



Cloud Infrastructure Strategies for the Entertainment and Media Sector: Building Resilient Platforms for Digital Entertainment Ecosystems

Archith Rapaka*

Atom Tickets, USA

* Corresponding Author Email: rsarchith@gmail.com - ORCID: 0000-0002-5247-6650

Article Info:

DOI: 10.22399/ijcesen.4503
Received : 15 October 2025
Revised : 05 December 2025
Accepted : 10 December 2025

Keywords

Cloud Infrastructure,
Entertainment Platforms,
Content Delivery Networks,
Fault Tolerance Systems,
Disaster Recovery Solutions

Abstract:

Cloud infrastructure has revolutionized the entertainment and media industries in their very essence, supporting the aspects of content distribution, audience interaction, and digital monetization at levels that were not possible earlier. Modern entertainment systems face novel operational issues associated with the extreme volatility of traffic, multi-format content delivery requirements, and system consistency during peak consumption periods. They must support real-time streaming with low processing latency, cater to many users at once in different parts of the world, and provide stable performance despite temporal variations in demand. Complex architectural patterns that include multi-region deployment, content delivery networks, and event-driven microservices have come to represent vital underpinnings for designing fault-resistant entertainment platforms. Advanced caching techniques running over multiple layers of architecture tune content delivery, while auto-scaling provides dynamic reallocation of resources with varying patterns of demand. Load balancing techniques provide traffic distribution across distributed pieces of infrastructure, ensuring efficient use of resources without compromising quality of service. Fault tolerance strategies based on temporary architecture configuration help maintain functional states in the event of component failure or unexpected disruptions. Cloud-based disaster recovery plans provide end-to-end backup, restore, and business continuity features that are essential to operations that generate revenue. System efficiency is dramatically improved by performance optimization methods such as machine learning-based content prediction, dynamic adaptive streaming protocols, and hierarchical resource utilization models. These technology innovations collectively allow entertainment platforms to offer dependable, high-performance digital experiences while addressing infrastructure complexity and operational expense constraints effectively.

1. Introduction

Entertainment and media have undergone a significant transformation in the digital age, with cloud infrastructure being the basis for contemporary systems of content dissemination, audience interaction, and monetization. From streaming services distributing content to millions of viewers simultaneously to ticketing systems processing enormous simultaneous transactions at the time of major event releases, success in operations in the sector relies upon a cloud architecture that is reliable, scalable, and available. This evolution has precipitated unprecedented technical issues, such as coping with traffic volatility, effectively and seamlessly delivering experiences to audiences in various content forms,

and keeping systems running smoothly during periods of heightened consumer usage.

The entertainment sector presents distinct operational characteristics that require specialized cloud infrastructure solutions fundamentally different from traditional enterprise applications. Real-time streaming systems exemplify the difficulty of managing traffic capture in dynamic environments. According to research on network traffic capture systems, modern solutions must handle processing with minimal delays - Libtins exhibits one of the lowest processing delays compared to other libraries, ranking second only to Libpcap [1]. These systems require capturing all traffic, including potentially malicious packets, with processing kept to a minimum to enable real-time capture. The volume of packets transferred on

high-speed networks necessitates storing log information in streaming systems rather than traditional files, as the latter would become excessively large [1]. Apache Kafka, "an open-source distributed event streaming platform used by thousands of companies for high-performance data pipelines, streaming analytics, data integration, and mission-critical applications," supports over one million messages per second for producers [1].

The global entertainment streaming market has experienced immense growth, with Content Delivery Networks (CDNs) becoming critical infrastructure. Research indicates that by 2050, 70% of the global population is predicted to live in cities, creating unprecedented demand for content delivery services [2]. CDNs consist of proxy servers built in numerous locations throughout the world, with duplicate copies of content created on origin servers and kept on server proxies. Studies show that traditional CDN systems face difficulties in scaling due to high deployment costs, leading to the emergence of cloud CDNs that lease resources such as memory and bandwidth [2]. Performance metrics from hierarchical optimized resource utilization models demonstrate significant improvements: overall memory usage reduced to 80%, CPU usage reduced to 20%, response time minimized to 2 seconds, and total congestion cost with respect to network load level reduced to 100 [2].

Multi-region cloud deployments have become standard practice for managing global content delivery. Research on CDN architectures reveals that streaming content caching in cooperative edge-cloud computing architectures, which combine edge computing and cloud computing, significantly reduces workload on streaming media service systems [2]. Edge servers utilize portions of storage space to cache specific information, particularly for video stream files with strict latency requirements. The cooperative cache domain formed by pooling cache spaces of all edge servers enables the sharing of cached contents. Studies demonstrate that when cache takes 20% of the entire storage capacity, hit rates improve as traffic and request numbers rise, with the HORCP (Hierarchical Optimized Resource Utilization-based Content Placement) method outperforming standard models [2]. The segregation of cache space allows storage of both dynamic and static data in the cache region, initially showing low hit rates when cache space is first used, but increasing as the number of requests rises and space is filled by previously processed requests [2].

This article outlines fundamental infrastructure practices that enable entertainment and media organizations to deliver reliable digital experiences

with predictable performance against the myriad of technical challenges faced. The analysis examines technologies such as Apache Kafka and ksqlDB, which enable real-time pattern matching using queries for network traffic analysis [1], and hierarchical probability routing models that enable reliable end-to-end data transmission with optimized routing paths [2]. These systems demonstrate the benefits and effectiveness of employing modern technologies for network traffic analysis and management in entertainment infrastructure, where real-time processing capabilities prove essential for monitoring, alerting, and generating live insights from streaming data.

2. Industry-Specific Infrastructure Requirements

The infrastructure issues within the entertainment and media sector are different from the issues associated with any other sector of the economy in that there is extensive fluctuation in traffic with a global scope and high-performance demands. According to research on large-scale video stream concurrent transmission, the video cloud market in China reached USD 5.04 billion in the second half of 2021 and is expected to reach USD 31.4 billion by 2025. Cisco's forecast indicates that video streaming will account for 80% of the total Internet traffic in 2023, demonstrating the dramatic increase in the proportion of video streaming of total Internet traffic [3].

Traffic patterns with entertainment clients differ significantly from traditional business applications. In smart education application scenarios involving cloud recording and live streaming systems, there is a high demand for supporting 200 edge recorded broadcast and live broadcast classrooms with 5000 concurrent logins and online viewing capabilities. These systems must handle continuous audio and video data transmitted quickly and frequently between edge cameras, edge servers, network switches, and storage devices [3]. The transmission performance requirements are particularly stringent, as experimental data shows that when transmitting 64B small stream packets, link aggregation transmission schemes can obtain 3.6 times the actual bandwidth compared to single network port benchmark schemes [3].

Content delivery requirements add significant complexity to infrastructure design. In edge computing clusters for video services, the system must process packets of varying sizes - from 64B to 1024B - with studies showing that Q-learning scheduling algorithms can reduce packet loss rates by up to 0.48% compared to traditional DPDK algorithms. When processing 100,000,000 packets,

optimized systems demonstrate that packet loss rates decrease gradually as the number of cores and queues increase, with performance continuing to improve when queue numbers increase to 24, ultimately reducing average time delay by 21% [3]. The global nature of the industry requires infrastructure delivering consistent performance across various regions and network conditions. Data centers have become critical infrastructure across Earth's five continents, underground, underwater, and in space, with territories from Singapore to Iceland, and from Cape Town through Chile to Northern Ireland being envisioned, planned, and zoned for new data center construction. The incidental failure of just one facility has proven capable of knocking out large parts of national services, transnational online platforms, and entertainment media, interrupting the pace of data transmission and digital service provision across several continents for weeks [4].

Regulatory compliance and content protection concerns further complicate existing infrastructure. Edge recorded broadcast and live broadcast systems typically include core business modules handling audio and video capture, pushing flow, transcoding, storage, distribution, and playback. These systems require load-balancing service clusters that allow users to share the load across different physical servers through control policies. Experimental results demonstrate that when using Q-learning load balancing, the utilization of each network port becomes relatively smooth, reducing CPU load rate by 18% compared to RSS scheduling [3].

The cyclical nature of entertainment consumption requires sophisticated resource management. Research shows that when the number of receiving ports is close to the number of sending ports, the average load ratio of all cores is the smallest, and performance is optimal. Systems must dynamically adjust based on load, with studies showing that Q-learning algorithms can achieve 3 times the throughput at 0 packet loss rate compared to benchmark UDP transmission algorithms, and 5 times the throughput at 0.01 packet loss rate. When compared to RSS algorithms, there is a 37.5% improvement in throughput at 0 packet loss rate and 23.25% improvement at 0.01 packet loss rate [3].

These requirements have led to the adoption of sophisticated monitoring and observability systems. Data center operations now depend on maintaining "uptime" - a measure of the time an online service is available - which has become a key metric in the data economy. The temporalities of data infrastructure and labor time get entangled through narratives of economic development, with data centers requiring continuous calibration to the tempos and rhythms of global data flow. Speed

differences of even a millisecond might put competing organizations out of business in sectors like high-frequency algorithmic trading, leading to careful planning and strategic placement of data centers to optimize traffic speeds [4].

The infrastructural temporalities of data centers create what researchers call "power-chronographies" that shape differential lived experiences of time through the biopolitics of labor. Workers in data centers must temporarily recalibrate their bodies to operate in stillness and silence, often alone, providing companionship and care for machines. This bodily recalibration machine time becomes part of the labor of maintaining machine uptime [4]. The result is infrastructure that must accommodate both the technical demands of content delivery and the human temporalities of those who maintain these systems.

3. Core Technology Strategies and Implementation

Layered architectural techniques are used in modern entertainment platforms with emphasis placed on scalability, reliability, and optimization of performance. These solutions use multi-region deployments, which are now critical in organizations that are developing resilient, high-performance apps with a global presence. AWS Global Accelerator can reduce time to first byte by as much as 60% compared to standard internet routing, particularly for users in geographically remote locations or regions with less developed internet infrastructure [5]. By distributing workloads and data across geographically diverse locations, businesses can enhance availability, reduce latency, ensure regulatory compliance, and strengthen disaster recovery capabilities while serving content from geographically proximate locations rather than requiring all traffic to traverse back to a single region [5].

Content delivery networks represent fundamental components of entertainment ecosystems, serving as primary mechanisms for delivering media content to worldwide audiences. AWS CloudFront accelerates content delivery by caching content at edge locations worldwide, reducing latency for global users while decreasing the load on origin servers [5]. CloudFront's integration with AWS WAF and Shield provides security at the edge, protecting multi-region infrastructure from DDoS attacks and malicious traffic. For applications requiring field-level encryption, CloudFront protects sensitive data throughout the entire application delivery path, which is particularly valuable for multi-region deployments where data

may traverse multiple networks and processing environments [5]. The service's origin failover capability automatically redirects requests to a secondary origin when the primary origin is unavailable, complementing Global Accelerator's health-aware routing to provide multiple layers of redundancy.

Event-driven architectures have emerged as essential patterns for managing complex workflows in entertainment systems, with research showing that event-driven e-commerce platforms can handle considerably more concurrent users than their monolithic counterparts [6]. During the COVID-19 pandemic, global e-commerce transactions increased significantly compared to pre-pandemic levels, while average e-commerce sites experienced substantially higher traffic during pandemic peaks, revealing the limitations of traditional monolithic architectures [6]. Event-driven systems achieve lower latency compared to traditional request-response models, with cloud-based event processing systems capable of handling substantial event volumes per second, providing the throughput necessary for high-volume operations [6]. The asynchronous nature of event-driven communication reduces system coupling by as much as 60%, allowing services to evolve independently without creating complex interdependencies [6].

Microservices patterns combined with containerization technologies allow entertainment platforms to scale individual components independently. Organizations implementing microservices report faster time-to-market for new features and capabilities, with deployment frequency improving significantly after microservices adoption, allowing businesses to rapidly respond to market changes and customer needs [6]. The horizontal scaling capabilities of microservices improve resource utilization compared to monolithic systems, with high-demand services such as product catalog and checkout scaling independently without requiring the entire system to scale, optimizing infrastructure costs while maintaining performance [6]. System recovery time reduces significantly with a distributed architecture that incorporates event sourcing principles, minimizing downtime during service disruptions [6].

API gateway implementations provide centralized management of service interactions through AWS Transit Gateway, which acts as a regional network hub connecting VPCs and on-premises networks through a central point. AWS Cloud WAN simplifies connectivity between cloud and on-premises environments by providing a single dashboard to define network policies that are

automatically applied across the global network, enabling connectivity for thousands of AWS VPCs and on-premises locations through a unified global network spanning multiple AWS regions [5]. Cloud WAN's core network policies can be customized to create network segments that adhere to specific regulatory requirements, while centralized policy management ensures consistent security controls across the global network [5].

AWS Global Accelerator leverages the AWS global network infrastructure to optimize the network path from users to applications, using Anycast IP addresses to route user traffic to the optimal AWS endpoint based on factors including geographic proximity, endpoint health, and routing policies [5]. Connection persistence is preserved on endpoint failures in the service, and this is critical to ensuring a smooth user experience when a region is unable to connect. Global Accelerator continuously monitors the health of application endpoints across regions and automatically redirects traffic away from unhealthy endpoints within seconds, ensuring that user traffic is only directed to healthy application instances [5].

Integration of microservices with event-driven patterns creates enhanced modularity levels, with enterprises reporting improved system stability after microservices adoption, particularly during peak traffic periods [6]. The separation of concerns inherent in microservices architecture reduces development complexity for large applications, with each service focusing on a specific business domain, allowing specialized optimization without the risk of unintended consequences across the broader system [6]. When combined with Command Query Responsibility Segregation (CQRS), event sourcing allows systems to optimize for both write-heavy and read-heavy operations, enabling specialized optimization for different access patterns while maintaining data consistency [6].

4. Performance Optimization and Scalability Solutions

Optimizing entertainment infrastructure performance requires sophisticated caching strategies operating at multiple architectural layers. According to research on advanced video streaming services, wireless caching that allows distributed entities with limited storage to cache popular content has demonstrated significant potential, with studies showing that online video services generate the most significant portion of global data traffic, where a small number of contents dominate traffic patterns [7]. The research indicates that cache hit rates can be substantially improved by storing

popular content expected to be requested by nearby users in caching nodes, immediately providing content from these nodes when users' desired content is proactively cached. Video files can be encoded into multiple quality versions (720p, 1080p, and 4K), creating a critical tradeoff between video quality and diversity in caching strategies [7]. Advanced caching mechanisms for video streaming have evolved to support Dynamic Adaptive Streaming over HTTP (DASH), which adaptively chooses the bitrate of transmitting video file partitions depending on time-varying environments to avoid playback stalls.

Auto-scaling capabilities represent essential functions for managing traffic variability characteristic of entertainment platforms. Research on dynamic, scalable auto-scaling models demonstrates that systems can maintain higher resource utilization while reducing energy costs through intelligent provisioning algorithms [8]. The proposed auto-scaling algorithm operates based on configurable thresholds, where new virtual machines are provisioned when all VMs reach specified upper thresholds for network bandwidth and active sessions. CloudSim simulations showed CPU utilization recommendations of 66% maximum for single-region instances and 47% for multi-region instances during normal operations, with 24-hour aggregate peaks reaching 91% [8]. The auto-scaling system demonstrated the capability to handle sudden load demands while automatically terminating idle virtual machines when resource utilization drops below predetermined lower thresholds, effectively balancing performance with operational costs.

Load balancing strategies in entertainment infrastructure must accommodate both traditional web traffic and specialized media streaming protocols. Research indicates that Apache HTTP Load-Balancer can effectively distribute web applications across front-end systems, with cloud systems automatically adding new web servers to virtual clusters by modifying load balancer configurations [8]. The study identified three types of AWS load balancers: Network Load Balancer (operating at Layer 4, capable of processing millions of requests simultaneously using hash algorithms for IP address and port-based routing), Application Load Balancer (handling HTTP/HTTPS traffic at Layer 7 with Target Groups-based routing), and Classic Load Balancer [8]. These load balancers integrate with Auto Scaling Groups (ASGs) that manage servers automatically based on scaling policies, including Target Tracking Scaling, Step Scaling, and Simple Scaling configurations.

Recent advances in video caching optimization demonstrate significant performance improvements through partition-based approaches. Research shows that video files typically consist of chunks responsible for seconds of playback time, with each chunk potentially having different bitrates in dynamic streaming systems [7]. Scalable Video Coding (SVC) enables files to be divided into a base layer essential for video playback at the lowest quality, plus enhancement layers for quality improvement. Studies on 360° video streaming revealed that tile-based caching strategies can reduce bandwidth costs and delivery delays by minimizing errors between users' requested resolution and cached Field of View (FoV) tile resolution [7]. Multi-view streaming technology divides entire scenes into tile units processed by different cameras, with 3D video services requiring substantially higher data traffic than traditional 2D content.

The integration of machine learning techniques with caching systems has shown promising results. Deep learning-based approaches using LSTM models predict content IDs users will request shortly, while reinforcement learning algorithms, including Q-learning, Deep Q-Network (DQN), and Deep Deterministic Policy Gradient (DDPG), can optimize caching policies without prior knowledge of popularity distributions [7]. Research on super-resolution technology demonstrates potential for enhancing video chunk quality through deep learning, enabling more aggressive video transcoding at transmitter sides while maintaining quality standards. Studies indicate that joint optimization of caching, transcoding, and delivery can minimize end-to-end delay, with recent work showing that DDPG-based video streaming successfully updates cache status of roadside units while considering user mobility patterns [7]. Performance monitoring data from cloud environments reveals critical operational metrics. Disk I/O measurements showed varying throughput patterns based on logical-to-physical I/O transformation through share settings and software RAID configurations [8]. Network traffic analysis demonstrated packet transmission patterns across multiple network interfaces, with incoming and outgoing traffic metrics providing insights into instance-level network utilization. The research confirmed that auto-scaling mechanisms are crucial for improving resource utilization and lowering infrastructure costs, with managed instance groups allowing automatic instance addition or removal based on load changes, enabling applications to handle traffic surges while reducing costs during lower demand periods [8].

5. Fault Tolerance and Disaster Recovery

The critical infrastructure in the entertainment industry requires advanced fault tolerance facilities, especially due to the complex interdependence in the modern systems-of-systems (SoS) architecture. The ReViTA framework, introduced for fault-tolerant SoS design, addresses the challenge where constituent systems can produce disturbances at the SoS level through failures, implementation changes, or deviations from their roles due to competing objectives [9]. These disturbances particularly impact entertainment platforms where multiple independent systems must coordinate seamlessly during high-traffic events.

Circuit breaker patterns and redundancy strategies align with ReViTA's concept of Transient Architectural Configurations (TAC), which enable systems to temporarily reconfigure their architecture to maintain functionality during disturbances [9]. The framework's evaluation with 14 professionals revealed that architectural reconfigurations involving modifying SoS composition and relationships provide more reasonable approaches to handling disturbances than traditional fault tolerance methods. Entertainment platforms implementing these principles can leverage the opportunistic nature of SoS design, where constituent systems operate independently but contribute toward common goals when interconnected [9].

Real-time monitoring requirements reflect findings from ReViTA's real-world evaluation at a large Brazilian public university, where four professionals applied the framework to respond to power outages across distributed institutes and campuses [9]. The study demonstrated that effective fault tolerance requires continuous SoS monitoring to detect situations where operational architectural configurations no longer meet mission requirements. For entertainment platforms, this translates to multi-level observability, detecting both individual component failures and system-wide performance degradation patterns that could impact content delivery during critical release periods. Disaster recovery planning complexity aligns with cloud-based solutions, offering significant advantages. According to IDG survey data, 72% of companies have disaster recovery plans, with 46% utilizing cloud-based disaster recovery solutions [10]. The Disaster Recovery Planning Council reports that 20% of organizations experienced disasters within the past five years, with 80% of those experiencing downtime or data loss. Financial impacts prove substantial, with the Ponemon Institute determining average business downtime costs of \$5,600 per minute [10].

Cloud-based disaster recovery adoption shows strong growth trajectories. The Disaster Recovery Journal survey found 38% of organizations consider cloud-based disaster recovery the most efficient recovery method [10]. Market forecasts project the cloud disaster recovery market to grow from \$12 billion in 2020 to a figure of \$82 billion by 2025 - a 29% compound annual growth rate (CAGR). 46% of organizations intend to increase their spending on cloud-based disaster recovery in the next year, and 72% of the organizations utilizing cloud-based disaster response have reported satisfaction with their deployed solutions [10].

Infrastructure as a Service (IaaS) provides critical disaster recovery capabilities for entertainment platforms, enabling organizations to replicate IT infrastructure to cloud service providers who offer computing resources, storage, and network facilities on demand [10]. This approach eliminates the need for additional hardware purchases while enabling rapid resource provisioning during emergencies. Disaster Recovery as a Service (DRaaS) extends these capabilities through comprehensive services including data protection via continuous backups to secure cloud-based data centers, regular testing and validation of DR environments, and automated recovery procedures restoring critical systems with minimal interruption [10].

Recovery time objectives vary significantly between service types. DRaaS typically measures recovery times in hours or minutes, while Backup as a Service (BaaS) recovery times extend to days or weeks [10]. This distinction proves critical for entertainment platforms where content availability directly impacts revenue. DRaaS aims to minimize data loss with recovery to recent states, while BaaS focuses on specific data sets with recovery points depending on backup frequency [10].

Testing protocols emerge as critical success factors. ReViTA's evaluation revealed that regular testing helps identify potential problems before business operations are impacted [9]. The framework emphasizes choosing the correct testing methods, ensuring all disaster recovery procedure elements receive adequate validation. For entertainment platforms, this includes validating both technical recovery capabilities and organizational response procedures under various failure scenarios.

Business continuity benefits extend beyond technical recovery. The ReViTA study demonstrated that fault tolerance frameworks improve stakeholder communication and enhance resource utilization [9]. Participants noted that employing such frameworks brings crucial insights into costs and planning essential for implementing fault-tolerance strategies. The framework facilitates

a comprehensive understanding of conflicts and weaknesses in constituent systems while fostering collaboration between domain experts and decision-makers [9].

Cost considerations influence disaster recovery strategy selection. DRaaS involves higher costs due to dedicated infrastructure and testing requirements, while BaaS offers lower costs using existing

infrastructure with less complexity [10]. Entertainment platforms must balance these cost factors against business impact considerations, where DRaaS proves critical for minimizing downtime and maintaining business continuity during high-revenue periods, while BaaS serves important roles in minimizing data loss for less time-critical content archives [10].

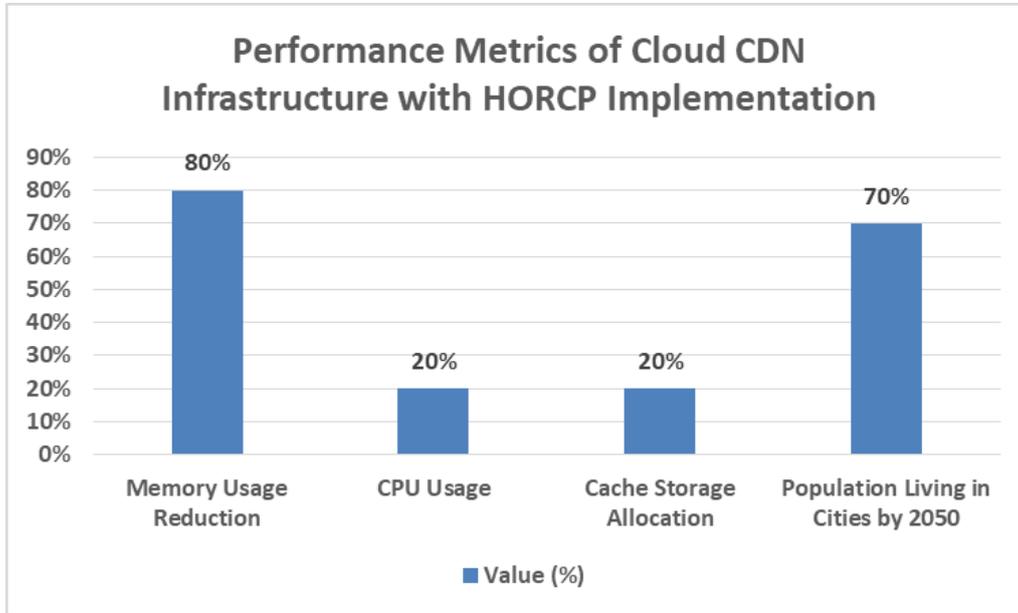


Figure 1: Performance Metrics of Cloud CDN Infrastructure with HORCP Implementation [1,2]

Table 1: Infrastructure optimization in video streaming systems [3,4]

Metric	Value
Video streaming % of Internet traffic (2023)	80%
China Video Cloud Market Projection 2025	\$31.4B
Bandwidth improvement (link aggregation vs single port)	3.6x
CPU load reduction (Q-learning vs RSS)	18%
Packet loss rate reduction (Q-learning vs DPDK)	0.48%
Average time delay reduction	21%
Throughput improvement at 0% packet loss (vs RSS)	37.5%

Table 2: Cloud Infrastructure Performance Metrics for Entertainment Platforms [7,8]

Performance Metric	Value
Single-region CPU utilization (maximum)	66%
Multi-region CPU utilization (maximum)	47%
24-hour aggregate CPU peak	91%
Network Load Balancer layer	Layer 4
Application Load Balancer layer	Layer 7
Video quality levels supported	3 (720p, 1080p, 4K)
ASG scaling policy types	3 types
Load balancer types (AWS)	3 types

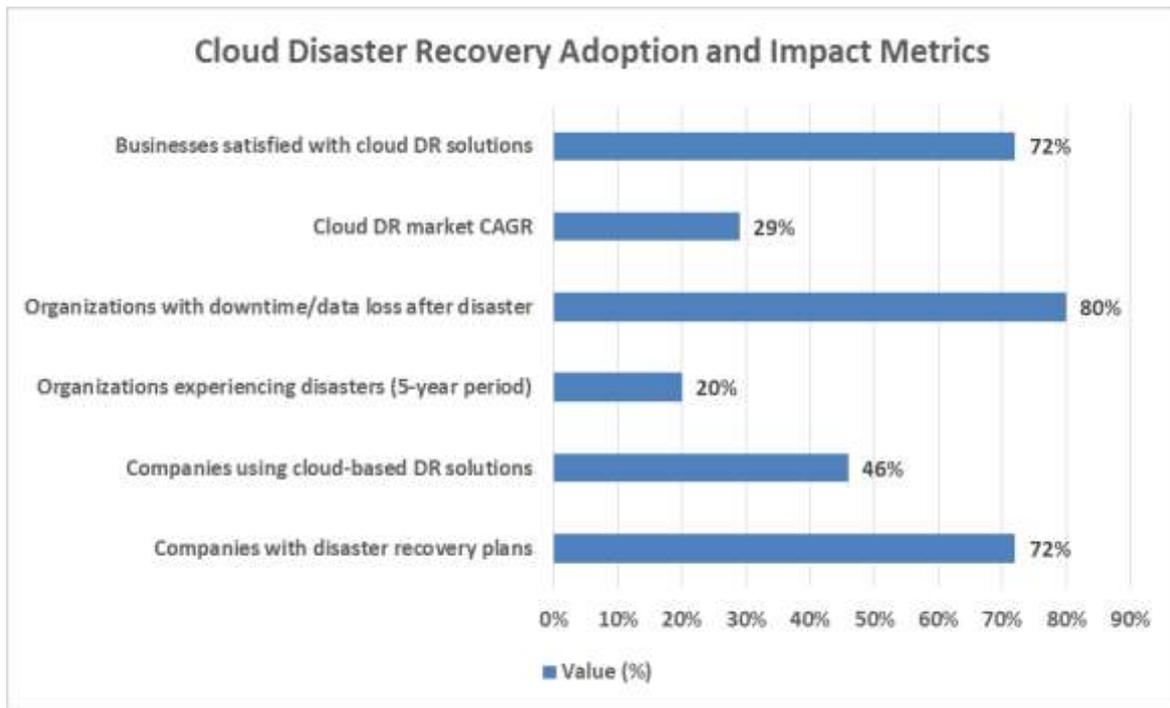


Figure 2: Cloud Disaster Recovery Adoption and Impact Metrics [9,10]

5. Conclusions

The use of cloud infrastructure has radically transformed the entertainment and media industries by developing novel patterns of content distribution, customer engagement, and platform resiliency. The integration of sophisticated features in architectural designs, including multi-region deployments, event-oriented microservices, and state-of-the-art caching mechanisms, has established robust bases to support global-scale entertainment situations. These technologies enable platforms to manage the high volatility of traffic, offer smooth content experiences across multiple formats, and maintain operational continuity during the peak event. Transient architectural designs and full-scale disaster recovery processes implemented guarantee system robustness against both technical outages and unforeseen interference, while scalability automation adjusts resources in accordance with varying demand patterns. The use of machine learning models along with legacy infrastructure management has improved content prediction accuracy, enhanced cache hit ratio, and cut delivery latency across distributed platforms. Hierarchical models of resource utilization and dynamic adaptive streaming protocols have optimized performance, thereby improving user experience and reducing operational expenditure. Cloud infrastructure strategies must evolve and adjust to the changing patterns of entertainment consumption to support immersive content formats, live participation, paradigms of interaction with the

audience, and other emerging technologies. Successful execution of such infrastructure strategies depends on prudent adaptation of technical requirements, business strategies, and operational constraints for providing entertaining platforms that can deliver reliable and high-performance digital experiences for meeting contemporary audience demands and financially enabling market competitiveness in the rapidly transforming digital entertainment landscape.

Author Statements:

- **Ethical approval:** The conducted research is not related to either human or animal use.
- **Conflict of interest:** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper
- **Acknowledgement:** The authors declare that they have nobody or no-company to acknowledge.
- **Author contributions:** The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- **Data availability statement:** The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.

References

- [1] Adrian-Tiberiu Costin et al., "A Real-Time Streaming System for Customized Network Traffic Capture", MDPI, 2023. Available: <https://www.mdpi.com/1424-8220/23/14/6467>
- [2] M. Sasikumar et al., "A Hierarchical Optimized Resource Utilization based Content Placement (HORCP) model for cloud Content Delivery Networks (CDNs)", Springer Open - Journal of Cloud Computing IEEE Cloud Computing, 2023. Available: <https://journalofcloudcomputing.springeropen.com/articles/10.1186/s13677-023-00519-2>
- [3] Haitao Liu et al., "An Optimization Method of Large-Scale Video Stream Concurrent Transmission for Edge Computing", MDPI, 2023. Available: <https://www.mdpi.com/2227-7390/11/12/2622>
- [4] Julia Velkova and Jean-Christophe Plantin, "Data centers and the infrastructural temporalities of digital media: An introduction", Sage Journals, 2023. Available: [Data centers and the infrastructural temporalities of digital media: An introduction - Julia Velkova, Jean-Christophe Plantin, 2023](#)
- [5] Piyush Dhar Diwan, "Multi-Region Networking and Global Traffic Management for AWS", International Journal of Network and Communication Research, Apr. 2025. Available: [Multi-Region-Networking.pdf](#)
- [6] George Thomas, "Microservices and event-driven architecture: Revolutionizing e-commerce systems", WJARR, May 2025. Available: https://journalwjarr.com/sites/default/files/fulltext_pdf/WJARR-2025-1663.pdf
- [7] Minseok Choi et al., "Caching, transcoding, delivery, and learning for advanced video streaming services", ScienceDirect, 2024. Available: <https://www.sciencedirect.com/science/article/pii/S2405959524000626>
- [8] Jvr Ravinda et al., "A Dynamic Scalable Auto-Scaling Model as a Load Balancer in the Cloud Computing Environment", ResearchGate, 2023. Available: https://www.researchgate.net/publication/372353970_A_Dynamic_Scalable_Auto-Scaling_Model_as_a_Load_Balancer_in_the_Cloud_Computing_Environment
- [9] Francisco Henrique Cerdeira Ferreira et al., "A framework for the design of fault-tolerant systems-of-systems", ScienceDirect, 2024. Available: <https://www.sciencedirect.com/science/article/abs/pii/S0164121224000530>
- [10] Aggidi Sathya and Bhuvana, "Cloud Disaster Recovery Management and Business Continuity", IRJMETS, 2024. Available: https://www.irjmets.com/uploadedfiles/paper/issue_5_may_2024/57003/final/fin_irjmets1716312021.pdf