# **JETRM** International Journal of Engineering Technology Research & Management Published By: <u>https://www.ijetrm.com/</u>

# INTELLIGENT EDGE ARCHITECTURES: AI AT THE BOUNDARY OF CLOUD AND DEVICE

# Srikanth Jonnakuti

Staff Software Engineer, Cloud Architect, Move Inc. operator of Realtor.com, Newscorp

#### ABSTRACT

Edge computing shifts intelligence closer to data sources, reducing latency and bandwidth by distributing ML inference between devices and cloud. This paper presents a concise study of hybrid edge–cloud architectures tailored for latency-sensitive retail and IoT applications. Building on paradigms such as fog computing, cloudlets, and MEC, we outline three architectural patterns—hierarchical edge–cloud, collaborative (split) inference, and on-device inference with cloud backup—and discuss orchestration strategies. Use cases in smart retail, industrial IoT, and smart cities illustrate real-world benefits of local inference with periodic cloud synchronization. Key challenges—resource constraints, heterogeneity, security, and management complexity—are analyzed alongside mitigation techniques. Finally, we survey emerging trends: advanced edge AI accelerators, 5G convergence, federated learning, autonomous edge management, and standards for interoperability. By summarizing state-of-the-art developments up to August 2022 in a condensed form, this article offers a 20% reduction in word count while retaining all 25 original references, enabling practitioners and researchers to grasp the essentials of intelligent edge architectures.

#### Keywords

Edge Computing  $\cdot$  Edge AI  $\cdot$  Hybrid Cloud  $\cdot$  Internet of Things  $\cdot$  Low Latency  $\cdot$  Machine Learning Inference

#### **1 INTRODUCTION**

Centralized cloud platforms offer scalable storage and model training, but skyrocketing data from billions of IoT devices has strained bandwidth and introduced unacceptable delays for real-time applications (Gartner, 2018). Latency from round-trip communication can breach the stringent timing requirements of autonomous systems, industrial control loops, and interactive retail services (Shi et al., 2016; Shi & Dustdar, 2016). Edge computing mitigates these issues by executing computation on—or near—data-generating devices, forming an "intelligent edge" that processes latency-critical inference tasks locally while leveraging the cloud for heavy aggregation and model refinement (Bonomi et al., 2014; Verbelen et al., 2012). Edge architectures span from microcontrollers executing TinyML models to on-premises gateways running complex analytics. This continuum demands strategies to partition workloads optimally, orchestrate tasks dynamically, and secure a vastly enlarged attack surface. In this article, we distill key architectural patterns, illustrate applications in retail and IoT domains, analyze practical challenges, and highlight future directions, all within a concise framework.

#### **2 RELATED WORKS**

#### 2.1 Evolution of Edge and Fog Computing

Early CDN servers in the 1990s foreshadowed today's edge by caching content closer to users (Davis et al., 2004). Cisco's 2012 fog computing model introduced a layered continuum—device, fog node, cloud—to support latency-sensitive IoT applications, culminating in the OpenFog reference architecture (Bonomi et al., 2014; Dolui & Datta, 2017). Cloudlets—trusted micro-data centers one hop from mobile clients—demonstrated significant latency reductions for AR and video analytics (Verbelen et al., 2012). Meanwhile, ETSI's Mobile Edge Computing (now Multi-access Edge Computing) standardized cloud capabilities at cellular base stations, enabling URLLC for 5G use cases like connected vehicles (ETSI, 2015; Mao et al., 2017). Surveys by Shi et al. (2016) and Abbas et al. (2018) have cataloged these paradigms, underlining persistent challenges in connectivity, energy, and orchestration.

#### 2.2 Edge Intelligence and On-Device AI

The push to embed ML at the edge—edge AI—addresses immediacy and privacy by processing sensitive data locally (Cao et al., 2015; Hassan et al., 2018). Techniques such as quantization, pruning, and knowledge distillation compress DNNs to run on constrained hardware, while TinyML frameworks enable inference on

**International Journal of Engineering Technology Research & Management** 

Published By:

https://www.ijetrm.com/

microcontrollers. Collaborative inference (split computing) further reduces bandwidth by executing initial model layers on devices and offloading intermediate features to the cloud for final classification (Kang et al., 2017; Teerapittayanon et al., 2017). Runtime frameworks like TensorFlow Lite and hardware like Google's Edge TPU and NVIDIA Jetson series have made on-device inference practical by 2020 (Li et al., 2018; Mohammadi et al., 2018). Together, these advances lay the groundwork for hybrid architectures combining edge autonomy with cloud-scale coordination.

# **3 PROPOSED ARCHITECTURES**

We outline three hybrid patterns that balance latency, bandwidth, and compute resources:

- 1. Hierarchical Edge-Cloud (Three-Tier)
  - Devices  $\rightarrow$  Edge/Fog Nodes  $\rightarrow$  Cloud.

Edge nodes filter, aggregate, and run inference on raw IoT streams, sending only critical summaries to the cloud for storage and retraining (Shi & Dustdar, 2016). In smart factories, on-site servers detect anomalies locally, triggering immediate alerts and logging to the cloud for model updates (El-Sayed et al., 2017).

### 2. Collaborative Inference (Split Computing)

Partition DNNs across device and cloud: devices execute early feature-extractor layers, transmitting compact feature maps to cloud servers for final layers, drastically cutting data transfer and improving privacy (Kang et al., 2017). Optimal cut-layer selection can adapt to runtime conditions—bandwidth, device load—via profiling or dynamic algorithms (Teerapittayanon et al., 2017).

### 3. On-Device Inference with Cloud Support

Edge devices host complete lightweight models for offline responsiveness. The cloud handles exceptional cases and periodic model updates. For instance, voice assistants process common commands locally and forward complex queries to the cloud, which also aggregates usage data to refine models (Li et al., 2018). Model caching strategies enable edge nodes to load context-specific models on demand (Mao et al., 2017).

**Orchestration** across these architectures requires intelligent runtimes that monitor latency, energy, and resource availability to decide task placement. Heuristic or learning-based schedulers can dynamically offload workloads to maintain SLAs under variable conditions (Mao et al., 2017; Abbas et al., 2018).

# 4 APPLICATIONS IN RETAIL AND IOT

#### 4.1 Smart Retail

Edge servers in stores process video feeds to track shelf inventory, analyze shopper behavior, and detect security events with minimal delay—alerts reach staff within seconds, preventing stockouts and theft (Scale Computing, 2020). Personalized digital signage leverages anonymous on-device vision models to tailor advertisements in real time, preserving privacy by discarding raw video after processing (Gartner, 2018). POS systems rely on local edge databases to continue transactions offline, synchronizing with the cloud once connectivity is restored (IEEE Innovation at Work, 2019).

#### 4.2 Industrial IoT

On-site gateways run predictive maintenance models on sensor streams to detect equipment faults before failures, triggering automatic shutdowns or alerts without cloud dependence (Cao et al., 2015). Cloud-aggregated data refines these models, which are then redistributed to factory edge nodes for continuous improvement (Mohammadi et al., 2018).

### 4.3 Smart Cities and Transportation

Edge nodes at intersections analyze traffic density and pedestrian flows, dynamically adjusting signal timings to ease congestion and enhance safety, with sub-10 ms responses unattainable by cloud-only systems (Abbas et al., 2018). Roadside MEC servers facilitate vehicle-to-edge communication for cooperative driving and hazard warnings (ETSI, 2015).

#### 4.4 Healthcare and Smart Homes

Wearable monitors and home hubs execute on-device arrhythmia detection or security camera person-detection models, issuing immediate alerts locally and forwarding summaries to cloud platforms for physician review, thereby safeguarding privacy and ensuring uninterrupted service during outages (Hassan et al., 2018).

International Journal of Engineering Technology Research & Management

Published By:

https://www.ijetrm.com/

# **5 CHALLENGES AND LIMITATIONS**

- **Resource Constraints**: Edge CPUs and microcontrollers impose strict limits on model size, compute, and energy. Compression and model-simplification techniques are essential, but often necessitate split computing or cloud assistance (Shi et al., 2016; Mohammadi et al., 2018).
- Heterogeneity & Scalability: Diverse hardware, OSes, and communication protocols complicate deployment. Lightweight container orchestration (e.g., K3s) and unified data-integration layers help, yet managing thousands of nodes remains a DevOps challenge (Abbas et al., 2018; Hamdan et al., 2020).
- Latency Trade-offs: While edge reduces end-to-end delay, over-offloading or multi-hop topologies can introduce bottlenecks. Network engineering and adaptive partitioning are required to confine critical loops to minimal hops (Shi & Dustdar, 2016; Mach & Becvar, 2017).
- Orchestration Complexity: Dynamic scheduling must consider latency, energy, and cost. Heuristics often replace optimal—but intractable—solutions, and building transparent, policy-driven runtimes is an ongoing research area (Abbas et al., 2018).
- Security & Privacy: Decentralized nodes broaden the threat surface. Secure boot, TEEs, and zero-trust models are vital to protect data at rest and in transit (El-Sayed et al., 2017; Teerapittayanon et al., 2017). Federated learning and differential privacy can mitigate raw data exposure (Li et al., 2018).
- **Maintenance Overhead**: OTA updates, remote diagnostics, and resilience to node failures demand robust EdgeOps platforms. Digital-twin abstractions and predictive maintenance of the edge infrastructure itself are emerging solutions (Merenda et al., 2019).
- **Model Staleness**: Edge-deployed models risk concept drift. Hybrid inference schemes can direct low-confidence cases to cloud models, but balancing accuracy, latency, and network usage requires sophisticated orchestration (Mao et al., 2017; Mohammadi et al., 2018).

# 6 FUTURE TRENDS

- Advanced Edge AI Hardware: Next-gen NPUs and energy-efficient SoCs (e.g., Edge TPU successors, Jetson upgrades) will shrink the performance gap with cloud GPUs, enabling larger models locally (Merenda et al., 2019).
- **5G–Edge Convergence**: URLLC and network slicing in 5G will embed MEC platforms within cellular networks, offering guaranteed low-latency slices for edge workloads (ETSI, 2015). Edge-as-a-Service offerings will let developers deploy workloads on telco-owned edge servers.
- Federated & Distributed Learning: FL will evolve to support heterogeneous devices synchronizing model updates without sharing raw data, benefiting healthcare, personalized services, and industrial deployments (Mohammadi et al., 2018). Peer-to-peer learning among fog nodes ("fog federation") promises faster propagation of local innovations (Samarah et al., 2018).
- Autonomous Edge Management: AI-driven orchestration, leveraging reinforcement learning and digital twins, will automate placement, scaling, and failure recovery across the edge continuum (Abbas et al., 2018; Hamdan et al., 2020).
- Standardization & Interoperability: Efforts by IEEE (e.g., IEEE 1934) and LF Edge will yield common APIs for discovery, model deployment, and telemetry, fostering an open edge ecosystem akin to cloud-native frameworks (Mach & Becvar, 2017).
- **Privacy-Enhancing Technologies**: Homomorphic encryption, secure multi-party computation, and differential privacy will mature enough for edge contexts, enabling encrypted aggregation of model updates or analytics outputs (Li et al., 2018).
- New Modalities: Edge genomics, on-orbit satellite processing, and AI-driven agricultural drones represent frontier domains where latency, autonomy, and intermittent connectivity necessitate advanced edge architectures (Scale Computing, 2020).

International Journal of Engineering Technology Research & Management

Published By:

https://www.ijetrm.com/

# 7 CONCLUSIONS

Intelligent edge computing bridges the gap between centralized cloud AI and resource-constrained devices, delivering low-latency, privacy-preserving services across retail, industry, and urban infrastructures. By surveying the evolution of fog, cloudlets, and MEC, we distilled three hybrid architectural patterns— hierarchical edge–cloud, split inference, and on-device models with cloud support—and examined their orchestration. Real-world use cases demonstrate significant gains in responsiveness, bandwidth savings, and resilience. Nonetheless, challenges in resource management, heterogeneity, security, and operational complexity remain active research areas. Emerging hardware accelerators, 5G integration, federated learning, and standardized platforms promise to enhance edge capabilities and simplify deployment. As edge devices grow more powerful and networks more reliable, the boundary between cloud and edge will blur, forming a seamless continuum where AI can execute optimally at any point—from microcontrollers to data centers—fulfilling the promise of real-time, ubiquitous intelligence.

### REFERENCES

[1] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge computing: Vision and challenges," IEEE Internet of Things Journal, vol. 3, no. 5, pp. 637–646, Oct. 2016, doi:10.1109/JIOT.2016.2579198.

[2] W. Shi and S. Dustdar, "The promise of edge computing," IEEE Computer, vol. 49, no. 5, pp. 78–81, May 2016, doi:10.1109/MC.2016.145.

[3] N. Abbas, Y. Zhang, A. Taherkordi, and T. Skeie, "Mobile edge computing: A survey," IEEE Internet of Things Journal, vol. 5, no. 1, pp. 450–465, Feb. 2018, doi:10.1109/JIOT.2017.2750180.

[4] K. Dolui and S. K. Datta, "Comparison of edge computing implementations: Fog computing, cloudlet and mobile edge computing," in Proc. Global Internet of Things Summit (GIoTS), Geneva, Switzerland, Jun. 2017, pp. 1–6, doi:10.1109/GIOTS.2017.8016213.

[5] F. Bonomi, R. Milito, P. Natarajan, and J. Zhu, "Fog computing: A platform for Internet of Things and analytics," in Big Data and Internet of Things: A Roadmap for Smart Environments, Springer, Cham, 2014, pp. 169–186, doi:10.1007/978-3-319-05029-4\_7.

[6] T. Verbelen, P. Simoens, F. De Turck, and B. Dhoedt, "Cloudlets: Bringing the cloud to the mobile user," in Proc. ACM Workshop on Mobile Cloud Computing & Services (MCS), Scottsdale, AZ, USA, Jul. 2012, pp. 29–36, doi:10.1145/2307849.2307858.

[7] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," IEEE Communications Surveys & Tutorials, vol. 19, no. 4, pp. 2322–2358, 2017, doi:10.1109/COMST.2017.2745201.

[8] M. Taleb, K. Samdanis, B. Mada, H. Flinck, S. Dutta, and D. Sabella, "On multi-access edge computing: A survey of the emerging 5G network edge architecture and orchestration," IEEE Communications Surveys & Tutorials, vol. 19, no. 3, pp. 1657–1681, 3rd qtr. 2017, doi:10.1109/COMST.2017.2692421.

[9] H. El-Sayed, M. Alam, Z. Kouki, M. Mohiuddin, and B. Baksi, "Edge of Things: The big picture on the integration of edge, IoT and cloud in a distributed computing environment," IEEE Access, vol. 6, pp. 1706–1717, 2018, doi:10.1109/ACCESS.2017.2786567.

[10] N. Hassan, S. Gillani, E. Ahmed, I. Yaqoob, and M. Imran, "The role of edge computing in Internet of Things," IEEE Communications Magazine, vol. 56, no. 11, pp. 110–115, Nov. 2018, doi:10.1109/MCOM.2018.1700906.

[11] Y. Cao, A. Marin-Perianu, I. Lombardi, E. Marin-Perianu, and P. J. Marin, "Distributed analytics and edge intelligence for the Internet of Things (IoT)," in Proc. IEEE International Workshop on Mobile Big Data, Shanghai, China, Jun. 2015, pp. 43–48, doi:10.1145/2757384.2757398.

[12] H. Li, K. Ota, and M. Dong, "Learning IoT in edge: Deep learning for the Internet of Things with edge computing," IEEE Network, vol. 32, no. 1, pp. 96–101, Jan. 2018, doi:10.1109/MNET.2018.1700202.
 [13] M. Mohammadi, A. Al Eugaha, S. Saraur, and M. Guizani, "Deep learning for IoT big data and streaming for IoT big data and streaming for IoT big data."

[13] M. Mohammadi, A. Al-Fuqaha, S. Sorour, and M. Guizani, "Deep learning for IoT big data and streaming analytics: A survey," IEEE Communications Surveys & Tutorials, vol. 20, no. 4, pp. 2923–2960, 2018, doi:10.1109/COMST.2018.2844341.

[14] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. N. Mudge, J. Mars, and L. Tang, "Neurosurgeon: Collaborative intelligence between the cloud and mobile edge," in Proc. ACM ASPLOS, Williamsburg, VA, USA, Apr. 2017, pp. 615–629, doi:10.1145/3037697.3037698.

**International Journal of Engineering Technology Research & Management** 

Published By:

https://www.ijetrm.com/

[15] S. Teerapittayanon, B. McDanel, and H. T. Kung, "Distributed deep neural networks over the cloud, the edge and end devices," in Proc. IEEE International Conference on Distributed Computing Systems (ICDCS), Atlanta, GA, USA, Jun. 2017, pp. 328–339, doi:10.1109/ICDCS.2017.226.

[16] A. Davis, J. Parikh, and W. E. Weihl, "Edge computing: Extending enterprise applications to the edge of the Internet," in Proc. 13th International World Wide Web Conference Alternate Track Papers & Posters, Budapest, Hungary, May 2004, pp. 180–187, doi:10.1145/1013367.1013397.

[17] R. Roman, J. Lopez, and M. Mambo, "Mobile edge computing, fog et al.: A survey and analysis of security threats and challenges," Future Generation Computer Systems, vol. 78, pp. 680–698, Jan. 2018, doi:10.1016/j.future.2016.11.031.

[18] M. Merenda, C. Porcaro, and D. Iero, "Edge machine learning for AI-enabled IoT devices: A review," Sensors, vol. 20, no. 9, Art. 2533, May 2020, doi:10.3390/s20092533.

[19] S. Samarah, M. G. H. Zamil, M. Rawashdeh, M. S. Hossain, G. Muhammad, and A. Alamri, "Transferring activity recognition models in fog computing architecture," Journal of Parallel and Distributed Computing, vol. 122, pp. 122–130, 2018, doi:10.1016/j.jpdc.2018.07.020.

[20] G. Premsankar, M. Di Francesco, and T. Taleb, "Edge computing for the Internet of Things: A case study," IEEE Internet of Things Journal, vol. 5, no. 2, pp. 1275–1284, Apr. 2018,

doi:10.1109/JIOT.2017.2783258.

[21] M. Satyanarayanan, P. Bahl, R. Caceres, and N. Davies, "The case for VM-based cloudlets in mobile computing," IEEE Pervasive Computing, vol. 8, no. 4, pp. 14–23, Oct.–Dec. 2009, doi:10.1109/MPRV.2009.82.
[22] S. Hamdan, M. Ayyash, and S. Almajali, "Edge-computing architectures for Internet of Things applications: A survey," Sensors, vol. 20, no. 22, Art. 6441, Nov. 2020, doi:10.3390/s20226441.

[23] P. Mach and Z. Becvar, "Mobile edge computing: A survey on architecture and computation offloading," IEEE Communications Surveys & Tutorials, vol. 19, no. 3, pp. 1628–1656, 2017, doi:10.1109/COMST.2017.2682318.

[24] J. Lin, W. Yu, N. Zhang, X. Yang, H. Zhang, and W. Zhao, "A survey on Internet of Things: Architecture, enabling technologies, security and privacy, and applications," IEEE Internet of Things Journal, vol. 4, no. 5, pp. 1125–1142, Oct. 2017, doi:10.1109/JIOT.2017.2694844.

[25] S. Wang, J. Wan, D. Zhang, D. Li, and C. Wu, "Edge computing: A survey of recent advances and future trends," Journal of Network and Computer Applications, vol. 154, p. 102620, Jan. 2020, doi:10.1016/j.jnca.2020.102620.