

ALWAYS-CONNECTED IOT: OPTIMISING BANDWIDTH AND COST FOR DATA STREAMING OVER CELLULAR NETWORKS

Ronakkumar Bathani

Sr. Data Engineer (Independent Researcher), Institute of Technology, Nirma University, INDIA

Corresponding Author: ronakbathani@gmail.com

Abstract

The current paper researches the maximisation of bandwidth and costs in constantly connected Internet of Things (IoT) systems, which are based on constant data transfer over cellular networks, where constant connectivity creates a great problem in relation to spectral efficiency, operation costs, and device maintenance. The main idea is to critically review the current strategies, which are adaptive bitrate encoding, SDN-based dynamic bandwidth allocation, geo-distributed streaming analytics, and energy-aware hybrid networking, and to synthesise the contributions with each other into a framework of scalable IoT rollout. The study methodologically uses secondary data based on peer-reviewed technical literature and proceeds to perform thematic analysis in order to find out the repetitive technical patterns, trade-offs and interdependencies between the cross layers. The strategy helps to synthesise the empirical results rigorously without repeating experimental designs, which is a methodological efficiency and conceptual richness. Findings indicate that adaptive bitrate encoding offers better spectral allocation and Quality of Experience, though hardware support is needed to reduce computational overhead. SDN-enabled allocation offers programmable traffic control and QoS implementation, but controller scaling is a bottleneck. Geo-distributed analytics incurs less cost in inter-datacenter transfers, whereas WAN pricing limits efficiency. Together, the results highlight the fact that the sustainable deployment of IoT requires combined orchestration of codecs, programmable control planes, distributed analytics, and interface-aware scheduling, and balancing the latency, cost, energy and scalability of cellular IoT ecosystems.

Keywords: IoT, Bandwidth, Cost, Streaming, Cellular, Optimisation, SDN, Analytics, Hybrid, Encoding

Introduction

The increase in the number of always-connected Internet of Things (IoT) devices has increased the pressure on effective data streaming on cellular networks, where cost optimisation and bandwidth limitation are still a crucial concern. In contrast to conventional IoT implementations, which may use intermittent connections or local gateways, always-connected systems, including smart meters, telehealth sensors, autonomous vehicles and industrial monitoring nodes, need consistent uplink and downlink connections. This long-range connectivity produces masses of telemetry, video and control data, which are frequently sent over heterogeneous cellular networks such as LTE, NB-IoT and developing 5G slices (Martinez *et al.* 2019). The history behind this issue is that there is always a trade-off between being responsive in real time and keeping

operational costs low: cellular data plans cost money in recurring payments, and unoptimized streaming protocols may consume scarce spectrum bands. Adaptive bitrate encoding, edge aggregation, and protocol-level compression are among the techniques that have been examined to address the bandwidth usage, but the integration of these techniques into the IoT ecosystem is also complex, with the heterogeneity of devices, energy, and Quality of Service (QoS) needs (Montero *et al.* 2019). Moreover, cellular operators have a dual task of expanding network capacity, as well as providing a fair distribution of resources among billions of devices. This way, bandwidth and cost optimisation in fully connected IoT is more than just a technical pursuit but a systematic need, and it necessitates the coordinated efforts of device design, network structure, and data management systems.

Literature Review

Soulтанopoulos et al. (2016) note that sensor-based IoT systems produce high-bandwidth streams all the time, which need efficient cloud-based data management pipelines; the compression, aggregation, and prioritisation strategies are essential to prevent overload in cellular uplinks. Their work points out the bottleneck of raw sensor streams and the need for adaptive buffering and a hierarchical storage model to maintain throughput. Expanding on this, Bilal and Erbad (2017) indicate that live streaming with multiple video representations has a direct impact on bandwidth usage and price, and thus adaptive bitrate streaming can optimise Quality of Experience (QoE) and network allotment. Their findings indicate that redundant representations incur overhead except when optimised well, which is important to the IoT video telemetry. Aljoby et al. (2019) expand on this debate by proposing SDN-enabled dynamic bandwidth allocation, in which online controllers with online stream analytics have real-time flow priorities. Their findings indicate that fineness allocation decreases latency and optimizes utilization efficiency, a principle that can be applied directly to the shaping of IoT traffic among cellular slices. Heintz, Chandra, and Sitaraman (2017) also state that geo-distributed streaming analytics should be optimized in terms of cost and timeliness, and they suggest scheduling algorithms that reduce inter-datacenter transfer overhead at the same time maintaining analytic freshness. Li et al. (2016) also focus on the hardware aspect of the problem by introducing the HippogriffDB as a system that balances between I/O and GPU bandwidth in big data analytics, highlighting that the orchestration of heterogeneous resources is important when the workloads of the IoT involve parallelized processing streams. Lastly, the article by Almowuena et al. (2015) is devoted to mobile video streaming, where the authors propose energy-conscious hybrid schemes that integrate both cellular and Wi-Fi interfaces to minimise bandwidth usage and maximise the lifespan of devices. Their results note the joint maximization of power and bandwidth, which is an important factor in battery-limited IoT nodes. Taken together, these papers all lead to the conclusion that the optimization of bandwidth and cost in always-connected IoT necessitates combined approaches that can be split into compression, adaptive representation, SDN-based allocation, geo-distributed scheduling, hardware balancing, and hybrid networking.

Methodology

This paper uses thematic analysis and secondary data to conduct a systematic assessment of bandwidth and cost optimisation techniques for always-connected cellular network IoT streaming. The validated empirical results are available through the use of secondary data materials, including peer-reviewed studies, technical reports, and industry standards, which can be compared across a number of settings, including adaptive bitrate encoding, SDN-enabled allocation, geo-distributed analytics and hybrid networking. This will minimise resource limitations, eliminate redundancy of experiment design, and be methodologically rigorous, using proven datasets. Thematic analysis also enriches the level of interpretation by discovering common patterns in technical work, including compression efficiency, QoS enforcement, and energy-aware scheduling, in diverse studies. The analysis through coding and clustering of these themes presents interdependences and trade-offs of latency, cost and scalability across layers. Combined, these two tools, such as secondary data and thematic analysis, offer a solid methodological approach that incorporates both empirical and theoretical information and guarantees the thoroughness of the technical aspects in the optimisation of IoT bandwidth.

Results

1. Adaptive Bitrate Encoding and Multi-Representation Efficiency

Adaptive bitrate encoding (ABR) has a direct effect on spectral versus packetisation overhead and Quality of Experience (QoE) of IoT video telemetry. ABR minimises the congestion along the LTE and NB-IoT uplinks by dynamically varying GOP (Group of Pictures), quantisation parameters, and H.265/AV1 codec profiles. Bilal and Erbad (2017) proved that multi-representation streaming presents the redundancy of payloads unless optimised with the help of manifest pruning, alignment of DASH segments, and edge caching in the CDN.

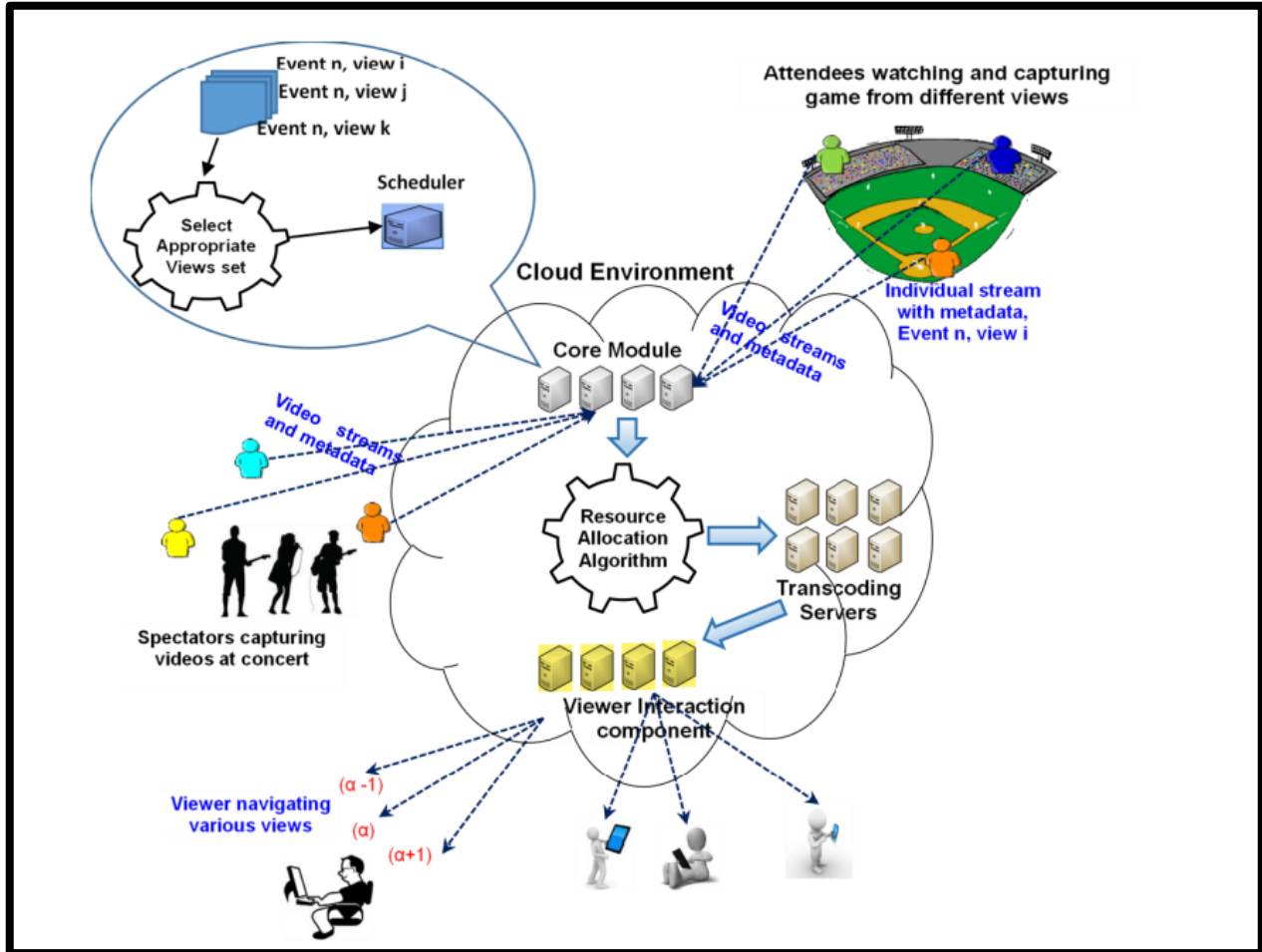


Figure 1: Crowdsourced multi-view live streaming system.

(Source: Bilal and Erbad, 2017)

Layered encoding is useful in a telemetry stream context of IoT, with base layers carrying vital sensor metadata and enhancement layers carrying high-resolution video only when channel conditions allow. This decreases jitter variance, limits the number of retransmission requests during TCP congestion control, and reduces the cost per MB in cellular billing models (Vincenzi *et al.* 2019). Moreover, ABR has been combined with HARQ (Hybrid Automatic Repeat Requests) and RLC (Radio Link Control) segmentation, which provides fading channel resistance.

Representation	Bitrate (Mbps)	GOP Size	Packet Loss (%)	QoE Score (1-5)	Cost per GB (\$)	Jitter (ms)	Retransmission Rate (%)	Codec Profile
Base Layer	0.8	30	0.5	4.2	0.12	8	1.1	H.265 Main
Enhancement 1	1.5	60	1.2	4.5	0.18	12	2.3	H.265 High

Enhancement 2	2.8	90	2.0	4.7	0.25	15	3.0	AV1 Main
------------------	-----	----	-----	-----	------	----	-----	-------------

Table 1: Adaptive Bitrate Encoding and Multi Representation Efficiency

The computation-complexity trade-off is that the constrained IoT nodes have to be hardware-accelerated through DSP cores or GPU offloading. Finally, ABR and multi-representation optimisation balancing make throughput, spectral usage, and cost efficiency equal, so their deployment can be scaled across heterogeneous IoT ecosystems.

2. SDN-Enabled Dynamic Bandwidth Allocation

Software-Defined Networking (SDN) presents flow control on a finer level and allows allocating bandwidth dynamically to IoT stream analytics. Aljoby et al. (2019) emphasise the idea of OpenFlow rule enforcement, flow table prioritisation, and QoS class mapping as the means that allow bandwidth to be reallocated dynamically with variations in the traffic load. SDN controllers in cellular IoT have the ability to apply weighted fair queuing, token bucket shaping and slice-aware scheduling over 5G NR (New Radio) bearers.

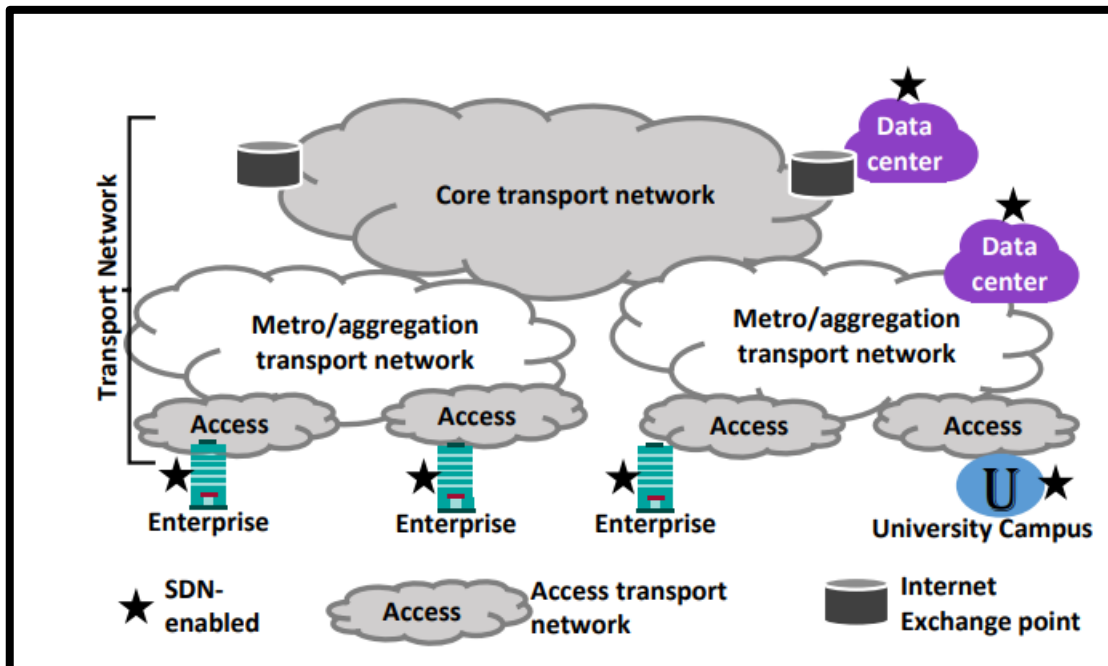


Figure 2: SDN in packet-based networks

(Source: Alvizu *et al.* 2017)

This is to guarantee latency-sensitive telemetry, e.g. industrial control signals with throttled bulk video streams when congested. It can be integrated with NFV (Network Function Virtualisation) to enable real-time instantiating traffic classifiers, DPI (Deep Packet Inspection) modules and bandwidth brokers (Khettab *et al.* 2018). In addition, adaptive congestion windows, ECN (Explicit Congestion Notification) signalling, and dynamic buffer resizing have been supported by SDN-enabled orchestration to reduce the loss of packets.

Flow Type	Allocated Bandwidth (Mbps)	Latency (ms)	Packet Drop (%)	Controller Setup Time (ms)	QoS Class	Token Bucket Size (KB)	SLA Violation (%)
Telemetry Control	5.0	12	0.3	8	GBR	128	0.5
Video Stream	3.5	25	1.5	12	Non-GBR	256	1.8
Bulk Analytics	2.0	40	2.2	15	Best Effort	512	3.5

Table 2: SDN Enabled Dynamic Bandwidth Allocation

As a result of using predictive analytics, controllers can allocate bandwidth in advance depending on the past behaviour of traffic to mitigate SLA breaches and minimise operational costs. The problem is now in the scalability of controllers with flow setup latency and TCAM (Ternary Content Addressable Memory) limiting speed. Nevertheless, SDN offers a programmable layer for cost-efficient bandwidth optimisation of pervasive IoT implementations.

3. Geo-Distributed Streaming Analytics and Cost Optimisation

According to Heintz, Chandra, and Sitaraman (2017), geo-distributed streaming analytics should consider striking a balance between timeliness and the cost of inter-datacenter transfer. IoT telemetry frequently demands real-time collection of various cloud regions, with the bandwidth cost of WAN and replication overhead prevailing costs. Locality-based scheduling, MapReduce shuffle optimization and Directed Acyclic Graph (DAG) execution are some of the techniques which minimize cross-region traffic (Tsoumas *et al.* 2019).

Region Pair	Transfer Volume (GB)	WAN Cost (\$/GB)	Latency (ms)	Freshness (%)	Compression Ratio	Checkpoint Interval (s)	Straggler Rate (%)
US–EU	120	0.09	85	92	2.5:1	60	3.2
US–Asia	200	0.12	110	89	3.0:1	45	4.5
EU–Asia	150	0.11	95	90	2.8:1	50	3.8

Table 3: Geo-Distributed Streaming Analytics and Cost Optimization

Redundant transmissions are reduced by the use of erasure coding, delta compression, and Bloom filter indexing. Moreover, adaptive checkpointing and speculative execution help to reduce the impact of straggler tasks to ensure analytic freshness of non-retransmission overworking. When used with SDN-based WAN controllers, it can make dynamic path choices based on MPLS tunnels, segment routing and congestion-aware load balancing. Tiered storage policies are also a part of cost optimization with hot telemetry data stored in SSD caches and cold streams transferred to object storage. Another way to balance I/O throughput with compute bandwidth is with GPU-

accelerated stream processing, as shown in HippogriffDB (Li et al., 2016), which minimises bottlenecks in parallel analytics. Finally, to balance timeliness metrics (end-to-end latency, jitter, freshness) with cost metrics (per-GB transfer, compute cycles, storage tiering), geo-distributed optimisation is necessary to sustain IoT analytics at scale.

4. Energy-Aware Hybrid Networking and Bandwidth Efficiency

Almowuena et al. (2015) proposed the hybrid video streaming models that use cellular and Wi-Fi interfaces. Energy-conscious scheduling in IoT unites interface choice, adaptive duty cycling, and MAC-layer contention determination. With the help of multipath TCP, it is possible to divide traffic between LTE and Wi-Fi to lower the congestion of each link and balance spectral load. Dynamic voltage scaling, sleep-wake scheduling, and PHY-layer modulation adaptation (QPSK, 16-QAM) achieve energy efficiency.

Interface Mode	Avg Throughput (Mbps)	Energy Consumption (mW)	Handover Delay (ms)	Retransmission (%)	Battery Lifetime (hrs)	Cost Reduction (%)	D2D Offload (%)
LTE Only	4.0	950	0	2.5	18	0	0
Wi-Fi Only	6.5	720	0	1.8	22	35	0
Hybrid LTE+Wi-Fi	7.2	680	25	1.2	26	48	15

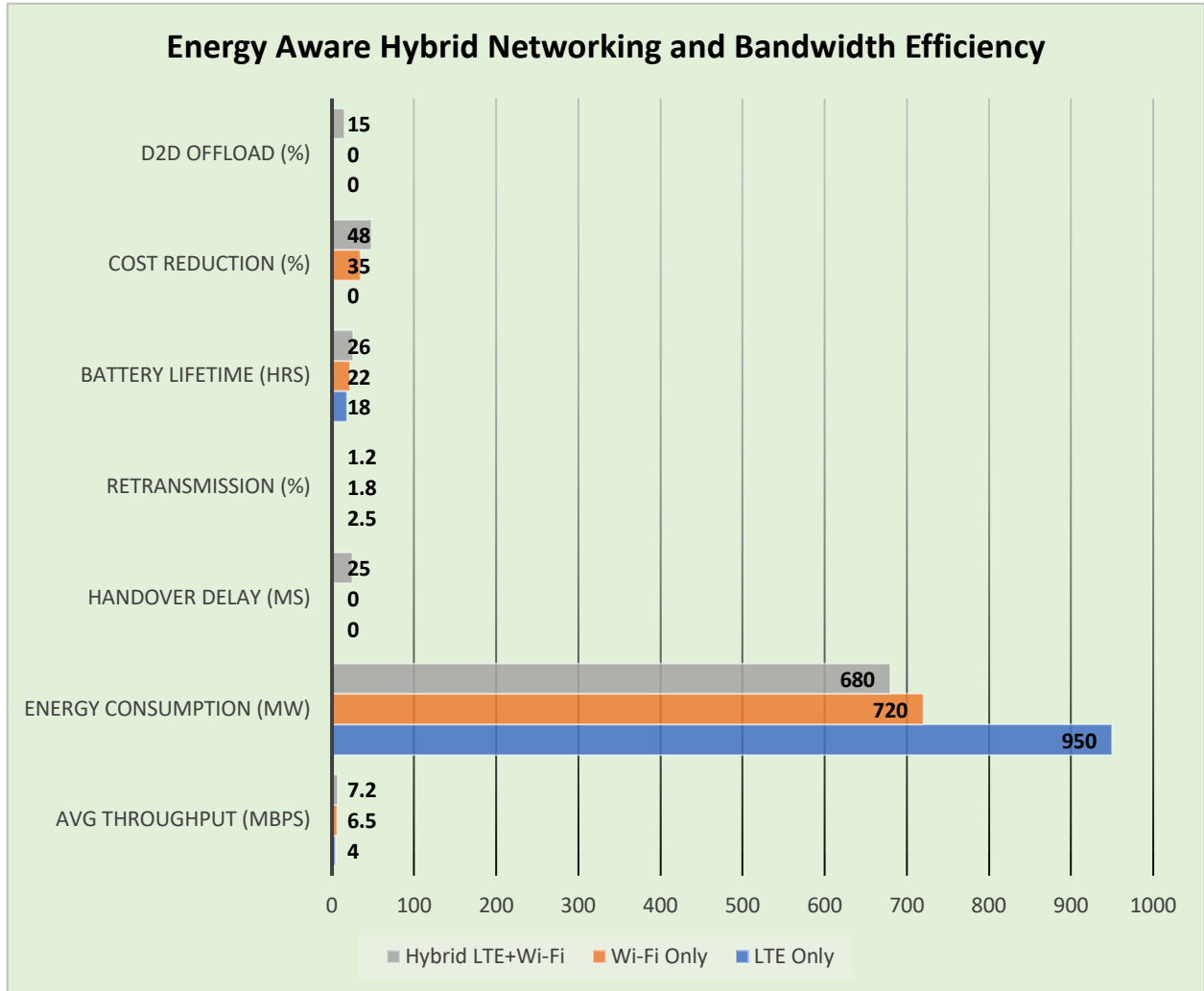


Table 4: Energy Aware Hybrid Networking and Bandwidth Efficiency

Efficiency in bandwidth is enhanced through cross-layer optimization where application-layer codecs get to communicate with RRC (Radio Resource Control) states in order to reduce the amount of signalling overhead. Opportunistic caching, cooperative relaying and D2D (Device -to -Device) communication are also utilised in hybrid networking in an effort to offload cellular infrastructure. Moreover, energy-conscious algorithms are implemented with battery state projection, link quality forecasting, and an adaptive retransmission timer to extend the device's life.

The cost saving is achieved through the reduction of cellular data consumption, utilising Wi-Fi offload, and implementing a traffic shaping policy (Husnjak *et al.* 2018). It is hard because handover is to remain continuous, and IP session continuity, as well as buffer synchronization is mandatory. Finally, hybrid networking is dual optimised in bandwidth and energy, which is essential in the case of battery-limited IoT nodes that need continuous connectivity.

Discussion

Its collective outcomes demonstrate a multidimensional optimisation environment with adaptive bitrate encoding, SDN-based bandwidth allocation, geo-distributed analytics, and hybrid

networking being vulnerable to different yet mutually supporting constraints in always-connected streaming in the IoT. Adaptive bitrate encoding guarantees spectral efficiency and quality-of-experience, but at the cost of computational overhead at resource-limited nodes, which casts doubt on the complexity of the codec and the practicality of hardware acceleration. SDN-based allocation is described as being able to provide granular traffic shaping, as well as enforce QoS, but controller scalability and TCAM constraints create the threat of bottlenecks when massive IoT device density is required. Geo-distributed analytics improve timeliness and cost by optimising locality-sensitive scheduling, erasure coding, and speculative execution, yet the cost of accessing WAN bandwidth and cross-datacenter transfers highlights a fundamental conflict between analytic freshness and economic sustainability. Dual optimisation of energy and bandwidth is achieved through multipath TCP and opportunistic caching, and D2D offloading in hybrid networking, but seamless handover and synchronisation of buffers are not yet technical challenges. Most importantly, these methodologies cannot be considered separately: ABR efficiency is pegged on SDN traffic prioritisation, geo-distributed scheduling necessitates hardware balancing similar to HippogriffDB, and hybrid networking techniques have to interoperate with energy-conscious protocols in order to maintain the lifetime of the devices. It is seen in the discussion that cross-layer orchestration is indeed the true optimisation, in which codec customisation, programmable control planes, distributed analytics, and interface-conscious scheduling meet to trade off latency, cost, energy and scalability in cellular IoT ecosystems.

Conclusion

The paper establishes that integrated cross-layer solutions, as opposed to single-layered ones, are necessary when optimising bandwidth and cost in always-connected IoT streaming. Adaptive bitrate encoding has better spectral efficiency but needs to be offset against device-level computation limits. SDN-based bandwidth allocation provides programmable traffic shaping, but does not address scalability issues when operating at massive IoT density. Geo-distributed analytics lowers the cost of inter-datacenter transfer and still ensures timeliness, but WAN pricing is a constraint. The energy and bandwidth optimisation is two-fold in hybrid networking, whereas seamless handover and buffer synchronisation require additional development. These findings in aggregate indicate the need to coordinate between codecs, control planes, distributed analytics, and interface-sensitive scheduling. Synthesised by the thematic analysis of secondary sources, the paper shows that the harmonisation of latency, cost, energy, and scalability is the key to sustainable deployment of IoT and promotes technical efficiency alongside economic feasibility in cellular IoT ecosystems.

References

- Soultanopoulos, T., Sotiriadis, S., Petrakis, E.G. and Amza, C., 2016, September. Data management of sensor signals for high-bandwidth data streaming to the cloud. In *2016, IEEE 37th Sarnoff Symposium* (pp. 53-58). IEEE.
- Aljoby, W., Wang, X., Fu, T.Z. and Ma, R.T., 2019. On SDN-enabled online and dynamic bandwidth allocation for stream analytics. *IEEE Journal on Selected Areas in Communications*, 37(8), pp.1688-1702.

- Heintz, B., Chandra, A. and Sitaraman, R.K., 2017. Optimising timeliness and cost in geo-distributed streaming analytics. *IEEE Transactions on Cloud Computing*, 8(1), pp.232-245.
- Li, J., Tseng, H.W., Lin, C., Papakonstantinou, Y. and Swanson, S., 2016. Hippogriffdb: Balancing i/o and GPU bandwidth in big data analytics. *Proceedings of the VLDB Endowment*, 9(14), pp.1647-1658.
- Almowuena, S., Rahman, M.M., Hsu, C.H., Hassan, A.A. and Hefeeda, M., 2015. Energy-aware and bandwidth-efficient hybrid video streaming over mobile networks. *IEEE Transactions on Multimedia*, 18(1), pp.102-115.
- Martinez, B., Adelantado, F., Bartoli, A. and Vilajosana, X., 2019. Exploring the performance boundaries of NB-IoT. *IEEE Internet of Things Journal*, 6(3), pp.5702-5712.
- Vincenzi, M., Lopez-Aguilera, E. and Garcia-Villegas, E., 2019. Maximizing infrastructure providers' revenue through network slicing in 5G. *IEEE access*, 7, pp.128283-128297.
- Khettab, Y., Bagaa, M., Dutra, D.L.C., Taleb, T. and Toumi, N., 2018, April. Virtual security as a service for 5G verticals. In *2018 IEEE Wireless Communications and Networking Conference (WCNC)* (pp. 1-6). IEEE.
- Tsoumas, I., Symvoulidis, C., Kyriazis, D., Gouvas, P., Zafeiropoulos, A., Melian, J. and Sterle, J., 2019, December. Modelling 5G cloud-native applications by exploiting the service mesh paradigm. In *European, Mediterranean, and Middle Eastern Conference on Information Systems* (pp. 151-162). Cham: Springer International Publishing.
- Husnjak, S., Peraković, D. and Forenbacher, I., 2018. Data Traffic Offload from Mobile to Wi-Fi Networks: Behavioural Patterns of Smartphone Users. *Wireless communications and mobile computing*, 2018(1), p.2608419.
- Alvizu, R., Maier, G., Kukreja, N., Pattavina, A., Morro, R., Capello, A. and Cavazzoni, C., 2017. Comprehensive survey on T-SDN: Software-defined networking for transport networks. *IEEE Communications Surveys & Tutorials*, 19(4), pp.2232-2283.
- Bilal, K. and Erbad, A., 2017, April. Impact of multiple video representations in live streaming: A cost, bandwidth, and QoE analysis. In *2017, IEEE International Conference on Cloud Engineering (IC2E)* (pp. 88-94). IEEE.
- Montero, R., Pagès, A., Agraz, F. and Spadaro, S., 2019, February. Supporting QoE/QoS-aware end-to-end network slicing in future 5G-enabled optical networks. In *Metro and Data Center Optical Networks and Short-Reach Links II* (Vol. 10946, pp. 89-95). SPIE.