

NICoLE: Are In-Network LLM-Based Agents Cost-Feasible for RTP Video Streaming?

Alireza Shirmarz ^{ufscar}, Fábio Luciano Verdi ^{ufscar}, Gyanesh Patra ^{ericsson}, Gergely Pongrácz ^{ericsson}
{ashirmarz, verdi}@ufscar.br, {gyanesh.patra, gergely.pongracz}@ericsson.com

Abstract—We present NICoLE (Network-Integrated Codec-LLM Engine), an in-network decision layer for real-time RTP video streaming that assists encoder adaptation without modifying end-to-end congestion control. NICoLE correlates RTP telemetry (frame size and packet pacing) with queue dynamics at the bottleneck to detect mismatch between encoded rate and available service capacity. Based on this, it selects discrete encoder profiles (resolution, frame rate, GoP) and applies short-lived ECN-based queue prioritization. The system is implemented using WebRTC (GCC) and GStreamer. Experiments under controlled bottleneck conditions show reduced queue buildup, lower stall ratio, and improved frame pacing stability at high frame rates (90–120 FPS), while gains remain limited at lower rates. To address reliability concerns, NICoLE constrains the LLM to structured inputs and a bounded action space, ensuring deterministic behavior and avoiding free-form reasoning. We show that inference latency (~ 383 ms) remains below GCC reaction time (600–1500 ms), enabling timely in-network intervention. These results indicate that LLM-assisted in-network control is feasible and beneficial under high-density, delay-sensitive workloads, but not universally required.

Index Terms—AI Agent, Real Time video Streaming, L4S, QoE

I. INTRODUCTION

RTP-based real-time video streaming underpins interactive applications such as cloud gaming and XR [1]–[4]. These systems operate under tight latency and jitter budgets, where transient queue buildup inflates motion-to-photon delay, disrupts frame pacing, and degrades Quality of Experience (QoE), particularly at high frame rates (60–120 FPS) [5], [6].

Recent work shows that improving real-time video performance requires both faster congestion feedback loops and better encoder adaptation. Tight-loop congestion control reduces reaction latency [7], while encoder-side designs improve rate–distortion behavior under varying bandwidth [8]. However, both remain end-host driven and do not leverage in-network observations. As a result, congestion is still detected after queue buildup, creating a mismatch between real-time network conditions and encoder configuration.

During transient congestion, queue buildup distorts packet departure timing before end-to-end control converges, degrading frame pacing and QoE, especially in high-FPS workloads [5]. However, intervening within this adaptation window is challenging: due to encryption and processing constraints, network devices observe only RTP-level metadata and queue state, without access to codec structure (e.g., resolution, FPS,

GoP). This creates a synchronization gap between congestion onset and encoder reconfiguration.

Challenge 1: Bitrate-only control is insufficient. WebRTC’s Google Congestion Control (GCC) [9] continuously adjusts bitrate based on delayed receiver feedback, while encoder parameters (resolution, FPS, GoP) evolve on slower timescales. During transient congestion, bitrate reduction alone does not immediately align the stream with instantaneous service capacity, leading to suboptimal operation before encoder adaptation converges.

Challenge 2: Encoder reconfiguration incurs transient mismatch. Even when adaptation is triggered, encoder reconfiguration and signaling require non-negligible time. During this interval, queue buildup and pacing distortion persist. While AQM mechanisms such as CoDel [10] or ECN-enabled L4S DualPI2 [11], [12] reduce queue delay, they operate at flow granularity and do not support short-lived prioritization aligned with encoder adaptation windows.

To address these challenges, the network must infer encoder pacing and service capacity directly at the bottleneck. RTP headers expose frame boundaries, timestamps, and packet sizes, enabling extraction of inter-frame gap (IFG) and frame size (FS) without payload inspection [4]. Combined with queue telemetry, this enables early detection of rate mismatch.

By correlating encoded rate with observed service capacity, the network can trigger profile-level adaptation at congestion onset. During the adaptation interval, temporary ECT(1) marking can be applied to leverage L4S prioritization for short-timescale mitigation. Unlike persistent L4S marking, this approach is transient and bounded.

This multi-dimensional decision problem—mapping RTP telemetry and queue state to discrete adaptation actions—extends beyond simple threshold-based control. Rule-based logic requires careful tuning across heterogeneous encoder behaviors, while conventional ML models rely on fixed feature mappings. In contrast, compact LLMs can map structured telemetry inputs to a bounded action space without manual threshold tuning [13]–[15]. This design avoids free-form reasoning and ensures predictable behavior suitable for network control.

Building on this direction, we propose *NICoLE*, an in-network LLM-based agent that assists real-time RTP streaming through (i) profile-aware encoder adaptation and (ii) transient L4S steering, while preserving end-to-end congestion control semantics.

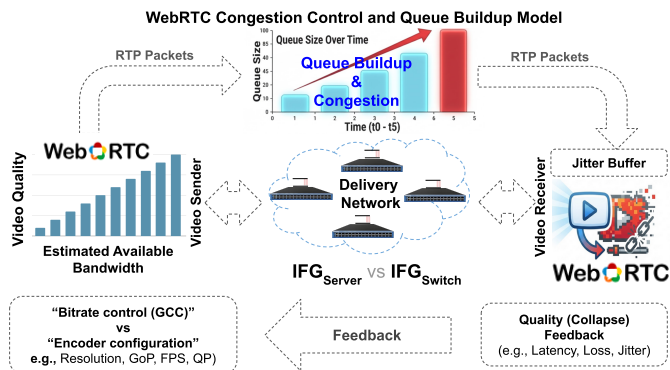


Fig. 1: WebRTC real-time streaming control loop.

NICoLE operates as a bounded decision layer that correlates RTP telemetry with queue dynamics. Upon transient congestion, it activates short-lived prioritization and selects a discrete encoder profile aligned with observed service capacity, complementing GCC’s bitrate control.

This work presents a cost-aware feasibility study of network-assisted encoder configuration using compact LLM agents. We evaluate whether such in-network intelligence delivers measurable QoE gains while respecting inference latency and deployment constraints.

Our contributions are as follows:

- **In-Network LLM Control Framework:** A bounded decision architecture linking RTP telemetry and queue state to encoder-level adaptation without modifying end-to-end congestion control;
- **Transient L4S Steering:** Short-timescale prioritization during encoder adaptation;
- **Profile-Aware Encoder Assistance:** Network-informed selection of discrete encoder profiles (resolution, FPS, GoP);
- **Cost-Benefit Evaluation:** CPU-based evaluation of QoE, queue behavior, and inference feasibility.

This work does not claim universal superiority of LLM-based control, but demonstrates its feasibility and benefit under high-density, delay-sensitive workloads.

II. BACKGROUND AND MOTIVATION

A. RTP Streaming and the WebRTC Control Loop

Interactive applications such as cloud gaming and XR commonly employ RTP over UDP to minimize latency and enable fine-grained playout control [1], [2], [16]. Encoded video frames (e.g., H.264) are packetized into RTP packets using standardized formats [17].

Fig. 1 shows the WebRTC control loop: RTP packets traverse shared network devices where queue buildup may occur; the receiver jitter buffer absorbs delay variation; and RTCP feedback drives end-to-end bitrate adaptation.

B. GCC Controls Bitrate; Encoder Controls Video Parameters

WebRTC commonly employs Google Congestion Control (GCC) [9], which estimates available bandwidth from receiver

feedback (e.g., TWCC delay trends [18]) and adjusts the sender bitrate.

Transport control and codec behavior are decoupled: *GCC regulates bitrate*, while the *encoder determines how bitrate is distributed* across resolution, frame rate (FPS), and quantization. Under congestion, bitrate reduction typically precedes profile adaptation, which occurs on slower timescales. If the encoded rate exceeds instantaneous service capacity, queue buildup and delay variation persist—an effect amplified at high frame rates.

C. Queue Buildup, L4S, and In-Network Control

Congestion manifests within network queues, yet control remains reactive: the sender detects congestion only after delay inflation is observed at the receiver via RTCP feedback [9], [18], [19].

L4S reduces persistent queue buildup through ECN-based dual-queue AQM (e.g., DualPI2) [11], [12]. However, it is typically applied at flow granularity via persistent ECT(1) marking, rather than as a short-lived mechanism aligned with transient encoder adaptation.

Recent learning-based approaches optimize queue management [13], [20], [21], but do not explicitly coordinate network observations with encoder-level adaptation.

In contrast, NICoLE uses L4S as a *transient* prioritization mechanism: ECT(1) marking is applied conditionally during encoder adaptation intervals, rather than persistently across the flow lifetime.

Motivation. These limitations motivate a bounded in-network mechanism that (i) infers encoder pacing and service capacity from RTP telemetry and queue state, (ii) assists discrete profile adaptation without modifying GCC semantics, and (iii) applies short-lived prioritization to suppress queue buildup during reconfiguration.

III. SYSTEM MODEL AND NICoLE ARCHITECTURE

A. System Overview

Fig. 2 shows NICoLE. A control-plane agent periodically processes RTP and queue telemetry and updates a per-flow action table. The dataplane enforces two bounded actions without per-packet inference: transient ECT(1)-based L4S steering during Δt and network-assisted encoder profile signaling.

B. Deployment Model

The system consists of a sender, a receiver, and NICoLE deployed in-network (Fig. 2). The sender performs video encoding and GCC-based bitrate adaptation, while the receiver maintains a jitter buffer and provides RTCP feedback.

NICoLE operates alongside GCC without modifying its control loop. The dataplane performs flow identification, RTP parsing, ECN/DSCP marking, queue classification, and telemetry extraction. Aggregated features are forwarded to the control plane, where the LLM is invoked periodically to update per-flow actions. Encoder profile updates are signaled via a WebRTC DataChannel.

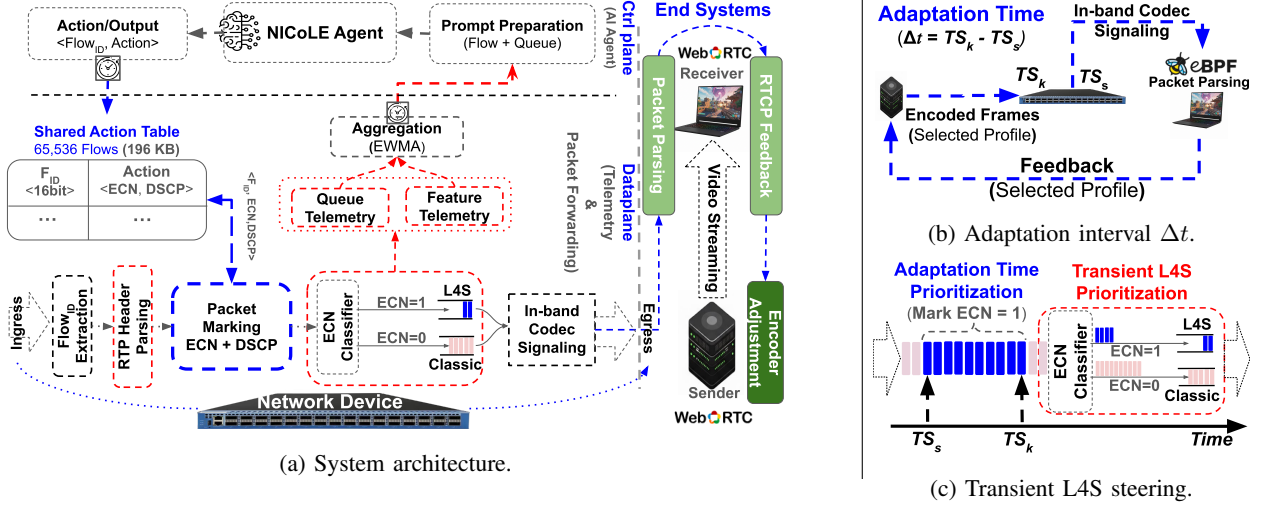


Fig. 2: NICoLE architecture. (a) Control/data-plane interaction. (b) Adaptation interval Δt . (c) Transient ECT(1)-based L4S steering during Δt .

C. Observation Space

NICoLE observes RTP telemetry and queue state:

- **RTP**: sequence number, timestamp, marker bit, packet size;
- **Pacing**: IFG_{server} and IFG_{switch} ;
- **Frame size**: mean RTP frame size;
- **Queue**: occupancy and ECN state.

Pacing distortion is defined as

$$\Delta_{IFG} = IFG_{switch} - IFG_{server}.$$

The encoded rate is approximated as

$$R_{enc} \approx \frac{FS}{IFG_{server}} \approx FS \times FPS.$$

Congestion arises when $R_{enc} > \mu$, where μ is the estimated service rate. Unlike GCC, NICoLE directly correlates encoder rate with in-network service capacity.

D. Action Space

NICoLE executes bounded actions:

- **Transient L4S**: conditionally mark packets ECT(1) during Δt to enable temporary low-latency queue admission, then revert after convergence;
- **Profile selection**: select the closest encoder profile aligned with μ , prioritizing FPS;
- **Signaling**: encode adaptation hints via DSCP toward the receiver.

At the receiver, DSCP markings are extracted (e.g., via eBPF) and forwarded to the sender through a WebRTC DataChannel.

E. Control Philosophy

Conventional WebRTC control reacts after congestion is observed at the receiver. NICoLE instead detects mismatch in-network and applies early profile adaptation with transient prioritization, stabilizing frame pacing during the adaptation interval.

F. Cost and Reliability

NICoLE uses an LLM as a bounded decision engine. Inputs are structured numeric telemetry, and outputs are restricted to a discrete action tuple (profile ID, L4S flag, DSCP). This eliminates free-form reasoning, prevents hallucination, and ensures deterministic behavior.

Inference is periodic and triggered only under sustained congestion, enabling cost-aware deployment in latency-sensitive environments.

IV. CONTROL ALGORITHM

NICoLE performs periodic control based on aggregated RTP and queue telemetry, detecting service-rate mismatch and triggering bounded profile adaptation without modifying GCC.

At each control interval:

- 1) **Telemetry aggregation**: RTP frame size (FS), sender pacing (IFG_{server}), switch inter-arrival gap (IFG_{switch}), and queue occupancy are smoothed using EWMA ($\alpha = 0.1$);
- 2) **Congestion detection**: congestion is declared when encoded rate R_{enc} exceeds service rate μ , or pacing distortion Δ_{IFG} exceeds threshold θ ;
- 3) **LLM decision**: a bounded LLM selects an encoder profile aligned with μ and determines whether transient L4S activation is required;
- 4) **Enforcement**: the dataplane applies ECT(1) marking during adaptation window Δt , updates DSCP hints, and maintains actions until convergence or congestion clears.

Profile transitions are rate-limited to prevent oscillation. NICoLE adjusts encoder operating points within GCC's bitrate budget.

A. Congestion Detection

For each flow f :

$$\Delta_{IFG}^{(f)} = IFG_{switch}^{(f)} - IFG_{server}^{(f)} \quad (1)$$

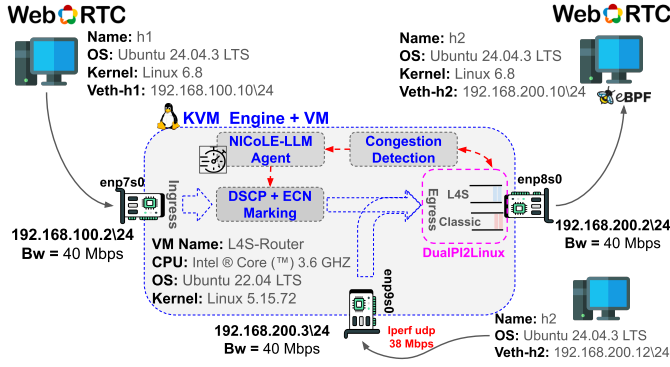


Fig. 3: Experimental testbed with an L4S-enabled router. The bottleneck link is shaped to 40 Mbps using `tc` with DualPI2.

Under uncongested conditions:

$$IFG_{switch} \approx IFG_{server}, \quad \Delta_{IFG} \approx 0$$

When $R_{enc} > \mu$:

$$IFG_{switch} > IFG_{server}, \quad \Delta_{IFG} > 0$$

Congestion onset is detected when:

$$Q^{(f)} > \theta_Q \quad \text{and} \quad \Delta_{IFG}^{(f)} > \theta_{IFG} \quad (2)$$

This joint condition distinguishes sustained congestion from transient jitter.

B. Profile Adaptation

Upon congestion, NICoLE selects a lower-capacity profile P_{k-1} from predefined set \mathcal{P} , prioritizing FPS and GoP stability while adapting spatial resolution to match service capacity. The selected profile remains consistent with GCC’s bitrate constraint.

When congestion subsides ($Q^{(f)} < \theta_{low}$), the profile is gradually restored under stable pacing.

C. Transient L4S Steering

Packets are conditionally marked ECT(1) during transient congestion to enable temporary L4S queue admission [4], [11], [12]. Unlike persistent flow-level L4S marking, this mechanism is short-lived and aligned with the encoder adaptation interval.

D. Control Loop and LLM Design

NICoLE operates with decision interval $T \approx 320$ ms. RTP and queue features (FS , IFG_{server} , IFG_{switch} , Δ_{IFG}) are aggregated into structured inputs, and the LLM outputs a bounded action tuple $(P_k, L4S_flag, DSCP_k)$ stored per flow and enforced in the dataplane.

Inputs are numeric and fixed-field; outputs are restricted to discrete actions. This design avoids free-form reasoning, ensures deterministic behavior, and prevents hallucination. Telemetry aggregation (~ 500 ms) aligns observation and inference latency.

TABLE I: Encoder profiles. FPS and GoP are fixed; only resolution is adapted.

Profile	Resolution	FPS	GoP
P0	3840×2160	Fixed	Fixed
P1	1920×1080	Fixed	Fixed
P2	1280×720	Fixed	Fixed
P3	640×360	Fixed	Fixed

V. EXPERIMENTAL SETUP

A. Testbed Architecture

Fig. 3 shows the topology. Two WebRTC endpoints communicate through a KVM-based router acting as the in-network device. The router runs Ubuntu 22.04 (kernel 5.15) with DualPI2 [12], traffic shaping via `tc`, and the NICoLE control plane.

Both endpoints use a GStreamer-based WebRTC stack for RTP/UDP control and DSCP signaling. At the receiver, DSCP values are extracted via eBPF and relayed to the sender through a WebRTC DataChannel, enabling network-assisted adaptation while preserving GCC semantics.

The bottleneck capacity is fixed at 40 Mbps.

B. LLM Model Selection

We evaluate compact LLMs (0.6B–8B parameters) under QLoRA fine-tuning and 4-bit quantization (GGUF) using `llama.cpp` on an 8-core 3.6 GHz CPU (8 GB RAM). Metrics include inference latency, throughput, and control accuracy.

As shown in Fig. 4, larger models (e.g., 8B) incur excessive latency (> 1 s), while smaller models reduce accuracy. LLaMA-1B provides the best trade-off, achieving sub-400 ms latency with high control accuracy, consistent with the control interval.

C. Traffic Scenario

A controlled single-flow setup is used to isolate transient congestion. RTP streams are generated at fixed frame rates (30–120 FPS) over a 40 Mbps bottleneck.

Background UDP traffic (38 Mbps via `iperf`) is injected for 60 s within a 120 s run to induce deterministic oversubscription. Higher frame rates increase encoded demand and amplify congestion sensitivity.

D. Profiles and Encoder Configuration

The encoder supports predefined spatial profiles (Table I). FPS and GoP are fixed, and NICoLE adapts only spatial resolution, while GCC regulates bitrate.

E. Metrics

We evaluate perceptual, transport, and network performance:

- **Perceptual:** VMAF;
- **Temporal:** FPS stability;
- **Rate:** bitrate dynamics;
- **Queue:** occupancy and ECN behavior;
- **Buffer:** jitter-buffer occupancy and drops;

- **Pacing:** Δ_{IFG} .

We compare **GCC-only** with **NICoLE + GCC + L4S**.

To contextualize the LLM, a rule-based baseline using threshold logic can be defined. While effective in static settings, such rules require manual tuning and generalize poorly; NICoLE replaces this with a bounded adaptive decision model.

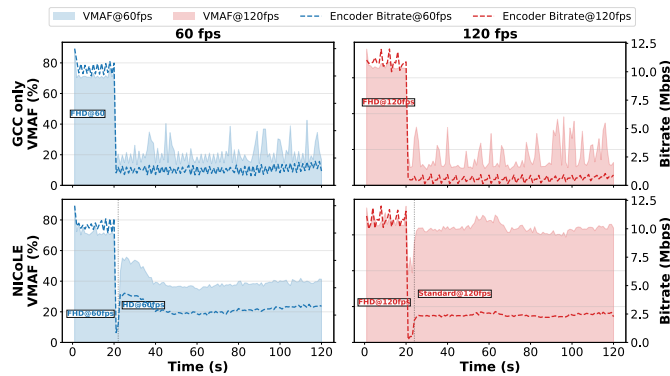
VI. EXPERIMENTAL RESULTS

A. LLM Inference Feasibility

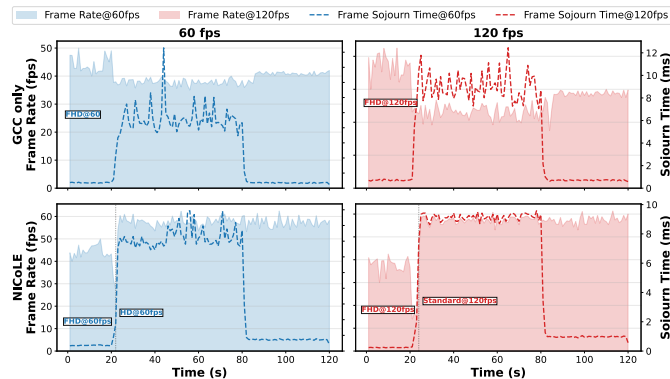
We evaluate whether in-network LLM control satisfies real-time constraints on commodity CPU hardware. Candidate models are tested under 4-bit quantization (GGUF) using `llama.cpp` on an 8-core 3.6 GHz CPU VM (8 GB RAM). Metrics include inference throughput, latency, and control accuracy.

Fig. 4 shows that smaller models provide higher throughput but lower accuracy, while larger models (e.g., 8B) incur excessive latency (>1 s). LLaMA-1B achieves 2.61 decisions/s with 383 ms latency and high accuracy, satisfying the 500–1000 ms control interval.

B. Flow Congestion Dynamics



(a) VMAF and bitrate (60 FPS, 120 FPS)



(b) Frame rate and queue delay

Fig. 5: Time-series comparison under a 40 Mbps bottleneck. NICoLE stabilizes queue, pacing, and quality.

We evaluate transient behavior during a 60 s congestion window within a 120 s run under a 40 Mbps bottleneck. Experiments span 30–120 FPS. Fig. 5 shows representative time-

TABLE II: Performance during congestion (mean \pm std over five runs).

FPS	Avg VMAF		Peak Delay (ms)	
	GCC	NICoLE	GCC	NICoLE
30	29.35 \pm 2.3	25.68 \pm 1.75	7.42 \pm 0.5	3.06 \pm 0.28
60	20.83 \pm 2.08	40.15 \pm 2.12	18.57 \pm 1.67	7.15 \pm 0.62
90	38.94 \pm 4.1	46.19 \pm 3.5	12.32 \pm 1.84	9.60 \pm 0.76
120	19.86 \pm 2.38	63.56 \pm 5.08	12.45 \pm 1.87	11.49 \pm 0.79

series for 60 and 120 FPS, while Table II reports mean \pm std over five runs.

Under GCC-only control, congestion is detected reactively via RTCP feedback, leading to delayed response, queue buildup, and bitrate oscillation. These effects intensify at higher frame rates due to shorter inter-frame intervals.

NICoLE detects congestion in-network using pacing distortion and queue state, enabling earlier intervention through profile adaptation and transient L4S steering.

Table II shows that benefits increase with temporal density. At 60 FPS, NICoLE reduces peak queue delay by 61% (18.57 \rightarrow 7.15 ms) and improves VMAF (20.83 \rightarrow 40.15). At 90 FPS, quality improves by \approx 18% with reduced delay variability. At 120 FPS, GCC collapses to 19.86 VMAF, while NICoLE maintains 63.56 VMAF with similar delay (11.49 vs 12.45 ms). Fig. 5 confirms that NICoLE stabilizes quality and pacing, whereas GCC exhibits oscillation and quantization collapse.

At 30 FPS, longer inter-frame intervals reduce congestion sensitivity. In this regime, GCC maintains stable performance and NICoLE provides limited benefit, reflecting a trade-off favoring delay stability over spatial fidelity.

C. QoE Stability: Stall and Jitter

We evaluate stall ratio, frame loss, and jitter-buffer stability across 30–120 FPS.

Stall ratio increases sharply with FPS under GCC, exceeding 50% at 120 FPS due to reactive bitrate control. NICoLE significantly reduces stalls at high frame rates by preserving temporal pacing through early profile adaptation.

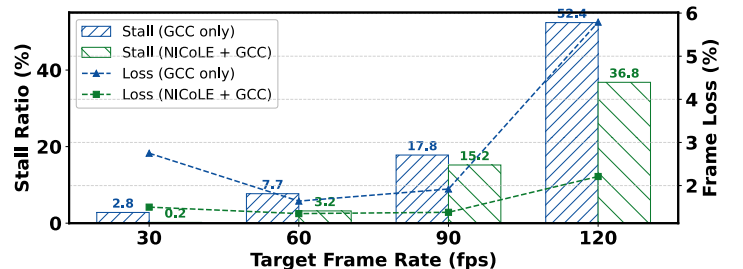


Fig. 6: Stall ratio and frame loss under congestion.

Fig. 7 shows that GCC exhibits large playback variance at high FPS, while NICoLE maintains higher median FPS with lower dispersion. Improvements are most pronounced at 90 and 120 FPS.

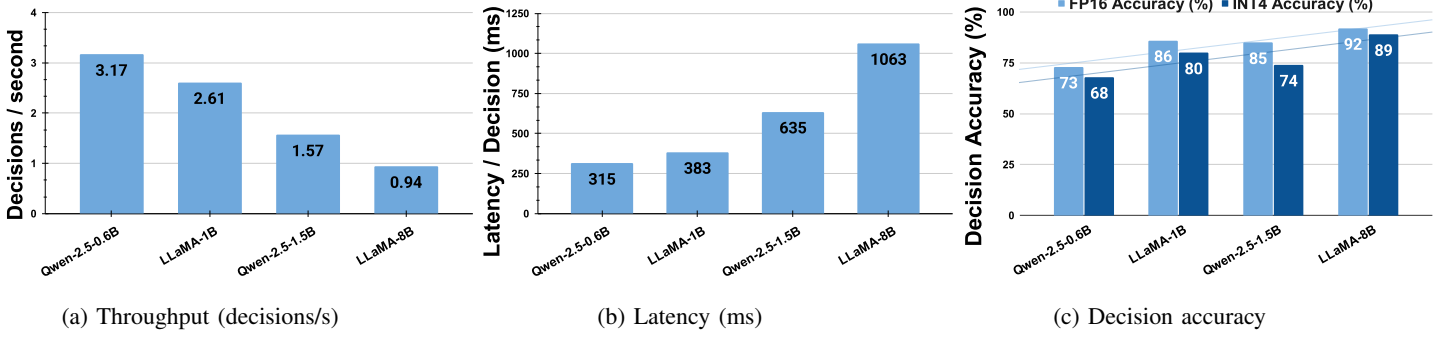


Fig. 4: CPU-only LLM evaluation under 4-bit quantization. LLaMA-1B achieves the best latency–accuracy trade-off.

These results confirm that NiCoLE improves both perceptual quality and temporal stability, particularly under high-density workloads.

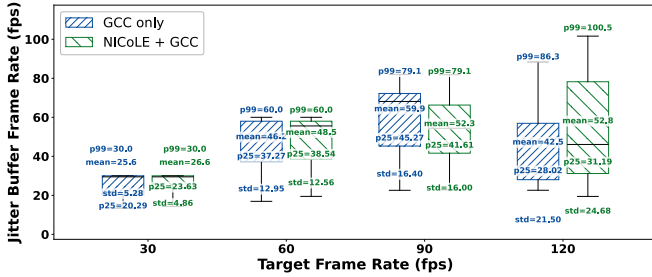


Fig. 7: Jitter buffer playback stability.

VII. DISCUSSION: COST–GAIN TRADE-OFF

An important question is when in-network LLM-assisted control justifies its computational cost. A simple feasibility condition is:

$$T_{inf} < T_{queue}$$

where T_{inf} is inference latency and T_{queue} is the effective GCC reaction time, including RTCP feedback and encoder response.

From Sections VI-A and V-B, $T_{inf} \approx 383$ ms for LLaMA-1B, while T_{queue} ranges from 600 to 1500 ms [9], [18], [19]. This enables proactive intervention before significant queue buildup.

Empirical results show that benefit depends on temporal density. At 30 FPS, GCC alone maintains stable performance. At 90–120 FPS, shorter inter-frame intervals amplify queue instability, and early profile adaptation significantly improves quality and delay stability, indicating a transition region between 60 and 90 FPS.

LLM-based control should be interpreted cautiously. The decision problem is low-dimensional and can be approximated using rule-based or lightweight ML methods. The advantage of LLMs lies in flexibility across heterogeneous conditions without manual tuning. NiCoLE therefore demonstrates conditional benefit rather than universal superiority.

A. Architectural Implications and Future Directions

Network-assisted codec adaptation improves temporal stability primarily under high-density workloads, but requires workload-aware deployment due to computational overhead.

Extending NiCoLE to multi-flow environments introduces fairness and resource-sharing challenges. Executing transformer inference directly in programmable dataplanes remains infeasible at line rate, though emerging systems such as FENIX [22] and Taurus [23] suggest potential paths toward in-network ML acceleration.

B. Limitations

The evaluation focuses on a single-flow scenario to isolate transient dynamics. In multi-flow settings, interactions between L4S and Classic queues may affect fairness and stability. The controlled bandwidth environment simplifies telemetry patterns.

The prototype uses CPU-based inference in a Linux VM; performance under higher concurrency and heterogeneous traffic remains future work.

VIII. CONCLUSION

This paper presented NiCoLE, an in-network agent that assists real-time RTP video streaming through proactive spatial adaptation and transient L4S steering. Results show reduced queue buildup, stall ratio, and playback instability under congestion, with the largest gains at high frame rates (90–120 FPS), where reactive bitrate-only control becomes insufficient.

The results highlight that LLM-assisted in-network control is conditionally beneficial: gains increase with temporal density and congestion sensitivity, while remaining limited in low-FPS or over-provisioned scenarios. This establishes a workload-aware model for deploying network-assisted streaming. Future work will extend NiCoLE to multi-flow environments and explore tighter integration with programmable dataplane architectures.

ACKNOWLEDGMENT

This work was supported by Ericsson Telecomunicações Ltda., and by the Sao Paulo Research Foundation (FAPESP), grant 2021/00199–8, CPE SMARTNESS .

REFERENCES

- [1] H. Schulzrinne *et al.*, “RTP: A Transport Protocol for Real-Time Applications,” 2003.
- [2] J. Uberti *et al.*, “WebRTC Media Transport and Use of RTP,” 2021.
- [3] A. Shirmarz *et al.*, “From Pixels to Packets: Traffic Classification of Augmented Reality and Cloud Gaming,” in *2024 IEEE 10th International Conference on Network Softwarization (NetSoft)*, 2024, pp. 195–203.
- [4] A. Shirmarz *et al.*, “In-Network AR/CG Traffic Classification Entirely Deployed in the Programmable Data Plane: Unlocking RTP Features and L4S Integration,” in *Proc. IEEE Int. Conf. Network Softwarization (NetSoft)*. IEEE, 2025, pp. 477–485.
- [5] “Understanding Latency Tolerance in Cloud-Based VR Gaming,” in *Proc. of the ACM Multimedia Conference*. ACM, 2025.
- [6] F. Monaco *et al.*, “Latency Sensitivity in Cloud Gaming: An Empirical Study,” *Computer Networks*, 2025.
- [7] P. Karimi *et al.*, “Tight Loops, Smooth Streams: Responsive Congestion Control for Real-Time Video,” vol. 139, pp. 9:1–9:0, 2026.
- [8] H. Meng *et al.*, “MAE: More Adaptive Video Encoder for Consistent Low Latency in High-Quality Real-Time Communication,” in *23rd USENIX Symposium on Networked Systems Design and Implementation (NSDI 26)*. Santa Clara, CA: USENIX Association, Apr. 2026.
- [9] H. Alvestrand *et al.*, “WebRTC Congestion Control: Requirements and Candidate Algorithms,” 2020.
- [10] K. Nichols *et al.*, “Controlling Queue Delay,” *Communications of the ACM*, vol. 55, no. 7, pp. 42–50, 2012.
- [11] IETF, “Low latency, low loss, scalable throughput (l4s) architecture,” 2023.
- [12] K. D. Schepper *et al.*, “Dual-Queue Coupled Active Queue Management (AQM) for Low Latency, Low Loss, and Scalable Throughput (L4S),” RFC 9332, Jan. 2023.
- [13] S. R. Pokhrel *et al.*, “Distilling Large Language Models for Network Active Queue Management,” 2025.
- [14] Ericsson, “AI Agents in the Telecommunication Network Architecture,” Ericsson AB, Tech. Rep., October 2025.
- [15] ETSI, “AI in the Evolution of Autonomous Networks,” November 2025.
- [16] G. Fairhurst *et al.*, “Operational Considerations for Streaming Media,” RFC 9317, Oct. 2022.
- [17] S. Wenger *et al.*, “RTP Payload Format for H.264 Video,” 2011.
- [18] S. Holmer *et al.*, “A Transport-Wide Congestion Control Mechanism for RTP,” 2020.
- [19] R. Jesup *et al.*, “Congestion Control Requirements for Interactive Real-Time Media,” RFC 8836, Jan. 2021.
- [20] D. Wu *et al.*, “NetLLM: Adapting Large Language Models for Networking,” in *Proc. ACM SIGCOMM 2024*. Association for Computing Machinery, Inc, Aug. 2024, pp. 661–678.
- [21] F. Liu *et al.*, “Themis: Toward Stable Near-Zero Queuing Delay in Congestion Control for Low-Latency Interactive Video Streaming,” in *Proc. of the 33rd ACM International Conference on Multimedia*, 2025, pp. 12 131–12 139.
- [22] X. Gao *et al.*, “FENIX: Enabling In-Network DNN Inference with FPGA-Enhanced Programmable Switches,” in *23rd USENIX Symposium on Networked Systems Design and Implementation (NSDI 26)*, ser. NSDI ’26, 2026.
- [23] T. Swamy *et al.*, “Taurus: A data plane architecture for per-packet ml,” in *Proc. Int. Conf. Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. Association for Computing Machinery, Feb. 2022, pp. 1099–1114.