

A Survey on Virtual Network Functions for Media Streaming: Solutions and Future Challenges

ROBERTO VIOLA, ÁNGEL MARTÍN, and MIKEL ZORRILLA, Fundación Vicomtech, Basque Research and Technology Alliance, Spain

JON MONTALBÁN, Department of Electronic Technology, University of the Basque Country, Spain

PABLO ANGUEIRA, Department of Communications Engineering, University of the Basque Country, Spain

GABRIEL-MIRO MUNTEAN, Performance Engineering Laboratory, School of Electronic Engineering, Dublin City University (DCU), Ireland

Media services must ensure an enhanced user's perceived quality during content playback to attract and retain audiences, especially while the streams are distributed remotely via networks. Thus, media streaming services rely heavily on good and predictable network performance when delivered to a large number of people. Furthermore, as the quality of media content gets high, the network performance demands are also increasing, and meeting them is challenging. Network functions devoted to improving media streaming services become essential to cope with the high dynamics of network performance and user mobility. Furthermore, new networking paradigms and architectures under the 5G networks umbrella are bringing new possibilities to deploy smart network functions, which monitor the media streaming services through live and objective metrics and boost them in real-time. This survey overviews the state-of-the-art technologies and solutions proposed to apply new network functions for enhancing media streaming.

CCS Concepts: • **Information systems** → **Multimedia streaming**; **Computing platforms**; • **Networks** → **Network management**.

Additional Key Words and Phrases: Media streaming, network functions, network virtualization

ACM Reference Format:

Roberto Viola, Ángel Martín, Mikel Zorrilla, Jon Montalbán, Pablo Angueira, and Gabriel-Miro Muntean. 2022. A Survey on Virtual Network Functions for Media Streaming: Solutions and Future Challenges. 1, 1 (October 2022), 35 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

In recent years, media streaming traffic has been experiencing a growing trend. Wireless and mobile devices are becoming the primary producers and consumers of rich media content. 5G networks must cope with these new traffic demands by supporting higher bandwidth and reduced latency. It is estimated that 5G connections will handle nearly three times more traffic than current LTE connections by 2023 [61]. New applications involving video streams are gaining relevance and are attracting an increased audience, including vertical industries where media has a residual

Authors' addresses: **Roberto Viola**, rviola@vicomtech.org; **Ángel Martín**, amartin@vicomtech.org; **Mikel Zorrilla**, mzorrilla@vicomtech.org, Fundación Vicomtech, Basque Research and Technology Alliance, Paseo Mikeletegi 57, San Sebastián, Spain, 20009; **Jon Montalbán**, jon.montalban@ehu.es, Department of Electronic Technology, University of the Basque Country, Bilbao, Spain, 48013; **Pablo Angueira**, pablo.angueira@ehu.es, Department of Communications Engineering, University of the Basque Country, Bilbao, Spain, 48013; **Gabriel-Miro Muntean**, gabriel.muntean@dcu.ie, Performance Engineering Laboratory, School of Electronic Engineering, Dublin City University (DCU), Dublin, Ireland.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

53 presence. Examples of application areas that can benefit from advanced media streaming include the Industrial Internet of
54 Things (IIoT), medical equipment, and connected and autonomous vehicles. Moreover, 360-degree and 3D video formats
55 enable support for new services beyond entertainment, i.e., professional applications. They build novel interaction
56 experiences and data navigation on top of several technologies, such as eXtended Reality (XR), Virtual Reality (VR), and
57 Augmented Reality (AR) [82]. Finally, online gaming and video conferencing are also highly popular, especially in the
58 last period. These services have increasing demands in terms of network support. However, although the networks
59 have growing capabilities, there is a significant increase in rich media streaming traffic, fueled mainly by the global
60 COVID-19 pandemic. This pandemic is transforming users' habits to access the Internet [86, 140] and media content
61 consumption [83, 126]. The Broadband Commission for Sustainable Development, a joint initiative of the International
62 Telecommunication Union (ITU) and the United Nations Educational, Scientific and Cultural Organization (UNESCO),
63 is also concerned about these user habit changes. Therefore, it is implementing an Agenda for Action to push an
64 emergency response to the pandemic, aiming at Internet access extension and boosting its capacity [5, 6].

65 All the factors mentioned above inevitably influence the evolution of all services, especially affecting the rich
66 media ones. It is therefore evident that there is a need for new network-related solutions to support high Quality of
67 Service (QoS) for these applications. The current traffic crosses networks working on a best-effort basis where no
68 details regarding packet delivery (e.g., time) are guaranteed. Therefore, best-effort networked-transmitted media traffic
69 may result in a lower user's Quality of Experience (QoE). A practical example of this QoE degradation is stalls or
70 artifacts during media playback on player devices. Employing Content Delivery Networks (CDNs) is the most common
71 solution to prevent adverse quality effects and make video delivery more efficient. CDNs are geographically distributed
72 hierarchical systems that cache and store video streams to foster efficiency and increase service coverage. CDN price is
73 decreasing, but the overall cost for the content provider is increasing, as the traffic from/to CDN is growing. [172].

74 Beyond CDNs, more advanced solutions such as load balancers, transcoders, and transraters based on Network
75 Function Virtualization (NFV) technologies [105] are investigated to support media streaming services. NFV allows
76 the deployment of Virtual Network Functions (VNF) devoted to empowering network abilities when delivering media
77 streaming traffic in an optimized and cost-effective manner [68, 117].

78 In order to take advantage of using VNFs in a media streaming context, three significant aspects must be considered
79 during their design. First, VNFs should monitor objective operational parameters of the network, such as throughput or
80 latency, representative of QoS of the media streaming dataflows, which directly influence user's satisfaction. However,
81 QoS metrics do not perfectly map to user experience, as users' perceived quality is highly subjective. VNFs should also
82 consider QoE, which compiles subjective evaluation elements, including rewards for playback quality and smoothness
83 and penalties for image freezes and unstable or low quality [28, 125]. Secondly, VNFs should give the CP more control
84 over the network. CP should be able to shape the network traffic and allocate resources by establishing business rules
85 for VNF deployment and life cycle management. These rules allow balancing the trade-offs between network resources
86 and business costs [107], i.e., Capital Expenditure (CAPEX) and Operational expenditures (OPEX), so they are highly
87 relevant. Last, VNFs should also consider energy efficiency. The volume, complexity, and real-time nature of the media
88 streaming traffic have an evident impact on the energy consumption of the network and devices managing the media
89 content. An optimized streaming delivery should reduce energy consumption by turning on or off the VNFs, depending
90 on the demand for network resources, at any time. By putting together real-time QoS and QoE metrics, business policies,
91 and energy efficiency constraints, the user and the CP can benefit from optimized VNFs for media streaming. VNFs can
92 provide the user with a target QoE and reduce business costs and energy consumption to what is strictly necessary to
93 achieve it.

In this context, the main contributions of this survey are:

- The paper provides an extensive overview of different technologies and protocols in the context of VNFs for media streaming;
- The review analyses state-of-the-art network media functions by classifying them into major categories, i.e., media casting, transcoding, and content caching, and providing a comparison among them;
- The paper presents technologies considered by the telecommunications industry as key enablers for the next generation networks and discusses remaining challenges, including emerging aspects such as data security, energy efficiency, and business models for new network assets;
- The survey provides an overview of international initiatives in the media streaming field, including research activities and projects.

The rest of the paper is structured as follows. First, Section 2 presents the objective of this work in the context of related surveys. Section 3 contains an overview of media streaming technologies and protocols. Section 4 identifies and describes the VNFs employed to date to enhance the performance of media streaming services. The employment of VNFs in media streaming has been growing in the last few years as the attention increases on media distribution over the newly deployed 5G networks [123]. VNFs are intrinsically designed to follow the principles of modularity, interoperability, scalability, and flexibility. Media streaming can leverage VNFs to enable higher network capacity and stability, media traffic optimization, and other performance-related advantages. In this sense, several network solutions to enable performance-driven management of the resources are already being employed and/or investigated with the final aim of increasing media streaming performance. These solutions include efficient use of network resources and end device capabilities [80] during their involvement in the streaming service. To better describe these performance-driven VNFs for media streaming, we classify them into major categories and provide a comparison among them. Nevertheless, even if in the current deployment of 5G networks, the VNFs have a significant role, as 5G aims to have a fully virtualized network deployment, there are still several open issues and challenges that need to be addressed in the future. Thus, Section 5 presents the current challenges in the virtualization process of network functions inside 5G and beyond to assess the open issues and scientific research directions. It discusses the future of VNFs to enable an improved media streaming process and enhanced user experience. Finally, we highlight some valuable international initiatives in Section 6 and assert our conclusions in Section 7.

2 PAPER OBJECTIVES IN THE CONTEXT OF RELATED SURVEYS

Table 1. Summary of Previous Surveys on Virtual Network Functions and Media Streaming.

Survey	Scope and topics	Network virtualization	Domain	Research focus	Year
Skorin-Kapov et al. [190]	QoE assessment and management for HAS, MEC monitoring	SDN/NFV, MEC	Media streaming	QoE-driven architecture	2018
Barakabitze et al. [33]	QoE assessment and management in SDN/NFV, QoE-driven HAS over MEC	SDN/NFV, MEC, Cloud/Fog	Media streaming	SDN routing solutions	2019
Zhang et al. [223]	VNF design considerations	VNF, Cloud/Edge	Agnostic	Virtualization solutions	2019
Fei et al. [85]	NFV/VNF current limitations and future directions	NFV, VNF	Agnostic	Virtualization concepts	2020
This work	NFV architecture, performance-driven VNF	NFV, VNF, MEC	Media streaming	Network media functions	2022

157 This survey aims to perform an extensive literature review on the proposed solutions in the realm of VNFs applied
158 to the field of media streaming. The paper also addresses future challenges in this research area.

159 Several surveys on network virtualization have been published in the last few years, in line with the increased
160 interest in virtualization. Zhang et al. [223] provide generic considerations when designing VNFs, while Fei et al. [85]
161 analyze the limitations and identify future research directions. Both surveys focus on the VNF design and utilization to
162 improve current networks, but they do not consider media streaming use cases as they remain domain-agnostic.

163 Limited to media streaming domain, Skorin-Kapov et al. [190] and Barakabitze et al. [33] discuss using virtualized
164 solutions to assess and improve the QoE. Nevertheless, the former is limited to theoretical concepts for achieving a
165 QoE-driven network architecture and does not provide evidence of implementations. The latter is focused only on
166 routing solutions based on Software-defined networking (SDN) and does not consider VNF approaches.

167 A comparison of our work with other surveys is shown in Table 1. Our survey describes the NFV architecture and
168 reviews VNFs solutions for the media streaming domain. The VNFs solutions are analyzed and categorized into major
169 categories. Therefore, our survey addresses a more specific scope, as it discusses the relation between the VNFs and
170 media streaming and provides an extensive review of state-of-art solutions concerning media-specific operations, such
171 as media casting, transcoding, and content caching.

172 A list of acronyms used throughout the paper is presented in Table 2.

173 3 MEDIA STREAMING OVERVIEW

174 Before analyzing systems and functions to enhance the media streaming services, it is important to understand the
175 technologies involved in media streaming services to design and implement better VNFs. Media streaming refers to
176 delivering media content (e.g., live television, video clip, etc.) from a streaming server to a streaming client over a
177 particular network infrastructure. The media source can be either live or pre-recorded. In some cases, the CP is also the
178 infrastructure owner employed to stream the content. However, diverse providers and operators have recently entered
179 the market with different roles in the media streaming process, e.g., Akamai, Netflix, etc. Some of them have their own
180 proprietary media streaming solutions. However, the first solutions were based on the Real-time Transport Protocol
181 (RTP) [184] on top of the User Datagram Protocol (UDP) [170], where the Real-time Transport Control Protocol (RTCP)
182 [184] was employed to monitor network metrics and update the rate control. The choice of UDP was based on its lower
183 latency than the Transmission Control Protocol (TCP) [171], even if it does not guarantee reliability when delivering
184 packets, i.e., lost packets are not re-transmitted when employing UDP. The later explosion of Over-the-top (OTT)
185 services, e.g., Netflix and Hulu, pushed the search for new solutions to deliver Video-on-Demand (VOD) contents, where
186 latency was not a concern, but scalability to cover the increasing demand for content. In OTT services, the CP streams
187 its content over a public network, and an Internet service provider (ISP) controls the actual content delivery. HTTP
188 adaptive streaming (HAS) [186] technologies were introduced to deliver OTT contents, where the use of TCP and HTTP
189 made them attractive since these protocols are ubiquitous. Additionally, almost every device or User Equipment (UE)
190 can establish HTTP-based communications. The HAS-based design has the following advantages over RTP/UDP-based
191 solutions:
192
193
194
195
196
197
198
199
200
201
202
203

- 204 • Traverse networks: HAS communications are performed on the HTTP/TCP stack and use pull-based streaming
205 protocols. Thus, they cross current network infrastructure components, such as Network Address Translation
206 (NAT) and firewall devices [169];
207

Table 2. List of Acronyms used in the paper.

209				
210				
211	3GPP	3 rd Generation Partnership Project	NFVI	NFV Infrastructure
212	5G	Fifth Generation	NFVO	NFV Orchestrator
213	6G	Sixth Generation	NS	Network Service
214	AES	Advanced Encryption Standard	O-RAN	Open RAN
215	AES-CBC	AES block cipher mode	ONAP	Open Network Automation Platform
216	AES-CTR	AES counter mode	OPEX	Operational Expenditure
217	ANN	Artificial Neural Network	OSM	Open Source MANO
218	API	Application Programming Interface	OTT	Over-the-top
219	AR	Augmented Reality	PoP	Point of presence
220	C-RAN	Cloud-RAN	QoE	Quality of Experience
221	CAPEX	Capital Expenditure	QoS	Quality of Service
222	CDN	Content Delivery Network	RAN	Radio Access Network
223	CMAF	Common Media Application Format	RNI	Radio Network Information
224	CN	Core Network	RNIS	RNI Service
225	COTS	Commercial off-the-shelf	RTCP	Real-time Transport Control Protocol
226	CP	Content Provider	RTMP	Real-time Messaging Protocol
227	DASH	Dynamic Adaptive Streaming over HTTP	RTSP	Real-time Transport Protocol
228	DNS	Domain Name System	RTSP	Real Time Streaming Protocol
229	ETSI	European Telecommunications Standards Institute	SCTP	Stream Control Transmission Protocol
230	FeMBMS	Further enhanced MBMS	SDN	Software-defined networking
231	FTRL	Follow The Regularized Leader	SDR	Software-defined radio
232	HAS	HTTP Adaptive Streaming	SLA	Service Level Agreement
233	HLS	HTTP Live Streaming	SON	Self-Organizing Network
234	HTTP	HyperText Transfer Protocol	SRT	Secure Reliable Transport
235	IaaS	Infrastructure as a Service	STUN	Session Traversal Utilities for NAT
236	IBN	Intent-Based Network	SVA	Streaming Video Alliance
237	IIoT	Industrial Internet of Things	TCP	Transmission Control Protocol
238	ISP	Internet Service Provider	TURN	Traversal Using Relays around NAT
239	ITU	International Telecommunication Union	UAV	Unmanned Aerial Vehicle
240	KPI	Key Performance Indicator	UDP	User Datagram Protocol
241	L1	Physical layer	UE	User Equipment
242	L2	Data link layer	UHD	Ultra-High-Definition
243	L3	Network layer	UNESCO	United Nations Educational, Scientific and Cultural Organization
244	L4	Transport layer	VIM	Virtual Infrastructure Manager
245	L7	Application layer	VNF	Virtual Network Function
246	LL CMAF	Low Latency CMAF	VNF-CC	VNF Chain Composition
247	LL-DASH	Low Latency DASH	VNF-FG	VNF Forwarding Graph
248	LL-HLS	Low Latency HLS	VNF-FGE	VNF Forwarding Graph Embedding
249	LTE	Long-Term Evolution	VNF-PC	VNF Placement and Chaining
250	M3U8	M3U UTF-8 Playlist File	VNF-SCH	VNF Scheduling
251	MANO	Management and Orchestration	VNFI	VNF Instance
252	MBMS	Multimedia Broadcast/Multicast Service	VNFM	VNF Manager
253	MEC	Multi-access Edge Computing	VOD	Video-on-Demand
254	ML	Machine Learning	VR	Virtual Reality
255	MPD	Media Presentation Description	vRAN	Virtual RAN
256	MPTCP	Multipath TCP	WebRTC	Web Real-Time Communication
257	Multi-RAT	Multiple Radio Access Technology	XR	eXtended Reality
258	NAT	Network Address Translation		
259	NFV	Network Function Virtualization		
260				

- Reuse and scalability: HAS-based media services can reuse existing CDN systems and caching infrastructures without modifications to reach broad audiences;
- User mobility and device heterogeneity: The dynamic content adaptation-enabled player mechanism is accommodated by all the latest heterogeneous UEs, i.e., smartphones and tablets, which support user mobility.

Figure 1 illustrates the HAS-based adaptive streaming principle. HAS works in pull mode, meaning the client pulls the data from a standard HTTP server, which hosts the media content. To reduce the effect of network fluctuations on the playback, HAS employs a dynamic content adaptation to provide a seamless streaming experience. The original media content is encoded at multiple representations, which differ in bitrate and/or resolution and are split into segments of

Table 3. Features of streaming technologies.

Tech.	Transport	Manifest file	Common issues	Latency	Available bitrate	Bitrate adaptation	CDN compatible	Encryption
RTP	UDP	no	packets lost & artifacts	very low ($\leq 1\text{sec}$)	RTCP	encoder	no	no
RTSP	UDP	SDP	packets lost & artifacts	very low ($\leq 1\text{sec}$)	RTCP	encoder	no	no
RTMP	TCP	no	packets lost & artifacts	low (1-3secs)	RTMP control messages	encoder	no	AES-128 CBC
SRT	UDP	no	packets lost & artifacts	very low ($\leq 1\text{sec}$)	SRT control messages	encoder	no	AES-128 / 265 CTR
WebRTC	UDP, QUIC-ready	SDP	packets lost & artifacts	very low ($\leq 1\text{sec}$)	RTCP	encoder	no	AES-128 CTR
HLS	HTTP 1.X / 2.0 over TCP	M3U8	segment buffering & quality switch	high (5-30secs)	representation player		yes	AES-128 CBC
DASH	HTTP 1.X / 2.0 over TCP, QUIC-ready	MPD	segment buffering & quality switch	high (5-30secs)	representation player		yes	AES-128 CBC / CTR
LL-HLS	HTTP 2.0 over TCP	M3U8	chunks buffering & quality switch	low (1-3secs)	representation player		yes	AES-128 CBC
LL-DASH	HTTP 1.1 Chunked over TCP	MPD	chunks buffering & quality switch	low (1-3secs)	representation player		yes	AES-128 CBC / CTR

fixed time duration (i.e., a segment is usually between 2 and 10 seconds). A manifest file is also generated and stored at the server, which contains information on the available representations, including HTTP URLs indicating where to download the segments of each representation. During a typical HAS session, the client constantly measures specific parameters, such as available network bandwidth and playback buffer level. When it requests content, the client receives the manifest file, which is examined. Then, following an internal adaptation algorithm that processes the monitored performance parameters' values and takes decisions according to the desired adaptation policy, the client requests to download the segment of an appropriate representation from the server.

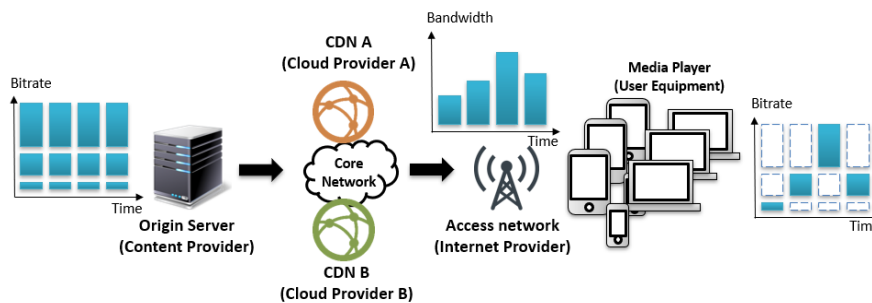


Fig. 1. HAS-based media streaming principle

Besides transport-layer protocols such as RTP, application-layer protocols are also employed in media streaming. Table 3 summarizes different aspects of interest when identifying the best protocol candidates for video streaming.

RTP and Real Time Streaming Protocol (RTSP) [185] perform low latency communications compatible with multicast media streaming. TCP-based Real-time Messaging Protocol (RTMP) [198] enables higher reliability than RTSP, but with higher latency. Secure Reliable Transport (SRT) [188] simplifies the delivery by enabling both push and pull modes of operation. Web Real-Time Communication (WebRTC) [109] enables media streaming through a web browser by exploiting Session Traversal Utilities for NAT (STUN) [215], and Traversal Using Relays around NAT (TURN) [142] protocols provided by third-party servers. SRT and WebRTC increase security by including mandatory encryption support, which is not always required for RTMP. HTTP Live Streaming (HLS) [166] and Dynamic Adaptive Streaming over HTTP (DASH) [191] increase latency due to an internal buffering to overcome network dynamics. In any case, violations of delivery timing could cause stalls and image freezes during the playback if the internal buffer gets empty. To minimize such issues, HAS allows dynamic adaptation mechanisms to track the variability of the network and select the appropriate bitrate. Thus, sudden networking problems are prevented by an alternative bitrate selection from the manifest. Common Media Application Format (CMAF) [114] was a proposal to merge major streaming formats around HLS and DASH. Moreover, its Low Latency mode (LL CMAF) aims to reduce the latency by enabling HTTP chunked/push mode. Thus, the latency can be reduced and get closer to UDP-based streaming technologies. In practice, CMAF did not achieve the integration of HLS and DASH streaming formats since the implementations of Low Latency HLS (LL-HLS) [71] and Low Latency DASH (LL-DASH) [49] still present some differences. Thus, LL-HLS and LL-DASH employ different approaches for HTTP transport and encryption schemes. For instance, a common feature of most HTTP-based solutions is the security by design where different encryption standards protect communications, such as Advanced Encryption Standard (AES) [60] with Cipher Block Chaining (AES-128 CBC) or Counter mode (AES-128 CTR).

Finally, even if most existing media streaming solutions employ UDP and/or TCP, some of them, such as DASH [40] and WebRTC [208], are already evolving and/or being tested with QUIC, a new transport protocol that is expected to substitute TCP when HTTP/3 will replace the current HTTP/2. QUIC lays on top of UDP to provide reduced latency, but with a connection control mechanism to guarantee the same reliability as TCP [133]. There are also proposals to use HAS-based media streaming with protocols such as Stream Control Transmission Protocol (SCTP) [162] and Multipath TCP (MPTCP) [87], which support multihoming and are very important in recent heterogeneous network environments. Noteworthy is that MPTCP is backward compatible with the vanilla TCP, which is very useful for service deployment. Finally, efforts are already being made to develop a multipath QUIC [67] protocol to combine the benefits of these approaches. However, no HAS-based media delivery solution has used it so far.

4 PERFORMANCE-DRIVEN NETWORK FUNCTIONS

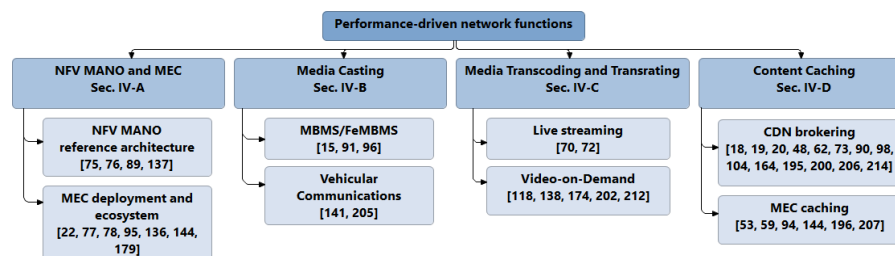


Fig. 2. Performance-driven Network Functions.

VNF brings the architecture, program model, and Application Programming Interfaces (APIs) to deploy specialized media functions as micro-services, exploiting cloud technologies for a smart and agile enhancement of media streaming sessions. This section presents an overview of VNF-based solutions designed to improve media streaming performance. These solutions employ knowledge from network studies and data acquired from live monitoring network traffic. Figure 2 illustrates the significant avenues that performance-driven VNF solutions take. First, we introduce NFV Management and Orchestration (NFV MANO) and Multi-access Edge Computing (MEC), as VNFs rely on these architectures introduced by European Telecommunications Standards Institute (ETSI) and embraced by 5G networks. Then, we discuss the state-of-the-art relevant media-related functions such as media casting, media transcoding, and content caching.

4.1 NFV Management and Orchestration and Multi-access Edge Computing

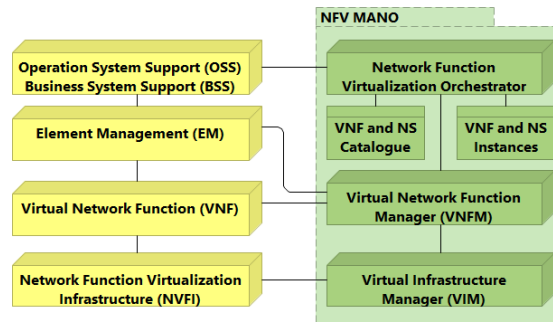


Fig. 3. ETSI NFV MANO architecture.

Apart from the performance leaps on Key Performance Indicators (KPI) in terms of speed, capacity, mobility, and reliability, brought by 5G radio technologies, the network core is also fully engaged in a revolution involving its digital transformation. The concept that one network fits all is over. It is time to adapt the network according to applicable resource efficiency and delivery performance trade-offs. The goal is to allow the network management system to coordinate the network in an agile, programmable and efficient way. This vision is being fueled by the transformation of network functions into dynamically controllable and configurable software components. Network functions are virtualized by exploiting cloud technologies and their scalable mechanisms, and their orchestration is done on top of standardized software solutions. Catalyzed by the network slices concept, the network would also connect groups of virtualized functions devoted to specific data flows or groups of users of specific services. The network would independently handle Service Level Agreements (SLAs) of multiple points of presence (PoPs) over a common bare-metal infrastructure.

To achieve it, 5G network embraces NFV and VNF [137] concepts and comes with an NFV MANO architecture [76], standardized by ETSI. NFV brings the primary virtualization step, providing computing, memory, storage, and network resources from a bare-metal infrastructure (NFV Infrastructure or NFVI). The utilization of NFV contributes to deploying a network by providing hardware and software decoupling. Thus, general purpose hardware, referred to as commercial off-the-shelf (COTS) hardware, can be used to run every network function having a software implementation (VNF). Cloud vendors mainly employ this architecture to provide Infrastructure as a Service (IaaS) solutions. Thus, hosting for systems on top of hardware and connectivity setup is performed on demand. VNFs go a step further in virtualization, deploying specific network functions on top of NFVI. VNFs can be deployed, configured, started, or stopped in a programmable manner. Thus, VNFs are intended to enable modularity, interoperability, scalability, and flexibility when a media streaming service is managed, and the generated traffic is delivered.

NFVI and VNFs are managed and orchestrated by NFV MANO, whose reference architecture is shown in Figure 3. Its functional blocks are:

- Virtual Infrastructure Manager (VIM): It manages and controls physical and virtual resources (compute, storage, and networking resources). Once a VNF is instantiated (VNF Instance or VNFI), it provides the VNFI with the resources it requires;
- VNF Manager (VNFM): It is responsible for managing the life cycle of VNFI through the resources provided by the VIM;
- NFV Orchestrator (NFVO): It combines more than one VNF to create end-to-end services. Several VNFs could share VIM resources and be meant to be used for the deployment of a unique Network Service (NS), e.g., one VNF deploys the back-end and another one the front-end; the combination of the two VNFs constitutes the NS.

Since 5G architecture allows for public and private network deployment, existing NFV MANO-compliant solutions encompass commercial and open-source alternatives for each of the three components. Some examples are Open Source MANO (OSM) [75], whose development is promoted by ETSI, and Open Network Automation Platform (ONAP) [89], supported by Linux Foundation.

All the described technologies that turn network functions into virtualized software systems facilitate a high level of automation and orchestration by network management systems. This trend is being deeply explored and investigated in the current generation of mobile networks (5G), and it will be a key pillar for the next ones (beyond 5G), and MEC infrastructures [179]. MEC architectures enable context-aware applications. It opens computing infrastructures co-located with the base stations to host services close to the mobile users. It aims to exploit the capillary distribution of cloud computing infrastructures at the edge of the cellular Radio Access Network (RAN).

The application of NFV and VNF technologies at the edge and the evolution of the RAN towards software components boosted by open-source software, such as OpenAirInterface [159] or srsLTE [99], ease the integration of MEC services with RAN systems. These open-source solutions implement the Mobile Packet Core (Evolved Packet Core for LTE, 5G Core for 5G) and the RAN on top of open-source hardware. They enable the deployment, management, and orchestration through NFV MANO of both the mobile packet core [156, 161] and RAN [91]. A RAN deployment through NFV and VNF is usually referred to as virtual RAN (vRAN). Furthermore, vRAN is also evolving towards the concept of Open RAN (O-RAN) [219], having open interfaces and network intelligence as key enablers to manage and tailor the network based on vendors' and operators' requirements. O-RAN enables multi-vendor vRAN deployments, resulting in a more competitive ecosystem [93]. In this context, MEC is an NFV MANO-compliant platform that also comes with a specific API to access Radio Network Information (RNI) [77].

Figure 4 shows how network core and edge leveraging virtualization technologies are monitored and orchestrated according to business and technical policies. Different NFVIs can be managed by a unique NFV MANO system and interact with each other through Software Defined Networking (SDN)-enabled communications [46]. The NFV MANO is in charge of deploying VNFs at different NFVIs and managing their life cycle, while the SDN controller configures the forwarding rules of the network to forward the packets between VNFs. Then, a monitoring system collects the performance metrics coming from the different network sections (network core and edge) and processes them to check that the established policies are effective. When considering media streaming, any policy changes to regulate the trade-off between OPEX and QoE involve adjusting the network configuration. It means acting on the NFV MANO system to manage the life cycle of the VNFs and on the SDN controller to update the forwarding rules.

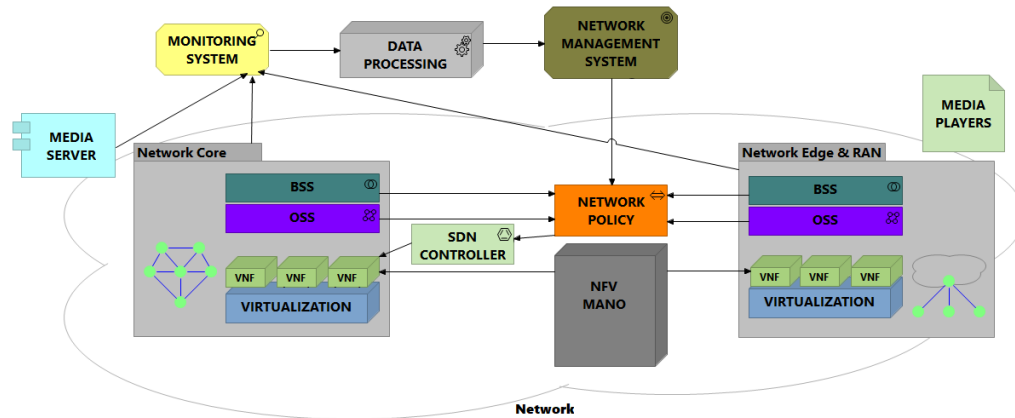


Fig. 4. ETSI NFV architecture applied to media streaming services.

While focusing on the edge architecture, Figure 5 illustrates the MEC components and their interactions with the rest of the RAN and Core Network (CN) building blocks. The MEC host manages the User-plane, while the Control-plane communication is managed by the CN (LTE Evolved Packet Core or 5G Core). Depending on whether the deployment is within an LTE or 5G network, the MEC host is equipped with User-plane Serving and Packet Gateways (SGW-U and PGW-U) or User Plane Function (UPF), respectively. These components are connected directly to the base station (eNB for LTE or gNB for 5G) and provide access to the Internet. Inside the MEC Host, the RNI Service (RNIS) oversees collecting RAN information which is later consumed by the application VNFs. Specifically, VNFs can be designed to exploit such information to increase the overall system performance.

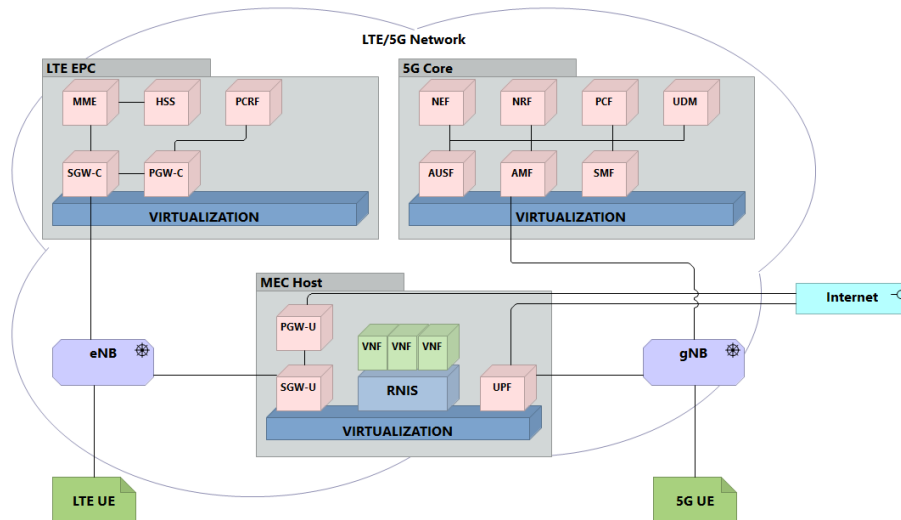


Fig. 5. MEC architecture and connection with RAN and CN.

VNF is also applicable for media-specific network functions beyond the 5G core and RAN, involving:

- media casting, in order to perform massive delivery of live data flows;
- media transcoding, such that streaming rate matches network available bandwidth, resulting in higher quality at destination;

- content caching, including storing popular data to help improve high traffic conditions and managing alternative endpoints to balance the data requests.

These network functions perform specialized operations within the media applications to improve network efficiency, reduce bandwidth overheads, favor idle resource allocation to other network flows, and enhance QoE with enforced KPIs according to SLAs.

ETSI includes several media streaming use cases to be considered for MEC deployment [78] to empower traditional media streaming applications. Typical applications are based on the interaction between the remote server (origin server or CDN) and the client, as shown in Figure 6. In this scenario, the MEC platform can host diverse VNFs, which exploit RNI to get a more comprehensive view of the local conditions to enhance media streaming service. In this line, some solutions, such as [95, 136, 144], exploit standard RAN interfaces and data reports to conclude better decisions for media applications.

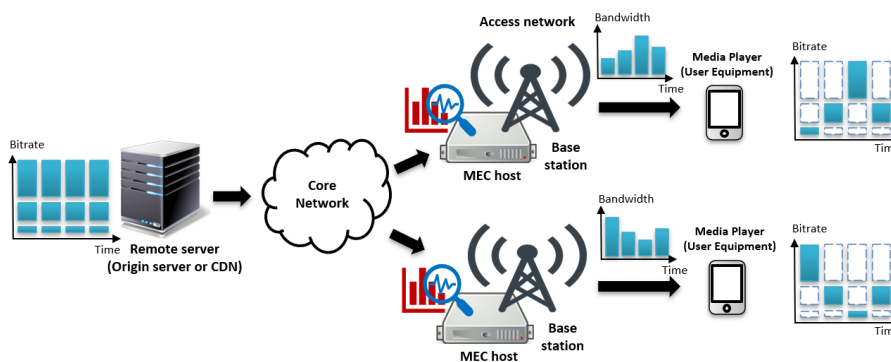


Fig. 6. MEC-powered media streaming.

The following sections analyze how the described core technologies of 5G are applied to expand the network functions with core components for improved media stream delivery, resulting in benefits in terms of enhanced quality and efficient resource utilization. Accordingly, Table 4 compiles and classifies all the research activities exploiting 5G to support performance-aware networking. The classification highlights the main features implemented and secondary aspects, as sometimes the same approach is applicable to more than one solution. Some proposals are limited to architecture design and do not achieve real implementation and experimentation. The implemented ones differ in activation and processing approach, as they could operate in a reactive or proactive manner and, in some cases, embed a processing algorithm using classic data analysis or Artificial Neural Network (ANN)-based approaches.

In any case, all the proposed solutions directly impact the performance of the media streaming systems. Performance increase can be provided in terms of QoS and QoE enhancement, more effective business costs (CAPEX and OPEX), and energy saving. However, most of the proposed solutions do not provide specific validation tests for all the performance aspects (QoS/QoE, CAPEX/OPEX, and energy). They mainly limit their analysis to HAS-centric QoE metrics and do not provide insights on applicable cost models, including business aspects or energy footprint evidence. Furthermore, when considering business aspects only, providing cost models is also difficult due to the lack of business models to use new network assets, i.e., NFVI and MEC. Clear business models are required from network operators to improve the evaluation of business costs [22].

Table 4. Performance driven networking for media streams using 5G technologies.

Main feature	Second feature	Activation	Processing approach	Reference	Network feature	Description	Limitations
Casting	-	Not applicable	Not applicable	[96]	FeMBMS	Design of 3GPP architecture for media multicast	To be validated in field test
Casting	-	Reactive	Not applicable	[91]	FeMBMS, VNF, SDR	Virtualization of FeMBMS with SDR setup	No optimization or overhead evaluation
Transcoding	-	Not applicable	Not applicable	[70]	NFV, VNF, 5G Core	Design of centralized virtual transcoder solution at 5G Core	Lack of mathematical analysis
Transcoding	-	Reactive	ANN	[72]	VNF, MEC	On-the-fly transcoder at the network edge	Cost performance trade-off not conducted
Transcoding	Caching	Proactive	Classic	[174]	L1 MC-NOMA, MEC	Solution empowered by multicarrier non-orthogonal multiple access	Slicing of computation without practical evaluation
Transcoding	Caching	Reactive	Classic	[138]	MEC, VNF	Transcoding and cache location in virtualized edge infrastructures	Latency factor not considered
Transcoding	Caching	Reactive / Proactive	Classic	[202]	MEC, VNF	Transcoding and cache location when content popularity is known (proactive) or not (reactive)	User experience factor not considered
Transcoding	Caching	Proactive	Classic	[118, 212]	MEC, VNF	Transcoding and cache location based on known content popularity	Lack of mobility
CDN Brokering	-	Reactive	Not applicable	[18–20]	L7	Proprietary solution for selection of CDN vendor at startup	Dynamic costs and experience excluded
CDN Brokering	-	Reactive	Classic	[98, 164, 200]	L3 DNS	Performance-driven solution based on DNS resolution	Prevention schemes not discussed
CDN Brokering	-	Not applicable	Not applicable	[48, 90, 195, 214]	L3 DNS	Design of CDN-ISP collaborative solutions	No practical evaluation
CDN Brokering	-	Proactive	ANN	[206]	L7, L3	Solution for proactive CDN selection employing ANN algorithm to forecast network metrics	Small real world field test
CDN Brokering	-	Reactive	Not applicable	[62, 73, 104]	L7	Cloud solution for cost-effective CDN switching	Fairness among users not discussed
CDN Caching	CDN Brokering	Reactive / Proactive	Classic	[207]	L7, L3, MEC	Statistical solution for CDN selection (reactive) and content caching (proactive)	Popularity of content not considered
Caching	-	Not applicable	Not applicable	[53]	VNF, Orchestration	Design of virtual CDNs for media distribution	No evaluation or comparison
Caching	-	Proactive	Classic	[196]	MEC, SDR	Solution at edge exploiting radio network information	Lack of mobility
Caching	Fair QoE	Reactive	Classic	[144]	MEC, SDR	Solution at edge exploiting radio network information	Delay factor not evaluated
Caching	Fair QoE	Reactive	Classic	[94]	MEC, SDR	Solution at edge exploiting radio network information and content popularity	User experience factor not considered
Caching	-	Proactive	ANN	[59]	MEC	Solution for proactive caching employing ANN technologies to predict popularity	Computational workload and time not considered

4.2 Media Casting

For massive delivery of common data synchronously at once, the broadcast is still much more efficient than unicast communications widely employed by cellular networks. That is why 3rd Generation Partnership Project (3GPP) introduced the Multimedia Broadcast/Multicast Service (MBMS) specification in Long-Term Evolution (LTE) Release 9. Later, it evolved towards further enhanced MBMS (FeMBMS) in Release 14 to enable higher per cell bandwidth

625 for MBMS services and simultaneous reception of both unicast and multicast services [96]. Furthermore, Release 16
626 includes feedback for increased reliability [15].

627 As this technology is tied to the RAN system, it has sense in some use cases as firmware/software updates, clock
628 synchronization, alarms, and massive media contents to be turned in the network edge from unicast communications
629 to broadcast signals. It would need support from MEC systems which will turn popular streams into broadcast flows
630 to expand the capacity of a cell. This is feasible as manifests of HAS technologies, such as HLS or DASH, keep the
631 manifests unencrypted even for encrypted contents, allowing simple processing to parse them by intermediaries, such
632 as CDNs or MEC systems, for efficient and intelligent media delivery.

633 This architecture brings three significant benefits by attracting all the ongoing live sessions to consume the broadcast
634 dataflow instead of establishing concurrent unicast sessions:

- 635 (1) Efficiency at the radio link, as the broadcast stream reduces radio link usage. Data traffic is independent of the
636 volume of users since everyone is consuming the same broadcast signal;
- 637 (2) Optimal fidelity, as the network can deliver the maximum resolution (bitrate representation) to all the audience;
- 638 (3) Enhanced QoE, as the media players sharing the radio link do not have to struggle with independent adaptive
639 mechanisms executed in each player competing for the available bandwidth. It means no bitrate or resolution
640 changes to track time-varying network conditions and no freezes to refill the buffer.

641 This approach is possible thanks to virtualization and softwarization paradigms to RAN technologies, making vRAN
642 and the containerization of some RAN network functions such as FeMBMS feasible [91].

643 Specifically, broadcast communications are gaining relevance in the vehicular communications field as they allow
644 synchronous provisioning of common awareness to vehicles, pedestrians, and Road-Side Units (RSU) in a surrounding
645 area. Common awareness can be essential for Cooperative, Connected and Automated Mobility (CCAM) applications
646 related to the safety of autonomous driving [141]. In these applications, media flows are important as the vehicles get
647 fitted with more camera-like sensors capturing the environment and exchanging the raw/compressed data or processed
648 insights/summaries from onboard computer vision systems [205].

657 4.3 Media Transcoding and Transrating

658 As summarized in Table 3, HAS technologies, such as DASH or HLS, are widely employed and need the provision
659 of several representations meaning different resolutions and bitrates [138]. Thus, VNF-based transcoders are being
660 developed under international funding initiatives to empower different use cases, e.g., live 3D media streaming [70] or
661 automotive [16]. Here, the generation of representations at edge servers is gaining relevance to get higher efficiency by
662 distributing the higher fidelity through the core and generating lower bitrate variants (transrating) at the edge. It would
663 reduce overheads in the core to send all the possible media variants. To this end, the media transcoding at the edge is
664 essential [202], stressing the fronthaul capacity and requiring Cloud-RANs (C-RANs) or MEC systems to minimize the
665 network delivery cost. Furthermore, the capillarity of the MEC systems brings a better adaptation to the local needs
666 when transcoding to produce variants.

667 However, transcoding is a heavy process to be performed at a resource-constrained MEC server. Then, it means
668 a challenge for delay-sensitive services. Here, different works deal with the optimal position of transcoding systems
669 in different edge hosts to respond to a distributed demand more efficiently and quickly, where players use a specific
670 base station as a gateway linked to the edge hosts. To overcome this challenge, a mechanism for optimal request
671 forwarding which respects the resource limitations and minimizes serving latency is required [174]. In [212], different
672

677 short/long-term decisions are concluded to deal with the time-varying conditions in terms of demand and network
678 dynamics.

679 Beyond the planning of such a transcoding process, other approaches consider different algorithms for reactive or
680 proactive planning [202]. In this case, the dynamics significantly impact the reaction time and forecast range. These
681 aspects are minimized using a segment duration in the HAS stream, which favors steady short-term conditions as
682 changes come on a segment duration basis.

683 These works focus on enhancing QoS metrics while managing the capacity of each processing asset. However, they
684 do not consider heterogeneous SLAs and cost penalties to apply trade-off policies. It is a primary feature to evaluate, as
685 the required GPU assets for HW-accelerated processing ensure parallelization of transcoding threads and significantly
686 impact infrastructure costs.

687 Finally, it is essential to underline that these solutions are often linked to caching strategies, as both transcoding and
688 caching can be executed at the edge to better match the local conditions, patterns, and demand features. Therefore, joint
689 transcoding processing and caching strategies are designed [118, 138, 174, 202, 212]. In [72], the authors only transcode
690 the content on-the-fly if the content is not cached.

696 4.4 Content Caching

697 *4.4.1 CDN brokering.* Caching is the most employed network function to improve performance when accessing online
698 content, particularly media streaming. A CDN is the most popular network solution to provide caching capabilities
699 in this context. It consists of a geographically distributed network of proxy servers and data centers to provide high
700 availability of the contents. Caching mechanisms are key inside a CDN, as CDN proxy servers work by selectively
701 storing the content so that users can quickly access it from nearby locations. The employment of CDN services by
702 the CPs increased in the last few years as the number of CDN vendors increased. Furthermore, major CPs also moved
703 to multi-CDN strategies to provide more reliable service while streaming their content. Thus, an improved service
704 also generates more satisfaction among the customers. Nevertheless, the different CDNs employed can differ from one
705 CP to another. Static selection of the CDN when a streaming session starts is the easiest and most widely employed
706 solution among CPs. In 2012, this strategy was used by Netflix [20] and Hulu [18], with big similarities [19]. In both
707 cases, they were using three different CDN vendors. They used to map the player device to a CDN depending on its
708 location or the subscriber when the streaming session starts. Moreover, the CDN never changes during the streaming
709 session, even when the performances decrease. Other solutions include client-side CDN selection [164] or Domain
710 Name System (DNS)-based solutions [200]. The client has a privileged position to measure end-to-end QoS metrics
711 (network bandwidth and latency) when choosing the CDN. However, this approach has the disadvantage of producing
712 an uncoordinated decision as each client selects the CDN independently from the others. A DNS-based solution means
713 resolving a fixed hostname owned by the CP into different IP addresses referring to several CDNs. Depending on the
714 DNS resolution, the client receives the content from the appropriate CDN. In any case, a sub-optimal CDN server
715 selection could decrease performance [98], affecting the user's satisfaction.

721 In recent years, other network caching solutions are also raising to empower delivery. The same Netflix changed its
722 streaming strategies. It developed and deployed an in-house CDN, called Open Connect [48], to reduce the dependency
723 on CDN vendors and streaming costs. Moreover, Open Connect is meant to be also run inside the ISP infrastructure, i.e.,
724 closer to the user, to guarantee better performances in terms of network bandwidth and latency [69]. Open Connect
725 also helps Netflix and other CPs have in-house solutions to control better the resources enabled for the streaming
726

729 session and reduce the costs. Anyway, it requires a significant investment to have such a solution that could not be
 730 affordable by small CPs.

731 The Streaming Video Alliance (SVA) is a joint initiative that works on different aspects of media streaming and aims
 732 to standardize the employed protocols and technologies. Its membership includes some major world-wide content
 733 production and delivery agents. Among its activities, the SVA Open Caching Working Group [195] oversees identifying
 734 the critical components of a non-proprietary caching system and establishing the basic guidelines for its implementation
 735 inside the ISP infrastructure. Thus, it wants to promote an architecture similar to Netflix’s Open Connect but with the
 736 advantage of being standardized.
 737
 738

739 Other collaborations between CDN and ISP are proposed in the literature. In [90], ISP provides the CDN provider
 740 with information concerning geographical user distribution and allows the CDN provider to allocate server resources
 741 inside the ISP network. The authors of [214] use a redirection center instance inside the ISP network, which intercepts
 742 the client requests and selects the appropriate CDN server. The process is transparent to the client as the redirection
 743 center employs a CDN surrogate to store the content and instructs an OpenFlow controller to migrate the traffic to the
 744 CDN surrogate. Beyond the employment of multi-CDN solutions, there are still possibilities for improvements. CDN
 745 Brokering [42] is proposed to make CDN utilization in a multi-CDN environment more effective. It redirects clients
 746 dynamically between two or more CDNs.
 747
 748

749 CDN brokers work as switching services that dynamically and seamlessly select the optimal CDN to use at any time.
 750 To achieve this, CDN brokers collect and analyze in real-time the performance of the available CDNs to select the best
 751 one. Thus, network analytics have a prominent role in CDN selection, in contrast with traditional multi-CDN strategies
 752 where the same CDN is kept during the streaming session. The approach from [206] applies ANN technologies to
 753 forecast dynamic demand and changeable performance to make decisions, including cost-performance trade-offs. In
 754 this context, a representative example is Eurovision Flow [73], proposed by the European Broadcasting Union (EBU).
 755 Similar solutions are also provided by Citrix [62] and Haivision [104].
 756
 757

758 **4.4.2 Edge caching.** In [207], a MEC proxy retrieves media streaming metrics of video players at the access point
 759 and CDNs performance metrics to enhance DASH media streaming. The MEC proxy evaluates the performance of
 760 different CDNs and switches players’ sessions when a CDN is underperforming and cannot support the demanded
 761 traffic. Moreover, it features local edge caching to reduce network traffic. Recurrent content is downloaded and cached
 762 once for every player. In [53], a similar MEC cache is proposed for empowering the delivery.
 763
 764

765 With a deeper integration with RAN interfaces, in [196] and [144], the MEC cache is improved by exploiting RNI.
 766 The media segments and representations are selectively cached depending on the network state. In [94], both RNI and
 767 knowledge of segment popularity are employed to decide the segments to cache. Moving from a reactive to a proactive
 768 approach, the authors of [59] empower the edge cache with neural collaborative filtering to predict content popularity.
 769 The predictions are exploited to proactively cache the content at the MEC, as more content popularity means a higher
 770 probability of being requested by the users.
 771
 772

773 Table 5. Comparison of performance driven network function categories.
 774

775 Network media function	775 Scalability	775 Video adaptation	775 Processing resources
776 Media Casting	High	No	Low
777 Transcoding	Low	Yes	High
778 CDN Brokering	Medium	Limited to HAS	Low
779 Edge Caching	High	Limited to HAS	Low

To sum up, the proposed performance-driven network functions for media streaming belong to four main categories, corresponding to specific actions on the content: Media Casting, Transcoding, CDN Brokering, and Edge Caching. Table 5 compares the four categories regarding scalability, video adaptation (resolution and bitrate), and required processing resources. Media Casting and Edge Caching are the best-performing solutions in terms of scalability. They are designed to receive the content from a remote server or CDN and then serve it to the UE(s). As the content is served from RAN or Edge Network, the traffic inside the Network Core is reduced (i.e., there is no redundant traffic). Another advantage is their low processing resource consumption, as they only have to forward packets and do not process the content. CDN Brokering also has low resource processing, like Casting and Edge caching, but its scalability is limited, as it does not reduce network traffic. CDN Brokering selects the remote CDN that the UE has to connect to at any time only. Therefore, the redundancy in Network Core traffic remains, and all the workload has to be absorbed by the remote CDN. Finally, when considering video adaptation (resolution and bitrate), Transcoding solutions are the best ones, as they allow to generate a personalized video representation in real-time. Each client can receive a representation according to its display capabilities and the network state. Nevertheless, this flexibility in video adaptation comes at the cost of using increased processing resources for the transcoding operation. Video adaptation is also available with CDN Brokering and Edge Caching. However, it is possible only when using HAS-based technologies and is limited to the pre-encoded video representations of the content at the CDN.

5 CHALLENGES OF VIRTUAL NETWORK FUNCTIONS FOR MEDIA STREAMING

VNF solutions play a significant role in successfully deploying 5G networks. It is backed by evidence, especially for supporting rich media applications such as multimedia streaming, as described in Section 4. However, VNF applications still require some challenges and open issues to be addressed, as shown in Figure 7. This section discusses and classifies these challenges around some key features studied concerning 5G networks and presents the open issues in the context of the 6G networks roadmap.

5.1 Self-Organizing Networks

Table 6. SON categories and use cases.

Self-configuration	Self-optimization	Self-healing
<ul style="list-style-type: none"> • IP address & connectivity • neighbour & context discovery • radio access parameters • policy management 	<ul style="list-style-type: none"> • load balancing • resource selection • caching infrastructure • coverage & capacity • radio interference management • mobility & handover 	<ul style="list-style-type: none"> • fault detection • fault classification • countermeasures operations

Agile deployment and life cycle management of VNFs exploiting an NFV MANO architecture are essential features to satisfy the expectations of smart 5G networks, but further research is still ongoing to increase network automation. In this context, the Self-Organizing Network (SON) paradigm [79] represents a next step to achieve a fully virtualized and automated network. SON empowers the network with specialized decision-making algorithms which monitor network resources and traffic patterns, and autonomously take actions to enforce or optimize network operations [26]. SON capabilities were initially meant to be included as add-on features of LTE, as 3GPP Release 8 started defining LTE

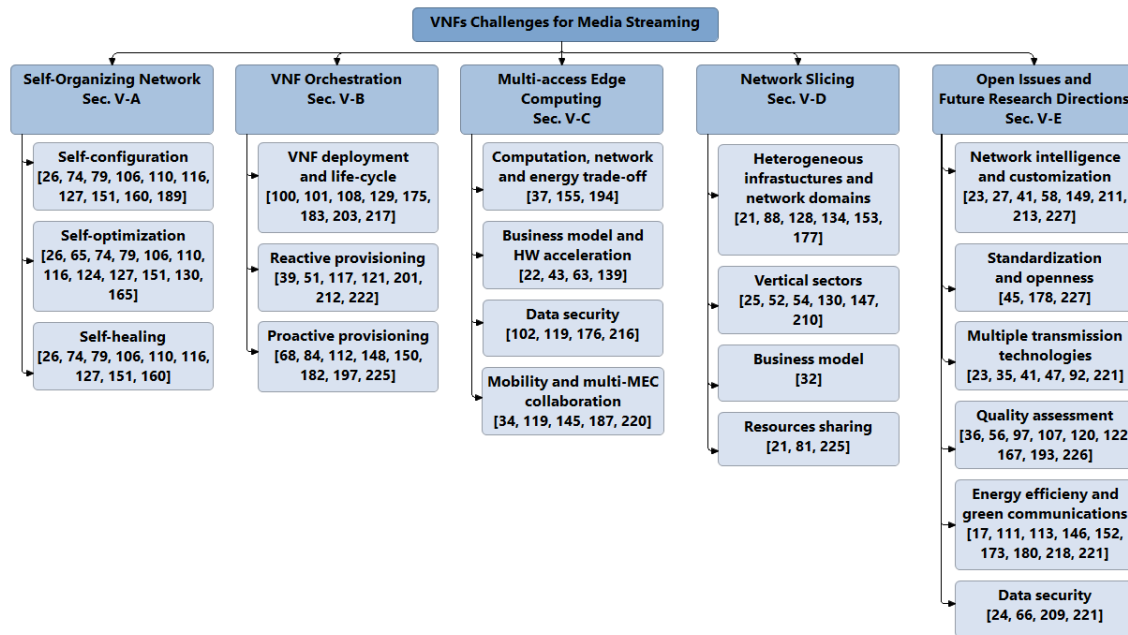


Fig. 7. Virtual Network Functions Challenges for Media Streaming.

and already set the basis for SON concepts and requirements [110]. However, SON is expected to enhance 5G network management by providing automation to cope with increasing network complexity [151].

Specifically, in the media streaming context, SON should provide the necessary network resources and guarantee target QoS or QoE scores when delivering media streams. More generally, SON turns static networks into dynamic ones by configuring network parameters, optimizing the allocated resources, and fixing or preventing issues in real-time.

A SON-enabled system can accomplish tasks belonging to three categories: self-configuration, self-optimization and self-healing [26]. Self-configuration techniques adjust network operational parameters to change network behavior and rules according to specific business policies and node neighborhood context. Self-optimization strategies are dynamically applied to ensure network performance is near optimal. They include real-time network monitoring and performance metrics processing to apply enhanced operational parameters proactively. Self-optimization techniques can be applied in many areas: load balancing, resource selection, caching infrastructure, coverage and capacity, radio interference management, mobility, and handover. Last, self-healing is necessary to generate a prompt reaction when network faults, failures, or any operational range violations occur. The objective is to continuously monitor the system and ensure a fast and seamless recovery, whatever reason causes the failure. In case of a failure event, self-healing functions detect (fault detection) and diagnose (fault classification) it. Then, according to applicable policies and the current setup, the appropriate countermeasure is applied to reestablish the desired network performance.

All these SON flavors need actionable data to process decision-making algorithms. Therefore, it is crucial to collect and exploit network data. Current networks are ready to probe and provide a considerable amount of data. However, specialized intelligence must be deployed within the network to infer valuable and helpful information from the collected data [116]. Such information helps take automatic actions to reach, recover or even improve the network performance. In the context of media streaming, it means that the SON paradigm has the potential to increase the QoS/QoE while decreasing the business costs and energy consumption to maintain the network. In this context, ML

885 techniques will become prominent, even if selecting the suitable algorithm is not trivial and depends on the considered
886 use case [127, 151]. Table 6 shows the most common use cases belonging to the three SON categories, as seen from the
887 network operator’s perspective. Some SON applications are already provided by network vendors and included in their
888 commercial hardware equipment. Some examples are HCL’s SON [106], Nokia’s EdenNet [160] and Ericsson’s SON
889 Optimization Manager [74].
890

891 In any case, SON systems need to have a more comprehensive view of the delivered traffic beyond the metrics
892 from the network functions and including the service domain. Operated SON policies are usually steered by network
893 statistics rather than application characteristics. In [143], the authors analyze network metrics of selected network
894 topologies, such as bandwidth and latency, when the network nodes are generating media streaming traffic. It proposes
895 an ML-powered algorithm to select the appropriate network topology that copes with the demanded network resources
896 for the media traffic. Nevertheless, the communication dynamics of applications delivered on top of the network have an
897 impact on network performance. Thus, the authors of [165] propose to design an application-driven SON to widen the
898 view with both network performance and the user’s QoE metrics. When considering media streaming applications, data
899 are available from network functions in the path and playback devices. Thus, data exploitation inside a SON-enabled
900 system needs further investigation, as the multi-domain data exploitation is still underexplored. In [131], the authors
901 classify the network traffic into four classes, including one for interactive games and video telephony, and propose to take
902 into account the QoS requirements of each class to enable a flexible mobility management (handover) scheme. In [189],
903 the authors introduce a self-organizing Unmanned Aerial Vehicle (UAV)-based communication framework for media
904 streaming. It aims at dynamically selecting the network configuration to achieve better network capabilities in terms of
905 throughput, end-to-end transmission delay, packet re-transmission, and packet loss during video transmission. UAV-
906 based communication is also considered in [65] to design an ML-based scheduling solution to react to the changeable
907 networking conditions and make the best decisions to stream Live Ultra-High-Definition (UHD) Video Streaming from
908 UAV to ground users. Finally, the authors of [124] propose a SON-enabled video transcoder to be deployed within the
909 network to let the video propagate over the network by exploiting knowledge about the state of the network.
910
911
912
913
914
915

916 5.2 VNF Orchestration

917 The orchestration on top of SDN programs and coordinates different network functions to support specific applications
918 and services further. The specialized network functions in the media context include servers, load balancers, caches,
919 transcoders, transraters, and encryptors. These functions are containerized to make media functions ready to be
920 deployed on top of cloud infrastructures and instantiated depending on resource policies. Thus a media system gains
921 virtualization and flexible scaling of cloud resources available as NS at the network core and edge. This evolution is
922 essential to reduce superfluous OPEX and lead to shorter time-to-market and lower CAPEX. In this context, DevOps
923 and micro-services platforms gain relevance to accomplishing the Function-as-a-Service (FaaS) paradigm. Here, the
924 ambition is to meet media-intensive use cases, where an on-demand VNF instantiation and deployment enables an
925 elastic infrastructure that enforces the performance and restrains costs [203].
926
927
928

929 ETSI NFV MANO architecture brings several functionalities at different layers operating different systems. For
930 example, MANO manages the allocation of virtualized resources of the NFVI layer using the VIM system, the VNFs
931 live-cycle at the VNF layer driving the VNFM, and orchestration scaling or deleting NS instances steering the NFVO
932 [100]. MANO is expected to support optimal media delivery in flexible networks toward handling target QoE and cost
933 trade-offs. However, this has not been fully achieved yet. The main problems of orchestration, still open and under the
934 focus of different research lines, are related to resource allocation and the instances placement.
935
936

937 First, deploying a VNF over a distributed platform requires allocating network and computing assets to be provisioned
938 to host the VNF. Network and computing resource allocation is a challenging feature whose interest is raising and
939 focusing on VNFs deployment, and life cycle management [108]. In this context, the NVFO is in charge of selecting the
940 appropriate resources, among the available ones at the NVFI, when deploying a VNF, which mainly includes: VNF Chain
941 Composition (VNF-CC) and VNF Placement and Chaining (VNF-PC), VNF Forwarding Graph embedding (VNF-FGE)
942 and VNF Scheduling (VNF-SCH). VNF-CC deals with the composition of several VNFs to be deployed jointly by the
943 NFVO. VNF-PC aims for the optimal placement and the required instances of VNFs needed to deploy Service Function
944 Chains (SFCs) while optimizing the cost of resource provisioning. How the traffic flows between VNFs is also described
945 through the definition of VNF Forwarding Graphs (VNF-FG). Thus, any network service can be considered composed of
946 a set of VNFs and VNF-FGs. Each VNF executes a small function of the entire application or service [217]. It aims to find
947 appropriate resources and locations to allocate the VNFs in NFVI. At this stage, resource selection and optimization
948 must be accomplished concerning the specific constraints defined by SLA [183]. Finally, VNF-SCH determines how to
949 schedule the processing operations of the deployed VNFs [175].
950

951 Specifically, when media VNFs come into play, the dynamic allocation and provisioning of VNF resources are essential
952 to mitigate or even prevent QoS/QoE violations when bottlenecks at some network path arise. When the VNFs are
953 already deployed and running, the required resources may vary during their life cycle, as they depend on the user
954 demand of the running function provided by the VNFs. Allocated resources could be optimized to fit the user's variable
955 demand. Increasing or decreasing the allocated resources means the VNFs also need to scale up/down dynamically.
956 Then, an efficient orchestration and automation of the VNFs require supporting this dynamic allocation of resources.
957 Changes in resource allocation should be applied according to real-time network traffic and service demands. Dynamic
958 resource allocation can be performed in reactive or proactive manners. Here, transcoder [129], and transrating [101]
959 systems are instantiated with VNF technologies triggered by a changeable demand applying Machine Learning (ML)
960 techniques to automate the monitoring and actuation.
961

962 Second, the placement of caches, transcoders, and transraters has been profoundly studied to solve this NP-hard
963 problem. Here, [51] analyzes the QoE feedback to on-the-fly conclude an optimal placement of transcoding VNFs. Others
964 such as [201, 212] go beyond and optimize at the same time caching, transcoding, and transrating in two timescales,
965 short-term and long-term, to favor steadiness.
966

967 However, most of the proposed solutions include two main limitations, the lack of a cost function to add effectiveness
968 to the equation and the inability to prevent issues firing actions to avoid SLA violations dynamically.
969

970 Thus, the related literature often employs an over-provisioning strategy, where the allocated resources for each
971 VNF are larger than required. This approach is inefficient in terms of OPEX and energy consumption generated by the
972 allocated resources which are not employed. It means that this approach is not cost-effective, as it is clear that adjusting
973 the resources allocated for the VNF to the actual demand would avoid over-provisioning and reduce costs. Some
974 approaches optimize the cache and transcoding, including simple cost models concurrently [39, 121]. They formulate a
975 constrained optimization problem to minimize each user request's total caching, computing, and bandwidth utilization.
976 In [222], the authors also include the hardware acceleration costs. Moreover, [117] meets the placement and chaining of
977 VNFs for media cache, including VNF instantiation, migration, hosting, and routing costs.
978

979 Furthermore, the solutions widely involve a reactive provisioning and allocation approach, which means changing
980 the allocated resources to react when traffic and/or demand change. In [68], the authors design a proactive algorithm
981 for VNF placement and allocation of caching and protection VNFs. They also aim to minimize the OPEX by considering
982 the bandwidth and host resource consumption trade-offs under diverse workload variations. Most of the literature
983

989 on generic VNF orchestration driven by traffic demand prediction employs ANN algorithms [182, 197]. However, the
990 application of such algorithms in practical solutions is limited, being mostly theoretical. Among the most innovative
991 solutions proposed, [84] describes a novel Follow The Regularized Leader (FTRL) online algorithm for VNF provisioning,
992 which handles workload fluctuations. The solution in [148] employs an ANN algorithm to predict future resource
993 requirements for each VNF contributing to a network service. The authors of [225] propose the POLAR algorithm,
994 combining online learning and optimization to proactively provision resources with VNFs provisioning. In contrast, the
995 VNFs chaining in a network service is ignored. In [112], proactive failure recovery is proposed when considering VNF
996 deployed at distributed edge computing nodes. Finally, the authors of [150] propose a multi-layer resource allocation
997 solution, which aims to proactively provide resources to the VNFs deployed in several VIMs and network resources
998 between VIMs.
999
1000
1001

1002 5.3 Multi-access Edge Computing (MEC) 1003

1004 MEC represents a novel technological solution integrated into 5G networks to bring computation closer to the user.
1005 MEC infrastructures create new potential revenue flows to network operators opening their edge infrastructures to
1006 host specialized services at the network edge. There are many aspects that require investigation to achieve a complete
1007 integration of MEC into the current network architecture and services. However, some avenues are already seen as
1008 highly beneficial for MEC deployment and use. For instance, media streaming is a crucial application of MEC solutions,
1009 as ETSI considers it one of MEC core use cases [78]. MEC platforms can host edge services to empower media streaming
1010 applications traditionally based on server-client communications. As explained in the previous section, MEC and VNFs
1011 enable the deployment of innovative media-related services such as media casting, transcoding, and content caching.
1012
1013

1014 More specifically, MEC resources are exploited by both the server and clients to offload computation tasks [37, 155].
1015 Offloading server tasks reduces network traffic and latency, as the processing is performed close to UEs. Computation
1016 offloading at the MEC is necessary to enable use cases where the remote server has a very high delay and/or the
1017 client has not enough computing capabilities. 6G potential use cases include resource-expensive and delay-sensitive
1018 applications such as augmented and virtual reality [154]. In any case, it is important to note that MEC resources are
1019 shared between different service providers, but how the resources are distributed among different service providers
1020 is still undefined. The authors of [155] propose to allocate MEC resources proportionally to each service provider's
1021 demanded resources and payment. On the other side, if it is the UE that offloads its tasks to the MEC host, it reduces not
1022 only the device computation load but also the power consumption on the device, as computing-intensive tasks heavily
1023 impact the battery duration. In [37], a video telephony application employs MEC to encode the content. It reduces
1024 processing operations at the UE but increases network traffic since uncompressed raw content is sent to the base station.
1025 The authors focus on power consumption but do not consider operational costs generated by using the MEC platform.
1026 In general, how to balance network traffic, power consumption, and operational costs trade-off needs to be studied. In
1027 [194], optimization of the allocation of both computing and network resources is discussed while considering energy
1028 efficiency. Even in this case, operational costs are not considered in the optimization problem. In general, business
1029 aspects raise complex discussions due to the lack of a clear business model [22]. MEC needs a business model equivalent
1030 to the one applicable in cloud computing infrastructures. However, unlike cloud computing, the decentralized location
1031 and utilization of shared resources between services make the cost model more complex. Resource accounting and
1032 monitoring must also be determined to create a complete business model. The debate on the business model is even
1033 more intricate if we consider hardware-acceleration assets, such as GPUs, required to accomplish critical tasks where
1034 general-purpose hardware (CPU) has limitations [139]. Some works suggest employing Field-Programmable Gate Array
1035
1036
1037
1038
1039
1040

(FPGA) approaches instead of GPU solutions due to their reduced price and power consumption [43, 63], but this possibility is again underexplored.

Regarding accessible information at MEC, the API to communicate with RNIS [77] has been recently standardized, and its development is ongoing [30, 199]. It means that services running at the MEC host cannot be further optimized. When RNIS implementations are available, edge services will embed more complex and precise algorithms (classic or ML models) to exploit RNI to improve their operations and overall system performance. However, improved capabilities due to RNI exploitation raise some security concerns regarding managing information at MEC hosts, an aspect that needs further investigation [102, 119]. In order to exploit a MEC decentralized approach, the deployment of location-aware services is necessary. Thus, mechanisms for user privacy protection and anonymity are needed. Moreover, modification of the networks to introduce MEC capabilities opens the door for potential attacks, including DDoS attacks, malware injection, authentication, and authorization attacks [176, 216]. Use cases like surveillance and CCAM may also include videos containing sensitive data. Thus, video and data stream anonymization [44] is an important matter to consider for further improving MEC-based solutions.

Mobility remains another primary concern and is becoming critical, as the explosion in availability and type of mobile devices (e.g., smartphones and tablets) involves an increasing number of UEs to be served. In the same way that connectivity is guaranteed when moving from a cell to another in a cellular network, migration support for MEC services is also required. Consequently, the investigation on multi-MEC cooperation should be addressed to guarantee seamless session migration across MEC servers [119, 187]. Moreover, this seamless migration becomes critical as delays may raise security concerns in some use cases. An example is represented by CCAM use cases, where the video is streamed between two self-driving cars that make real-time decisions [34].

From the perspective of media services, user QoE plays an important role and a wide MEC deployment definitely should target it, especially as transcoding and caching capabilities would be provided closer to UEs. Balancing the cost of MEC-based caching and transcoding and provision of high user QoE is an essential direction for future research [119]. Moreover, finding suitable locations where MEC instances should be deployed becomes relevant, as it may affect the fulfillment of the demanded requirements. It is especially true for low latency multimedia services, where the distance between the MEC host and UE affects the overall delay [145]. Finally, content caching mechanisms in the network have been studied both at the core and at the edge, but a convergent solution has not been identified yet. Caching solutions that integrate both core and edge caching could improve network performance in terms of energy consumption, network throughput, latency, and user QoE [220].

5.4 Network Slicing

Network slicing [21, 88] is introduced in 5G networks as a solution involving several virtual/logical networks (slices) on top of a shared physical network, where each virtual/logical network delivers the traffic generated by a specific service [153, 177]. To achieve it, it provides network and computing resources across different networks. It can be considered that a network slice is associated with a set of network resources and VNFs, which that slice can provide. In this context, NFV, together with SDN, plays an important role, especially in deploying and managing network slices [163, 224]. NFV enables life cycle management and orchestration of the VNFs, while SDN allows for the configuration and control of the routing and forwarding planes of the underlying network infrastructure, providing communication between the deployed VNFs. This results in a logical network of resources and VNFs built over a common underlying physical infrastructure, separated into diverse network slices. The introduction of network slicing is necessary since a best-effort network cannot guarantee that appropriate network resources are offered in each use case. With network slicing,

1093 different sets of proprieties can be established, each one associated with specific network resources and supporting
1094 relevant use cases. For instance, use cases requiring high throughput, i.e., on-demand video streaming, are logically
1095 separated from use cases requiring low latency, i.e., real-time video streaming or mission-critical services, and are
1096 supported by different slices. Each network slice provides support for service in terms of end-to-end connectivity,
1097 meaning that network slicing provisioning refers to three different aspects: at the air interface, in the RAN and in the
1098 CN [128, 134].
1099

1100 Network slicing at the air interface refers to partitioning physical radio resources (physical layer or L1) into subsets
1101 of several physical resources. Each subset associated with a different network slice is mapped into logical resources to
1102 be provided to the Medium Access Control (MAC) sublayer at the datalink layer (or L2) and higher layers.
1103

1104 In the RAN, network slicing changes RAN operations, including MEC-operated ones, such as device association and
1105 access control, from a cell-specific perspective to a slice-specific one. Thus, the RAN operations are service-oriented
1106 instead of physical cell-oriented. Configuration of control and user planes is tailored and/or tuned considering the
1107 requirements of each slice individually. Then, factors such as QoS requirements, traffic load, or type of service/traffic
1108 are prominent when operating the RAN.
1109

1110 Finally, network slicing in the CN enables the definition of vertical networks, where each one aims to support a
1111 service belonging to a specific vertical industry. NFV and SDN have a higher impact in this aspect of the network,
1112 where each vertical industry should be able to run its VNF-specific solutions. CN needs flexible management to enable
1113 resource scalability and migration when required by the network traffic associated with a service.
1114

1115 A videoconferencing system is deployed in [25] through the deployment of two different slices to split audio and
1116 video transmissions, as they have different requirements in terms of network throughput. In [210], the authors focus
1117 on the eHealth vertical, where services are typically media-rich and mission-critical and are high QoS demanding.
1118 Then, a MEC-based application, empowered with end-to-end network slicing, is designed and developed to enable
1119 in-ambulance applications. The application is accessed by paramedics in the ambulance and sends audiovisual data
1120 to the hospital/doctor. The same vertical is addressed by [54] to enable real-time communication between hospital
1121 staff and patients. In [52] and [147], applications of network slicing for Vehicle-to-Everything (V2X) services are
1122 investigated. Different use cases are considered in a vehicle, including those related to safety and traffic efficiency,
1123 autonomous or teleoperated driving, media & entertainment, and remote diagnostics. Each use case requires different
1124 latency, throughput, and communication reliability requirements. Consequently, different network slices with different
1125 configurations are required on the same physical network of resources. The authors of [130] present several use cases
1126 belonging to different verticals, such as protection and smart metering in the smart grid sector, car and passenger data
1127 exchange in an intelligent transportation system and best-effort data delivery in a multimedia system. Each use case
1128 and vertical sector requires different capabilities in terms of latency and throughput. The different types of traffic are
1129 prioritized by splitting them into specialized network slices.
1130

1131 Network slicing-related research has increased importance in the current 5G network context. Ongoing challenges
1132 include solutions to allow wide employment and operation of slices for different industry verticals. Most slicing
1133 operations relate to the exploitation of resources the network operator provides. However, the effects of changes in
1134 network operators' business models for operating network slicing are unknown [32]. The increase in the number of
1135 devices belonging to different verticals and their mobility management in the presence of different technologies (LTE,
1136 5G, Wi-Fi) also needs further investigation [225]. An end-to-end network slice implies that slice segments potentially
1137 stretch across different administrative domains. There are two requirements to achieve a unified control of the network
1138 slice. First, an exchange point that performs the resource negotiation between different administrative domains is
1139
1140
1141
1142
1143
1144

necessary to enable multi-domain slices. Then, standardized APIs should make transparent the underlying domains and simplify the negotiations to provide control on the slice [21]. Finally, network slicing leverages algorithms to accommodate applications with diverse requirements over the same physical network. Thus, complex algorithms are necessary for deciding how to allocate efficiently, manage, and control the physical resources shared across diverse slices [204]. Concerning these algorithms, the application of ML in network systems is capturing increased research attention lately, and this trend is expected to continue in the future [81].

5.5 Open Issues and Future Research Directions

The benefits of virtualization for media streaming communications will become increasingly evident in the next few years, as the 5G coverage will be extended. Complementary technologies such as MEC, SON, and network slicing are still not fully integrated. Further efforts in integrating all these new paradigms and/or architectures are envisioned to provide a more efficient and intelligent network [115].

ML-powered network intelligence to manage NFV and VNFs is only partially achieved in 5G networks, but it will also be a key factor for the future 6G networks [227]. The concept of Intent-Based Network (IBN) [213] means employing ML solutions to transform business intents into a network configuration, operation, and maintenance strategies. In order to meet the massive service demands and overcome limitations due to time-varying network traffic, the network can continuously learn and adapt to the time-varying network environment based on the massive collected network data in real time. An intelligent-native network exploits ML algorithms to improve its capabilities and reduce the business costs for service deployment and management [58, 149]. The advantages of an intelligent-native network are two-fold. First, the network can analyze the user's behavior in real-time and autonomously learn his needs to predict his future behavior. Then, the user's information can be employed for network customization to achieve a user-centric network [211]. Second, the network can meet changing requirements of a network service during its life cycle by autonomously matching the requirements to the corresponding network communication, computing, and caching assets. This is also valid for new emerging services. Holographic (AR and VR) and haptic communications are meant to be wider available thanks to the future 6G network [41]. Moreover, the global COVID-19 pandemic is accelerating the digital transformation of multiple and heterogeneous verticals, such as the development of new services for smart cities and innovation in the eHealth, including telemedicine, medical and thermal imaging, and robotics for medicine practice [23, 27].

Openness is also essential to achieve flexible networks and services [227]. An open network platform and interfaces (O-RAN, NFV MANO, SDN, etc.) allow interconnection and interoperability of different vendors, which is essential for sharing physical infrastructure. Thus, agents of diverse vertical industries may deploy their private physical infrastructure and manage it through NFV MANO solutions and SDN controllers independent from public networks operated by mobile network operators [178]. The standardization process will continue in the following years to fulfill the remaining gaps and guarantee interoperability of heterogeneous implementations of open network solutions [45].

The cooperation of different physical networks will also attract attention. Multiple Radio Access Technology (multi-RAT) aims to employ different access networks to improve the overall connectivity [92]. Its application to improve media streaming is already being investigated [35, 47]. However, new transmission solutions based on space, UAV-based and underwater communications will be integrated with terrestrial ones [23, 41]. Flexibility to operate the network at any level (spectrum/band, physical and MAC, etc.), despite the different involved technologies, will be imperative [221].

In a media streaming context, performance assessment focuses on evaluating the system employed to stream the content. The performance assessment is done by employing metrics and collecting measurements and/or estimations of

1197 such metrics. These involve quantifiable values to track and monitor a streaming session, i.e., video resolution and
1198 bitrate, or a related factor that can influence it, i.e., network bandwidth and latency, at any network device. Moreover,
1199 there exists lots of different metrics to describe the media system under different points of view, e.g., ISP, CP or user's
1200 view, including both objective (QoS) [56, 120] and subjective ones (QoE) [122, 193, 226]. In the context of the current
1201 5G networks and devices, metrics collection is not so easy, as the heterogeneity in devices and virtualization solutions
1202 increases the complexity. In a virtualized and heterogeneous network environment, media and network-related metrics
1203 to describe the media system's performance should be further expanded to characterize this new context where the
1204 media streaming solution is deployed. Monitoring systems represent an essential part of a virtualized architecture.
1205 Monitoring the resources means defining metrics that apply to the virtualization environment and describing the state of
1206 the NFVI and the running VNFI. It allows to manage the physical and virtual resources more efficiently, including those
1207 allocated for encoding operations [36, 167] or content caching [97]. Furthermore, the network flexibility, guaranteed by
1208 NFV and SDN, raises the discussion on business costs monitoring [107]. The CP may be able to enforce its business
1209 policy to balance the QoS/QoE and business costs trade-off. Thus, the definition of business metrics is necessary to
1210 allow CP's business decisions when running each component of its media streaming solution.
1211

1212 Energy efficiency and green communications [113] are envisioned to enable more sustainable networking [173].
1213 Energy efficiency concerns are also relevant for media streaming services [17, 146]. 6G network will be developed by
1214 taking into account self-sustainability devices and solutions [111]. Here, low-power wireless devices could harvest energy
1215 from the available high-power radio waves [221]. Thus, battery-free implementations will be an interesting topic to be
1216 further explored in different use cases. Clearly, IoT applications will benefit the most from battery-free implementations
1217 [218]. However, the high presence of video content traffic also suggests their use for media communications to support
1218 a self-sustainable 6G network [152, 180].
1219

1220 Finally, the growth of network and media traffic will have consequences on security. Critical media use cases, e.g.,
1221 eHealth applications [24], and autonomous driving systems [66], need to be secured with security mechanisms that
1222 will complement the conventional cryptography-based ones. Increasing security will be assured with the design of
1223 cross-layer algorithms to protect the transferred information [209, 221].
1224

1225 6 INTERNATIONAL INITIATIVES

1226 Employing VNFs for media streaming is a research topic that has attracted the attention of international organizations
1227 and international funding programs for many years. Recently, the European Commission has funded numerous research
1228 projects to develop and implement VNFs for different research scenarios and vertical industries. Table 7 summarizes
1229 the most relevant actions. The project list includes initiatives targeting generic architectural design (i.e., CogNet [7],
1230 SELFNET [11] and SliceNet [11]), activities building testbed environments and pilot environments for use case definition
1231 and testing (FLAME [14], SoftFIRE [13] and 5GTango [4]), projects targeting specific application verticals and developing
1232 required functionalities (5G-Media [1], 5Growth [3], 5GCity [2]) and finally international software communities to
1233 provide open-source platforms (OpenAirInterface Software Alliance [10], Mosaic5G [8] and O-RAN Alliance [9]).
1234

1235 Regarding architectural definition, SELFNET H2020 project designed and tested an autonomous network management
1236 framework capable of the automatic detection and mitigation of common failures in the network [55]. Among others, it
1237 proposed the smart integration of state-of-the-art technologies in NFV. One of the outcomes is presented in [157], where
1238 the SELFNET framework preserves the network's health maximizing the QoE and minimizing the end-to-end energy
1239 consumption. SliceNet project addressed both management and control planes of network slicing to leverage QoS for
1240 sliced services [57]. The project proposed an integrated network management, control and orchestration framework
1241

Table 7. Major SDN/NFV related research activities.

Project	Period	Area	References
CogNet (Building an Intelligent System of Insights and Action for 5G Network Management)	2015-2018	Architecture	[7, 31, 38]
SELFNET (Framework for Self-Organized Network Management in Virtualized and Software Defined Networks)	2015-2018	Architecture	[11, 55]
SliceNet (End-to-End Cognitive Network Slicing and Slice Management Framework in Virtualised Multi-Domain, Multi-Tenant 5G Networks)	2017-2020	Architecture	[12, 181, 210]
SoftFIRE (Software Defined Networks and Network Function Virtualization Testbed within FIRE+)	2016-2018	TestBeds	[13, 132]
FLAME (Facility for Large-scale Adaptive Media Experimentation)	2017-2020	TestBeds	[14, 103]
5GTango (5G Development and validation platform for global industry-specific network services and Apps)	2017-2020	TestBeds	[4, 168, 192]
5G-Media (Programmable edge-to-cloud virtualization fabric for the 5G Media industry)	2017-2020	Application Verticals	[1, 29, 50]
5Growth (5G-enabled Growth in Vertical Industries)	2019-2021	Application Verticals	[3, 135]
5GCity (A Distributed Cloud and Radio Platform for 5G Neutral Hosts)	2017-2020	Application Verticals	[2, 64]
OpenAirInterface Software Alliance	2014-	Development Platforms	[10, 159]
Mosaic5G	2016 -	Development Platforms	[8, 158]
O-RAN Alliance	2018-	Development Platforms	[10, 219]

and applied the concept to various use cases. One of those cases related to multimedia health services is described in [210], where demanding QoS requirements (i.e., latency) must be fulfilled. The network intelligence topic is tackled by CogNet, a project focused on realizing the well-known control loop MAPE (Monitor, Analyze, Plan and Execute) with ML techniques and policy-based mechanisms for a vision of softwarized 5G networks. CogNet validated its vision in different use cases that include SLA Enforcement and Mobile Quality Predictors [31], [38].

A second group of projects aimed at creating platforms and testbed environments where specific use cases, applications, algorithms, and interoperability solutions could be designed and validated. FLAME stands out in this area as a facility for large-scale experiments in the Adaptive Media field. Since 2017, FLAME has hosted different proposals [103] to offload video content proactively to the edge of the network on an SDN/NFV environment. FLAME tests include AR applications and smart video surveillance for aiding impaired citizens. SoftFIRE is another testbed environment to experiment with VNF services and applications in SDN/NFV. SoftFIRE aims to assess the maturity level of solutions in programmability, interoperability, and security and show how they can support the full potential of these properties in a real-world case [132]. Finally, 5GTango puts the focus on network flexible programmability [168] by providing software development kits (SDKs) [192]. This project included qualification and verification mechanisms as well as a modular service platform to bridge the gap between business needs and network operational management systems. 5GTango was demonstrated in two verticals through specific pilots: advanced manufacturing and immersive media [168].

The third category encompasses some examples of projects designing the required building blocks that enable the applications for specific vertical sectors. 5GCity was an H2020 project aiming at designing, implementing, and demonstrating a distributed cloud and radio platform for municipalities and infrastructures with neutral hosting

1301 capabilities. One of the project's primary outcomes was the 5GCity Orchestration Platform, which supported the NFV
1302 MANO model. In [64], the authors demonstrate that the virtualized platform could address different use cases related
1303 to media streaming, such as real-time video acquisition and production at the edge, UHD Video Distribution, and
1304 immersive services or mobile real-time transmission. 5G-Media [1] exploits the principles of NFV and SDN to facilitate
1305 the development, deployment, and operation of VNF-based media services on 5G networks. Key in this project is the
1306 development of a platform for service virtualization that provides an advanced cognitive management environment for
1307 the provisioning of network services and media applications [29]. The use cases include teleimmersive gaming, mobile
1308 journalism, and UHD content distribution [50]. 5Growth [3] supports diverse industry verticals by developing the tools
1309 for interfacing those verticals with the 5G end-to-end platforms. The system provides the creation of network slices
1310 with closed-loop automation and SLA life cycle service control. ML-driven solutions are also part of the project targets
1311 to optimize access, transport, core and cloud, edge and fog resources across multiple technologies and domains [135].
1312

1313 Finally, OpenAirInterface Software Alliance [10], Mosaic5G [8], and O-RAN Alliance [9] are mixed academic and
1314 industrial communities to create ecosystems of open-source projects for studying, building, and sustaining open flexible
1315 and integrated 5G network. OpenAirInterface Software Alliance [10] provides 5G network tools extensively used by
1316 researchers from both industry and academia. This initiative gathers developers from around the world, who work
1317 together to build wireless cellular RAN and CN technologies [159]. Mosaic5G [8] develops a set of 5G software solutions
1318 and has already hosted experiments targeting low latency MEC services, orchestration solutions, and programmable
1319 RANs [158]. O-RAN Alliance [9] is pushing the standardization and the development of the O-RAN. RAN industry is
1320 moving towards open, intelligent, virtualized, and fully interoperable RAN [219].
1321

1322 7 CONCLUSIONS

1323 The popularity of media streaming services is constantly growing due to an increasing number of users, the diversity of
1324 rich media experiences, e.g., online gaming, VR/AR applications, and the utilization beyond entertainment services,
1325 i.e., in professional application domains. The latest smart mobile devices also have an essential role in the success of
1326 media streaming, as their processing and rendering capabilities support streaming content at very high resolutions, e.g.,
1327 UHD or 4K. Consequently, media streaming traffic represents a substantial share of the total Internet traffic and, more
1328 importantly, an increasing one.

1329 To cope with this increasing media traffic and high dynamics of network performance and user mobility, improved
1330 network capabilities are required to maintain high QoS and QoE performance while achieving the best trade-off with
1331 business costs and energy efficiency. 5G networking is bringing new possibilities to deploy intelligent network functions,
1332 which monitor the media streaming service through live and objective metrics and boost it in real-time. Under the 5G
1333 umbrella, NFV will have a prominent role in virtualizing network functions and their management and orchestration.

1334 In this context, this work provided a state-of-the-art on VNFs applied to media streaming. To this end, we considered
1335 the factors that concur with the design and implementation of a stable VNF. VNF-based media streaming functions
1336 monitors and collects performance metrics to tune their configuration to enhance the processed media streaming
1337 sessions. Beyond that base features limited to session quality, advanced solutions adapt the VNF life cycle deployment
1338 and management, dealing with cost-performance trade-offs and the balanced setup when infrastructures with several
1339 VNFs come into play. Moreover, network traffic monitoring and analysis allow the creation of models to approximate the
1340 behavior of the network and predict future network events to take actions proactively. Thus, any network malfunction
1341 or issue that affects the media steaming session can be prevented.

Several VNF solutions to improve media streaming are presented. Solutions, including media casting, transcoding, and content caching, can be employed at any network segment. Thanks to the NFV MANO architecture, the deployment of VNFs is interoperable and can be automatically operated and orchestrated across a virtualized infrastructure. However, the exploitation of media streaming VNFs is not limited to the Network Core, but they can also be run at MEC hosts. Capillarity of the MEC allows computing operations close to the base stations and reduces the latency when dealing with live streaming services.

Finally, research challenges and open issues have been presented in the realm of VNFs applied to media streaming services. The main venues where the research will focus in the next few years are the achievement of dynamic VNF deployment and orchestration, complete MEC integration, and network slicing. Long-term research will also address the strong employment of ML to foster network capabilities and the utilization of open network solutions and/or new access technologies, also combining them to increase capacity. Green communications and security will also be significant concerns, as future networks should reduce their environmental impact and guarantee the security of the processed information. In conclusion, VNFs represent an essential enabler in improving media streaming services. However, despite the research done under international initiatives pushing 5G and network virtualization, several research challenges still exist and provide opportunities for further research activities.

ACKNOWLEDGMENTS

This work was fully supported by the Open-VERSO project (Red Cervera programme, Spanish government's Centre for the Development of Industrial Technology) and by the European Union's Horizon 2020 Research and Innovation programme under Grant Agreement 870610 for the TRACTION project. The support of the Science Foundation Ireland (SFI) Research Centres Programme grant 12/RC/2289_P2 (Insight SFI Research Centre for Data Analytics) and EJ/GV Grant IT1436-22 are also gratefully acknowledged.

REFERENCES

- [1] [n.d.]. 5G-Media (Programmable Edge-to-Cloud Virtualization Fabric for the 5G Media Industry). <http://www.5gmedia.eu/>
- [2] [n.d.]. 5GCity – A distributed cloud & radio platform for 5G Neutral Hosts. <https://www.5gcity.eu/>
- [3] [n.d.]. 5Growth (5G-enabled Growth in Vertical Industries). <https://5growth.eu/>
- [4] [n.d.]. 5GTango (5G Development and validation platform for global industry-specific network services and Apps). <https://5gtango.eu/>
- [5] [n.d.]. *The affordability of ICT services 2020*. https://www.itu.int/en/ITU-D/Statistics/Documents/publications/prices2020/ITU_A4AI_Price_Briefing_2020.pdf
- [6] [n.d.]. *Broadband Commission Agenda for Action for Faster and Better Recovery*. <https://www.broadbandcommission.org/COVID19/Pages/default.aspx>
- [7] [n.d.]. CogNet (Building an Intelligent System of Insights and Action for 5G Network Management). <http://www.cognet.5g-ppp.eu/>
- [8] [n.d.]. Mosaic5G. <https://mosaic5g.io/>
- [9] [n.d.]. O-RAN Alliance. <https://www.o-ran.org/>
- [10] [n.d.]. OpenAirInterface – 5G software alliance for democratising wireless innovation. <https://openairinterface.org/>
- [11] [n.d.]. SELFNET (Framework for Self-Organized Network Management in Virtualized and Software Defined Networks). <https://selfnet-5g.eu/>
- [12] [n.d.]. SliceNet (End-to-End Cognitive Network Slicing and Slice Management Framework in Virtualised Multi-Domain, Multi-Tenant 5G Networks). <https://slicenet.eu/>
- [13] [n.d.]. SoftFIRE (Software Defined Networks and Network Function Virtualization Testbed within FIRE+). <https://www.softfire.eu/>
- [14] 2017. FLAME (Facility for Large-scale Adaptive Media Experimentation). <https://www.ict-flame.eu/>
- [15] 3GPP. 2020. Overall Description of LTE-Based 5G Broadcast; Version 16.0.0; Technical Report (TR) 36.976. *3rd Generation Partnership Project (3GPP)* (2020).
- [16] 5GINFIRE. [n.d.]. *D2.2-5GinFIRE Experimental Infrastructure Architecture and 5G Automotive Use Case(update)*. https://bscw.5g-ppp.eu/pub/bscw.cgi/d302168/D2-2-5GINFIRE_Experimental_Infrastructure_Architecture_and_5G_Automotive_Use_Case-v1-0.pdf
- [17] Hatem Abou-Zeid and Hossam S Hassanein. 2013. Predictive green wireless access: Exploiting mobility and application information. *IEEE wireless communications* 20, 5 (2013), 92–99.

- 1405 [18] Vijay Kumar Adhikari, Yang Guo, Fang Hao, Volker Hilt, and Zhi-Li Zhang. 2012. A tale of three CDNs: An active measurement study of Hulu and
1406 its CDNs. In *2012 Proceedings IEEE INFOCOM Workshops*. IEEE, 7–12.
- 1407 [19] Vijay K Adhikari, Yang Guo, Fang Hao, Volker Hilt, Zhi-Li Zhang, Matteo Varvello, and Moritz Steiner. 2014. Measurement study of Netflix, Hulu,
1408 and a tale of three CDNs. *IEEE/ACM Transactions on Networking* 23, 6 (2014), 1984–1997.
- 1409 [20] Vijay Kumar Adhikari, Yang Guo, Fang Hao, Matteo Varvello, Volker Hilt, Moritz Steiner, and Zhi-Li Zhang. 2012. Unreeling netflix: Understanding
1410 and improving multi-cdn movie delivery. In *2012 Proceedings IEEE INFOCOM*. IEEE, 1620–1628.
- 1411 [21] Ibrahim Afolabi, Tarik Taleb, Konstantinos Samdanis, Adlen Ksentini, and Hannu Flinck. 2018. Network slicing and softwarization: A survey on
1412 principles, enabling technologies, and solutions. *IEEE Communications Surveys & Tutorials* 20, 3 (2018), 2429–2453.
- 1413 [22] Ejaz Ahmed, Arif Ahmed, Ibrar Yaqoob, Junaid Shuja, Abdullah Gani, Muhammad Imran, and Muhammad Shoaib. 2017. Bringing computation
1414 closer toward the user network: Is edge computing the solution? *IEEE Communications Magazine* 55, 11 (2017), 138–144.
- 1415 [23] Muhammad Waseem Akhtar, Syed Ali Hassan, Rizwan Ghaffar, Haejoon Jung, Sahil Garg, and M Shamim Hossain. 2020. The shift to 6G
1416 communications: vision and requirements. *Human-centric Computing and Information Sciences* 10, 1 (2020), 1–27.
- 1417 [24] Yazan Al-Issa, Mohammad Ashraf Ottom, and Ahmed Tamrawi. 2019. eHealth cloud security challenges: a survey. *Journal of healthcare engineering*
1418 2019 (2019).
- 1419 [25] Pol Alemany, L Juan, Ana Pol, Anton Roman, Panagiotis Trakadas, Panagiotis Karkazis, Marios Touloupou, Evgenia Kapassa, Dimosthenis Kyriazis,
1420 Thomas Soenen, et al. 2019. Network slicing over a packet/optical network for vertical applications applied to multimedia real-time communications.
1421 In *2019 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*. IEEE, 1–2.
- 1422 [26] Osianoh Glenn Aliu, Ali Imran, Muhammad Ali Imran, and Barry Evans. 2012. A survey of self organisation in future cellular networks. *IEEE*
1423 *Communications Surveys & Tutorials* 15, 1 (2012), 336–361.
- 1424 [27] Zaheer Allam and David S Jones. 2021. Future (post-COVID) digital, smart and sustainable cities in the wake of 6G: Digital twins, immersive
1425 realities and new urban economies. *Land Use Policy* 101 (2021), 105201.
- 1426 [28] Mohammed Alreshoodi and John Woods. 2013. Survey on QoE/QoS correlation models for multimedia services. *arXiv preprint arXiv:1306.0221*
1427 (2013).
- 1428 [29] Federico Alvarez, David Breitgand, David Griffin, Pasquale Andriani, Stamatia Rizou, Nikolaos Zioulis, Francesca Moscatelli, Javier Serrano,
1429 Madeleine Keltch, Panagiotis Trakadas, T. Khoa Phan, Avi Weit, Ugur Acar, Oscar Prieto, Francesco Iadanza, Gino Carrozzo, Harilaos Koumaras,
1430 Dimitrios Zarpalas, and David Jimenez. 2019. An Edge-to-Cloud Virtualized Multimedia Service Platform for 5G Networks. *IEEE Transactions on*
1431 *Broadcasting* 65, 2 (June 2019), 369–380. <https://doi.org/10.1109/TBC.2019.2901400>
- 1432 [30] Sagar Arora, Pantelis A Frangoudis, and Adlen Ksentini. 2019. Exposing radio network information in a MEC-in-NFV environment: the RNISaaS
1433 concept. In *2019 IEEE Conference on Network Softwarization (NetSoft)*. IEEE, 306–310.
- 1434 [31] Haytham Assem, Lei Xu, Teodora Sandra Buda, and Declan O’Sullivan. 2016. Machine learning as a service for enabling Internet of Things and
1435 People. *Personal and Ubiquitous Computing* 20, 6 (Nov. 2016), 899–914. <https://doi.org/10.1007/s00779-016-0963-3>
- 1436 [32] Alcardo Alex Barakabitze, Arslan Ahmad, Rashid Mijumbi, and Andrew Hines. 2020. 5G network slicing using SDN and NFV: A survey of
1437 taxonomy, architectures and future challenges. *Computer Networks* 167 (2020), 106984.
- 1438 [33] Alcardo Alex Barakabitze, Nabajeet Barman, Arslan Ahmad, Saman Zadtootaghaj, Lingfen Sun, Maria G Martini, and Luigi Atzori. 2019. QoE
1439 management of multimedia streaming services in future networks: a tutorial and survey. *IEEE Communications Surveys & Tutorials* 22, 1 (2019),
1440 526–565.
- 1441 [34] Hamid R Barzegar, Nabil El Ioini, Claus Pahl, et al. 2020. Service continuity for ccam platform in 5g-carmen. In *2020 International Wireless*
1442 *Communications and Mobile Computing (IWCMC)*. IEEE, 1764–1769.
- 1443 [35] Pavlos Basaras, George Iosifidis, Stepan Kucera, and Holger Clausen. 2020. Multicast Optimization for Video Delivery in Multi-RAT Networks.
1444 *IEEE Transactions on Communications* 68, 8 (2020), 4973–4985.
- 1445 [36] Sreejata Basu. 2020. *Cloud Video Encoding vs On-Premise: Pros, Cons and Beyond*. <https://www.muvi.com/blogs/cloud-video-encoding-vs-on-premise.html>
- 1446 [37] Michael Till Beck, Sebastian Feld, Andreas Fichtner, Claudia Linnhoff-Popien, and Thomas Schimper. 2015. ME-VoLTE: Network functions for
1447 energy-efficient video transcoding at the mobile edge. In *2015 18th International Conference on Intelligence in Next Generation Networks*. IEEE,
1448 38–44.
- 1449 [38] Imen Grida Ben Yahia, Jaafar Bendriss, Alassane Samba, and Philippe Dooze. 2017. CogNitive 5G networks: Comprehensive operator use cases
1450 with machine learning for management operations. In *2017 20th Conference on Innovations in Clouds, Internet and Networks (ICIN)*. IEEE, Paris,
1451 252–259. <https://doi.org/10.1109/ICIN.2017.7899421>
- 1452 [39] Ilias Benkacem, Tarik Taleb, Miloud Bagaa, and Hannu Flinck. 2018. Performance benchmark of transcoding as a virtual network function in CDN
1453 as a service slicing. In *2018 IEEE Wireless Communications and Networking Conference (WCNC)*. 1–6. <https://doi.org/10.1109/WCNC.2018.8377402>
- 1454 [40] Divyashri Bhat, Amr Rizk, and Michael Zink. 2017. Not so QUIC: A performance study of DASH over QUIC. In *Proceedings of the 27th workshop on*
1455 *network and operating systems support for digital audio and video*. 13–18.
- 1456 [41] Jagadeesha R Bhat and Salman A Alqahtani. 2021. 6G Ecosystem: Current Status and Future Perspective. *IEEE Access* 9 (2021), 43134–43167.
- [42] Alexandros Biliris, Chuck Cranor, Fred Douglass, Michael Rabinovich, Sandeep Sibal, Oliver Spatscheck, and Walter Sturm. 2002. CDN brokering.
Computer Communications 25, 4 (2002), 393–402.

- [43] Saman Biokhaghazadeh, Ming Zhao, and Fengbo Ren. 2018. Are fpgas suitable for edge computing?. In *{USENIX} Workshop on Hot Topics in Edge Computing (HotEdge 18)*.
- [44] Pascal Birnstill, Daoyuan Ren, and Jürgen Beyerer. 2015. A user study on anonymization techniques for smart video surveillance. In *2015 12th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE, 1–6.
- [45] Leonardo Bonati, Michele Polese, Salvatore D’Oro, Stefano Basagni, and Tommaso Melodia. 2020. Open, Programmable, and Virtualized 5G Networks: State-of-the-Art and the Road Ahead. *Computer Networks* 182 (2020), 107516. <https://doi.org/10.1016/j.comnet.2020.107516>
- [46] Michel S Bonfim, Kelvin L Dias, and Stenio FL Fernandes. 2019. Integrated NFV/SDN architectures: A systematic literature review. *ACM Computing Surveys (CSUR)* 51, 6 (2019), 1–39.
- [47] Sem Borst, Aliye Özge Kaya, Doru Calin, and Harish Viswanathan. 2017. Dynamic path selection in 5G multi-RAT wireless networks. In *IEEE INFOCOM 2017-IEEE Conference on Computer Communications*. IEEE, 1–9.
- [48] Timm Böttger, Felix Cuadrado, Gareth Tyson, Ignacio Castro, and Steve Uhlig. 2018. Open connect everywhere: A glimpse at the internet ecosystem through the lens of the netflix cdn. *ACM SIGCOMM Computer Communication Review* 48, 1 (2018), 28–34.
- [49] Nassima Bouzakaria, Cyril Concolato, and Jean Le Feuvre. 2014. Overhead and performance of low latency live streaming using MPEG-DASH. In *IISA 2014, The 5th International Conference on Information, Intelligence, Systems and Applications*. IEEE, 92–97.
- [50] David Breitgand, Avi Weit, Stamatia Rizou, David Griffin, Ugur Acar, Gino Carrozzo, Nikolaos Zioulis, Pasquale Andriani, and Francesco Iadanza. 2018. Towards Serverless NFV for 5G Media Applications. In *Proceedings of the 11th ACM International Systems and Storage Conference (SYSTOR ’18)*. Association for Computing Machinery, Haifa, Israel, 118. <https://doi.org/10.1145/3211890.3211916>
- [51] Utku Bulkan, Muddesar Iqbal, and Tasos Dagiuklas. 2018. Load-Balancing for Edge QoE-Based VNF Placement for OTT Video Streaming. In *2018 IEEE Globecom Workshops (GC Wkshps)*. 1–6. <https://doi.org/10.1109/GLOCOMW.2018.8644214>
- [52] Claudia Campolo, Antonella Molinaro, Antonio Iera, and Francesco Menichella. 2017. 5G network slicing for vehicle-to-everything services. *IEEE Wireless Communications* 24, 6 (2017), 38–45.
- [53] Gino Carrozzo, Francesca Moscatelli, Gabriel Solsona, Oscar Prieto Gordo, Madeleine Keltsch, and Martin Schmalohr. 2018. Virtual CDNs over 5G networks: scenarios and requirements for ultra-high definition media distribution. In *2018 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*. IEEE, 1–5.
- [54] Alberto Huertas Celdrán, Manuel Gil Pérez, Félix J García Clemente, Fabrizio Ippoliti, and Gregorio Martínez Pérez. 2019. Dynamic network slicing management of multimedia scenarios for future remote healthcare. *Multimedia Tools and Applications* 78, 17 (2019), 24707–24737.
- [55] Alberto Huertas Celdran, Manuel Gil Perez, Felix J. Garcia Clemente, and Gregorio Martinez Perez. 2017. Enabling Highly Dynamic Mobile Scenarios with Software Defined Networking. *IEEE Communications Magazine* 55, 4 (April 2017), 108–113. <https://doi.org/10.1109/MCOM.2017.1600117CM>
- [56] Dan Chalmers and Morris Sloman. 1999. A survey of quality of service in mobile computing environments. *IEEE Communications surveys* 2, 2 (1999), 2–10.
- [57] Chia-Yu Chang, Navid Nikaein, Osama Arouk, Kostas Katsalis, Adlen Ksentini, Thierry Turetletti, and Konstantinos Samdanis. 2018. Slice Orchestration for Multi-Service Disaggregated Ultra-Dense RANs. *IEEE Communications Magazine* 56, 8 (Aug. 2018), 70–77. <https://doi.org/10.1109/MCOM.2018.1701044>
- [58] Dawei Chen, Yin-Chen Liu, BaekGyu Kim, Jiang Xie, Choong Seon Hong, and Zhu Han. 2020. Edge computing resources reservation in vehicular networks: A meta-learning approach. *IEEE Transactions on Vehicular Technology* 69, 5 (2020), 5634–5646.
- [59] Yu Chen, Yong Liu, Jingya Zhao, and Qinghua Zhu. 2020. Mobile edge cache strategy based on neural collaborative filtering. *IEEE Access* 8 (2020), 18475–18482.
- [60] Pete Chown. 2002. *Advanced encryption standard (AES) ciphersuites for transport layer security (TLS)*. Technical Report. RFC 3268, June.
- [61] Cisco. 2020. *Cisco Annual Internet Report (2018–2023) White Paper*. <https://www.cisco.com/c/en/us/solutions/collateral/executive-perspectives/annual-internet-report/white-paper-c11-741490.html>
- [62] Citrix. [n.d.]. *Intelligent Traffic Management*. <https://www.citrix.com/products/citrix-intelligent-traffic-management/>
- [63] Ian Colbert, Jake Daly, Ken Kreutz-Delgado, and Srinjoy Das. 2021. A Competitive Edge: Can FPGAs Beat GPUs at DCNN Inference Acceleration in Resource-Limited Edge Computing Applications? *arXiv preprint arXiv:2102.00294* (2021).
- [64] Carlos Colman-Meixner, Hamzeh Khalili, Konstantinos Antoniou, Muhammad Shuaib Siddiqui, Apostolos Papageorgiou, Antonino Albanese, Paolo Cruschelli, Gino Carrozzo, Luca Vignaroli, Alexandre Ulisses, Pedro Santos, Jordi Colom, Ioannis Neokosmidis, David Pujals, Rita Spada, Antonio Garcia, Sergi Figerola, Reza Nejabati, and Dimitra Simeonidou. 2019. Deploying a Novel 5G-Enabled Architecture on City Infrastructure for Ultra-High Definition and Immersive Media Production and Broadcasting. *IEEE Transactions on Broadcasting* 65, 2 (June 2019), 392–403. <https://doi.org/10.1109/TBC.2019.2901387>
- [65] Ioan-Sorin Comşa, Gabriel-Miro Muntean, and Ramona Trestian. 2020. An innovative machine-learning-based scheduling solution for improving live UHD video streaming quality in highly dynamic network environments. *IEEE Transactions on Broadcasting* 67, 1 (2020), 212–224.
- [66] Jin Cui, Lin Shen Liew, Giedre Sabaliauskaite, and Fengjun Zhou. 2019. A review on safety failures, security attacks, and available countermeasures for autonomous vehicles. *Ad Hoc Networks* 90 (2019), 101823.
- [67] Quentin De Coninck and Olivier Bonaventure. 2020. Multipath Extensions for QUIC (MP-QUIC). *draft-deconinck-quick-multipath-06* (2020).
- [68] Mouhamad Dieye, Shohreh Ahvar, Jagruti Sahoo, Ehsan Ahvar, Roch Glitho, Halima Elbiaze, and Noel Crespi. 2018. CPVNF: Cost-efficient proactive VNF placement and chaining for value-added services in content delivery networks. *IEEE Transactions on Network and Service Management* 15, 2 (2018), 774–786.

- [69] Trinh Viet Doan, Vaibhav Bajpai, and Sam Crawford. 2020. A Longitudinal View of Netflix: Content Delivery over IPv6 and Content Cache Deployments. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*. IEEE, 1073–1082.
- [70] Alexandros Doumanoglou, Nikolaos Zioulis, David Griffin, Javier Serrano, Truong Khoa Phan, David Jiménez, Dimitrios Zarpalas, Federico Alvarez, Miguel Rio, and Petros Daras. 2018. A system architecture for live immersive 3D-media transcoding over 5G networks. In *2018 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*. IEEE, 11–15.
- [71] Kerem Durak, Mehmet N Akcay, Yigit K Erinc, Boran Pekel, and Ali C Begen. 2020. Evaluating the Performance of Apple’s Low-Latency HLS. In *2020 IEEE 22nd International Workshop on Multimedia Signal Processing (MMSp)*. IEEE, 1–6.
- [72] Sunny Dutta, Tarik Taleb, Pantelis A Frangoudis, and Adlen Ksentini. 2016. On-the-fly QoE-aware transcoding in the mobile edge. In *2016 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 1–6.
- [73] EBU. [n.d.]. *Eurovision FLOW*. <https://tech.ebu.ch/docs/groups/flow/Eurovision%20Flow%20Brochure.pdf>
- [74] Ericsson. [n.d.]. *SON Optimization Manager*. <https://www.ericsson.com/en/portfolio/digital-services/automated-network-operations/network-management/son-optimization-manager>
- [75] ETSI. [n.d.]. *Open Source MANO (OSM)*. <https://osm.etsi.org/>
- [76] ETSI. 2014. *ETSI GS NFV-MAN 001: Network Functions Virtualisation (NFV); Management and Orchestration*. https://www.etsi.org/deliver/etsi_gs/nfv-man/001_099/001/01.01.01_60/gs_nfv-man001v010101p.pdf
- [77] ETSI. 2017. *ETSI GS MEC 012: Mobile Edge Computing (MEC); Radio Network Information API*. https://www.etsi.org/deliver/etsi_gs/MEC/001_099/012/01.01.01_60/gs_MEC012v010101p.pdf
- [78] ETSI. 2018. *ETSI GS MEC 002: Multi-access Edge Computing (MEC): Phase 2: Use Cases and Requirements*. https://www.etsi.org/deliver/etsi_gs/MEC/001_099/002/02.01.01_60/gs_MEC002v020101p.pdf
- [79] ETSI. 2020. *ETSI TS 128 313: 5G; Self-Organizing Networks (SON) for 5G networks*. https://www.etsi.org/deliver/etsi_ts/128300_128399/128313/16.00.00_60/ts_128313v160000p.pdf
- [80] Kristian Evensen, Tomas Kupka, Haakon Riiser, Pengpeng Ni, Ragnhild Eg, Carsten Griwodz, and Pål Halvorsen. 2014. Adaptive media streaming to mobile devices: challenges, enhancements, and recommendations. *Advances in Multimedia 2014* (2014).
- [81] Zubair Md Fadlullah, Fengxiao Tang, Bomim Mao, Nei Kato, Osamu Akashi, Takeru Inoue, and Kimihiro Mizutani. 2017. State-of-the-art deep learning: Evolving machine intelligence toward tomorrow’s intelligent network traffic control systems. *IEEE Communications Surveys & Tutorials* 19, 4 (2017), 2432–2455.
- [82] Ching-Ling Fan, Wen-Chih Lo, Yu-Tung Pai, and Cheng-Hsin Hsu. 2019. A survey on 360 video streaming: Acquisition, transmission, and display. *ACM Computing Surveys (CSUR)* 52, 4 (2019), 1–36.
- [83] Thomas Favale, Francesca Soro, Martino Trevisan, Idilio Drago, and Marco Mellia. 2020. Campus traffic and e-Learning during COVID-19 pandemic. *Computer Networks* 176 (2020), 107290.
- [84] Xincan Fei, Fangming Liu, Hong Xu, and Hai Jin. 2018. Adaptive VNF scaling and flow routing with proactive demand prediction. In *IEEE INFOCOM 2018-IEEE Conference on Computer Communications*. IEEE, 486–494.
- [85] Xincan Fei, Fangming Liu, Qixia Zhang, Hai Jin, and Hongxin Hu. 2020. Paving the way for NFV acceleration: A taxonomy, survey and future directions. *ACM Computing Surveys (CSUR)* 53, 4 (2020), 1–42.
- [86] Anja Feldmann, Oliver Gasser, Franziska Lichtblau, Enric Pujol, Ingmar Poesse, Christoph Dietzel, Daniel Wagner, Matthias Wichtlhuber, Juan Tapiador, Narseo Vallina-Rodriguez, Oliver Hohlfeld, and Georgios Smaragdakis. 2020. The Lockdown Effect: Implications of the COVID-19 Pandemic on Internet Traffic. *Proceedings of the ACM Internet Measurement Conference (2020)*, 1–18. <https://doi.org/10.1145/3419394.3423658>
- [87] Alan Ford, Costin Raiciu, Mark Handley, Sebastien Barre, Janardhan Iyengar, et al. 2011. Architectural guidelines for multipath TCP development. *IETF, Informational RFC 6182* (2011), 2070–1721.
- [88] Xenofon Foukas, Georgios Patounas, Ahmed Elmokashfi, and Mahesh K Marina. 2017. Network slicing in 5G: Survey and challenges. *IEEE Communications Magazine* 55, 5 (2017), 94–100.
- [89] Linux Foundation. [n.d.]. *Open Network Automation Platform (ONAP)*. <https://www.onap.org/>
- [90] Benjamin Frank, Ingmar Poesse, Yin Lin, Georgios Smaragdakis, Anja Feldmann, Bruce Maggs, Jannis Rake, Steve Uhlig, and Rick Weber. 2013. Pushing CDN-ISP collaboration to the limit. *ACM SIGCOMM Computer Communication Review* 43, 3 (2013), 34–44.
- [91] Alvaro Gabilondo, Javier Morgade, Roberto Viola, Pablo Angueira, and Jon Montalbán. 2020. Realising a vRAN based FeMBMS Management and Orchestration Framework. In *2020 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*. IEEE, 1–7.
- [92] Olga Galinina, Alexander Pyattaev, Sergey Andreev, Mischa Dohler, and Yevgeni Koucheryavy. 2015. 5G multi-RAT LTE-WiFi ultra-dense small cells: Performance dynamics, architecture, and trends. *IEEE Journal on Selected Areas in Communications* 33, 6 (2015), 1224–1240.
- [93] Liljana Gavrilovska, Valentin Rakovic, and Daniel Denkovski. 2020. From Cloud RAN to Open RAN. *Wireless Personal Communications* (2020), 1–17.
- [94] Chang Ge, Ning Wang, Severin Skillman, Gerry Foster, and Yue Cao. 2016. QoE-driven DASH video caching and adaptation at 5G mobile edge. In *Proceedings of the 3rd ACM Conference on Information-Centric Networking*. 237–242.
- [95] Francesco Giannone, Pantelis A Frangoudis, Adlen Ksentini, and Luca Valcarenghi. 2020. Orchestrating heterogeneous MEC-based applications for connected vehicles. *Computer Networks* 180 (2020), 107402.
- [96] Jordi Joan Gimenez, Jose Luis Carcel, Manuel Fuentes, Eduardo Garro, Simon Elliott, David Vargas, Christian Menzel, and David Gomez-Barquero. 2019. 5G new radio for terrestrial broadcast: A forward-looking approach for NR-MBMS. *IEEE Transactions on Broadcasting* 65, 2 (2019), 356–368.

- [97] Kostas Giotis, Yiannos Kryftis, and Vasilis Maglaris. 2015. Policy-based orchestration of NFV services in software-defined networks. In *Proceedings of the 2015 1st IEEE conference on network softwarization (netsoft)*. IEEE, 1–5.
- [98] Utkarsh Goel, Mike P Wittie, and Moritz Steiner. 2015. Faster web through client-assisted CDN server selection. In *2015 24th International conference on computer communication and networks (ICCCN)*. IEEE, 1–10.
- [99] Ismael Gomez-Miguel, Andres Garcia-Saavedra, Paul D Sutton, Pablo Serrano, Cristina Cano, and Doug J Leith. 2016. srsLTE: An open-source platform for LTE evolution and experimentation. In *Proceedings of the Tenth ACM International Workshop on Wireless Network Testbeds, Experimental Evaluation, and Characterization*. 25–32.
- [100] Andres J. Gonzalez, Gianfranco Nencioni, Andrzej Kamisiński, Bjarne E. Helvik, and Poul E. Heegaard. 2018. Dependability of the NFV Orchestrator: State of the Art and Research Challenges. *IEEE Communications Surveys Tutorials* 20, 4 (2018), 3307–3329. <https://doi.org/10.1109/COMST.2018.2830648>
- [101] Yashuang Guo, F. Richard Yu, Jianping An, Kai Yang, Chuqiao Yu, and Victor C. M. Leung. 2020. Adaptive Bitrate Streaming in Wireless Networks With Transcoding at Network Edge Using Deep Reinforcement Learning. *IEEE Transactions on Vehicular Technology* 69, 4 (2020), 3879–3892. <https://doi.org/10.1109/TVT.2020.2968498>
- [102] Gürkan Gür, Pawani Porabage, and Madhusanka Liyanage. 2020. Convergence of ICN and MEC for 5G: Opportunities and Challenges. *IEEE Communications Standards Magazine* 4, 4 (2020), 64–71.
- [103] Kay Haensge, Dirk Trossen, Sebastian Robitzsch, Michael Boniface, and Stephen Phillips. 2019. Cloud-Native 5G Service Delivery Platform. In *2019 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*. 1–7. <https://doi.org/10.1109/NFV-SDN47374.2019.9040042>
- [104] Haivision. [n.d.]. *Haivision Lightflow MultiCDN*. <https://www.haivision.com/products/haivision-lightflow-multicdn/>
- [105] Bo Han, Vijay Gopalakrishnan, Lusheng Ji, and Seungjoon Lee. 2015. Network function virtualization: Challenges and opportunities for innovations. *IEEE Communications Magazine* 53, 2 (2015), 90–97.
- [106] HCL. [n.d.]. *SON*. <https://www.hcltech.com/ERX/Telecom-and-5G/SON>
- [107] Enrique Hernandez-Valencia, Steven Izzo, and Beth Polonsky. 2015. How will NFV/SDN transform service provider opex? *IEEE Network* 29, 3 (2015), 60–67.
- [108] Juliver Gil Herrera and Juan Felipe Botero. 2016. Resource allocation in NFV: A comprehensive survey. *IEEE Transactions on Network and Service Management* 13, 3 (2016), 518–532.
- [109] Christer Holmberg, Stefan Hakansson, and G Eriksson. 2015. Web real-time communication use cases and requirements. *Request for Comments (RFC) 7478* (2015).
- [110] Honglin Hu, Jian Zhang, Xiaoying Zheng, Yang Yang, and Ping Wu. 2010. Self-configuration and self-optimization for LTE networks. *IEEE Communications Magazine* 48, 2 (2010), 94–100.
- [111] Jie Hu, Qing Wang, and Kun Yang. 2020. Energy self-sustainability in full-spectrum 6G. *IEEE Wireless Communications* 28, 1 (2020), 104–111.
- [112] Huawei Huang and Song Guo. 2019. Proactive failure recovery for NFV in distributed edge computing. *IEEE Communications Magazine* 57, 5 (2019), 131–137.
- [113] Tongyi Huang, Wu Yang, Jun Wu, Jin Ma, Xiaofei Zhang, and Daoyin Zhang. 2019. A survey on green 6G network: Architecture and technologies. *IEEE Access* 7 (2019), 175758–175768.
- [114] K Hughes and D Singer. 2017. Information technology–Multimedia application format (MPEG-A)–Part 19: Common media application format (CMAF) for segmented media. *ISO/IEC 19* (2017), 23000.
- [115] Chih-Lin I, Slawomir Kuklinski, and Tao Chen. 2020. A Perspective of O-RAN Integration with MEC, SON, and Network Slicing in the 5G Era. *IEEE Network* 34, 6 (2020), 3–4. <https://doi.org/10.1109/MNET.2020.9277891>
- [116] Ali Imran, Ahmed Zoha, and Adnan Abu-Dayya. 2014. Challenges in 5G: how to empower SON with big data for enabling 5G. *IEEE network* 28, 6 (2014), 27–33.
- [117] Narjes Tahghigh Jahromi, Somayeh Kianpisheh, and Roch H Glitho. 2018. Online VNF placement and chaining for value-added services in content delivery networks. In *2018 IEEE International Symposium on Local and Metropolitan Area Networks (LANMAN)*. IEEE, 19–24.
- [118] Qingmin Jia, Renchao Xie, Haijun Lu, Wenbin Zheng, and Hong Luo. 2019. Joint Optimization Scheme for Caching, Transcoding and Bandwidth in 5G Networks with Mobile Edge Computing. In *2019 IEEE 5th International Conference on Computer and Communications (ICCC)*. IEEE, 999–1004.
- [119] Xiantao Jiang, F Richard Yu, Tian Song, and Victor CM Leung. 2021. A Survey on Multi-Access Edge Computing Applied to Video Streaming: Some Research Issues and Challenges. *IEEE Communications Surveys & Tutorials* (2021).
- [120] Jingwen Jin and Klara Nahrstedt. 2004. QoS specification languages for distributed multimedia applications: A survey and taxonomy. *IEEE multimedia* 11, 3 (2004), 74–87.
- [121] Yichao Jin, Yonggang Wen, and Cedric Westphal. 2015. Optimal Transcoding and Caching for Adaptive Streaming in Media Cloud: an Analytical Approach. *IEEE Transactions on Circuits and Systems for Video Technology* 25, 12 (2015), 1914–1925. <https://doi.org/10.1109/TCSVT.2015.2402892>
- [122] Parikshit Juluri, Venkatesh Tamarapalli, and Deep Medhi. 2015. Measurement of quality of experience of video-on-demand services: A survey. *IEEE Communications Surveys & Tutorials* 18, 1 (2015), 401–418.
- [123] Madeleine Keltch, Sebastian Prokesch, Oscar Prieto Gordo, Javier Serrano, Truong Khoa Phan, and Igor Fritsch. 2018. Remote production and mobile contribution over 5G networks: scenarios, requirements and approaches for broadcast quality media streaming. In *2018 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*. IEEE, 1–7.

- [124] Javad I Khan, Seung Su Yang, Qiong Gu, Darsan Patel, Patrick Mail, Oleg Komogortsev, Wansik Oh, and Zhong Guo. 2001. Resource adaptive netcentric systems: A case study with sonet-A self-organizing network embedded transcoder. In *Proceedings of the ninth ACM international conference on Multimedia*. 617–620.
- [125] Hyun Jong Kim and Seong Gon Choi. 2010. A study on a QoS/QoE correlation model for QoE evaluation on IPTV service. In *2010 The 12th International Conference on Advanced Communication Technology (ICACT)*, Vol. 2. IEEE, 1377–1382.
- [126] Daniel L King, Paul H Delfabbro, Joel Billieux, and Marc N Potenza. 2020. Problematic online gaming and the COVID-19 pandemic. *Journal of Behavioral Addictions* 9, 2 (2020), 184–186.
- [127] Paulo Valente Klaine, Muhammad Ali Imran, Oluwakayode Onireti, and Richard Demo Souza. 2017. A survey of machine learning techniques applied to self-organizing cellular networks. *IEEE Communications Surveys & Tutorials* 19, 4 (2017), 2392–2431.
- [128] Zbigniew Kotulski, Tomasz Nowak, Mariusz Sepczuk, Marcin Tunia, Rafal Artych, Krzysztof Bocianiak, Tomasz Osko, and Jean-Philippe Wary. 2017. On end-to-end approach for slice isolation in 5G networks. Fundamental challenges. In *2017 Federated conference on computer science and information systems (FedCSIS)*. IEEE, 783–792.
- [129] Harilaos Koumaras, Christos Sakkas, Michail Alexandros Kourtis, Christos Xilouris, Vaios Koumaras, and Georgios Gardikis. 2016. Enabling agile video transcoding over SDN/NFV-enabled networks. In *2016 International Conference on Telecommunications and Multimedia (TEMU)*. 1–5. <https://doi.org/10.1109/TEMU.2016.7551916>
- [130] Fabian Kurtz, Caner Bektaş, Nils Dorsch, and Christian Wietfeld. 2018. Network slicing for critical communications in shared 5G infrastructures—an empirical evaluation. In *2018 4th IEEE Conference on Network Softwarization and Workshops (NetSoft)*. IEEE, 393–399.
- [131] Yeunwoong Kyung and Tae-Kook Kim. 2020. QoS-Aware Flexible Handover Management in Software-Defined Mobile Networks. *Applied Sciences* 10, 12 (2020), 4264.
- [132] David Lake, Gerry Foster, Serdar Vural, Yogarathnam Rahulan, Bong-Hwan Oh, Ning Wang, and Rahim Tafazolli. 2017. Virtualizing and orchestrating a 5G evolved packet core network. In *2017 IEEE Conference on Network Softwarization (NetSoft)*. 1–5. <https://doi.org/10.1109/NETSOFT.2017.8004215>
- [133] Adam Langley, Alistair Riddoch, Alyssa Wilk, Antonio Vicente, Charles Krasic, Dan Zhang, Fan Yang, Fedor Kouranov, Ian Swett, Janardhan Iyengar, et al. 2017. The quic transport protocol: Design and internet-scale deployment. In *Proceedings of the conference of the ACM special interest group on data communication*. 183–196.
- [134] Qian Li, Geng Wu, Apostolos Papathanassiou, and Udayan Mukherjee. 2016. An end-to-end network slicing framework for 5G wireless communication systems. *arXiv preprint arXiv:1608.00572* (2016).
- [135] Xi Li, Andres Garcia-Saavedra, Xavier Costa-Perez, Carlos J Bernardos, Carlos Guimarães, Kiril Antevski, Josep Mangués-Bafalluy, Jorge Baranda, Engin Zeydan, Daniel Corujo, et al. 2021. 5Growth: An End-to-End Service Platform for Automated Deployment and Management of Vertical Services over 5G Networks. *IEEE Communications Magazine* 59, 3 (2021), 84–90.
- [136] Yue Li, Pantelis A Frangoudis, Yassine Hadjadj-Aoul, and Philippe Bertin. 2016. A mobile edge computing-based architecture for improved adaptive HTTP video delivery. In *2016 IEEE Conference on Standards for Communications and Networking (CSCN)*. IEEE, 1–6.
- [137] Chengchao Liang, F Richard Yu, and Xi Zhang. 2015. Information-centric network function virtualization over 5G mobile wireless networks. *IEEE network* 29, 3 (2015), 68–74.
- [138] Chunyu Liu, Heli Zhang, Hong Ji, and Xi Li. 2021. MEC-assisted flexible transcoding strategy for adaptive bitrate video streaming in small cell networks. *China Communications* 18, 2 (2021), 200–214.
- [139] Shaoshan Liu, Liangkai Liu, Jie Tang, Bo Yu, Yifan Wang, and Weisong Shi. 2019. Edge computing for autonomous driving: Opportunities and challenges. *Proc. IEEE* 107, 8 (2019), 1697–1716.
- [140] Andra Lutu, Diego Perino, Marcelo Bagnulo, Enrique Frias-Martinez, and Javad Khangosstar. 2020. A Characterization of the COVID-19 Pandemic Impact on a Mobile Network Operator Traffic. In *Proceedings of the ACM Internet Measurement Conference*. 19–33.
- [141] Huisheng Ma, Shufang Li, Erqing Zhang, Zhengnan Lv, Jing Hu, and Xinlei Wei. 2020. Cooperative Autonomous Driving Oriented MEC-Aided 5G-V2X: Prototype System Design, Field Tests and AI-Based Optimization Tools. *IEEE Access* 8 (2020), 54288–54302.
- [142] Rohan Mahy, Philip Matthews, and Jonathan Rosenberg. 2010. *Traversal using relays around nat (turn): Relay extensions to session traversal utilities for nat (stun)*. Technical Report. RFC 5766 (Proposed Standard), Internet Engineering Task Force.
- [143] Angel Martin, Jon Egaña, Julián Flórez, Jon Montalbán, Igor G Olaizola, Marco Quartulli, Roberto Viola, and Mikel Zorrilla. 2018. Network resource allocation system for QoE-aware delivery of media services in 5G networks. *IEEE Transactions on Broadcasting* 64, 2 (2018), 561–574.
- [144] Angel Martin, Roberto Viola, Mikel Zorrilla, Julián Flórez, Pablo Angueira, and Jon Montalbán. 2019. MEC for Fair, Reliable and Efficient Media Streaming in Mobile Networks. *IEEE Transactions on Broadcasting* 66, 2 (2019), 264–278.
- [145] Jorge Martín-Pérez, Luca Cominardi, Carlos J Bernardos, Antonio de la Oliva, and Arturo Azcorra. 2019. Modeling mobile edge computing deployments for low latency multimedia services. *IEEE Transactions on Broadcasting* 65, 2 (2019), 464–474.
- [146] Abbas Mehrabi, Matti Siekkinen, and Antti Ylä-Jääski. 2019. Energy-aware QoE and backhaul traffic optimization in green edge adaptive mobile video streaming. *IEEE Transactions on Green Communications and Networking* 3, 3 (2019), 828–839.
- [147] Jie Mei, Xianbin Wang, and Kan Zheng. 2019. Intelligent network slicing for V2X services toward 5G. *IEEE Network* 33, 6 (2019), 196–204.
- [148] Rashid Mijumbi, Sidhant Hasija, Steven Davy, Alan Davy, Brendan Jennings, and Raouf Boutaba. 2016. A connectionist approach to dynamic resource management for virtualised network functions. In *2016 12th International Conference on Network and Service Management (CNSM)*. IEEE, 1–9.

- [149] Jean-Baptiste Monteil, Jernej Hribar, Pieter Barnard, Yong Li, and Luiz A DaSilva. 2020. Resource Reservation within Sliced 5G Networks: A Cost-Reduction Strategy for Service Providers. In *2020 IEEE International Conference on Communications Workshops (ICC Workshops)*. IEEE, 1–6.
- [150] FJ Moreno-Muro, M Garrich, C San-Nicolás-Martínez, M Hernández-Bastida, P Pavón-Mariño, A Bravalheri, AS Muqaddas, N Uniyal, R Nejabati, D Simeonidou, et al. 2019. Joint VNF and multi-layer resource allocation with an open-source optimization-as-a-service integration. In *45th European Conference on Optical Communication (ECOC 2019)*. IET, 1–4.
- [151] Jessica Moysen and Lorenza Giupponi. 2018. From 4G to 5G: Self-organized network management meets machine learning. *Computer Communications* 129 (2018), 248–268.
- [152] Saman Naderiparizi, Mehrdad Hessar, Vamsi Talla, Shyamath Gollakota, and Joshua R Smith. 2018. Towards battery-free {HD} video streaming. In *15th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 18)*. 233–247.
- [153] Akihiro Nakao, Ping Du, Yoshiaki Kiriha, Fabrizio Granelli, Anteneh Atumo Gebremariam, Tarik Taleb, and Miloud Bagaa. 2017. End-to-end network slicing for 5G mobile networks. *Journal of Information Processing* 25 (2017), 153–163.
- [154] Sabuzima Nayak and Ripon Patgiri. 2022. 6G Communication: A Vision on the Potential Applications. In *Edge Analytics*. Springer, 203–218.
- [155] Anselme Ndikumana, Saeed Ullah, Tuan LeAnh, Nguyen H Tran, and Choong Seon Hong. 2017. Collaborative cache allocation and computation offloading in mobile edge computing. In *2017 19th Asia-Pacific Network Operations and Management Symposium (APNOMS)*. IEEE, 366–369.
- [156] Van-Giang Nguyen, Anna Brunstrom, Karl-Johan Grinnemo, and Javid Taheri. 2017. SDN/NFV-based mobile packet core network architectures: A survey. *IEEE Communications Surveys & Tutorials* 19, 3 (2017), 1567–1602.
- [157] James Nightingale, Qi Wang, Jose M Alcaraz Calero, Enrique Chirivella-Perez, Marian Ulbricht, Jesus A Alonso-Lopez, Ricardo Preto, Tiago Batista, Tiago Teixeira, Maria Joao Barros, et al. 2016. QoE-Driven, Energy-Aware Video Adaptation in 5G Networks: The SELFNET Self-Optimisation Use Case. *IJDSN* 12, 1 (2016), 7829305–1.
- [158] Navid Nikaein, Chia-Yu Chang, and Konstantinos Alexandris. 2018. Mosaic5G: agile and flexible service platforms for 5G research. *ACM SIGCOMM Computer Communication Review* 48, 3 (Sept. 2018), 29–34. <https://doi.org/10.1145/3276799.3276803>
- [159] Navid Nikaein, Mahesh K. Marina, Saravana Manickam, Alex Dawson, Raymond Knopp, and Christian Bonnet. 2014. OpenAirInterface: A Flexible Platform for 5G Research. *ACM SIGCOMM Computer Communication Review* 44, 5 (Oct. 2014), 33–38. <https://doi.org/10.1145/2677046.2677053>
- [160] Nokia. [n.d.]. *EdenNet*. <https://www.nokia.com/networks/portfolio/self-organizing-networks>
- [161] Andres F Ocampo, Thomas Dreiholz, Mah-Rukh Fida, Ahmed Elmokashfi, and Haakon Bryhni. 2020. Integrating Cloud-RAN with Packet Core as VNF Using Open Source MANO and OpenAirInterface. In *Proceedings of the 45th IEEE Conference on Local Computer Networks (LCN), Sydney, New South Wales/Australia (November 2020)*.
- [162] Lyndon Ong, John Yoakum, et al. 2002. *An introduction to the stream control transmission protocol (SCTP)*. Technical Report. RFC 3286 (Informational), May.
- [163] Jose Ordonez-Lucena, Pablo Ameigeiras, Diego Lopez, Juan J Ramos-Munoz, Javier Lorca, and Jesus Folgueira. 2017. Network slicing for 5G with SDN/NFV: Concepts, architectures, and challenges. *IEEE Communications Magazine* 55, 5 (2017), 80–87.
- [164] John S Otto, Mario A Sánchez, John P Rula, Ted Stein, and Fabián E Bustamante. 2012. namehelp: Intelligent client-side DNS resolution. In *Proceedings of the ACM SIGCOMM 2012 conference on Applications, technologies, architectures, and protocols for computer communication*. 287–288.
- [165] Ye Ouyang, Zhongyuan Li, Le Su, Wenyuan Lu, and Zhenyi Lin. 2018. Application behaviors driven self-organizing network (SON) for 4G LTE networks. *IEEE Transactions on Network Science and Engineering* 7, 1 (2018), 3–14.
- [166] Roger Pantos and William May. 2017. HTTP Live Streaming. *rfc 8216, August (2017)*.
- [167] Alain Pellen. 2020. *Cost Comparison: On-Premises Vs Cloud Computing*. <https://www.harmonicinc.com/insights/blog/on-prem-vs-cloud>
- [168] Manuel Peuster, Stefan Schneider, Mengxuan Zhao, George Xilouris, Panagiotis Trakadas, Felipe Vicens, Wouter Tavernier, Thomas Soenen, Ricard Vilalta, George Andreou, Dimosthenis Kyriazis, and Holger Karl. 2019. Introducing Automated Verification and Validation for Virtualized Network Functions and Services. *IEEE Communications Magazine* 57, 5 (May 2019), 96–102. <https://doi.org/10.1109/MCOM.2019.1800873>
- [169] Lucian Popa, Ali Ghodsi, and Ion Stoica. 2010. HTTP as the Narrow Waist of the Future Internet. In *Proceedings of the 9th ACM SIGCOMM Workshop on Hot Topics in Networks*. 1–6.
- [170] Jon Postel et al. 1980. User datagram protocol. *STD 6, RFC 768, August (1980)*.
- [171] Jon Postel et al. 1981. Transmission control protocol. *STD 7, RFC 793, September (1981)*.
- [172] Dan Rayburn. 2020. *CDN/Media Pricing See’s Big Drop for Largest Customers: Pricing Down to \$0.0006*. <https://www.streamingmediablog.com/2020/05/q1-cdn-pricing.html>
- [173] Daniela Renga and Michela Meo. 2018. From self-sustainable green mobile networks to enhanced interaction with the smart grid. In *2018 30th International Teletraffic Congress (ITC 30)*, Vol. 1. IEEE, 129–134.
- [174] Sepehr Rezvani, Saeedeh Parsaeefard, Nader Mokari, Mohammad R Javan, and Halim Yanikomeroglu. 2019. Cooperative multi-bitrate video caching and transcoding in multicarrier NOMA-assisted heterogeneous virtualized MEC networks. *IEEE Access* 7 (2019), 93511–93536.
- [175] Jordi Ferrer Riera, Eduard Escalona, Josep Batalle, Eduard Grasa, and Joan A Garcia-Espin. 2014. Virtual network function scheduling: Concept and challenges. In *2014 international conference on smart communications in network technologies (SaCoNeT)*. IEEE, 1–5.
- [176] Rodrigo Roman, Javier Lopez, and Masahiro Mambo. 2018. Mobile edge computing. fog et al.: A survey and analysis of security threats and challenges. *Future Generation Computer Systems* 78 (2018), 680–698.
- [177] Peter Rost, Christian Mannweiler, Diomidis S Michalopoulos, Cinzia Sartori, Vincenzo Sciancalepore, Nishanth Sastry, Oliver Holland, Shreya Tayade, Bin Han, Dario Bega, et al. 2017. Network slicing to enable scalability and flexibility in 5G mobile networks. *IEEE Communications*

- 1717 *magazine* 55, 5 (2017), 72–79.
- 1718 [178] Ahmad Rostami. 2019. Private 5G networks for vertical industries: Deployment and operation models. In *2019 IEEE 2nd 5G World Forum (5GWF)*.
1719 IEEE, 433–439.
- 1720 [179] Dario Sabella, Vadim Sukhomlinov, Linh Trang, Sami Kekki, Pietro Paglierani, Ralf Rossbach, Xinhui Li, Yonggang Fang, Dan Druta, Fabio Giust,
1721 et al. 2019. Developing software for multi-access edge computing. *ETSI white paper* 20 (2019), 1–38.
- 1722 [180] Ali Saffari, Mehrdad Hesar, Saman Naderiparizi, and Joshua R Smith. 2019. Battery-free wireless video streaming camera system. In *2019 IEEE*
1723 *International Conference on RFID (RFID)*. IEEE, 1–8.
- 1724 [181] Pablo Salva-Garcia, Jose M. Alcaraz-Calero, Qi Wang, Miguel Arevalillo-Herraez, and Jorge Bernal Bernabe. 2020. Scalable Virtual Network
1725 Video-Optimizer for Adaptive Real-Time Video Transmission in 5G Networks. *IEEE Transactions on Network and Service Management* 17, 2 (June
2020), 1068–1081. <https://doi.org/10.1109/TNSM.2020.2978975>
- 1726 [182] P Sandhir and K Mitchell. 2008. A neural network demand prediction scheme for resource allocation in cellular wireless systems. In *2008 IEEE*
1727 *Region 5 Conference*. IEEE, 1–6.
- 1728 [183] Frederico Schardong, Ingrid Nunes, and Alberto Schaeffer-Filho. 2020. NFV resource allocation: A systematic review and taxonomy of VNF
1729 forwarding graph embedding. *Computer Networks* (2020), 107726.
- 1730 [184] Henning Schulzrinne, Stephen Casner, Ron Frederick, Van Jacobson, et al. 1996. RTP: A transport protocol for real-time applications. *rfc* 1889,
1731 *January* (1996).
- 1732 [185] Henning Schulzrinne, Anup Rao, and Robert Lanphier. 1998. Real time streaming protocol (RTSP). *rfc* 2326, *April* (1998).
- 1733 [186] Michael Seufert, Sebastian Egger, Martin Slanina, Thomas Zinner, Tobias Hoßfeld, and Phuoc Tran-Gia. 2014. A survey on quality of experience of
1734 HTTP adaptive streaming. *IEEE Communications Surveys & Tutorials* 17, 1 (2014), 469–492.
- 1735 [187] Sonia Shahzadi, Muddesar Iqbal, Tasos Dagiuklas, and Zia Ul Qayyum. 2017. Multi-access edge computing: open issues, challenges and future
1736 perspectives. *Journal of Cloud Computing* 6, 1 (2017), 1–13.
- 1737 [188] Maria Sharabayko, Maxim Sharabayko, Jean Dube, Joonwoong Kim, and Jeongseok Kim. 2021. The SRT Protocol. *draft-sharabayko-srt-00* (2021).
- 1738 [189] Chetna Singhal and BN Chandana. 2021. Aerial-SON: UAV-based Self-Organizing Network for Video Streaming in Dense Urban Scenario. In *2021*
1739 *International Conference on COMmunication Systems & NETWORKS (COMSNETS)*. IEEE, 7–12.
- 1740 [190] Lea Skorin-Kapov, Martín Varela, Tobias Hoßfeld, and Kuan-Ta Chen. 2018. A survey of emerging concepts and challenges for QoE management
1741 of multimedia services. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 14, 2s (2018), 1–29.
- 1742 [191] Iraj Sodagar. 2011. The mpeg-dash standard for multimedia streaming over the internet. *IEEE multimedia* 18, 4 (2011), 62–67.
- 1743 [192] Thomas Soenen, Wouter Tavernier, Manuel Peuster, Felipe Vicens, George Xilouris, Stavros Kolometsos, Michail-Alexandros Kourtis, and Didier
1744 Colle. 2019. Empowering network service developers: enhanced nfv devops and programmable mano. *IEEE Communications Magazine* (May 2019).
1745 <https://doi.org/10.1109/MCOM.2019.1800810>
- 1746 [193] Guan-Ming Su, Xiao Su, Yan Bai, Mea Wang, Athanasios V Vasilakos, and Haohong Wang. 2016. QoE in video streaming over wireless networks:
1747 perspectives and research challenges. *Wireless networks* 22, 5 (2016), 1571–1593.
- 1748 [194] Yang Sun, Tingting Wei, Huixin Li, Yanhua Zhang, and Wenjun Wu. 2020. Energy-efficient multimedia task assignment and computing offloading
1749 for mobile edge computing networks. *IEEE Access* 8 (2020), 36702–36713.
- 1750 [195] SVA. [n.d.]. *Open Caching*. <https://www.streamingvideoalliance.org/working-group/open-caching/>
- 1751 [196] Yiming Tan, Ce Han, Ming Luo, Xiang Zhou, and Xing Zhang. 2018. Radio network-aware edge caching for video delivery in MEC-enabled cellular
1752 networks. In *2018 IEEE Wireless Communications and Networking Conference Workshops (WCNCW)*. IEEE, 179–184.
- 1753 [197] David Alexander Tedjopurnomo, Zhifeng Bao, Baihua Zheng, Farhana Choudhury, and AK Qin. 2020. A survey on modern deep neural network
1754 for traffic prediction: Trends, methods and challenges. *IEEE Transactions on Knowledge and Data Engineering* (2020).
- 1755 [198] Michael Thornburgh. 2014. Adobe’s RTMFP Profile for Flash Communication. *Internet Engineering Task Force (IETF)* (2014).
- 1756 [199] Lechosław Tomaszewski, Sławomir Kukliński, and Robert Kołakowski. 2020. A new approach to 5G and MEC integration. In *IFIP International*
1757 *Conference on Artificial Intelligence Applications and Innovations*. Springer, 15–24.
- 1758 [200] Ruben Torres, Alessandro Finamore, Jin Ryong Kim, Marco Mellia, Maurizio M Munafo, and Sanjay Rao. 2011. Dissecting video server selection
1759 strategies in the youtube cdn. In *2011 31st International Conference on Distributed Computing Systems*. IEEE, 248–257.
- 1760 [201] Tuyen X. Tran, Parul Pandey, Abolfazl Hajisami, and Dario Pompili. 2016. Collaborative Multi-bitrate Video Caching and Processing in Mobile-Edge
1761 Computing Networks. [arXiv:1612.01436 \[cs.NI\]](https://arxiv.org/abs/1612.01436)
- 1762 [202] Tuyen X Tran and Dario Pompili. 2018. Adaptive bitrate video caching and processing in mobile-edge computing networks. *IEEE Transactions on*
1763 *Mobile Computing* 18, 9 (2018), 1965–1978.
- 1764 [203] Refik Fatih Ustok, Ugur Acar, Selcuk Keskin, David Breitgand, Avi Weit, Petros Drakoulis, Alexandros Doumanoglou, Nikolaos Zioulis, Dimitrios
1765 Zarpalas, Petros Daras, Francesco Iadanza, Francesca Moscatelli, and Giacomo Bernini. 2020. Service Development Kit for Media-Type Virtualized
1766 Network Services in 5G Networks. *IEEE Communications Magazine* 58, 7 (2020), 51–57. <https://doi.org/10.1109/MCOM.001.1900613>
- 1767 [204] Spyridon Vassilaras, Lazaros Gkatzikis, Nikolaos Liakopoulos, Ioannis N Stiakogiannakis, Meiyu Qi, Lei Shi, Liu Liu, Merouane Debbah, and
1768 Georgios S Paschos. 2017. The algorithmic aspects of network slicing. *IEEE Communications Magazine* 55, 8 (2017), 112–119.
- 1769 [205] Gorka Velez, Ángel Martín, Giancarlo Pastor, and Edward Mutafungwa. 2020. 5G Beyond 3GPP Release 15 for Connected Automated Mobility in
1770 Cross-Border Contexts. *Sensors* 20, 22 (2020), 6622.

- 1769 [206] Roberto Viola, Angel Martin, Javier Morgade, Stefano Masneri, Mikel Zorrilla, Pablo Angueira, and Jon Montalbán. 2020. Predictive CDN Selection
1770 for Video Delivery Based on LSTM Network Performance Forecasts and Cost-Effective Trade-Offs. *IEEE Transactions on Broadcasting* (2020).
- 1771 [207] Roberto Viola, Angel Martin, Mikel Zorrilla, and Jon Montalbán. 2018. MEC proxy for efficient cache and reliable multi-CDN video distribution. In
1772 *2018 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*. IEEE, 1–7.
- 1773 [208] W3C. 2020. *QUIC API for Peer-to-peer Connections*. <https://w3c.github.io/webrtc-quick/>
- 1774 [209] Minghao Wang, Tianqing Zhu, Tao Zhang, Jun Zhang, Shui Yu, and Wanlei Zhou. 2020. Security and privacy in 6G networks: New areas and new
1775 challenges. *Digital Communications and Networks* 6, 3 (2020), 281–291.
- 1776 [210] Qi Wang, Jose Alcaraz-Calero, Ruben Ricart-Sanchez, Maria Barros Weiss, Anastasius Gavras, Navid Nikaein, Xenofon Vasilakos, Bernini
1777 Giacomo, Giardina Pietro, Mark Roddy, Michael Healy, Paul Walsh, Thuy Truong, Zdravko Bozakov, Konstantinos Koutsopoulos, Pedro Neves,
1778 Cristian Patachia-Sultanoiu, Marius Iordache, Elena Oproiu, Imen Grida Ben Yahia, Ciriaco Angelo, Cosimo Zotti, Giuseppe Celozzi, Donal
1779 Morris, Ricardo Figueiredo, Dean Lorenz, Salvatore Spadaro, George Agapiou, Ana Aleixo, and Cipriano Lomba. 2019. Enable Advanced
1780 QoS-Aware Network Slicing in 5G Networks for Slice-Based Media Use Cases. *IEEE Transactions on Broadcasting* 65, 2 (June 2019), 444–453.
<https://doi.org/10.1109/TBC.2019.2901402>
- 1781 [211] Xinwei Wang, Jiandong Li, Lingxia Wang, Chungang Yang, and Zhu Han. 2019. Intelligent user-centric network selection: A model-driven
1782 reinforcement learning framework. *IEEE Access* 7 (2019), 21645–21661.
- 1783 [212] Yanting Wang, Yan Zhang, Min Sheng, and Kun Guo. 2019. On the Interaction of Video Caching and Retrieving in Multi-Server Mobile-Edge
1784 Computing Systems. *IEEE Wireless Communications Letters* 8, 5 (2019), 1444–1447. <https://doi.org/10.1109/LWC.2019.2921759>
- 1785 [213] Yiming Wei, Mugen Peng, and Yaqiong Liu. 2020. Intent-based networks for 6G: Insights and challenges. *Digital Communications and Networks* 6,
1786 3 (2020), 270–280.
- 1787 [214] Matthias Wichtlhuber, Robert Reinecke, and David Hausheer. 2015. An SDN-based CDN/ISP collaboration architecture for managing high-volume
1788 flows. *IEEE Transactions on Network and Service Management* 12, 1 (2015), 48–60.
- 1789 [215] Dan Wing, Philip Matthews, Rohan Mahy, and Jonathan Rosenberg. 2008. Session traversal utilities for NAT (STUN). *RFC5389, October* (2008).
- 1790 [216] Yin hao Xiao, Yizhen Jia, Chunchi Liu, Xiuzhen Cheng, Jiguo Yu, and Weifeng Lv. 2019. Edge computing security: State of the art and challenges.
1791 *Proc. IEEE* 107, 8 (2019), 1608–1631.
- 1792 [217] Yanghao Xie, Zhixiang Liu, Sheng Wang, and Yuxiu Wang. 2016. Service function chaining resource allocation: A survey. *arXiv preprint*
arXiv:1608.00095 (2016).
- 1793 [218] Chenren Xu, Lei Yang, and Pengyu Zhang. 2018. Practical backscatter communication systems for battery-free Internet of Things: A tutorial and
1794 survey of recent research. *IEEE Signal Processing Magazine* 35, 5 (2018), 16–27.
- 1795 [219] Mao Yang, Yong Li, Depeng Jin, Li Su, Shaowu Ma, and Lieguang Zeng. 2013. OpenRAN: A Software-Defined Ran Architecture via Virtualization.
1796 *SIGCOMM Comput. Commun. Rev.* 43, 4 (Aug. 2013), 549–550. <https://doi.org/10.1145/2534169.2491732>
- 1797 [220] Jingjing Yao, Tao Han, and Nirwan Ansari. 2019. On mobile edge caching. *IEEE Communications Surveys & Tutorials* 21, 3 (2019), 2525–2553.
- 1798 [221] Ahmet Yazar, Seda Doğan Tusha, and Huseyin Arslan. 2020. 6G VISION: AN ULTRA-FLEXIBLE PERSPECTIVE. *ITU Journal on Future and Evolving*
Technologies 1, 1 (2020).
- 1799 [222] Jongwon Yoon and Suman Banerjee. 2020. Hardware-Assisted, Low-Cost Video Transcoding Solution in Wireless Networks. *IEEE Transactions on*
1800 *Mobile Computing* 19, 3 (2020), 581–597. <https://doi.org/10.1109/TMC.2019.2898834>
- 1801 [223] Chuanji Zhang, Harshvardhan P Joshi, George F Riley, and Steven A Wright. 2019. Towards a virtual network function research agenda: A
1802 systematic literature review of vnf design considerations. *Journal of Network and Computer Applications* 146 (2019), 102417.
- 1803 [224] Shunliang Zhang. 2019. An overview of network slicing for 5G. *IEEE Wireless Communications* 26, 3 (2019), 111–117.
- 1804 [225] Xiaoxi Zhang, Chuan Wu, Zongpeng Li, and Francis CM Lau. 2017. Proactive VNF provisioning with multi-timescale cloud resources: Fusing
1805 online learning and online optimization. In *IEEE INFOCOM 2017-IEEE Conference on Computer Communications*. IEEE, 1–9.
- 1806 [226] Tiesong Zhao, Qian Liu, and Chang Wen Chen. 2016. QoE in video transmission: A user experience-driven strategy. *IEEE Communications Surveys*
& Tutorials 19, 1 (2016), 285–302.
- 1807 [227] Yiqing Zhou, Ling Liu, Lu Wang, Ning Hui, Xinyu Cui, Jie Wu, Yan Peng, Yanli Qi, and Chengwen Xing. 2020. Service aware 6G: an intelligent and
1808 open network based on convergence of communication, computing and caching. *Digital Communications and Networks* (2020).
- 1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820