Beyond the Cloud: Quantifying Edge Computing's Real-Time Gains in Healthcare, Industry, and Autonomous Mobility

Prathmesh S Kolte¹, Palak Rajput², Kalyani Gond³

^{1,3}. Student, Bachelor of Engineering, Information and Technology, Mauli Group of Institution's College of Engineering and Technology, Shegaon

²Student, Electronics and Telecommunication, Bharati Vidyapeeth College of Engineering for Womens,

Pune

Abstract-Edge computing extends computation and storage from remote clouds to the network periphery, enabling data to be processed where it is produced. This paper provides an in-depth examination of edge computing's theoretical foundations, lavered architecture and performance advantages, supported by three domain-diverse case studies-remote patient monitoring, industrial predictive maintenance and autonomous vehicles. Quantitative results show latency reductions of up to 90 %, 40 % declines in unplanned downtime and a 15× improvement in vehicle reaction times. A comparative discussion highlights recurring benefits (low latency, bandwidth relief, privacy enhancement) and persistent barriers (resource constraints, heterogeneity, large-scale orchestration). Finally, we propose future research directions, including lightweight AI models, standards-driven interoperability and cloud-edge co-design methodologies. The findings demonstrate that edge computing can deliver transformative value across healthcare, manufacturing and transportation, provided that open challenges are addressed through coordinated research and engineering effort.

 Index
 Terms—Edge
 computing
 · Latency

 · Internet of Things
 · Industrial IoT
 · Autonomous

 vehicles
 · Predictive maintenance
 · Real-time analytics

I. INTRODUCTION

The proliferation of the Internet of Things (IoT), real-time analytics and bandwidth-hungry applications has exposed inherent deficiencies of the centralised cloud model—chiefly high round-trip latency, privacy concerns and network congestion. Edge computing mitigates these issues by relocating computation, storage and intelligence from distant data centres to resources "one network hop" away from data sources [1]. Figure 1 illustrates the continuum from central cloud to fog to edge, emphasising that the edge layer resides closest to sensors and actuators, thereby shortening the control loop.

Edge's ability to process data in situ unlocks new classes of latency-critical workloads—robotic control loops (<10 ms), telesurgery (<5 ms) and collaborative autonomous driving (<20 ms)—that are impractical over long-haul networks. This paper analyses edge computing holistically: Section 2 synthesises related work; Section 3 details the research methodology; Section 4 formalises a reference architecture; Section 5 presents and quantifies three real-world case studies; Section 6 discusses cross-cutting insights and unresolved challenges; Section 7 suggests future research trajectories; and Section 8 concludes.



Figure 1: Cloud-Fog-Edge Continuum.

Three-tier latency pyramid showing distance from data source (bottom: device; middle: edge; top: cloud) (Source: geeksforgeeks)

II. BACKGROUND AND RELATED WORK

Shi *et al.* [1] introduced the seminal edge computing vision, framing it as a resource-rich successor to mobile cloud offload. Subsequent studies contrasted edge and fog paradigms [2] and surveyed security models [3]. Industrial consortia (ETSI MEC, LF Edge) have published reference stacks, yet empirical, cross-sector performance evidence remains scarce. Prior papers often isolate a single domain or overlook deployment-scale orchestration. By triangulating results from healthcare, manufacturing and transportation, this study provides a broader empirical foundation.

III. METHODOLOGY

We performed a systematic literature survey (2018–2024) across IEEE Xplore, ACM DL and SpringerLink using the query "*edge computing*" *AND* "*deployment*" *AND* "*performance*". 213 records were screened; 37 met inclusion criteria (detailed quantitative results plus architectural transparency). Three high-maturity deployments—remote patient monitoring, industrial predictive maintenance and autonomous vehicles—were selected as case studies because they:

- 1. Operate in safety-critical contexts,
- 2. Provide pre- and post-deployment metrics, and
- 3. Represent distinct latency envelopes (10–200 ms).

We reproduced published benchmarking procedures (e.g., round-trip ping, throughput counters, fault logs) on open-access datasets or vendor white-papers to verify reported gains. Performance gains are normalised to each domain's cloud-baseline.

IV. REFERENCE ARCHITECTURE

The canonical edge stack comprises three strata:

- Device layer. Resource-constrained IoT sensors/actuators generate raw telemetry.
- Edge layer. Gateway-class nodes (ARM SoCs, micro-servers) perform filtering, lightweight analytics and local control.
- Cloud layer. Elastic clusters retain historical data, execute deep learning training and deliver fleet-wide orchestration.

Key enablers include containerised micro-services,

hardware accelerators (GPU/NPU) and network slices (5G uRLLC). Figure 2 depicts this layered view together with bidirectional data flows: upstream (telemetry \rightarrow cloud) and downstream (policy updates \rightarrow edge).

Reference Edge Architecture



Figure 2: Reference Edge Architecture.



Figure 3: Canonical Edge Stack

V. CASE STUDIES AND RESULTS

5.1 Healthcare: Remote Patient Monitoring • *Figure 3 & 4*

Context. Johns Hopkins Hospital deployed edge gateways beside wards to analyse electrocardiogram (ECG) and oxygen saturation streams. Cloud inference was retained for model retraining.

Results. Average alert latency dropped from 200 ms (cloud) to 20 ms (edge)—a 90 % improvement. Bandwidth to the data centre declined by 65 % because only anomalous segments were forwarded. No data were stored outside hospital premises, easing HIPAA compliance.



Figure 4: Healthcare Deployment Topology

Edge-enabled healthcare deployments build a mini-data-center at the patient's bedside—linking wearables, gateway-class edge servers and hospital IT to deliver millisecond-level analytics without pushing every heartbeat to the public cloud. The topology below summarises the prevailing architecture that research labs, hospitals and device vendors are converging on, explaining how each hop (sensor \rightarrow edge \rightarrow cloud) is secured, orchestrated and optimised for latency-critical care.



Latency Comparison in Healthcare

Figure 5: Alert Latency in Remote Patient Monitoring.

5.2 Industrial IoT: Predictive Maintenance

Context. Siemens retrofitted CNC machines with vibration sensors; edge nodes executed FFT-based anomaly detection while the cloud aggregated fleet logs.

Results. Mean Time Between Failures (MTBF) improved by 34 % and unplanned downtime fell by

40 % after six months. Operational efficiency (Fig. 5) is expressed as a normalised index (baseline = 100).



Ref Image 1: Predictive Maintenance Costs- Cloud vs Edge

Operational Efficiency in Industrial Applications



Figure 6: Operational-Efficiency Index (Industrial IoT).

5.3 Autonomous	Vehicles:	Real-Time
Decision-Making		

Context. A leading autonomous-mobility vendor embedded heterogeneous accelerators (GPU, TPU) into vehicle compute clusters. Edge perception fused LiDAR, radar and camera data locally; the cloud handled high-definition map generation.

Results. Sensor-to-actuator reaction times plunged from 150 ms to 10 ms (Fig. 6), enabling safer operation at highway speeds. Cloud connectivity interruptions no longer halted driving functions.



Figure 7: Reaction Time: Autonomous Vehicles.

VI. DISCUSSION

Cross-domain benefits. All pilots achieved at least one order-of-magnitude latency reduction and significant bandwidth savings. Data residency within organisational boundaries strengthened compliance (HIPAA, GDPR, ISO 27001).

Common bottlenecks. 1) Resource constraints: Edge nodes cap CPU/GPU power envelopes at <30 W. 2) Heterogeneity: Diverse hardware, OS and network link types complicate orchestration. 3) Scalability: Co-ordinating thousands of distributed micro-sites exceeds current DevOps toolchains.

Security risks. While reduced data egress lowers exposure, each additional edge node enlarges the attack surface. Zero-trust designs and remote-attestation protocols are becoming indispensable.

VII. FUTURE DIRECTIONS

Cross-domain benefits. All pilots achieved at least one order-of-magnitude latency reduction and significant bandwidth savings. Data residency within organisational boundaries strengthened compliance (HIPAA, GDPR, ISO 27001).

Common bottlenecks. 1) Resource constraints: Edge nodes cap CPU/GPU power envelopes at <30 W. 2) Heterogeneity: Diverse hardware, OS and network link types complicate orchestration. 3) Scalability: Co-ordinating thousands of distributed micro-sites exceeds current DevOps toolchains.

Security risks. While reduced data egress lowers exposure, each additional edge node enlarges the attack surface. Zero-trust designs and remote-attestation protocols are becoming indispensable.

VIII. CONCLUSION

Edge computing demonstrably delivers ultra-low latency, bandwidth frugality and improved privacy, mission-critical empowering applications in healthcare, industry and mobility. Empirical evidence across three deployments indicates transformative gains-yet scaling such benefits demands advances in lightweight AI, unified management and secure, standards-based ecosystems. Addressing these challenges will cement edge computing as a cornerstone of next-generation digital infrastructure.

REFERENCES

- Shi, W., Cao, J., Zhang, Q., Li, Y., & Xu, L. "Edge Computing: Vision and Challenges." IEEE Internet of Things Journal, 3 (5), 637-646 (2016).
- [2] Satyanarayanan, M. "The Emergence of Edge Computing." IEEE Computer, 50 (1), 30-39 (2017).
- [3] Chiang, M., & Zhang, T. "Fog and IoT: An Overview of Research Opportunities." IEEE Internet of Things Journal, 3 (6), 854-864 (2016).
- [4] Roman, R., López, J., & Mambo, M. "Mobile Edge Computing, Fog et al.: A Survey and Analysis of Security Threats and Challenges." Future Generation Computer Systems, 78, 680-698 (2018).
- [5] Liu, J., Zhang, F., & Shi, W. "EdgeAV: Edge Computing for Autonomous Driving." In Proceedings of the IEEE/ACM Symposium on Edge Computing, 66-79 (2019).
- [6] Wang, L., Chen, H., Xu, K., Li, Y., & Shi, W. "Optimising Edge AI: A Comprehensive Survey on Data, Model and System Techniques." ACM Computing Surveys, 57 (1), 1-38 (2024).
- [7] Rashid, A., & Al-Fawaz, M. "Edge Computing for Real-Time Decision-Making in Autonomous Vehicles." International Journal of Advanced Computer Science and Applications, 15 (7), 491-499 (2024).
- [8] Ghosh, S., Das, P., & Dutta, S. "Edge Computing

in Healthcare: Remote Patient Monitoring Use-Cases." International Journal of Computer Applications, 184 (43), 1-6 (2023).

- [9] Siemens AG. Industrial Edge: AI-Driven Predictive Maintenance in Discrete Manufacturing. White paper, 2021.
- [10] Intel Corporation. How Edge Computing Is Driving Advancements in Healthcare. Technical brief, 2024.