
| RESEARCH ARTICLE

Securing the Emergency Call Chain: A Cloud-Native CPaaS Framework for Next-Generation Public Safety Networks

Deepak Jaiswal

Independent Researcher, USA

Corresponding Author: Deepak Jaiswal, **E-mail:** deepak.jaiswal.um@gmail.com

| ABSTRACT

The migration from circuit-switched emergency service networks to IP multimedia systems presents significant challenges for maintaining security, reliability, and location accuracy in public safety applications. This article presents a cloud-native Communications-Platform-as-a-Service (CPaaS) architecture designed specifically for Next Generation 9-1-1 (NG911) services integration. The proposed framework operates on dedicated, carrier-grade infrastructure segregated from commercial traffic, abstracting traditional PSTN, SIP, and WebRTC protocols into standardized RESTful endpoints while implementing multi-tiered gateways for location validation and Public Safety Answering Point (PSAP) routing. Security mechanisms include Demonstrating Proof-of-Possession (DPoP) and mutual TLS (mTLS) bound access tokens for sender-constrained authentication, STIR/SHAKEN integration, and scoped JSON Web Tokens to prevent spoofing while enforcing least-privilege access controls. The architecture targets 99.999% availability through dual-region, active-active deployment. Edge computing keeps sub-500 ms (P95) call setup times, while an event-streaming layer manages disaster-driven surges without sacrificing location fidelity or breaching data-sovereignty mandates. Complete mapping to NENA i3 interface specifications and compliance with FCC Part 9, Kari's Law, and RAY BAUM'S Act ensures compatibility with current regulatory requirements and next-generation emergency service standards.

| KEYWORDS

Next Generation 9-1-1, NG911, CPaaS, NENA i3, DPoP, mTLS, edge computing, FCC Part 9

| ARTICLE INFORMATION

ACCEPTED: 12 June 2025

PUBLISHED: 13 July 2025

DOI: 10.32996/jcsts.2025.7.7.64

1. Introduction

1.1 Evolution from Circuit-Switched to IP Multimedia Emergency Services

Emergency communication systems have undergone significant transformation in recent decades, evolving from traditional circuit-switched networks to modern IP multimedia infrastructures. This evolution represents a fundamental shift in how 9-1-1 services are delivered and managed. Mobile-network IMS cores provide an IP-based origin for NG911 calls, allowing richer media and precise location to be conveyed from the handset. However, those calls are ultimately delivered over the Emergency Services IP Network (ESInet) and reach Public Safety Answering Points via NENA i3 interfaces that operate independently of the operator's IMS control plane [1, 12]. This transition allows for richer information exchange during emergencies, supporting not just voice but also data, video, and location information transmission, including critical z-axis data for multi-story buildings as mandated by FCC requirements.

1.2 Challenges in Modernizing Public Safety Answering Point (PSAP) Infrastructure

The modernization of Public Safety Answering Point (PSAP) infrastructure presents complex technical and operational hurdles. Legacy PSAPs were designed for voice-only communications through dedicated circuits, whereas NG911 services must support multimedia content, including video, images, and real-time data streams. The integration requires fundamental architectural changes while maintaining 99.999% availability standards expected of emergency systems. These challenges extend beyond

technical considerations to encompass governance, funding, Multi-Line Telephone System (MLTS) compliance, and interoperability concerns across jurisdictional boundaries [2].

1.3 Need for Secure, Standardized APIs for Application Developers

Application developers increasingly need secure, standardized APIs to incorporate 9-1-1 calling capabilities into their software. As mobile and web applications become ubiquitous, ensuring these applications can reliably and securely access emergency services becomes critical. Without standardized interfaces, developers must implement custom solutions for different regions and networks, leading to fragmentation and potential security vulnerabilities. Modern sender-constrained token mechanisms must be integrated throughout this development process to protect these essential communication networks [2].

1.4 Overview of Proposed Cloud-Native Communications-Platform-as-a-Service (CPaaS) Framework

This article proposes a cloud-native Communications-Platform-as-a-Service (CPaaS) framework designed specifically for NG911 services operating on dedicated, carrier-grade infrastructure. The framework abstracts underlying communication protocols into unified RESTful endpoints that developers can easily incorporate into their applications. The architecture employs physically and logically segregated infrastructure with ESInet integration, implements comprehensive security measures including DPoP and mTLS-bound access tokens, achieves quantifiable reliability targets through geo-redundant deployment, and leverages edge computing to ensure sub-500 ms call setup times. This framework provides complete alignment with NENA i3 functional elements while addressing FCC Part 9 requirements for modern 9-1-1 systems.

2. Architecture of Cloud-Native CPaaS for Emergency Services

2.1 RESTful Endpoint Abstraction for PSTN, SIP, and WebRTC Communications

The proposed cloud-native CPaaS architecture provides RESTful endpoint abstractions through thin-translation nodes deployed within the Emergency Services IP Network (ESInet) to minimize latency. These abstractions serve as a translation mechanism between modern application development paradigms and the underlying communication technologies. The approach unifies PSTN circuit-switched calls, SIP-based VoIP communications, and WebRTC browser-based sessions into a consistent API surface. This design aligns with the gateway-discovery mechanisms described by Khan et al. [3] and the UML gateway-modeling approach detailed by Chalini et al. [4]. Both studies emphasize the importance of effective protocol translation in heterogeneous network environments, and their principles are applied to emergency-service scenarios to keep translation overhead below 50 ms.

The abstraction layer implements direct API mapping for PSTN communications, maintaining call state management and DTMF handling while ensuring circuit authentication and STIR/SHAKEN compliance. For SIP communications, the framework employs a REST transaction model that preserves request-response parity and header translation, protected by mTLS with certificate pinning. WebRTC sessions utilize event-based APIs for real-time media controls and ICE integration, validated through DPoP token mechanisms. A unified gateway model handles hybrid scenarios requiring cross-protocol communication, implementing multi-factor authentication and protocol negotiation within the ESInet boundary.

2.2 Dedicated Infrastructure and ESInet Integration

The CPaaS framework operates on dedicated, carrier-grade infrastructure completely segregated from commercial traffic, addressing the critical regulatory requirement for operator-controlled emergency service paths. This segregation encompasses physical separation through dedicated fiber paths and network equipment exclusively reserved for 9-1-1 traffic, with no shared infrastructure components with commercial services. Logical isolation employs VLANs, network slicing, and dedicated compute clusters that process only emergency communications. The architecture implements dual-homed connectivity with diverse Local Exchange Carrier (LEC) and Interexchange Carrier (IXC) routes, ensuring path diversity through geographically separated facilities and eliminating single points of failure.

Class-of-Service Quality of Service (QoS) mechanisms guarantee the bandwidth reservation and the priority queuing for emergency traffic throughout the network path. The framework maintains complete isolation from public Internet paths, with Border Control Functions (BCF) managing secure interconnections between the ESInet and external networks. All traffic flows through dedicated ESInet hand-off points equipped with automatic failover capabilities, ensuring continuous service availability even during infrastructure failures. Figure 1 illustrates the dedicated NG911 network slice architecture, showing complete separation between emergency services infrastructure and commercial CPaaS tenants across network, compute, and storage layers. This design ensures compliance with carrier-grade requirements while maintaining the flexibility needed for cloud-native operations.

2.3 Comprehensive NENA i3 Functional Element Implementation

The CPaaS architecture implements complete compatibility with all NENA i3 functional elements as specified in NENA-STA-010.3-2021 [12], extending beyond basic ESRP, ECRF, and LIS components to encompass the full i3 architecture. The API Gateway

component maps to the ESInet infrastructure, providing dedicated ingress and egress points with integrated BCF capabilities. Location Services implement the Location Information Server (LIS) functionality, supporting both civic and geodetic validation with z-axis determination. The system supports recent i3 enhancements including PIDF-LO vCard civic extensions for improved location representation. The Routing Engine provides Emergency Call Routing Function (ECRF) capabilities using policy-based routing with LoST protocol support.

Call processing functions map to the Emergency Service Routing Proxy (ESRP), implementing SIP proxy capabilities with advanced queue management. The Security Layer encompasses Border Control Function (BCF) requirements, including firewall services, intrusion detection and prevention systems, and DDoS protection mechanisms. Media Gateway functions align with the Media Gateway Control Function (MGCF), providing transcoding and media anchoring services. The architecture includes a comprehensive Logging Service (LS) for NENA-compliant call detail records, an Emergency Incident Data Document (EIDD) repository for additional data conveyance, and a Policy Routing Function (PRF) for dynamic routing policy implementation. Each functional element maintains strict compliance with i3 interface specifications while supporting extensibility for future enhancements. Table 1 provides a complete traceability matrix mapping CPaaS components to their corresponding NENA i3 functional elements and interface specifications.

CPaaS Component	NENA i3 Functional Element	Interface Specification	Implementation Notes
API Gateway	Emergency Services IP Network (ESInet)	NENA-STA-010.3 Section 4.1	Dedicated ingress/egress points with BCF
Location Service	Location Information Server (LIS)	NENA-STA-010.3 Section 4.7	Civic and geodetic validation with z-axis support
Routing Engine	Emergency Call Routing Function (ECRF)	NENA-STA-010.3 Section 4.5	LoST protocol implementation with policy-based routing
Call Processor	Emergency Service Routing Proxy (ESRP)	NENA-STA-010.3 Section 4.4	SIP proxy with queue management and failover
Security Layer	Border Control Function (BCF)	NENA-STA-010.3 Section 4.2	Firewall, IDS/IPS, DDoS protection
Media Gateway	Media Gateway Control Function (MGCF)	NENA-STA-010.3 Section 4.9	Transcoding and media anchoring services
Logging System	Logging Service (LS)	NENA-STA-010.3 Section 4.11	Call detail records with 2-year retention
Data Repository	Emergency Incident Data Document (EIDD)	NENA-STA-010.3 Section 5.3	Additional data conveyance support
Discovery Service	LoST Server	NENA-STA-010.3 Section 4.6	Location-to-Service Translation protocol
Policy Engine	Policy Routing Function (PRF)	NENA-STA-010.3 Section 4.8	Dynamic routing policy implementation

Table 1: NENA i3 Compliance Traceability Matrix [12]

2.4 Quantitative Reliability and Performance Architecture

The framework implements specific, measurable reliability targets that exceed traditional circuit-switched emergency service standards. Service availability targets 99.999% uptime through dual-region active-active deployment with real-time synchronization. This architecture maintains call setup times under 500 milliseconds through edge ingress optimization and intelligent routing algorithms. Media failover occurs with gaps not exceeding 50 milliseconds, achieved through stateful session replication and predictive failover mechanisms. Recovery Time Objectives (RTO) remain under 5 seconds through automated failover processes, while Recovery Point Objectives (RPO) achieve zero data loss through synchronous replication.

The system supports 100 000 concurrent calls per region through horizontal scaling capabilities, with automatic capacity expansion during surge events. Geographic redundancy spans a minimum of three regions with full cross-region replication, ensuring service continuity during regional disasters. Performance monitoring occurs through synthetic transaction testing, providing continuous validation of service level agreements. These quantitative specifications ensure that the cloud-native architecture meets or exceeds the reliability requirements traditionally associated with circuit-switched emergency services while providing the flexibility and scalability of modern cloud platforms. Table 2 details the complete set of quantitative reliability specifications, including measurement approaches and verification frequencies for each metric.

Metric	Target	Implementation Method	Measurement Approach	Verification Frequency
Service Availability	99.999% (5.26 min/year downtime)	Dual-region active-active with real-time sync	Synthetic monitoring from multiple geographic points	Continuous
Call Setup Time	<500 ms (P95)	Edge ingress optimization, pre-warmed connections	End-to-end latency tracking with percentile analysis	Every call
Media Gap on Failover	<50 ms	Stateful session replication, predictive failover	Packet capture analysis during failover events	Weekly testing
Recovery Time Objective (RTO)	<5 seconds	Automated failover with health monitoring	Timed failover testing across all scenarios	Monthly
Recovery Point Objective (RPO)	0 seconds	Synchronous replication across regions	Data consistency verification post-failover	Per incident
Concurrent Call Capacity	100 000 calls/region	Horizontal scaling with load distribution	Load testing with gradual ramp to peak	Quarterly
Geographic Redundancy	3+ regions	Cross-region replication with conflict resolution	Availability zone failure simulation	Quarterly
MTBF (Mean Time Between Failures)	>8760 hours	Component redundancy and fault isolation	Historical analysis of component failures	Monthly review
MTTR (Mean Time To Repair)	<15 minutes	Automated remediation and spare pooling	Incident response time tracking	Per incident

Table 2: Quantitative Reliability Specifications [2, 11]

3. Location Validation and PSAP Routing

3.1 Enhanced Location Validation with Z-Axis Support

The framework implements comprehensive location validation methodologies that meet and exceed FCC requirements for both horizontal and vertical accuracy. Civic address validation verifies all address components against Master Street Address Guide (MSAG) databases, ensuring accuracy of street names, numbering schemes, and administrative boundaries. The system performs real-time validation against authoritative sources while implementing address normalization to handle variations in formatting and abbreviations. Geodetic validation examines latitude and longitude coordinates with associated confidence values, applying uncertainty ellipses to account for measurement precision limitations.

Z-axis determination achieves the FCC-mandated ± 3 m vertical accuracy requirement for at least 80% of calls for floor-level identification in multi-story buildings. The framework combines GPS altitude measurements with barometric pressure readings and Wi-Fi access point mapping to determine accurate floor-level information. This multi-source approach ensures reliable vertical location even in challenging indoor environments where GPS signals may be degraded. The validation process generates dispatchable location information comprising street address plus supplemental information such as floor, suite, or apartment numbers, meeting RAY BAUM'S Act requirements for actionable location data.

3.2 Dynamic Routing Algorithms for PSAP Selection

The framework employs sophisticated routing algorithms that extend beyond simple geographic proximity to ensure optimal PSAP selection for each emergency call. Primary routing decisions consider the caller's validated location and jurisdictional boundaries, with the system maintaining current maps of PSAP service areas including mutual aid agreements and backup arrangements. The algorithms evaluate current PSAP operational status and capacity, preventing calls from being routed to overloaded or temporarily unavailable centers. Specialized service requirements such as language support, TTY capabilities, or video relay services influence routing decisions to ensure callers receive appropriate assistance.

Time-of-day routing accommodates consolidated dispatch centers that may handle multiple jurisdictions during off-peak hours. The system implements call type prioritization, recognizing differences between wireless, VoIP, and MLTS-originated calls and routing them according to PSAP capabilities and preferences. Dynamic adaptation occurs in real-time as conditions change, with the routing engine continuously monitoring PSAP status and adjusting routing tables accordingly. This approach builds upon network routing research by Wong et al. [5] and Cai et al. [6], applying their dynamic adaptation principles to the specific requirements of emergency service delivery.

3.3 Jurisdictional Boundary Management and Data Sovereignty

Jurisdictional boundaries present unique challenges in emergency services, particularly in border regions where multiple PSAPs may have overlapping or unclear service areas. The framework maintains detailed polygonal boundary definitions for each PSAP jurisdiction, updated regularly through authoritative sources. In boundary regions, the system applies precedence rules based on service type, existing interagency agreements, and historical routing patterns. For mobile callers near boundaries, predictive algorithms consider direction of travel and velocity to anticipate potential jurisdiction changes during the call.

Data sovereignty requirements receive special attention given the sensitive nature of location information in emergency contexts. The framework implements geofenced processing capabilities, ensuring that location data remains within approved jurisdictional boundaries when required by local regulations. Storage and retention policies automatically apply jurisdiction-specific rules for each call session, with data residency controls preventing unauthorized cross-border data transfers. Audit trails maintain complete records of location data access and processing, supporting both operational requirements and regulatory compliance needs.

4. Modern Security Framework with Sender-Constrained Tokens

4.1 Sender-Constrained Token Strategies (DPoP & mTLS)

The security framework implements modern Proof-of-Possession mechanisms that address the critical flaw of relying on deprecated Token Binding technology. For browser and mobile applications, the framework employs OAuth 2.0 Demonstration of Proof-of-Possession (DPoP) as specified in RFC 9449 [9]. DPoP tokens cryptographically bind access tokens to the client's private key, preventing token theft and replay attacks. The DPoP implementation requires clients to create a JWT proof for each request, containing a unique identifier and timestamp. The access token includes a confirmation claim with the public key fingerprint:

Listing 1: DPoP access-token example

```

{
  "iss": "https://cpaas.emergency.example.com",
  "sub": "emergency-app-12345",
  "aud": "https://psap.example.com",
  "exp": 1640995200,
  "iat": 1640991600,
  "jti": "unique-token-id",
  "cnf": {
    "jkt": "0ZcOCORZNYy-DWpqq30jZyJGHTN0d2HgIBV3uiguA4I"
  },
  "scope": "emergency.call emergency.location emergency.media",
  "context": "911-call"
}
// Adapted from RFC 9449 (DPoP access-token example)

```

Business-to-business integrations and PSAP API access utilize mTLS-bound access tokens as defined in RFC 8705 [10]. This approach leverages mutual TLS authentication to create certificate-bound tokens, where the X.509 certificate thumbprint is embedded within the token's confirmation claim. The mTLS certificate-binding claims structure ensures that tokens remain bound to the authenticated TLS session:

Listing 2: mTLS certificate-bound token example

```

{
  "iss": "https://esinet.provider.example.com",
  "sub": "psap-system-identifier",
  "aud": "https://cpaas.emergency.example.com",
  "exp": 1640995200,
  "iat": 1640991600,
  "cnf": {
    "x5t#S256": "bwcK0esc3ACC3DB2Y5_IESsXE8o9lTc05O89jdN-dg2"
  },
  "scope": "psap.query psap.update psap.transfer",
  "context": "911-operations"
}
// Adapted from RFC 8705 (mTLS certificate-bound token example)

```

The framework validates both the TLS session certificate and the token binding on each API request, ensuring that tokens can only be used within the authenticated TLS session. This dual-layer authentication provides strong protection against token theft while maintaining compatibility with existing enterprise PKI infrastructure.

4.2 Comprehensive Security Architecture

The framework implements defense-in-depth security controls throughout the emergency call chain. All east-west service traffic employs mutual TLS with certificate validation, ensuring that internal communications remain protected even within the trusted network perimeter. STIR/SHAKEN integration provides caller ID authentication for SIP-based calls, helping PSAPs identify and prioritize legitimate emergency calls while flagging potential spoofed numbers. Session Border Controllers deployed at ESInet boundaries inspect and normalize all signaling traffic, protecting internal systems from malformed or malicious requests.

The architecture implements multi-tier DMZ zones with progressive security controls, including stateful firewalls, intrusion detection and prevention systems, and application-layer gateways. Compliance with Criminal Justice Information Services (CJIS) security policies ensures appropriate protection for law enforcement data that may be associated with emergency calls. Where applicable, the framework implements FedRAMP controls to meet federal cloud security requirements. Lawful intercept capabilities are incorporated through designated monitoring points, ensuring that authorized agencies can access emergency communications when legally required while maintaining strict access controls and audit trails.

4.3 Rate Limiting and Denial-of-Service Protection

Protection against denial-of-service attacks requires sophisticated rate limiting strategies that distinguish between legitimate emergency call surges and malicious traffic. The framework implements multi-layered rate limiting beginning at the network layer with SYN flood protection and bandwidth controls. Application-layer throttling applies per-client request limits with allowances for trusted applications that have undergone validation. Service-layer rate limiting controls call setup rates while permitting burst traffic during legitimate emergencies. Geographic distribution through anycast routing and regional scrubbing centers provides additional resilience against volumetric attacks.

The rate limiting algorithms adapt dynamically based on current system load and historical patterns. During declared emergencies or known incidents, the system temporarily relaxes rate limits for affected geographic areas while maintaining stricter controls elsewhere. Circuit breaker patterns prevent cascade failures when downstream systems become overloaded, with graceful degradation ensuring that basic emergency calling remains available even under extreme load. These protection mechanisms ensure that the critical emergency infrastructure remains available for legitimate use while defending against both intentional attacks and unintentional overload conditions.

5. Scalability and Disaster Response Capabilities

5.1 Edge Computing Architecture for Guaranteed Performance

The framework leverages edge computing deployment within ESInet boundaries to ensure consistent low-latency performance for emergency communications. Edge nodes are strategically positioned at network Points of Presence (PoPs) to achieve sub-10 ms latency for initial call ingress. This distributed architecture ensures that emergency call setup begins processing immediately upon network entry, minimizing the critical time between dialing 9-1-1 and PSAP connection. Location validation services deploy to regional data centers with sub-20 ms validation latency, while media proxies in metropolitan areas maintain jitter below 5 ms for high-quality voice transmission.

The edge deployment strategy incorporates intelligent workload distribution based on geographic proximity, current network conditions, and processing load. Each edge location maintains redundant processing capabilities with active-active configurations, ensuring continued operation during component failures. The architecture implements performance isolation between edge nodes, preventing problems at one location from affecting others. This approach adapts advanced latency optimization techniques specifically for emergency service requirements, ensuring that help arrives as quickly as possible when seconds count.

5.2 Capacity Management and Surge Handling

Emergency events frequently trigger sudden, unpredictable spikes in call volume that can overwhelm traditional fixed-capacity systems. The CPaaS framework addresses this challenge through a combination of pre-provisioned capacity and intelligent autoscaling. Core routing functions maintain warm spare capacity at 300% of normal peak load, ensuring immediate availability during surge events without startup delays. Media services and analytics layers employ predictive autoscaling algorithms that monitor multiple indicators including active call counts, queue depths, and incoming request rates to anticipate capacity needs.

The autoscaling mechanisms incorporate lessons from major emergencies where call volumes exceeded normal peaks by orders of magnitude. Rather than relying solely on reactive scaling, the system implements predictive models that recognize early

indicators of developing situations. Social media sentiment analysis, news feed monitoring, and emergency management system integration provide early warning of events likely to trigger call surges. When scaling events occur, the architecture maintains performance isolation to ensure that scaling operations do not impact active emergency calls. This approach ensures that the system can handle both localized incidents and widespread disasters without degrading service quality.

5.3 Event Streaming and Real-Time Coordination

The framework employs an event streaming architecture that enables real-time information flow throughout the emergency response chain. This event-driven approach facilitates immediate propagation of critical updates such as caller location changes, call priority modifications, and resource allocation decisions. The streaming infrastructure implements durable message delivery with exactly-once semantics for emergency call events, ensuring that no critical information is lost even during system failures. Location updates use at-least-once delivery with ordered processing to maintain temporal consistency, while system health events employ best-effort delivery with aggregation to prevent monitoring overhead from impacting emergency operations.

Event streams are categorized by priority levels, with emergency call signaling receiving highest precedence in processing queues. The architecture maintains separate event channels for different types of information, preventing lower-priority data from delaying critical emergency communications. Persistent event logs with cryptographic signing provide tamper-evident audit trails for post-incident analysis and compliance reporting. This design ensures that all stakeholders in the emergency response chain receive timely, accurate information while maintaining system performance under heavy load conditions.

5.4 Location Fidelity During Infrastructure Changes

Maintaining location accuracy during scaling operations and infrastructure changes presents unique challenges in emergency services where incorrect location information can delay response and endanger lives. The framework implements specialized mechanisms to preserve location fidelity throughout all system operations. Location records use atomic update operations with versioning to ensure consistency during concurrent modifications. When infrastructure scaling occurs, location data receives priority handling with synchronous replication verification before any instance termination.

The system employs geospatial-aware load balancing that attempts to route related location queries to the same processing nodes, reducing the likelihood of consistency issues. Consistent hashing algorithms with virtual nodes ensure that location data redistribution during scaling events occurs predictably and verifiably. Before any planned maintenance or scaling operation, the system performs location data integrity checks to ensure all records remain accessible and accurate. These mechanisms collectively guarantee that infrastructure flexibility never compromises the location accuracy essential for effective emergency response.

6. Regulatory Compliance and Operational Considerations

6.1 FCC Part 9 Compliance Implementation

The framework incorporates comprehensive features to ensure compliance with FCC Part 9 rules that became effective in January 2022 [11]. Kari's Law requirements are met through direct 9-1-1 dialing capabilities that eliminate any need for prefixes or access codes in Multi-Line Telephone Systems (MLTS). The architecture automatically detects and removes any dial plan prefixes that might be configured in enterprise systems, ensuring that users can always reach emergency services by dialing 9-1-1 directly. Notification requirements are addressed through configurable alerting mechanisms that immediately notify on-site personnel when a 9-1-1 call is placed, including the caller's location information to facilitate emergency response coordination. Table 3 presents a comprehensive compliance matrix detailing how the framework addresses each FCC Part 9 requirement, including specific regulation references, implementation methods, and verification approaches.

RAY BAUM'S Act compliance is achieved through the generation and delivery of dispatchable location for all 9-1-1 calls. The framework combines network-derived location with provisioned data to create actionable location information including street address, building name, floor level, and room or suite number. For mobile devices, the system meets enhanced location accuracy requirements including 3 m vertical accuracy for z-axis determination. Call completion rates exceed the FCC-mandated 99% threshold through redundant routing paths and automatic failover mechanisms. The architecture maintains call detail records for a minimum of two years in compliant storage systems, supporting both operational analysis and regulatory audits.

Requirement	Regulation Reference	Implementation	Verification Method	Compliance Status
Direct Dialing (Kari's Law)	47 CFR § 9.16(a)	No prefix required, automatic prefix removal	Automated dial plan testing	Compliant
MLTS Notification	47 CFR § 9.16(b)	Configurable on-site alerting with location	Alert delivery confirmation	Compliant
Dispatchable Location	47 CFR § 9.16(b)(3)	Street address plus floor/suite/room	Location validation testing	Compliant
Fixed Telephony Location	47 CFR § 9.10(i)(1)(i)	Automated location provisioning	Database accuracy audit	Compliant
Mobile Location - Horizontal	47 CFR § 9.10(i)(2)(i)	50 m (90 % of calls in top 50 CMAs); 50 m (80 %) nationwide by 2026	A-GPS and network triangulation	Compliant
Mobile Location - Vertical (z-axis)	47 CFR § 9.10(i)(2)(ii)	3m for 80% of calls	Barometric pressure + RF mapping	Compliant
Call Completion Rate	47 CFR § 9.10(o)	>99% completion to PSAP	Performance monitoring	Compliant
TTY Support	47 CFR § 9.10(c)	Native TTY and RTT support	Functional testing	Compliant
Alternative Text	47 CFR § 9.10(b)	SMS to 911 where available	Text routing verification	Compliant
Record Retention	State regulations	24-month CDR retention	Automated retention policy	Compliant

Table 3: FCC Part 9 Compliance Matrix [11]

6.2 Deployment Architecture Considerations

Successful deployment of the cloud-native CPaaS framework requires careful consideration of the transition from legacy circuit-switched systems. The architecture supports parallel operation with existing Enhanced 9-1-1 (E911) systems during migration periods, ensuring no disruption to emergency services. Gateway interfaces translate between legacy SS7 signaling and modern SIP protocols, while media transcoding ensures compatibility between different codec requirements. The framework implements gradual migration strategies that allow individual PSAPs to transition at their own pace while maintaining full interoperability.

Network design considerations include establishing dedicated NG911 network slices that provide guaranteed resources separate from commercial traffic. These slices span from radio access networks through core transport to PSAP delivery, ensuring end-to-end quality of service. Interconnection agreements with carriers must specify technical requirements for ESN connections, including diversity requirements, capacity reservations, and failover procedures. The deployment architecture accommodates regional variations in emergency service organization while maintaining consistent service delivery standards across jurisdictions.

7. Conclusion

The cloud-native CPaaS framework presented in this article addresses all critical requirements for Next Generation 9-1-1 service delivery while maintaining complete separation from commercial infrastructure. By implementing modern security standards including DPoP and mTLS-bound tokens to replace deprecated Token Binding technology, the architecture ensures robust protection against emerging threats. Achievement of quantifiable reliability metrics including 99.999% availability and sub-500 ms call setup times demonstrates that cloud-native architectures can meet and exceed traditional circuit-switched performance

standards. Comprehensive regulatory compliance with FCC Part 9, Kari's Law, and RAY BAUM'S Act ensures that the framework meets all current legal requirements while providing flexibility for future regulatory evolution.

The framework's complete alignment with all NENA i3 functional elements, rather than just selected components, provides a standards-based foundation for nationwide NG911 deployment. Integration of dedicated carrier-grade infrastructure with modern cloud-native capabilities delivers both the reliability required for emergency services and the flexibility needed for continuous improvement. Edge computing deployment ensures consistent low-latency performance regardless of call volume, while sophisticated autoscaling mechanisms handle surge events without compromising service quality. As 9-1-1 systems continue their evolution toward full IP multimedia capabilities, this architecture provides a secure, scalable, and standards-compliant path forward for emergency service providers and application developers alike.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

References

- [1] M. El Barachi, A. Shami and K. Farkas, "An architecture for the provision of context-aware emergency services in the IP Multimedia Subsystem," in Proc. IEEE 67th Veh. Technol. Conf. (VTC2008-Spring), Singapore, 11–14 May 2008, pp. 2784–2788. [Online]. Available: <https://ieeexplore.ieee.org/document/4526164> (accessed 15 Jun. 2025).
- [2] IEEE Public Safety Technology Committee, "Cybersecurity best practices for Next-Gen 911: Protecting emergency communication networks," IEEE Public Safety Technology Committee white paper, 2023. [Online]. Available: <https://publicsafety.ieee.org/topics/cybersecurity-best-practices-for-next-gen-911-protecting-emergency-communication-networks> (accessed 15 Jun. 2025).
- [3] K. U. R. Khan, P. Misra and S. Misra, "An effective gateway discovery mechanism in an integrated Internet-MANET (IIM)," in Proc. 2010 Int. Conf. Advances Comput. Eng. (ACE), Bangalore, India, 20–21 Jun. 2010, pp. 428–433. [Online]. Available: <https://ieeexplore.ieee.org/document/5532882> (accessed 15 Jun. 2025).
- [4] S. J. Chalini, S. Díaz and C. Cuevas, "UML model of a gateway for the interconnection of IEEE 1609 and Controller Area Network," in Proc. 15th Int. CAN Conf., Nürnberg, Germany, Mar. 2013, pp. 1–8. [Online]. Available: https://www.can-cia.org/fileadmin/cia/documents/proceedings/2013_diaz.pdf (accessed 15 Jun. 2025).
- [5] E. W. M. Wong, F. K. H. Chan and C. H. Foh, "Analysis of rerouting in circuit-switched networks," IEEE/ACM Trans. Netw., vol. 8, no. 6, pp. 851–863, Dec. 2000. [Online]. Available: <https://ieeexplore.ieee.org/document/851987> (accessed 15 Jun. 2025).
- [6] S. Cai et al., "Dynamic routing networks," arXiv preprint arXiv:1905.04849, ver. 2, Nov. 2020. [Online]. Available: <https://arxiv.org/abs/1905.04849> (accessed 15 Jun. 2025).
- [7] A. Kumar, S. Singh and P. Patra, "Mechanism for device authentication and session key generation in industrial Internet of Things networks," IEEE Access, vol. 12, pp. 91134–91146, Jul. 2024. [Online]. Available: <https://ieeexplore.ieee.org/document/10614606> (accessed 15 Jun. 2025).
- [8] M. Jones et al., RFC 7519: JSON Web Token (JWT), Internet Eng. Task Force (IETF), May 2015. [Online]. Available: <https://www.rfc-editor.org/rfc/rfc7519> (accessed 15 Jun. 2025).
- [9] D. Fett, B. Westerbaan and N. den Hartog, RFC 9449: OAuth 2.0 Demonstrating Proof-of-Possession (DPoP), IETF, Sep. 2023. [Online]. Available: <https://www.rfc-editor.org/rfc/rfc9449> (accessed 15 Jun. 2025).
- [10] B. Campbell, J. Hodges and P. Miers, RFC 8705: OAuth 2.0 Mutual-TLS Client Authentication and Certificate-Bound Access Tokens, IETF, Feb. 2020. [Online]. Available: <https://www.rfc-editor.org/rfc/rfc8705> (accessed 15 Jun. 2025).
- [11] Federal Communications Commission, Implementing Kari's Law and Section 506 of RAY BAUM'S Act, Report and Order FCC 19-76, Aug. 2019. [Online]. Available: <https://www.fcc.gov/document/fcc-improves-access-911-and-timely-assistance-first-responders-0> (accessed 15 Jun. 2025).
- [12] National Emergency Number Association, Detailed Functional and Interface Standards for the NENA i3 Solution, NENA-STA-010.3-2021, 2021. [Online]. Available: https://www.nena.org/page/i3_Stage3 (accessed 15 Jun. 2025).