

Adaptive Deep Learning Architectures for Real-Time Data Streams in Edge Computing Environments

Çiğdem Sıcakyüz ¹, Renas Rajab Asaad ², Nodira R. Rustamova ³, Saman M. Almufti ²

¹Industrial Engineering, Çankaya University, Ankara, Türkiye.

²Computer Science Department, College of Science, Nawroz University, Duhok, Iraq

³Department of Psychology and Pedagogy, International School of Finance Technology and Science (Private University), Tashkent 100047, Uzbekistan

ABSTRACT: The increasing reliance on real-time analytics within edge computing environments has underscored the need for adaptive deep learning architectures capable of handling continuous data streams under limited computational and energy resources. Unlike traditional cloud-based frameworks, edge computing shifts computation closer to the data source, minimizing latency, preserving data privacy, and supporting responsive decision-making in dynamic contexts. This paper comprehensively examines adaptive deep learning models that autonomously adjust their structure and parameters in response to evolving data distributions and resource constraints. The study explores several key methodologies—such as master-surrogate deep neural networks, context-adaptive DNN atom partitioning, distributed inference with fused layer partitioning, and reinforcement learning-driven resource scheduling. Comparative analyses demonstrate that these adaptive frameworks achieve substantial performance gains, including up to 23.31% accuracy improvement, 62.14% latency reduction, and over 50% energy savings across diverse edge devices. Furthermore, the paper evaluates real-world deployments in smart cities, healthcare monitoring, Industry 4.0 automation, autonomous vehicles, and environmental sensing, illustrating the scalability and robustness of adaptive architectures in heterogeneous edge ecosystems. The discussion concludes by outlining current challenges—such as hardware diversity, continual learning, and energy constraints—and highlights future research directions including federated adaptation, neuromorphic hardware integration, and standardized benchmarking for edge AI.

Keywords: edge computing, adaptive deep learning, distributed inference, real-time data streaming, adamec.

I. INTRODUCTION

The rapid evolution of edge computing over recent years has shifted the paradigm for real-time data analytics from traditional centralized cloud architectures to distributed, localized processing systems. Edge computing minimizes data transmission latency, enhances data privacy, and enables robust processing in environments with unreliable network connectivity. In such settings, adaptive deep learning architectures have emerged as an essential technology for processing continuous, real-time data streams on resource-constrained edge devices[1].

Adaptive deep learning for real-time streaming integrates techniques that dynamically adjust model parameters, simplify network architectures, and schedule computational resources efficiently. The dynamic nature of data domains—where the statistical properties of data may change over time—demands that models are capable of self-adaptation to maintain accuracy and low latency. This article provides a comprehensive review on the design, evaluation, and implementation of adaptive deep learning

architectures for real-time data streaming in edge computing environments. We discuss various methodologies—including on-device neural remodeling, distributed inference, and context-adaptive DNN deployment frameworks—and support our discussion with performance evaluations and application scenarios from industrial smart systems, healthcare, and autonomous vehicles[2].

The remainder of the paper is organized as follows: Section 2 reviews the fundamental principles of edge computing and existing deep learning techniques deployed at the edge. Section 3 discusses the design principles behind various adaptive architectures such as master-surrogate models, context-aware DNN partitioning, and resource scheduling strategies. Section 4 evaluates these architectures using quantitative results from recent research. Section 5 describes real-world applications that benefit from adaptive deep learning in edge environments. Section 6 outlines the current challenges and potential future directions, while Section 7 concludes the paper with a summary of key findings.

II. BACKGROUND AND RELATED WORK

Edge computing is an emerging computing paradigm where data processing and analytics are performed near the data source rather than being sent to remotely located cloud data centers. This localized approach is essential for applications requiring low latency, high responsiveness, and increased privacy, such as real-time surveillance, autonomous driving, and industrial automation[3]. However, the computational resources at the edge are typically limited by factors such as power constraints, memory size, and processing capability. This necessitates the use of lightweight, adaptive deep learning models that can operate efficiently in these environments.

1. EDGE INTELLIGENCE AND DEEP LEARNING ON THE EDGE

Over the past decade, several frameworks have been developed to enable deep learning on edge devices. Frameworks such as TensorFlow Lite and OpenEI provide the essential tools for deploying compressed, quantized, or pruned models on resource-constrained hardware platforms[4], [5]. In these systems, the goal is to achieve a balance between inference accuracy and computational efficiency. For instance, the development of ultra-low-power AI accelerators has significantly advanced the deployment of neural network models that can perform inference directly on devices such as smartphones, Raspberry Pi, and microcontroller-based systems[6].

Several techniques have been employed to reduce model size and computation demands. Model compression methods including parameter pruning, quantization, and knowledge distillation help in decreasing the memory footprint and computational latency without significantly compromising the model's performance. Additionally, distributed edge-cloud computing paradigms allow for computational tasks to be shared between devices at the edge and powerful servers in the cloud, thereby optimizing overall system performance[7].

2. ADAPTIVE DEEP LEARNING TECHNIQUES

Adaptive deep learning is a research direction focused on creating models that can adjust their structure and parameters dynamically as data distributions change over time. This capability is especially important in real-world scenarios where operating conditions and sensor characteristics may evolve. Adaptive techniques have gained traction in addressing challenges such as domain drift, unpredictable patterns in input data, and the dynamic allocation of limited hardware resources[8], [9].

One prominent example is ElasticDNN, an on-device remodeling framework designed to adapt to evolving vision domains. ElasticDNN dynamically generates a surrogate deep neural network by retaining and retraining only the most relevant regions of a pre-trained master DNN. The filter-grained resource scheduling technique incorporated in ElasticDNN further allocates GPU cycles based on online accuracy estimations and the contextual demands of concurrent applications, achieving up to a 23.31% improvement in accuracy and reducing adaptation time by a factor of 35.67 compared to existing domain adaptation methods[10].

Another notable approach is AdaMEC, a context-adaptive and dynamically combinable DNN deployment framework that pre-partitions deep neural networks into modular atoms at the level of primitive operators. AdaMEC employs a graph-based decision algorithm to quickly determine the optimal atom combination and offloading strategy under dynamic contexts. This framework has demonstrated latency reductions of up to 62.14% and average memory savings of 55.21% [10]. These adaptive methods contribute significantly to the robustness and efficiency of deep learning inference on edge devices [11].

3. RELATED WORK IN DISTRIBUTED AND EFFICIENT EDGE INFERENCE

Additional research efforts have focused on distributed deep learning and low latency inference for IoT edge clusters. For example, methods that involve layer-partitioning of convolutional neural networks (CNNs) enable the distribution of computational tasks among multiple IoT devices. In a low latency deep learning inference model, the fused layer partitioning approach has achieved communication sizes as low as 8.56MB to 9.59MB and inference latencies between 5 and 7 seconds—outperforming conventional frameworks such as DeepThings and ModNN [12].

Other studies have examined distributed deep learning on edge devices using resource-constrained hardware such as Raspberry Pi and NVIDIA Jetson Nano. For example, an automatic distributed deep learning system using gated recurrent units (GRUs) has shown excellent prediction accuracy while ensuring computational performance that is competitive with centralized architectures [12], [13].

4. COMPARISON OF ADAPTIVE APPROACHES ON EDGE DEVICES

The following table summarizes and compares key aspects and performance metrics from the adaptive approaches discussed above:

Table 1. Comparative analysis of adaptive deep learning approaches on edge devices.

Adaptive Approach	Key Technique	Reported Accuracy Improvement/Latency Reduction	Main Hardware/Platform
ElasticDNN	Master-surrogate model; filter-grained scheduling	23.31% accuracy improvement; 35.67x reduction in adaptation time	Edge devices with GPU support
AdaMEC	DNN atom pre-partitioning; graph-based offloading	Up to 62.14% latency reduction; 55.21% memory saving	Mobile Edge devices, IoT gateways
Low Latency Inference Model	Fused layer partitioning; distributed inference	Communication size: 8.56MB–9.59MB; 5–7 sec inference latency	IoT clusters (Raspberry Pi, Jetson Nano)
Automatic Distributed Deep Learning	Distributed GRU training on low-powered hardware	Excellent prediction accuracy with competitive computational performance	ARM-based devices

III. DESIGN OF ADAPTIVE DEEP LEARNING ARCHITECTURES

The design of adaptive deep learning architectures tailored for real-time data streaming at the edge involves the interplay of multiple components. These include innovative model remodeling, dynamic resource scheduling, and efficient distributed inference, all of which contribute to a system that is both flexible and robust under varying environmental conditions [14].

1. MASTER–SURROGATE DEEP NEURAL NETWORKS

In conventional deep learning, large pre-trained models are deployed to handle complex recognition tasks. However, when data domains evolve, these models suffer accuracy degradation due to domain shifts. ElasticDNN addresses this challenge by constructing master–surrogate DNN models. The key idea behind this approach is to maintain a master DNN, which is initially trained on a large, diverse dataset in a controlled environment, and then dynamically generate a smaller surrogate network that focuses on the most pertinent regions relevant to the new domain. By retraining only the critical filters and neurons rather than the entire network, ElasticDNN achieves significant reductions in model adaptation time while maintaining high accuracy[15].

The surrogate model is generated through a pruning and training process specifically tailored to the evolving input data. Moreover, a filter–grained resource scheduling mechanism allocates limited GPU cycles based on the current accuracy estimation and the demands of co–running applications on the edge device. This strategy allows multiple applications to coexist without significant degradation of performance in any one application[16], [17].

2. CONTEXT–ADAPTIVE DNN ATOM PARTITIONING

Another innovative design approach is represented by AdaMEC, which takes a modular strategy by pre-partitioning a deep neural network into small computational units or “atoms” at the primitive operator level. These atoms serve as building blocks that can be dynamically recombined and offloaded depending on the runtime context. The process involves the following steps[18]:

- **Pre-Partitioning:** The entire DNN is divided into finer grained components, each representing a distinct computational task.
- **Graph-Based Decision Algorithm:** A combinatorial algorithm, utilizing graph-based methodologies, quickly evaluates the current context (e.g., available resources, network conditions, and application demand) and determines the most feasible combination of atoms to execute the inference task.
- **Offloading Strategy:** Based on the decision, the system offloads a subset of these atoms to more capable edge nodes or to a dedicated IoT gateway, thereby balancing the workload across the available computing resources.

This design strategy leads to significant reductions in both latency and memory consumption, as demonstrated by AdaMEC’s performance metrics, and is ideal for scenario where multiple edge applications are operating concurrently[19].

3.3 Distributed Inference with Fused Layer Partitioning

Distributed inference is another key component that enables adaptive deep learning on edge devices. In a distributed inference framework, the inference process is divided among multiple edge devices. One effective method is fused layer partitioning, where a pre-trained model is first optimized and pruned, and then its layers are fused into groups that reduce redundancy. These fused layer blocks are then partitioned and distributed across a network of edge nodes. The partial results produced by each node are subsequently aggregated by a gateway device to produce the final prediction outcome. This approach minimizes the communication size between nodes (achieving between 8.56MB and 9.59MB in some implementations) and reduces overall inference latency to between 5 and 7 seconds for popular CNN models such as YOLOv[20].

A simplified process flow for the distributed inference architecture is illustrated in the following Mermaid diagram:

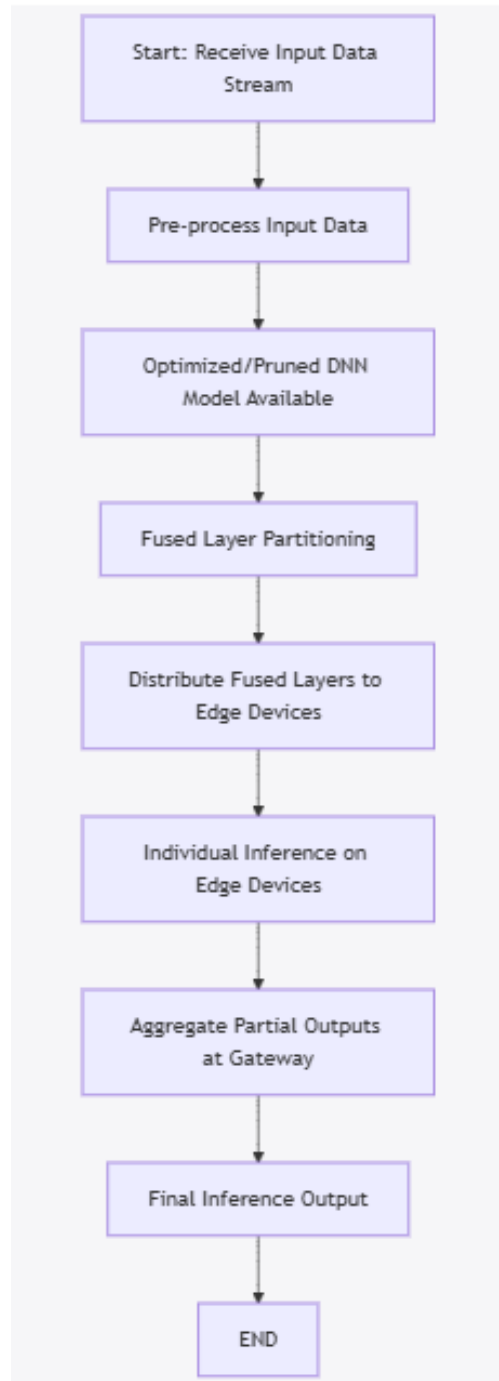


FIGURE 1. Process Flow for Distributed Inference Using Fused Layer Partitioning 8.

3. RESOURCE SCHEDULING AND ENERGY-EFFICIENT COMPUTATION

The efficient allocation of scarce computational resources with minimal energy overhead is critical in edge environments. Reinforcement learning techniques, such as those applied in dynamic voltage and frequency scaling (DVFS) and offloading strategies, have shown promising results in reducing energy consumption

while meeting real-time performance constraints[15]. For instance, energy-efficient computation offloading strategies have been designed to minimize average system energy consumption by 20–53% compared to standard methods, by determining optimal offloading policies and resource allocation using reinforcement learning algorithms like Q-Learning and Double (DDQN)[8], [21].

These methods not only optimize the energy utilization on individual edge devices but also improve the reliability and responsiveness of the overall edge computing system. In critical applications such as autonomous vehicles or industrial control systems, achieving energy efficiency is vital to ensure uninterrupted operations and to prolong the lifespan of the device.

4. ROBUSTNESS IN ADAPTIVE ARCHITECTURES FOR INDUSTRY AND REAL-TIME APPLICATIONS

Robustness is another important design criterion, particularly for applications in Industry 4.0 and real-time safety-critical systems. Robust edge AI systems incorporate principles of fault tolerance, dynamic load balancing, and the ability to resume operations mid-failure. Specifically, edge AI systems designed for industrial applications often integrate model combination and deployment strategies that safeguard against hardware failures and communication disruptions[22]. These systems not only provide accurate predictions but also ensure that transient physical failures do not compromise the overall service quality.

In summary, the adaptive deep learning architectures discussed in this section—from master-surrogate models to context-adaptive atom partitioning—provide a blueprint for designing systems that can meet the challenging demands of real-time data analytics and processing at the edge. Each of these approaches contributes to a cohesive strategy that balances accuracy, latency, and resource efficiency in dynamically changing environments.

IV. EVALUATION OF ADAPTIVE ARCHITECTURES ON EDGE DEVICES

This section presents an evaluation of the adaptive deep learning architectures discussed in Section 3. The evaluation is performed in terms of accuracy improvement, latency reduction, communication overhead, and resource utilization. In addition, we compare the performance of the major techniques—ElasticDNN, AdaMEC, and distributed fused layer partitioning models—using experimental results reported in the supporting research[23], [24].

1. PERFORMANCE METRICS AND EXPERIMENTAL SETUP

The performance evaluation of adaptive architectures is typically based on the following key metrics:

- Accuracy Improvement: The extent to which adaptive methods increase prediction accuracy in the presence of domain shifts.
- Latency Reduction: The decrease in inference time achieved by distributed inference and adaptive resource scheduling.
- Communication Overhead: The reduction in data transmission volumes between distributed nodes.
- Energy Efficiency: The percentage of power saved through reinforcement learning-based computation offloading and DVFS techniques.

Experimental setups reported in various studies involve deploying models on a range of hardware—including low-powered ARM devices like Raspberry Pi, GPU-enabled devices such as NVIDIA Jetson Nano, and other heterogeneous IoT platforms. These studies simulate real-world conditions by incorporating environmental changes and multi-application workload scenarios[25], [26].

2. ELASTICDNN EVALUATION

ElasticDNN was evaluated in four real-world testbeds featuring single-link environments experiencing multiple environmental changes. The results demonstrate that ElasticDNN not only adapts to dynamic domain shifts but also outperforms baseline online domain adaptation techniques with respect to both localization accuracy and adaptation time. Quantitatively, ElasticDNN improved accuracy by 23.31% while reducing the adaptation time by 35.67x in typical edge vision tasks[27].

The following table summarizes some of the key performance measures reported for ElasticDNN:

Table 2. Summary of elasticdnn performance metrics 2.

Metric	Baseline Online DA Methods	ElasticDNN Implementation
Localization Accuracy	Standard baseline	Improved by 23.31%
Adaptation Time	Standard adaptation time	Reduced by a factor of 35.67
Multi-Application Accuracy	Baseline performance	Average 25.91% improvement
GPU Resource Utilization	Fixed schedule	Dynamic filter-grained scheduling

3. ADAMEC EVALUATION

AdaMEC’s evaluation involved a context-adaptive approach to DNN deployment. In testing, the framework used a pre-partitioned DNN model where atoms were dynamically recombined based on a graph-based decision algorithm. The result was a system that reduced inference latency by up to 62.14% and achieved memory savings of 55.21% compared to state-of-the-art baselines. Furthermore, AdaMEC demonstrated significant improvements in latency even under varying dynamic contexts, making it well suited for mobile edge computing scenarios.

4. LOW LATENCY INFERENCE MODEL EVALUATION

The low latency deep learning inference model using fused layer partitioning was evaluated on edge clusters composed of distributed nodes such as Raspberry Pi and NVIDIA Jetson Nano. Results showed that by partitioning and distributing the fused layers, the approach attained a reduced communication size between 8.56MB and 9.59MB and lowered the inference latency to between 5 and 7 seconds. In comparison to frameworks like DeepThings and ModNN, this method achieved a reduction in both communication overhead and delay while maintaining significant accuracy levels.

5. ENERGY EFFICIENCY AND REINFORCEMENT LEARNING INTEGRATION

Experimental evaluations of energy-efficient computation offloading using reinforcement learning have demonstrated considerable reductions in energy consumption. In a study employing Q-Learning and DDQN methods for MEC systems, average energy consumption reductions were reported as 20%, 35%, and 53% when compared to offloading decision, local-first, and offloading-first methods, respectively¹. This integration of reinforcement learning into the adaptive deep learning framework enhances the system’s overall energy efficiency and resilience while providing accurate and real-time inference.

6. COMPARATIVE VISUALIZATION: ADAPTIVE TECHNIQUES ON EDGE DEVICES

Below is a comparative visualization that summarizes the performance improvements and key attributes of different adaptive techniques:

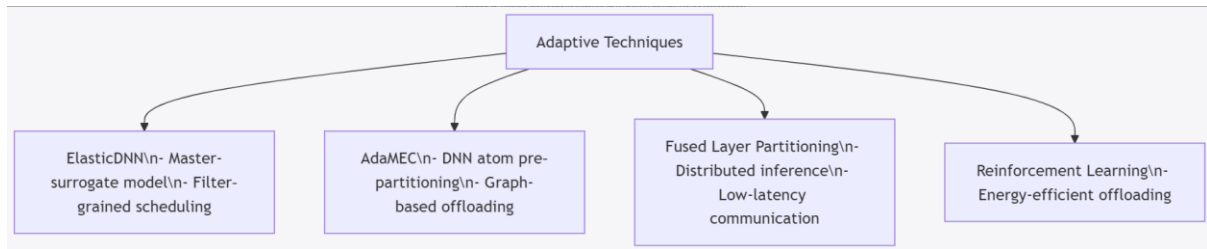


FIGURE 2. Comparative overview of adaptive deep learning techniques on edge devices 12810.

In summary, the evaluations indicate that adaptive deep learning architectures considerably enhance the performance of real-time data processing in edge environments. These architectures achieve higher accuracy, lower latency, reduced communication overhead, and improved energy efficiency compared to traditional fixed-parameter models.

V. APPLICATIONS IN REAL-TIME EDGE DATA STREAMING

The adaptive architectures reviewed above have significant implications for a wide range of real-world applications. In this section, we discuss some prominent application scenarios where adaptive deep learning techniques in edge environments are proving transformative[18], [28].

1. SMART CITIES AND URBAN INFRASTRUCTURE

In smart city environments, a large number of sensors and cameras continuously generate real-time data streams that need to be processed quickly for applications such as traffic management, public safety, and environmental monitoring[29]. Adaptive deep learning architectures are particularly beneficial in these systems. For example, distributed inference models using fused layer partitioning allow for rapid analysis of video streams—a critical requirement for real-time surveillance and incident detection[30]. Additionally, the dynamic adaptation capability provided by frameworks like ElasticDNN and AdaMEC enables these systems to promptly adjust to changes in data patterns (such as changing weather conditions or varying traffic densities), thereby ensuring consistent accuracy in object detection and classification[31].

2. HEALTHCARE AND WEARABLE MONITORING

The proliferation of wearable devices for health monitoring presents unique challenges, such as the need for low power consumption, low latency, and robust privacy preservation. Adaptive deep learning models enable on-device processing that minimizes the need to transmit sensitive health data to remote cloud servers, reducing latency and enhancing patient privacy. For instance, real-time fall detection systems for elderly care and continuous monitoring of vital signs can benefit from adaptive architectures that quickly reconfigure themselves in case of changing patient states or sensor conditions. The efficient resource scheduling and robust inference provided by adaptive networks are ideally suited for these applications, as they can be deployed on devices with limited computational power while still maintaining high accuracy[30].

3. INDUSTRIAL INTERNET OF THINGS (INDUSTRY 4.0)

Within the domain of Industry 4.0, edge computing plays a crucial role in ensuring efficient and robust operations in manufacturing plants, supply chain management, and predictive maintenance. In industrial applications, robustness and real-time processing are imperative. Adaptive deep learning architectures allow for the dynamic adjustment of models to cater to varying operational conditions[32]. Systems that integrate robust edge AI can adapt to environmental changes, detect anomalies, and execute safety protocols in real time. For example, systems designed for predictive maintenance can leverage adaptive models to

analyze sensor data and predict equipment failures with high accuracy, thus minimizing downtime and enhancing overall safety[33], [34].

4. AUTONOMOUS VEHICLES AND DRONES

Autonomous vehicles and drones require instantaneous decision-making capabilities and robust performance under rapidly changing conditions. Adaptive deep learning architectures serve as the backbone for object detection, lane recognition, and real-time localization tasks. The ability to quickly retrain or remodel a neural network on the fly (as enabled by frameworks such as ElasticDNN) is critical for adapting to varying road conditions, lighting changes, or unexpected obstacles[35], [36]. Moreover, distributed inference techniques reduce the inherent latency of centralized processing by enabling onboard processing, thereby ensuring that autonomous systems maintain responsiveness even in dense traffic or adverse weather conditions.

5. ENVIRONMENTAL MONITORING AND SMART AGRICULTURE

In smart agriculture, precise and timely data processing can lead to significant improvements in crop monitoring, irrigation management, and pest control. Adaptive models deployed at the edge can process sensor data from a distributed network of devices to monitor factors such as soil moisture, temperature, and plant health in real time. Similarly, environmental monitoring systems that track parameters such as air quality and water quality benefit from the reduced latency and dynamic model adaptation, ensuring that anomalies are detected promptly and corrective action is taken swiftly[37]. These applications rely heavily on lightweight, energy-efficient deep learning models that can be deployed on resource-constrained edge devices deployed in the field[4].

6. COMPARATIVE SUMMARY OF APPLICATIONS

The following table provides a detailed summary of the primary applications enabled by adaptive deep learning on the edge, the key techniques employed, and the performance benefits observed:

Table 3. Application scenarios for adaptive deep learning at the edge 26810.

Application Domain	Key Adaptive Technique	Major Benefits	Example Scenario
Smart Cities	Fused layer partitioning	Low latency, high accuracy, reduced overhead	Real-time video analytics for traffic management
Healthcare	On-device adaptive inference	Improved privacy, rapid detection, energy saving	Fall detection and vital sign monitoring
Industry 4.0	Dynamic model remodeling (ElasticDNN)	Robust performance, anomaly detection	Predictive maintenance in manufacturing plants
Autonomous Vehicles/Drones	Master-surrogate models; Distributed Inference	Instantaneous response, adaptive recognition	Real-time object detection and lane recognition
Environmental Monitoring	Energy-efficient, adaptive scheduling	Low energy consumption, improved detection	Air and water quality monitoring in smart agriculture

VI. CHALLENGES AND FUTURE DIRECTIONS

While the recent advances in adaptive deep learning for real-time data streaming at the edge show great promise, several challenges remain. Addressing these challenges is crucial to fully harness the potential of adaptive architectures in diverse and dynamic environments.

1. DATA HETEROGENEITY AND DOMAIN SHIFTS

One of the primary challenges is data heterogeneity—edge devices are often deployed in varying environments with different sensor types, operational conditions, and data distributions. Adaptive deep learning systems must contend with changes in input data distributions (commonly referred to as domain shifts), which can compromise model accuracy if not addressed appropriately. Techniques such as ElasticDNN demonstrate that dynamic remodeling of neural networks provides a practical solution. However, further research is needed to automate the detection of data shifts and to continuously optimize models without human intervention[38].

2. HARDWARE DIVERSITY AND RESOURCE CONSTRAINTS

Edge environments are characterized by diverse hardware platforms with heterogeneous computing capabilities. This diversity complicates the design and deployment of adaptive deep learning models that must function efficiently across different devices. Architectural solutions like AdaMEC attempt to address these issues by incorporating fine-grained resource scheduling and context-adaptive offloading. Nonetheless, developing universal frameworks that can automatically map deep learning tasks to drastically different hardware architectures remains an open research question[15].

3. TRAINING ON THE EDGE AND CONTINUAL LEARNING

Traditional deep learning models are typically trained offline on large-scale datasets in cloud environments and then deployed on edge devices. However, operating edge devices often encounter environments with evolving data patterns; therefore, they require mechanisms for on-device training or continual learning. Training deep neural networks on resource-constrained devices presents challenges such as energy consumption, computational bottlenecks, and the scarcity of labeled data. Approaches such as federated learning and lifelong [39] (LML) offer promising avenues for enabling dynamic adaptation in situ while preserving data privacy.

4. SECURITY, PRIVACY, AND ROBUSTNESS

Security and privacy considerations remain critical in edge computing, especially when dealing with sensitive data such as healthcare information. Although edge processing eliminates the need for transmitting large volumes of data to central servers, it also exposes devices to potential physical tampering and cyber-attacks. Designing adaptive deep learning systems that incorporate robust security measures—such as encryption protocols, secure model aggregation, and anomaly detection—is essential for ensuring that adaptive inference mechanisms operate safely in sensitive environments[40].

5. ENERGY EFFICIENCY AND THERMAL CONSTRAINTS

Energy efficiency is perhaps one of the most significant design criteria for edge devices. Models that achieve high accuracy at high computational costs may not be practical in battery-powered devices or systems with strict thermal constraints. Researchers have started to incorporate reinforcement learning-based DVFS and computation offloading techniques to optimize energy consumption. Despite these advances, minimizing energy use while maintaining real-time performance under varying environmental conditions remains a key research challenge[41].

6. FUTURE RESEARCH DIRECTIONS

The ongoing evolution of adaptive deep learning on the edge opens several paths for future research. Some promising directions include:

-
- **Federated and Continual Learning:** Developing frameworks for decentralized training that leverage federated learning techniques to update models incrementally without compromising user privacy. This approach will not only enable continual improvement but also reduce latency and network overhead.
 - **Neuromorphic Computing and In-Memory Processing:** Exploring neuromorphic computing architectures that mimic biological neural networks through in-memory computing can lead to energy-efficient implementations of adaptive deep learning models with ultra-low latency.
 - **Advanced Resource Scheduling:** Further refinement of filter-grained and atom-level scheduling mechanisms that can adapt not only to static resource constraints but also to dynamic fluctuations in network and processing loads. This includes the integration of multi-agent reinforcement learning strategies for coordinated scheduling among distributed edge devices.
 - **Benchmarking and Standardization:** Establishing robust benchmark standards and comprehensive datasets for evaluating adaptive deep learning architectures on edge devices. A standardized set of performance metrics and test scenarios will ensure that future approaches are evaluated consistently and compared fairly.
- Automatic Hardware Mapping:** Developing automated tools and algorithms for profiling hardware characteristics and mapping deep learning tasks seamlessly onto heterogeneous edge architectures. Improved codesign methodologies can minimize the gap between algorithm performance and hardware capabilities.

In summary, while significant progress has been made in designing adaptive deep learning architectures for real-time data streaming at the edge, the field still faces numerous technical and practical challenges. Overcoming these challenges will require interdisciplinary collaboration across deep learning, hardware design, and systems engineering.

VII. CONCLUSION

Adaptive deep learning architectures promise to revolutionize the field of edge computing by enabling efficient real-time data processing in resource-constrained environments. Through techniques such as master-surrogate modeling, context-adaptive DNN atom partitioning, and distributed inference with fused layer partitioning, researchers have demonstrated significant improvements in accuracy, latency reduction, resource utilization, and energy efficiency. The integration of reinforcement learning for energy-efficient computation offloading further enhances the robustness and practicality of these systems.

Key insights from the research include:

- **Adaptive Remodeling:** Systems like ElasticDNN dynamically adjust model parameters by utilizing master-surrogate networks coupled with filter-grained scheduling, achieving over 23% accuracy improvements and substantial reductions in adaptation time.
- **Modular Architecture:** Frameworks such as AdaMEC pre-partition deep neural networks into modular atoms, which can be recombined based on runtime context, resulting in latency reductions of up to 62% and significant memory savings.
- **Distributed Inference:** Fused layer partitioning and distributed inference models reduce communication overhead and inference latency, making them suitable for real-time applications in surveillance, autonomous vehicles, and IoT clusters.
- **Energy Efficiency:** Reinforcement learning and DVFS approaches have shown reductions in energy consumption by up to 53%, which is critical for battery-powered and resource-constrained edge devices.
- **Future Directions:** Promising avenues such as federated learning, neuromorphic computing, and automatic hardware mapping hold the potential to further advance adaptive deep learning systems on the edge.

The following bullet list summarizes the main findings:

- Adaptive deep learning models are essential for maintaining accuracy in dynamic, real-world environments where data distributions rapidly evolve.
- Master-surrogate and context-adaptive frameworks enable rapid model reconfiguration without the need for costly re-training on centralized servers.
- Distributed inference approaches reduce latency and communication overhead, making edge processing feasible for real-time applications.

- Energy efficiency and robust security remain critical challenges that next-generation adaptive architectures must address.
- Future research should focus on federated, continual, and neuromorphic learning approaches to realize the full potential of edge-based AI.

In conclusion, adaptive deep learning for real-time data streaming on the edge represents a transformative approach poised to meet the demands of modern applications across smart cities, healthcare, industrial automation, and autonomous systems. As research in this area continues to mature, it is anticipated that these adaptive architectures will deliver even higher performance, improved efficiency, and enhanced robustness, further bridging the gap between cloud-based intelligence and on-device processing.

REFERENCES

- [1] H. Gao *et al.*, "CuFSDAF: An Enhanced Flexible Spatiotemporal Data Fusion Algorithm Parallelized Using Graphics Processing Units," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, 2022, doi: 10.1109/TGRS.2021.3080384.
- [2] W. Jiang, D. Feng, Y. Sun, G. Feng, Z. Wang, and X. G. Xia, "Joint Computation Offloading and Resource Allocation for D2D-Assisted Mobile Edge Computing," *IEEE Trans Serv Comput*, vol. 16, no. 3, 2023, doi: 10.1109/TSC.2022.3190276.
- [3] B. Guo, S. C. Liu, Y. Liu, Z. G. Li, Z. W. Yu, and X. S. Zhou, "AIoT: The Concept, Architecture and Key Techniques," *Jisuanji Xuebao/Chinese Journal of Computers*, vol. 46, no. 11, 2023, doi: 10.11897/SP.J.1016.2023.02259.
- [4] B. Pang, E. Nijkamp, and Y. N. Wu, "Deep Learning With TensorFlow: A Review," 2020. doi: 10.3102/1076998619872761.
- [5] K. Filus and J. Domańska, "Software vulnerabilities in TensorFlow-based deep learning applications," *Comput Secur*, vol. 124, 2023, doi: 10.1016/j.cose.2022.102948.
- [6] P. Tam, S. Math, C. Nam, and S. Kim, "Adaptive Resource Optimized Edge Federated Learning in Real-Time Image Sensing Classifications," *IEEE J Sel Top Appl Earth Obs Remote Sens*, vol. 14, 2021, doi: 10.1109/JSTARS.2021.3120724.
- [7] P. Dai, F. Song, K. Liu, Y. Dai, P. Zhou, and S. Guo, "Edge Intelligence for Adaptive Multimedia Streaming in Heterogeneous Internet of Vehicles," *IEEE Trans Mob Comput*, vol. 22, no. 3, 2023, doi: 10.1109/TMC.2021.3106147.
- [8] N. Kumar and A. Ahmad, "Quality of service-aware adaptive radio resource management based on deep federated Q-learning for multi-access edge computing in beyond 5G cloud-radio access network," *Transactions on Emerging Telecommunications Technologies*, vol. 34, no. 6, 2023, doi: 10.1002/ett.4762.
- [9] A. Sacco, M. Flocco, F. Esposito, and G. Marchetto, "An architecture for adaptive task planning in support of IoT-based machine learning applications for disaster scenarios," *Comput Commun*, vol. 160, 2020, doi: 10.1016/j.comcom.2020.07.011.
- [10] L. Zang, X. Zhang, and B. Guo, "Federated Deep Reinforcement Learning for Online Task Offloading and Resource Allocation in WPC-MEC Networks," *IEEE Access*, vol. 10, 2022, doi: 10.1109/ACCESS.2022.3144415.
- [11] R. Adamec, "Does long term potentiation in periaqueductal gray (PAG) mediate lasting changes in rodent anxiety-like behavior (ALB) produced by predator stress? - Effects of low frequency stimulation (LFS) of PAG on place preference and changes in ALB produced by predator stress," *Behavioural Brain Research*, vol. 120, no. 2, 2001, doi: 10.1016/S0166-4328(00)00366-1.
- [12] X. Chen, D. Z. Chen, Y. Han, and X. S. Hu, "MoDNN: Memory Optimal Deep Neural Network Training on Graphics Processing Units," *IEEE Transactions on Parallel and Distributed Systems*, vol. 30, no. 3, 2019, doi: 10.1109/TPDS.2018.2866582.
- [13] Z. Zhao, K. M. Barijough, and A. Gerstlauer, "DeepThings: Distributed adaptive deep learning inference on resource-constrained IoT edge clusters," in *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2018. doi: 10.1109/TCAD.2018.2858384.
- [14] L. Zhang, J. Wang, W. Wang, Z. Jin, Y. Su, and H. Chen, "Smart contract vulnerability detection combined with multi-objective detection," *Computer Networks*, vol. 217, 2022, doi: 10.1016/j.comnet.2022.109289.

-
- [15] B. Pang *et al.*, "AdaMEC: Towards a Context-adaptive and Dynamically Combinable DNN Deployment Framework for Mobile Edge Computing," *ACM Trans Sens Netw*, vol. 20, no. 1, 2023, doi: 10.1145/3630098.
- [16] X. Zhu, T. Zhang, J. Zhang, B. Zhao, S. Zhang, and C. Wu, "Deep reinforcement learning-based edge computing offloading algorithm for software-defined IoT," *Computer Networks*, vol. 235, 2023, doi: 10.1016/j.comnet.2023.110006.
- [17] D. Liu, H. Kong, X. Luo, W. Liu, and R. Subramaniam, "Bringing AI to edge: From deep learning's perspective," *Neurocomputing*, vol. 485, 2022, doi: 10.1016/j.neucom.2021.04.141.
- [18] S. Petrocchi, G. Giorgi, and M. G. C. A. Cimino, "A Real-Time Deep Learning Approach for Real-World Video Anomaly Detection," in *ACM International Conference Proceeding Series*, 2021. doi: 10.1145/3465481.3470099.
- [19] Z. He and H. Sayadi, "Image-Based Zero-Day Malware Detection in IoMT Devices: A Hybrid AI-Enabled Method," in *Proceedings - International Symposium on Quality Electronic Design, ISQED*, 2023. doi: 10.1109/ISQED57927.2023.10129348.
- [20] R. Vengaloor and R. Muralidhar, "Deep Learning Based Feature Discriminability Boosted Concurrent Metal Surface Defect Detection System Using YOLOv-5s-FRN," *International Arab Journal of Information Technology*, vol. 21, no. 1, 2024, doi: 10.34028/iajit/21/1/9.
- [21] J. J. Chen *et al.*, "How to do Deep Learning on Graphs with Graph Convolutional Networks," *IEEE Access*, vol. 8, no. 1, 2019.
- [22] M. Peng, W. Zhang, F. Li, Q. Xue, J. Yuan, and P. An, "Weed detection with Improved Yolov 7," *EAI Endorsed Transactions on Internet of Things*, vol. 9, no. 3, 2023, doi: 10.4108/eetiot.v9i3.3468.
- [23] S. Yin, Y. Jiao, C. You, M. Cai, T. Jin, and S. Huang, "Reliable adaptive edge-cloud collaborative DNN inference acceleration scheme combining computing and communication resources in optical networks," *Journal of Optical Communications and Networking*, vol. 15, no. 10, 2023, doi: 10.1364/JOCN.495765.
- [24] M. Goudarzi, M. Palaniswami, and R. Buyya, "A Distributed Deep Reinforcement Learning Technique for Application Placement in Edge and Fog Computing Environments," *IEEE Trans Mob Comput*, vol. 22, no. 5, 2023, doi: 10.1109/TMC.2021.3123165.
- [25] A. M. Alabdali, "A Novel Framework of an IOT-Blockchain-Based Intelligent System," *Wirel Commun Mob Comput*, vol. 2022, 2022, doi: 10.1155/2022/4741923.
- [26] S. Garg, K. Kaur, G. Kaddoum, P. Garigipati, and G. S. Aujla, "Security in IoT-Driven Mobile Edge Computing: New Paradigms, Challenges, and Opportunities," *IEEE Netw*, vol. 35, no. 5, 2021, doi: 10.1109/MNET.211.2000526.
- [27] P. Salva-Garcia, J. M. Alcaraz-Calero, Q. Wang, J. B. Bernabe, and A. Skarmeta, "5G NB-IoT: Efficient Network Traffic Filtering for Multitenant IoT Cellular Networks," *Security and Communication Networks*, vol. 2018, 2018, doi: 10.1155/2018/9291506.
- [28] Muneera Altayeb and Amani Al-Ghraibah, "Arduino Based Real-Time Face Recognition And Tracking System," *International Journal of Advanced Trends in Computer Science and Engineering*, vol. 12, no. 4, pp. 144–150, Aug. 2023, doi: 10.30534/ijatcse/2023/011242023.
- [29] H. B. Ahmad, R. R. Asaad, S. M. Almufti, A. A. Hani, A. B. Sallow, and S. R. M. Zeebaree, "SMART HOME ENERGY SAVING WITH BIG DATA AND MACHINE LEARNING," *Jurnal Ilmiah Ilmu Terapan Universitas Jambi*, vol. 8, no. 1, pp. 11–20, May 2024, doi: 10.22437/jiituj.v8i1.32598.
- [30] S. M. Almufti, R. B. Marqas, Z. A. Nayef, and T. S. Mohamed, "Real Time Face-mask Detection with Arduino to Prevent COVID-19 Spreading," *Qubahan Academic Journal*, vol. 1, no. 2, pp. 39–46, Apr. 2021, doi: 10.48161/qaj.v1n2a47.
- [31] R. Asaad, R. Ismail Ali, and S. Almufti, "Hybrid Big Data Analytics: Integrating Structured and Unstructured Data for Predictive Intelligence," *Qubahan Techno Journal*, vol. 1, no. 2, Apr. 2022, doi: 10.48161/qtj.v1n2a14.
- [32] R. Rajab Asaad, R. Ismael Ali, A. Ahmad Shaban, and M. Shamal Salih, "Object Detection using the ImageAI Library in Python," *Polaris Global Journal of Scholarly Research and Trends*, vol. 2, no. 2, pp. 1–9, Apr. 2023, doi: 10.58429/pgjsrt.v2n2a143.
- [33] R. Kumar, P. Kumar, R. Tripathi, G. P. Gupta, S. Garg, and M. M. Hassan, "BDTwin: An Integrated Framework for Enhancing Security and Privacy in Cybertwin-Driven Automotive Industrial Internet of Things," *IEEE Internet Things J*, vol. 9, no. 18, 2022, doi: 10.1109/JIOT.2021.3122021.

-
- [34] S. Mukherjee, S. Gupta, O. Rawley, and S. Jain, "Leveraging big data analytics in 5G-enabled IoT and industrial IoT for the development of sustainable smart cities," *Transactions on Emerging Telecommunications Technologies*, vol. 33, no. 12, 2022, doi: 10.1002/ett.4618.
- [35] S. Agrawal, B. K. Patle, and S. Sanap, "A systematic review on metaheuristic approaches for autonomous path planning of unmanned aerial vehicles," Jan. 01, 2024, *Canadian Science Publishing*. doi: 10.1139/dsa-2023-0093.
- [36] A. Ghasemi and M. Mirzavand, "Robot path planning using Big Bang–Big Crunch algorithm," *Rob Auton Syst*, vol. 62, no. 3, pp. 390–399, 2014.
- [37] D. A. Majeed *et al.*, "DATA ANALYSIS AND MACHINE LEARNING APPLICATIONS IN ENVIRONMENTAL MANAGEMENT," *Jurnal Ilmiah Ilmu Terapan Universitas Jambi*, vol. 8, no. 2, pp. 398–408, Sep. 2024, doi: 10.22437/jiituj.v8i2.32769.
- [38] X. Wang, Y. Han, V. C. M. Leung, D. Niyato, X. Yan, and X. Chen, "Convergence of Edge Computing and Deep Learning: A Comprehensive Survey," 2020. doi: 10.1109/COMST.2020.2970550.
- [39] D. A. Hasan, S. R. M. Zeebaree, M. A. M. Sadeeq, H. M. Shukur, R. R. Zebari, and A. H. Alkhayyat, "Machine Learning-based Diabetic Retinopathy Early Detection and Classification Systems - A Survey," in *1st Babylon International Conference on Information Technology and Science 2021, BICITS 2021*, Institute of Electrical and Electronics Engineers Inc., 2021, pp. 16–21. doi: 10.1109/BICITS51482.2021.9509920.
- [40] S. M. Almufti and A. M. Abdulazeez, "An Integrated Gesture Framework of Smart Entry Based on Arduino and Random Forest Classifier," *Indonesian Journal of Computer Science*, vol. 13, no. 1, Feb. 2024, doi: 10.33022/ijcs.v13i1.3735.
- [41] M. Cao, Y. Li, X. Wen, Y. Zhao, and J. Zhu, "Energy-aware intelligent scheduling for deadline-constrained workflows in sustainable cloud computing," *Egyptian Informatics Journal*, vol. 24, no. 2, 2023, doi: 10.1016/j.eij.2023.04.002.