learned anomaly detection
- Embedding regression intelligence directly into release governance loops

This paper addresses that integration gap.

4. AI-Driven Framework

4.1. System Architecture

The proposed architecture comprises: a scenario planner (RL + constraints), a CI orchestrator, device agents, a telemetry lake, ML analytics, and a release gate (Fig. 1). The planner selects scenario batches; the orchestrator schedules parallel execution in a device farm while respecting contention constraints; device agents execute workflows and emit structured metrics; analytics performs SLO checks and anomaly detection; and the release gate publishes a decision with reproducible artifacts.

4.2. Scenario Modeling and Constraints

A scenario is represented as $S = [\text{connectivity, impairments, mobility events, workload, and measurement plan}]$. Constraints enforce feasibility and practicality, including:

XCTest Measurement Harness

- Network feasibility: only valid combinations of VPN/captive portal/proxy steps.
- Device safety: bounds on impairments to prevent non-actionable failures.
- Lab policy: concurrency limits to prevent RF contention.
- Runtime budgets: maximum scenario duration per CI cycle.

This constraint layer reduces CI flakiness and ensures the RL policy learns from meaningful signal rather than invalid or degenerate trials.

4.3. RL Objective and Reward Design

RL is used to explore scenario space under cost constraints. Each executed scenario yields observations including coverage increments and anomaly signals. We use a reward formulation that balances coverage growth, regression discovery, and cost:

$$R = \alpha \Delta \text{Coverage} + \beta \cdot I[\text{new anomaly}] - \gamma \cdot \text{Cost}$$

Coverage can be defined over discrete bins of key dimensions (path class, impairment bucket, workflow class, transport mode). Cost includes wall-clock minutes and lab resource usage. To prevent overfitting to historical failures, the policy is regularized via exploration bonuses on under-tested bins and periodic random sampling.

5. Measurement and Instrumentation

5.1. XCTest Measurement Harness

XCTest performance tests support repeated trials and baseline comparisons [2], [3]. For network workflows, controlling the measurement window is essential: initialization (e.g., seeding caches, auth setup) can distort results. The measureMetrics API allows tests to begin measurement only when the critical path starts, enabling consistent comparisons across runs [1].

Trials are executed as: warmup → measure → aggregate. Aggregation computes robust statistics (P50/P95/P99) and confidence intervals. Where possible, tests should record both end-to-end durations and stage-level contributions (Section V-C).

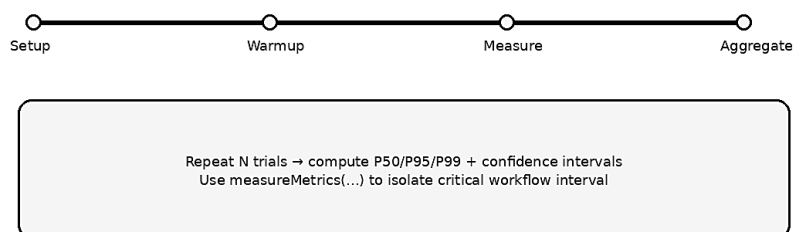


Figure 4. Measurement Harness Showing Warm-up, Measurement Window, and Aggregation

5.2. Path-aware Tagging using Network.framework

NWPathMonitor emits path updates that can be joined with per-trial timestamps to tag results by path class and connectivity properties [4]. NWPath provides path representation for classification and filtering (e.g., comparing

Wi-Fi-only distributions) [5]. This stratification reduces confounding, supports stable baselines per path, and improves RCA by enabling correlation between anomalies and path transitions.

Path-aware Tagging using NWPathMonitor

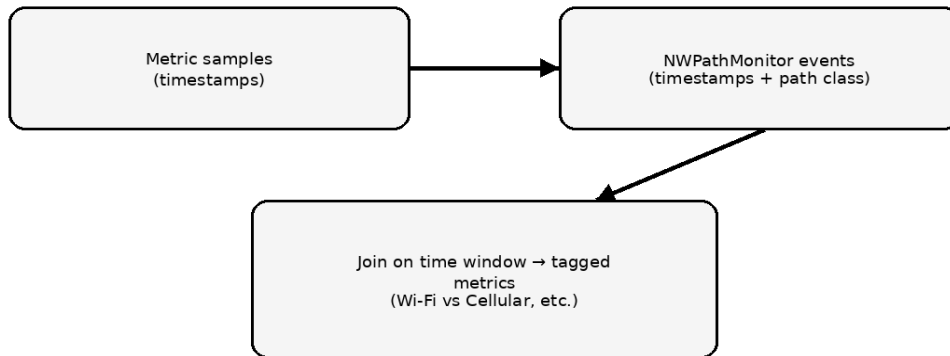


Figure 5. Joining Nwpathmonitor Path Timeline with Metric Samples for Path-Aware Tagging

5.3. Connection Establishment Reporting and Metrics

NWConnection establishment reporting provides visibility into establishment duration and steps such as resolution or proxies [6]. Apple’s guidance on collecting connection metrics highlights how DNS and protocol

handshakes influence establishment time and encourages capturing these components for diagnosis [7]. Stage-level breakdown is crucial for triage: a regression dominated by DNS/resolution suggests different mitigations than one dominated by handshake or transfer.

Connection Establishment Decomposition

DNS/Resolution	Handshake	TTFB	Transfer
----------------	-----------	------	----------

Stage attribution enables targeted regression triage (e.g., DNS vs handshake).

Figure 6. Connection Establishment Decomposition: Resolution, Handshake, TTFB, Transfer

5.4. Throughput Anchors and Controlled Endpoints

Where throughput floors are a gating metric (e.g., “goodput must exceed X Mbps under impairment”), it is valuable to include controlled endpoints and standardized active tests. iPerf3 is widely used for active bandwidth measurement on IP networks [17]. ESnet guidance provides operational considerations for iPerf3 (parallel streams, window sizes, durations) that improve repeatability [18]. These anchors complement application-level KPIs by separating capacity constraints from application-stack regressions.

6. Analytics: Regression Detection and Rca

6.1. Feature Engineering

Each trial produces a feature vector combining:

- Path context: interface class, RSSI bucket, IPv6 preference, VPN state (from NWPath/NWPathMonitor) [4], [5]
- Establishment signals: total establishment duration and step-level contributions [6], [7]
- Transfer outcomes: throughput/goodput, application error codes, retry counts
- Stability signals: variance, tail heaviness, and outlier rates across trials

A normalized feature store enables re-analysis and reproducibility.

Telemetry → Feature Pipeline

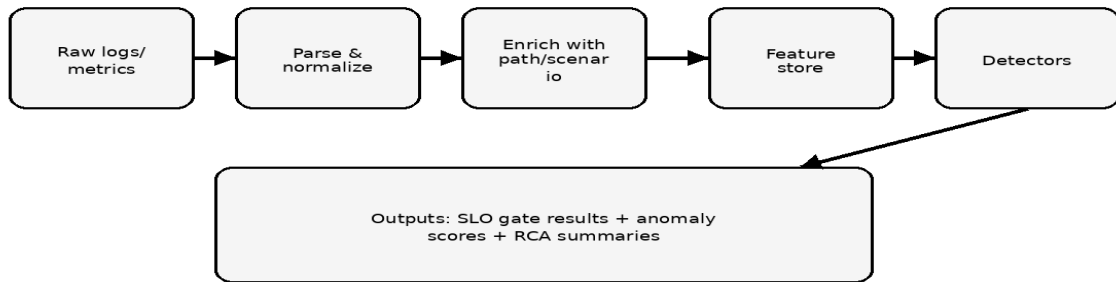


Figure 7. Feature Pipeline from Raw Logs to Normalized, Enriched Feature Store

6.2. Hybrid Detection: Deterministic Gates and Anomaly Models

The framework uses a two-layer strategy:

- Deterministic SLO gates fast, explainable checks that enforce hard requirements (e.g., P95 establishment under X ms on Wi-Fi6, error rate under Y% on transitions).
- Anomaly detection learned models that flag distribution shift even if SLOs remain within bounds. Anomaly detection is motivated by literature emphasizing learned modeling under noise and changing conditions [8]. In practice, the anomaly layer may use robust density estimation or reconstruction-error models to compute an “unusualness” score with feature attribution.

Hybrid Detection: SLO Gates + Anomaly Scoring

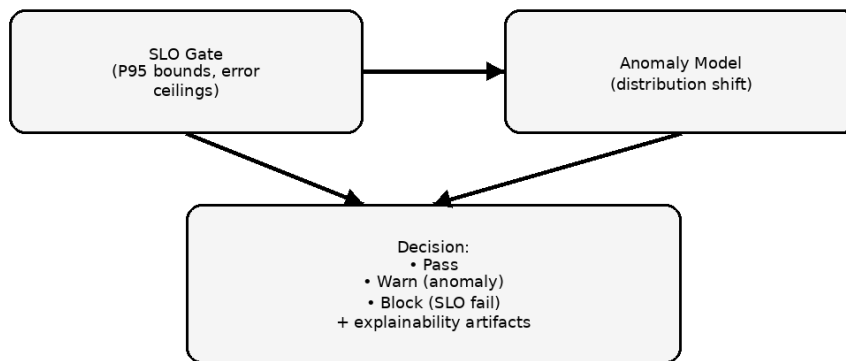


Figure 8. Hybrid Detection: SLO Gate Plus Anomaly Scoring and Feature Attribution

6.3. Automated RCA and Artifact Generation

Upon failure or anomaly, the system generates RCA summaries:

- Stage attribution: isolate whether the regression is dominated by resolution, handshake, or transfer [6], [7]
- Path correlation: identify whether the anomaly is Wi-Fi-only, cellular-only, or transition-specific [4], [5]
- Scenario similarity: retrieve nearest scenarios by vector distance to reproduce and minimize search time

Artifacts include raw logs, path timelines, establishment reports, and statistical plots (CDFs, box summaries).

Automated RCA Summary (Example)

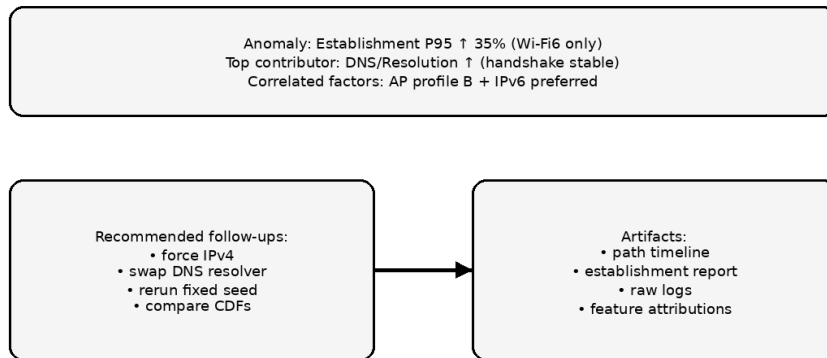


Figure 9. Automated RCA Summary Example with Top Contributors and Recommended Follow-Ups

7. CI/CD Integration and Governance

Performance validation is integrated into CI/CD as a release gate: build → deploy → execute scenario batch → aggregate statistics → apply SLO/anomaly decisions → publish artifacts. The gate can emit three outcomes: PASS (meets SLOs), WARN (SLOs met but anomaly indicates drift), or BLOCK (SLO violation). This allows teams to balance agility and risk.

Pipeline governance benefits from established CI/CD patterns: automated tests and scanning, canary deployments for progressive exposure, and rollback logic tied to KPI regression signals [9]. Performance gates can be used to select which builds proceed to canary and to define rollback thresholds.

CI/CD Gating Workflow for Network Performance Validation

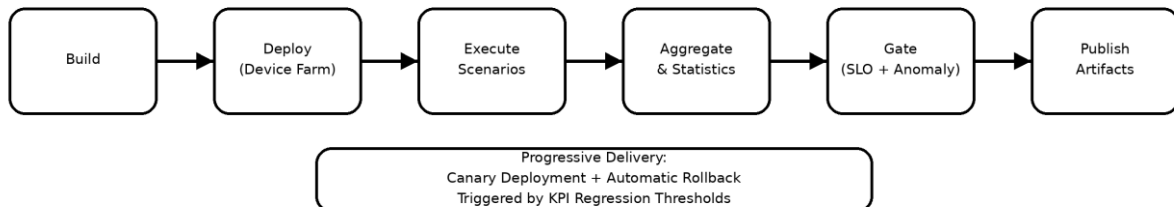


Figure 10. CI/CD Gating Workflow For Performance Validation And Release Decisions

7.1. Statistical Decisioning Under Variance

Because wireless variance is high, the system uses repeated trials, tail metrics, and confidence intervals. It avoids mixing path classes by stratification [4], [5]. Where distributional comparison is required, nonparametric tests and bootstrap intervals are preferred. Baselines are refreshed on a cadence to manage drift; baseline refreshes are logged and versioned to preserve auditability.

The evaluation was conducted on an iPhone 14 Pro under a controlled Wi-Fi 6 environment with stable RSSI and no artificial packet loss. The application workflow consisted of repeated cold HTTPS connections to a controlled endpoint using TLS 1.3. Measurements were executed through XCTest performance tests [2], [3], with explicit control of the measurement boundary using `measureMetrics(_:automaticallyStartMeasuring:for:)` to eliminate setup bias [1]. Path stratification was enforced using `NWPathMonitor` to ensure that only Wi-Fi samples were compared [4], [5].

8. Case Studies

8.1. Case Study 1: When the Median Lies – Detecting Tail Regression in Connection Establishment

Connection establishment latency is one of the most visible user-facing metrics in mobile applications. Even minor degradations in the tail can manifest as “random slowness” in the field. This case study illustrates how a regression that would have passed traditional median-based validation was detected through distributional analysis and phase-level attribution.

Fifty repeated trials were collected for both the baseline build and the candidate OS build. At first glance, the results appeared healthy. The median establishment time remained nearly identical between builds. Under conventional gating logic, the release would have passed. However, deeper distributional analysis revealed a significant divergence beginning at the 80th percentile. While the median remained

stable, the P95 latency increased by approximately 40%, and P99 exhibited similar amplification. The cumulative distribution functions (CDFs) clearly separated in the tail region, indicating a systematic shift rather than random noise. Network benchmarking guidance emphasizes comprehensive distributional reporting rather than reliance on single summary statistics [10], and this case demonstrated precisely why.

Anomaly detection further flagged the shift as statistically significant despite the median stability [8]. This early warning prompted phase-level investigation using Apple's establishment reporting and connection metrics APIs [6], [7]. Decomposition of the establishment duration revealed that TLS handshake and transfer times were unchanged. Instead, the regression was concentrated in the DNS/resolution phase,

where variance had increased substantially. This phase attribution dramatically narrowed the root-cause search space. Engineering analysis uncovered changes in resolver prioritization logic in the candidate build that increased fallback retries under specific timing conditions. After corrective adjustment, tail latency returned to baseline levels, and revalidation confirmed normalization of P95 and P99 distributions.

This case study highlights three important lessons. First, median stability does not guarantee user-perceived stability. Second, path-aware stratification is essential for valid comparison [4], [5]. Third, regression detection must be coupled with phase-level observability to remain actionable [6], [7]. Without this layered instrumentation, the issue would likely have escaped into production.

Case Study 1: Tail-Latency Shift (Conceptual CDF)

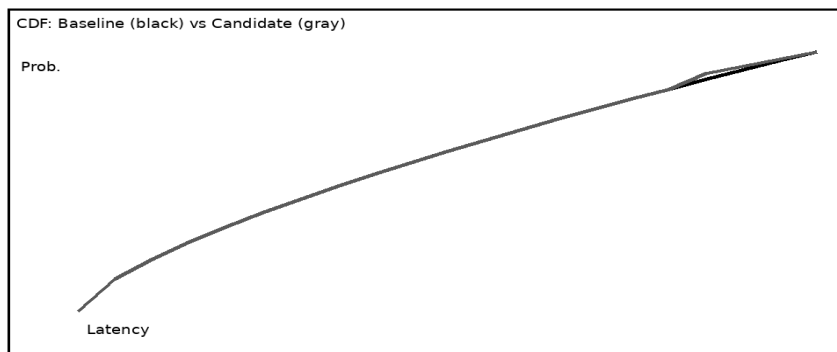


Figure 11. Case Study 1 Conceptual CDF Showing Divergence in Tail Percentiles

8.2. Case Study 2: Invisible Instability Transition-Induced Performance Drift

Mobile devices operate in inherently dynamic environments. Users move between Wi-Fi and cellular networks, often within seconds. While steady-state performance may appear stable, transitions can introduce subtle regressions that are difficult to detect through deterministic gating.

To evaluate transition resilience, the device was subjected to repeated Wi-Fi drop → cellular takeover → Wi-Fi restoration cycles. The workload consisted of a streaming start request combined with periodic API polling to simulate a realistic application session. Mild jitter and a 1% packet loss profile were introduced, and VPN was enabled to emulate enterprise environments. NWPathMonitor was used to timestamp and classify path transitions [4], ensuring that metrics could be correlated precisely with connectivity state changes. Each transition sequence was repeated forty times to capture distributional variability.

Under predefined service-level objectives (SLOs), the candidate build passed. Recovery time remained within threshold, error rates were below acceptable limits, and no request timeouts were recorded. A traditional gating system would have marked the build as safe. However, anomaly detection identified transient spikes in error rates and

increased variance in request latency during transition windows [8]. Although these deviations did not exceed static thresholds, their distribution shifted relative to the baseline build.

The reinforcement learning based scenario planner subsequently increased sampling density around this transition profile, recognizing elevated anomaly probability. This adaptive exploration confirmed reproducibility of the transient instability. Using establishment metrics [6], [7] and path-aware tagging [4], [5], engineers determined that during VPN-enabled transitions, path re-resolution introduced short-lived DNS retry amplification. Connection reuse logic was temporarily invalidated, increasing tail latency during the transition phase.

Given that the regression was transient and below hard SLO thresholds, the governance decision was not to block release but to proceed with a controlled canary deployment. Rollback triggers were configured based on live KPI monitoring, consistent with CI/CD progressive delivery practices [9].

Field telemetry confirmed the issue was minor but measurable. A targeted fix was scheduled for a subsequent release without disrupting delivery cadence. This case study illustrates that deterministic gating alone cannot capture

distributional instability during dynamic conditions. Hybrid detection combining SLO enforcement and anomaly modeling provides both safety and agility.

Case Study 2: Wi-Fi↔Cellular Transition Timeline

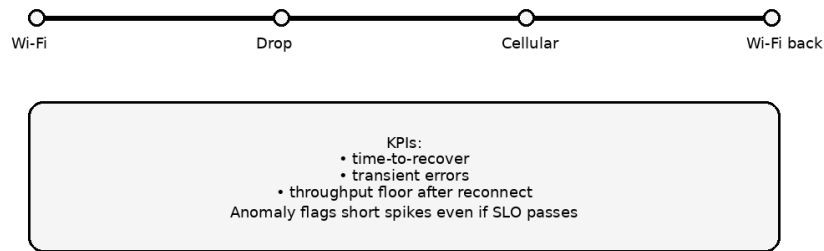


Figure 12. Case Study 2 Transition Timeline Showing Transient Errors and Recovery Behavior

8.3. Cross-Case Insights

Together, these case studies demonstrate that network regressions manifest in different forms:

- Sustained tail amplification under stable conditions
- Transient instability during mobility events

Both require:

- Repeated trials
- Distributional comparison
- Path-aware stratification
- Phase-level decomposition

More importantly, they show that validation must be tightly integrated with governance mechanisms. In one case, the correct decision was to block release. In the other, the appropriate response was controlled exposure via canary rollout [9]. The system's value lies not only in detection, but in enabling proportional, evidence-based release decisions.

9. Discussion and Limitations

The proposed approach improves coverage and reduces brittleness by combining Apple-native measurement and observability [1]–[7], benchmarking discipline [10], [11], throughput anchors [17], [18], and adaptive AI for selection and detection [8], [16]. Path-aware stratification is essential for credible gating, and stage-level attribution accelerates triage. Limitations include restricted access to certain radio-layer internals in third-party harnesses and the gap between lab models and field variability. MetricKit can complement lab tests by providing system-collected, longer-horizon metrics [13]–[15]. Finally, RL policies can overfit to historically high-yield failures; coverage rewards and periodic exploration mitigate this.

10. Conclusion

This paper presented an AI-driven automation framework for Apple device network performance validation integrating

RL-driven scenario exploration, constraints for validity, XCTest performance measurement [1]–[3], Network framework path awareness [4], [5], and establishment reporting/metrics for diagnosis [6], [7]. A hybrid SLO plus anomaly detection strategy provides both explainable gates and early warning signals aligned with anomaly detection literature [8]. CI/CD integration aligns with canary and rollback governance patterns [9].

References

- [1] Apple Developer Documentation, “measureMetrics(_:automaticallyStartMeasuring:for:),” XCTestCase API Reference.
- [2] Apple Developer Documentation, “Performance Tests,” XCTest Documentation.
- [3] Apple Developer Documentation, “Writing and running performance tests,” Xcode Documentation.
- [4] Apple Developer Documentation, “NWPathMonitor,” Network Framework Documentation.
- [5] Apple Developer Documentation, “NWPath,” Network Framework Documentation.
- [6] Apple Developer Documentation, “NWConnection.EstablishmentReport,” Network Documentation.
- [7] Apple Developer Documentation, “Collecting Network Connection Metrics,” Network Documentation.
- [8] B. Lindemann et al., “A Survey on Anomaly Detection for Technical Systems using LSTM Networks,” arXiv:2105.13810, 2021.
- [9] S. Palvai and V. Jain, “Developing End-to-End Concourse CI CD Pipelines With Automated Testing, Scanning, Canary Deployments, and Rollback Logic,” International Journal of Emerging Trends in Computer Science and Information Technology, vol. 7, no. 1, pp. 23–29, 2026, doi: 10.63282/3050-9246.IJETCSIT-V7I1P105.

- [10] S. Bradner and J. McQuaid, "Benchmarking Methodology for Network Interconnect Devices," IETF RFC 2544, 1999.
- [11] A. Morton et al., "Framework for TCP Throughput Testing," IETF RFC 6349, 2011.
- [12] Apple WWDC, "Introducing Network.framework: A modern alternative to sockets," WWDC 2018 Session 715.
- [13] Apple Developer Documentation, "MetricKit," MetricKit Documentation.
- [14] Apple Developer Documentation, "MXMetricManager," MetricKit API Reference.
- [15] Apple Developer Documentation, "MXNetworkTransferMetric," MetricKit API Reference.
- [16] M. Baqar and R. Khanda, "The Future of Software Testing: AI-Powered Test Case Generation and Validation," arXiv:2409.05808, 2024.
- [17] iPerf Project, "iPerf3," 2025.
- [18] ESnet, "iperf3 Documentation," 2025.
- [19] VIAVI Solutions, "RFC 6349 Testing – TrueSpeed Application Notes," 2025.
- [20] Apple Developer Documentation, "NWConnection," Network Documentation.