# End-to-End Resource Allocation Management Model in Next-Generation Network: Survey

## Azhar Hamza Abdulkadhim[a], Ali Saeed Alfoudi[b], Firas Hussean Maghool[c]

[a]University of Al-Qadisiyah/ College of computer science and Information Technology,  it.mast.23.1@qu.edu.iq, https://orcid.org/0009-0008-4192-7538

[b]University of Al-Qadisiyah/ College of computer science and Information Technology, https://orcid.org/0000-0002-1234-7795

[c] University of Al-Qadisiyah/ College of computer science and Information Technology, https://orcid.org/0000-0003-0829-9425A

A B S T R A C T

Network communication has grown rapidly with massive demands of services. Moreover, resource allocation in networking is a fundamental and crucial issue that cannot ensure the network's stability and efficiency with the myriad requirements of different services. Various vertical businesses may seek varied network services, particularly in the Fifth-Generation networks and Beyond (5G+).  The pros of Fifth-Generation communication networks are to outperform 4G in performance by having higher bandwidth, minimum latency, more capacity, and QoS (Quality of Service).  Software-Defined Network (SDN) and Network Function Virtualization (NFV) are two technologies that are combined in the next generation cellular network to provide improved network management. The primary idea behind resource allocation (RA) in the next generation network is the concept of network slicing where the network resources are virtually partitioning into many separate networks. Each separated network must satisfy the unique needs of the required service to achieve the required QoS.  In this survey, we focus on resource management issues related to network slicing and tackling the biggest obstacles in this field while offering a thorough and up-to-date overview of this field. Thus, thorough analysis of the allocation of resources on the access side and core side of the network communication was sought. Also, demonstrates how revolutionary techniques that are used to support the management of sliced networks which are based on Machine Learning (ML) and Artificial Intelligence (AI). Importantly, use appropriate ML techniques such as deep learning for predicting the network condition and Reinforcement learning to learn optimal allocation policy without depending on prior knowledge and other techniques such as classification and clustering to aggregate the similar needs of users into separate slices. This could help to enhance resource utilization by allocating a sufficient amount of resource as needed based on ML algorithms and optimal utilization of resources and reducing operational costs by real-time adjustment of it based on user demands and network conditions.

MSC..

## 1. Introduction

It has been expected that the fifth-generation subscribers will approach three billion and fifty million in  2026. In this regard, the use rate of data will be exceeded about 35 GB/ month for each user. The outcome of this estimation may result in 400 use cases in 70 industries (Ericsson 2021). Older generations such as 2G, 3G, and 4G were constructed for communication of human-based. Nevertheless, this design has a limitation of flexibility that is required for the recent generation of 5G and beyond (Ksentini and Nikaein 2017). Therefore, the separation of the network into slices has been known by Next Generation Mobile Network (NGMN) (Ksentini and Nikaein 2017).

∗Corresponding author

Email addresses: it.mast.23.1@qu.edu.iq

Communicated by 'sub etitor'

In 5G and future networks, a sliced network has growing traction as a potential option for the effective utilization of resources to meet this urgent need (Liu, DIng, and Liu 2021). Building a unique physical network for every business scenario would unavoidably result in issues like difficult network maintenance and operation, excessive expense, and limited scalability.

Thus, network slicing technology arises at the opportune time to enable many corporate applications with disparate performance requirements on one network's physical infrastructure. To fulfill the diverse business needs of 5G, operators can employ network-slicing technology to split actual network infrastructures into many simulated networks,  in accordance to the distinct business needs of 5G networks. This allows for customization and differentiation among users.

### 1.1. Network Slicing:

The main purpose of separating the network's physical infrastructure into various dependent networks is to obtain specific needs and get a higher service quality (Ordonez-Lucena et al. 2017). This is carried out by separating the networks virtually according to top priority in which many tenants can share computers, storage resources, and network. A particular slice will have a collection comprising functions of networks (NF) and their corresponding amount of resource allocations. This E2E is independently conceptualized to provide End-to-End (E2E) services based on demands for each job (Rost et al. 2017). There are many sorts of slices such as enhanced Mobile Broadband (eMBB), Ultra-Reliable and Low Latency Communications (URLLC), and massive Machine-Type Communications (mMTC). All these types are depending on the 3GPP TS (Ji et al. 2018). This technique enhances the network resources to be sorted logically based on service requirements to provide better performance and separation of resources (Pang and Zhang 2020). The pathway of network slices depends on (SDN) and (NFV) that enabled network programmability, flexibility, with modularity needs (Ordonez-Lucena et al. 2017).

### 1.2. The Technologies that Enabling of Software-defined 5G Networks:

### 1.2.1. Software-Defined Network (SDN)

The traditional IP network includes three interconnected layers: control plane, data plane, and management layer (Thirupathi et al. 2019). However, the formulation of policies is the responsibility of the plane of management , which are incorporated via the control planes that permit the flow of traffic through the enforcement of the rules that are set by the plane of control. This platform has numerous issues via limiting innovation. Thus, difficult to solve misconfiguration makes it costly to add new network capabilities. To solve these issues, SDN isolates network control from the introduced layer using network programmability, as seen in (Figure 1). This is performed for improvement the flexibility of networks, controllability, and softwarization (Bakhshi 2017). The benefits of Separated control plane from the data plane hardware are to enable the controller to construct network forwarding devices lighter and more cost-effective than conventional routers and switches. The distributed network (SDN) controller's responsibilities include network cognition, collecting and evaluating information from a wider network viewpoint, and regulation decision-making for automated optimization of network and administration (Asakipaam, Kponyo, and Gyasi 2023). However, SDN is defined as an adaptable resource management paradigm that offers intelligent controller on the network. By separating the control plane from the data plane, traffic flows can be modified to meet the QoS of an Industrial Internet of Things (IIoT) requirements (Bektas et al. 2019; Haque and Abu-Ghazaleh 2016). SDN transmission of packets relies on flow-based processing rather than individual packet handling. which sets it apart from traditional IP networking. Various network hardware technologies can be accommodated by the flow abstraction, which is not dependent on them. The Network Operating System (NOS), sometimes referred to as the logically centralized controller, monitors and controls the network with a comprehensive global point of view. Network programmability, together with adaptable network administration, reconfiguration, and protocol evolution, are made possible by the ability to write network applications on top of the Network Operating System (NOS). In particular, if  there's exist a needs to be add an additional feature to the network, it could be integrated by installing a supplementary application onto the existing network operating system (NOS). This eliminates the need to add supplementary, expensive hardware,   need to install additional, costly hardware, which is necessary in conventional networks (Haque and Abu-Ghazaleh 2016) .
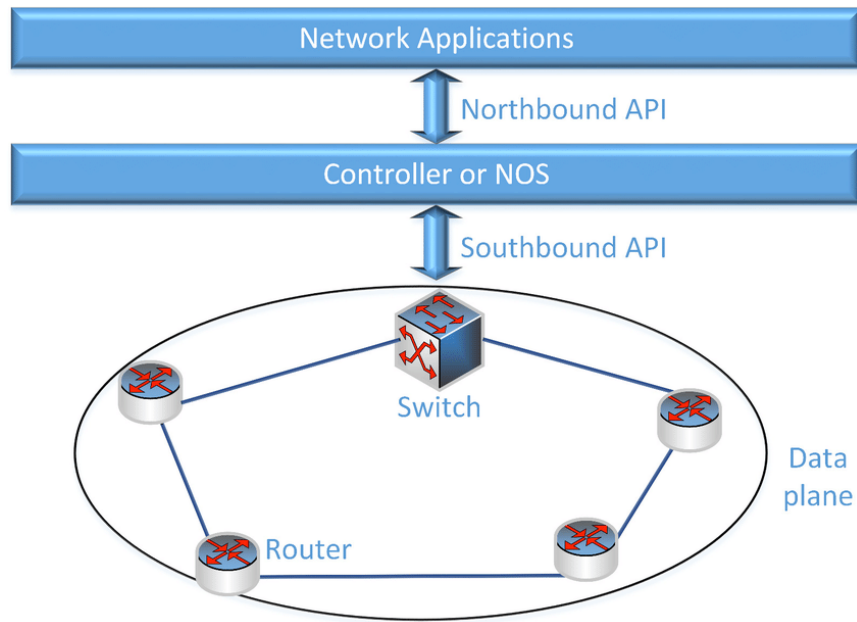
**Figure 1. the architecture of an SDN, Taken from (Haque and Abu-Ghazaleh 2016).**

The standards of an SDN establish the two Application Programming Interfaces (APIs) that connect the Control Plane, and Data Plane, with network applications. To be more precise, The southbound API facilitates the work of controller unit by transmitting the control rules to the data plane components, while the northbound API is established between network-based applications and the controller.

### 1.2.2. OpenFlow

The OpenFlow protocol demonstrates a standard interaction among the control and data planes. OpenFlow is a protocol that allows for the control of network resources by programming flow tables in supported switches, which are responsible for packet forwarding. the flow table is one of the main components of an Open Flow Switch, which is a collection of rules and related actions, an OpenFlow protocol that allows the controller to configure and update the information of the flow table, and a secure channel for communication. OpenFlow switches delegate decision-making authority for packet forwarding to the controller; they do not possess this capability (Oulahyane et al. 2024). These switches implement flow table-specified policies to perform data plane operations including packet forwarding and discarding. They are compatible with Ethernet switches that support OpenFlow and routers. A packet received by an OpenFlow-enabled switch is compared to the flow table that has been saved to carry out operations such as forwarding (directing it to a specified port, the controller, or the regular processing queue), dropping, or modifying. A control packet is delivered to the controller for the necessary actions if there is no match founded. Figure 2 depicts a network of commercial switches and access points with OpenFlow functionality. The OpenFlow Protocol permits a switch to be handled by two or more controllers for improved performance or resilience. In this example, all of the Flow Tables are managed by a single controller (Kreutz et al. 2008) .
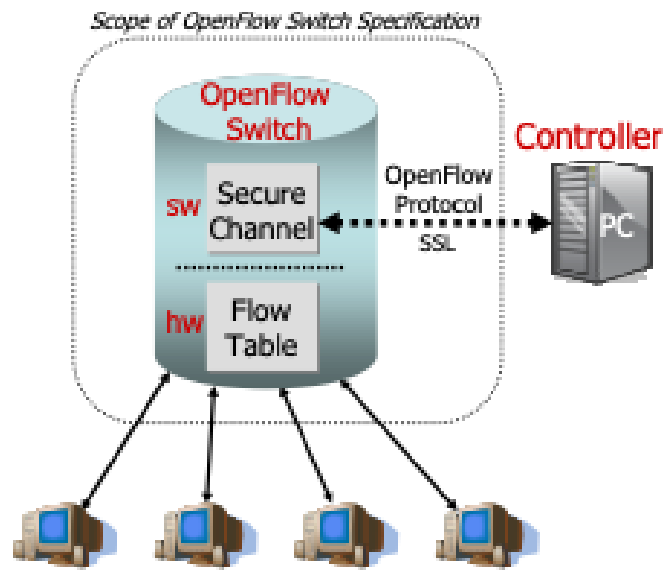
**Figure 2. An OpenFlow Switch. a remote controller managing the Table of Flow via the Secure Channel, taken from (Kreutz et al. 2008).**

### 1.2.3. Network Function Virtualization (NFV)

Traditionally, network services like firewalls, intrusion detection systems, and network optimizers have been carried out by service providers using specialized hardware middleboxes (Jin et al. 2013). However, With the extensive connectivity of 5G networks and their varied traffic demands, these specialized hardware middleboxes are costly and restrict deployment and administrative flexibility (Bega et al. 2020). By providing network services in a virtual environment as a service and enhancing resource consumption, application performance, and network resource usage, Network Function Virtualization (NFV) overcomes these drawbacks. Additionally, NFV makes network orchestration and administration more flexible (Van Rossem et al. 2016). This is achieved by employing VNFs, which operate on virtual machines (VMs) constructed on general-purpose hardware and managed by hypervisors, to implement each service in software (Nyanteh et al. 2021). Robust and customized network services can be produced by combining multiple VNFs, which can be hosted in big data centers as well as smaller locations near the edge of the network (ETSI 2017). Figure 3 demonstrates the idea of virtualizing network functions.
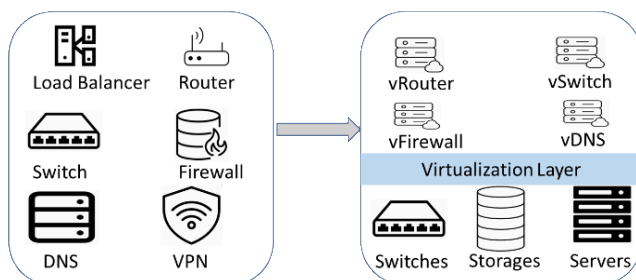


**Figure 3. Virtualization of Network Functions, taken from (Asakipaam, Kponyo, and Gyasi 2023)**

### 2. Literature Review

Reviewing relevant research about resource allocation in next-generation networks is the aim of this section

## 2.1.    Machine Learning-Based Resource Optimization Algorithms

Machine learning is becoming a fundamental element of our daily being alive since it is now widely embedded in many technologies and applications. Machine learning approaches can convert data into Suitable algorithm that salutes the platform's needs (Luu et al. 2020). It is a software that can be utilized for purposes of learning from data. To efficiently extract information, the type of algorithm and its corresponding job must have a thorough knowledge to match what we want to obtain from the data we have. Although there are varied machine learning methods with several categories, ML is basically categorized depending on the involvement of human monitoring in the learning process. These categories are supervised, unsupervised, and reinforcement learning. In addition, the emergence of deep learning (DL) provides the capability to simplify complicated optimization problems. The following briefly describes some machine learning categories. Additionally, we provide supplementary resources for those seeking a more deep understanding of each topic (Nurcahyani and Lee 2021).

### 2.1.1.    Reinforcement Learning (RL)

One kind of machine learning technique called reinforcement learning allows an agent to learn in an interactive setting by making mistakes and getting feedback from its experiences and actions. Reinforcement learning involves rewarding an agent based on its performance in a given scenario. The agent's objective in this kind of machine learning will be to maximize the long-term reward that it receives. While mapping between inputs and outputs is a common feature of both supervised and reinforcement learning, in contrast to supervised learning, rewards and penalties are utilized as signals to enhance the system's overall performance. The primary distinction between reinforcement learning (RL) and other machine learning techniques is that in RL, the agent makes decisions without referencing past information. In each state, the agent selects an action from a range of options, and it determines the quality of the action depending on input from the system. This makes it possible to handle complicated decision-making problems while supplying the minimal amount of information required to do so (Azimi et al. 2022; Cao et al. 2024) As proven in (Han, Feng, and Schotten 2018) , the Markov Decision Process (MDPs) can be recruited to provide an optimal policy for NS. This involves linking each unique "state" to an associated "action" of the system, which in turn produces a "reward, MDP doesn't need historical data and only depends on the current state. The reward function in NS resource handling issues is typically non-convex over a large space of policy. Because of this, professionals in this industry frequently decide to use Reinforcement Learning, which is renowned for its great efficacy and practical application in resolving Markovian choice issues. An innovative attempt was made to use RL to improve the network slicing mechanism (Bega et al. 2017), whereby the authors have shown that their Learning solution (QL) can considerably exceed randomized benchmark policies and effectively approach the ideal admission policy that optimizes the revenue of MNO's. The Q-learning approach can be implemented in an online learning way with a far more reasonable calculation cost than the value iteration method, which only slightly reduces revenue in exchange for achieving the optimum. Additionally, by carefully choosing the reward functions, RL algorithms can be created without a model, which greatly increases their resilience against inaccurate estimates of the slicing statistics, as previously shown in (Bega et al. 2017). authors of (Oladejo and Falowo 2018) tried to use reinforcement learning (RL) for congestion control within a slice, or resource allocation. With this goal in mind, they have put forth a system in which the MNO makes policy-based judgments on the availability of resources at any given time as well as the priorities of individual slices, and real-time slice elasticity is achieved based on requests for the grant of additional resources from each existing slice. This allows an admission-control-like mechanism to complete the cross-slice resource allocation task. It has been demonstrated that a Q-learning method significantly increases slice flexibility. (Bektas et al. 2019) suggested deep learning approach-based resource allocation techniques to offer effective resource management; however, the training process is computationally expensive or the training data are unavailable, making the proposed deep learning

approach unsuitable for large-scale systems and unable to achieve the requirements of dynamic slices and hence degraded QoS. In dynamic contexts, the reinforcement learning (RL) technique can adjust to the changes. It has therefore been used for resource scheduling (Messaoud et al. 2021), and optimization of assignments (Y. Zhou et al. 2020). By continuously interacting with the environment, which can be represented as (MDP), the RL agent can enhance its policies (Hernandez-Leal, Kartal, and Taylor 2019).

### 2.1.2. Supervised learning:

Using the provided data, the supervised learning aims to estimate the mapping. The goal data serves as the controller for the learning in supervised learning. A dataset with labels is comprised of the goal data. With this labeling, the supervisor can offer information if the machine makes a mistake while learning. In this manner, the algorithm fine-tunes itself to improve accuracy. Classification and regression are the two task classes in supervised learning, which are based on the output of the type of learning process. When obtaining a labeled dataset proves to be challenging, unlabeled data can be used in the learning process to facilitate the classification process. it refers to this type of learning as semi-supervised learning. Semi-supervised learning is the term for this type of learning. Unsupervised and supervised learning are combined to create semi-supervised learning. The goal of this kind of learning is to enhance the effectiveness of grouping or categorization. There are many unlabeled datasets and few labeled datasets used in this learning approach. In semi-supervised learning, supervision information from labeled datasets is used to improve clustering tasks by guiding which unlabeled datasets belong to the same class (van Engelen and Hoos 2020).

Several studies utilized AI technology to obtain knowledge about different types of services and classify network traffic. They devised strategies for allocating network slices based on these classification and prediction outcomes. They also employed dynamic scaling technology to automatically expand slices and monitor the usage status of each slice by evaluating the current network conditions. The study offered and built an AI-based traffic classifier and resource allocation mechanism for slicing, along with techniques for slicing allocation (Z. X. Wu et al. 2024). One branch of supervised learning is Artificial neural networks (ANNs), the most significant component of contemporary techniques that depend on artificial intelligence, which are well-known for their effectiveness in simulating non-linear systems. This can help improve reinforcement learning (RL) techniques into deep reinforcement learning (DRL) techniques, as the deep Q-Learning technique described in (Ye, Li, and Juang 2019) allocates DRL resources to provide a novel and promising policy in V2V communications. They integrated Q-learning, Deep Q Networks, and the Markov decision process (MDP) into RL. This algorithm draws comparisons between its performance and that of a neuronal network in a biological system. Artificial neural network (ANN) models that are mathematically based imitate the organic architecture of the human brain. In this sense, abstraction and generalization—two unique capacities of an organism—can be carried out by the ANN algorithm. In order to identify patterns in the input data and forecast the results of a fresh dataset that is comparable, the ANN algorithm goes through a learning process. ANN requires two fundamental building blocks: synapses/edges and neurons/nodes. The input, output, and hidden layers are the various layers that make up an artificial neural network. The output layer is responsible for forecasting the outcome of the learning process, whereas the input layer works directly with the input data. The heart of an ANN is the buried layer, where the processes of computing and learning take place. Neurones are found in every stratum. Neurones in one layer are connected to neurons in the next layer by means of edges that have a specific weight. Information from the input that may contribute to the generation or inhibition of the signal that is transmitted at each layer is contained in the weights on the edges. The deep neural network (DNN) or DL algorithm is built on the backbone of ANN. One of the machine learning subdomains known as deep learning (DL) is able to extract predictions from the input data and identify hidden patterns in the dataset. DL consists of interconnected input and output layers as well as several hidden layers between them (Nurcahyani and Lee 2021). Deep-Q Network was trained using training and

testing stage algorithms. Every agent can learn how to meet the V2V constraints with the use of a $\mathcal{E}$-greedy policy while limiting barriers to V2I communications; the proposed method's shortcoming is that it has a little amount of interference.

Shen et al. (Shen et al. 2021) proposed Uplink scheduling is improved by the DRL-Based Scheme for Resource Allocation at 5G Service-Oriented RoF-MmWave RAN, however it has an allocation fairness problem.

Xiong et al. (Xiong et al. 2019) DRL can be used to tackle important ideas linked to RA challenges that have been studied in 5G and beyond. They put into practice the DRL-based module for network slicing throughput optimization for 5G and beyond. They took into consideration a (state-actions) pair for the Q value, a greedy system for the best rewards from current actions, and a Q-learning list to keep the mapping of the system. They used mathematical formulas, which necessitated a powerful computer.

Li et al.(Sánchez, Casilimas, and Rendon 2022) provides a summary of Q-Learning and RL and explains the motivation behind creating Deep Q-Learning (DQL) using Q-Learning Network Slicing. They oversee resource distribution via DQL for radio resource slicing and the priority-based core network. Another option that DQL might offer is to serve people first by offering enhanced lucrative value in order to make better use of the computing resources and cut down on wait time. To guarantee QoE per user, DQL improves network slicing's effectiveness and adaptability. Its drawbacks include the slow learning rate and expensive GPU. Wang et al. (Z. Wang et al. 2020) proposed Although the DRL-MDP and DDPG algorithms improve the usefulness for MVNO by allocating resources at the edge network slicing, their complexity is a drawback. the study by Nguyen et al. (Nguyen et al. 2021) examined the most recent iterations of the resource management and machine learning techniques. These approaches are applied to cloud computing, edge computing, and fog computing in the 5G vehicular network. They fulfill many QoS requests. The authors suggested a vehicle communication system based on FDRL. UAVs and their role in assisting with vehicular communication.

### 2.1.3. Unsupervised learning

Unsupervised learning algorithms are typically employed for clustering tasks, where they are given a set of unlabeled data and asked to accurately anticipate the outcome. K-means, fuzzy C-means, Principal Component Analysis (PCA), Auto encoders (AEs), Self-Organizing Maps (SOMs), hidden Markov Model (HMM), and restricted Boltzmann machine (RBM) are examples of common unsupervised techniques. Additionally, Deep Learning (DL) techniques like Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) algorithms can perform better when unsupervised machine learning (ML) is applied (Zhang and Zhu 2020).

• K-means: This technique is used for grouping unlabeled input data into distinct clusters. Every new data point is assigned to a cluster by the K-means algorithm based on its distance from the closest related centroid. The process is continued until neither the data points nor the centroids are changed. The centroids are updated in accordance with the allocated data point. K, which stands for the intended number of clusters, has a significant effect on the algorithm's performance (Azimi et al. 2022)

• Self-organizing map (SOM): SOM is frequently applied to data grouping and dimensionality reduction. SOM is composed of two layers: an input layer and a map layer. Numerous neurons are present in each layer, and each neuron has a unique weight vector. Using an unsupervised competitive learning methodology, SOM constructs and rearranges the map throughout the training phase. Any new input vector is grouped in a cluster by the neuron that emerged victorious (van Engelen and Hoos 2020).

• Hidden Markov model (HMM): This method models a system by using a Markov process with unknown parameters. The goal of HMM is to extract unknown parameters from known ones. Markov models are helpful in situations where memory is lacking since the future state estimation sequence in them rely only on the current state. The system can gradually alter each state's probability distribution in an HMM (Zhang and Zhu 2020)

• Autoencoders: These are learning circuits that replicate inputs into outputs with the goal of producing the least amount of variance. To fine-tune the architecture, autoencoders are stacked and taught unsupervised bottom-up, with the top layer being trained by supervised learning. These designs have the potential to produce precise and useful outcomes for situations involving both regression and classification (Baldi 2012).

### 2.2.     Algorithms Based on Evolution:

evolutionary algorithms known as a significant class of machine learning methods that use statistical evolutions to learn random strategies from system feedback. An illustration of how they are used the allocation of network resources within per slice is provided by (D. Wu et al. 2019), in which  various people connected to various network slices ,Where had continually evolving network connection with constantly updating. Clients are categorized into separated group ,In order to ensure every user within group have the same requirements of service. hence, with a deconstructing of the complex model of management resources in sliced network, this approach facilitates the improvement the utilization of resources strategy. However, in the setting of slice admission control, as demonstrated in our earlier study (Haque and Abu-Ghazaleh 2016)  The efficiency with which genetic algorithms (GAs) work. with the implementation of the admission decision into binary sequences that represent a chromosome, GA could produce a modified policy that recursively produce  optimal value. Furthermore, the outperform of this technique in a first generation over the benchmark due to the manual integration of (randomly) benchmark policies. Additionally, it demonstrates strong resistance to changing surroundings.

### 3.   End-to-End communication networks.

As seen in Fig. 4, network slicing is defined as a logical end-to-end network that may dynamically supply one or more network services according to the slice requirements. Every network slice has certain processes in place to ensure that users' performance requirements are met (Azimi et al. 2022).

(Hassine 2017) defined the sliced network model which involves the slicing of the Radio Access Network (RAN) with the Core Network (CN) as End-to-End slicing, which was proposed by Next Generation Mobile Network (NGMN). It is the process of running and managing several logically separate virtual communication networks atop a common physical infrastructure with the goal of improving future network services such as the various requirements, flexibility, adaptively scaling, and security. This necessitates highly countable, divisible, and isolatable network resources and functions, which are achievable with today's network function virtualization technology.
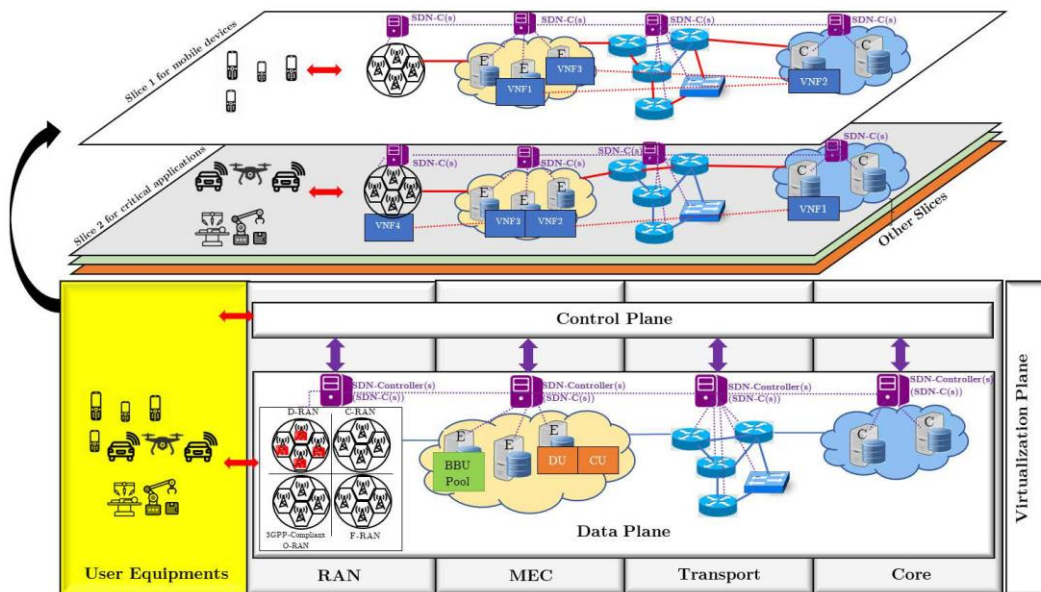
**Figure 4: Using slicing enablers in network slicing(Azimi et al. 2022).**

(Khamse-Ashari et al. 2022) suggested a novel approach to provide End-to-End resources that avoids the difficulties of resolving the MILP problem and maximizes tenant satisfaction and InP operational expenses. With the goal avoid the necessity for access points or centralized data for sharing the capacity of their resources, the authors developed a distributed privacy-saving approach, however, this increased the burden of Signaling. Because of the multiple slice managers that were used, a less-than-ideal global solution might have been produced.

(Wen et al. 2019) developed a heuristic strategy using variable locality searching to examine forceful technologies of network slicing to mitigate disquiet about vulnerabilities in VNFs and fluctuations of the traffic requirement. however, this approach resulted in a large signal overhead. In reality, only the user's requirements are known; the authors assumed that each slice would have a constant number of VNFs. As a result, testing the virtual network embedding procedure repeatedly continues until the right number of VNFs is found. In related studies, (Luu et al. 2020; Luu, Kerboeuf, and Kieffer 2021) without any prior information on virtual network topology or resources, they examined NS determination and embedding and Proposed a heuristic approach using user distribution and demands. However, the authors neglected to take into account the reconfiguration and reallocation of resources in a dynamic setting and optimize the sharing of physical resources to meet the Service Level Agreement for different slices' needs, On the other hand, While evaluating the suggested solutions, We must employ the 5G service-oriented architecture and network Key Performance Indicators. (Kazemifard and Shah-Mansouri 2021) created a heuristic technique for resource dimensioning to figure out how much virtual and computational power is To optimize time to reaction for specific needs, the Network Slicing Orchestration System (NSOS) can be used. To keep response times within the delay requirement, the suggested method only permits the NSOS to adjust its resources in response to future workloads. The system blocks excessive Slice Orchestration Requests (SORs) as service demand rises, and it reserves certain resources for idle time when demand falls. The delay that SOR encountered due to various NSO entities was the only one that the authors assessed. Nevertheless, iterative techniques for resource estimates for every slice require a lot of computing power, as mentioned in (Kazemifard and Shah-Mansouri 2021). Furthermore, resource distribution inside a slice isn't always fair to all users. (Alfoudi et al. 2019) suggested a technique for

Resource Management in network slices that could be used to allocate resources in a mobile LTE network slicing. This method guarantees isolation between users and fairly allocated bandwidth. Nevertheless, an exponential smoothed method was proposed, that had only two possible effects: a consistent allocation of resources and a set interval redistribution of physical resources. Due to the unpredictability of demands, the fixed interval redistribution method may result in either excessive or insufficient allocation of resources. This is because the model does not include any additional resources available in each period.

## 3.1. **Resource allocation at Radio access network**

The literature has recently given dynamic RAN radio resource allocation algorithms more attention. The following is a summary of some of the pertinent research studies: (Kamel, Le, and Girard 2014) researched the resource sharing for an effective LTE network slicing into numerous virtual networks (VNs) to provide services to various service providers.

(Gharehgoli et al. 2023) tackles the problem of unknown information in network slicing for 5G+ communication. It focuses on the unpredictability of channel status information (CSI) and demand (user requests, bandwidth, and workloads associated with virtual network functions). an objective is to use the algorithms of deep reinforcement learning (DRL) for distributed of resources in an end-to-end NS, In order to maximize the utility of the infrastructure provider. The researchers present a recurrent deterministic policy gradient (RDPG) algorithm that outperforms other algorithms. To solve the problem of resource allocation as a non-convex mixed-integer non-linear programming problem. According to the simulation results, the RDPG method beats the SAC approach by an average of 70%, the RDPG algorithm's action selection computational complexity is O($H^2$), making it practical to use. In (L. Zhou et al. 2020), a method for allocating radio resources to maximize service adoption is suggested to optimize the service level agreement (SLA) also maintaining isolation among slices. In fact, such problem is evaluated utilizing the dual Lagrangian algorithm technique. In (Tang, Shim, and Quek 2019), the researchers presented the cloud RAN (C-RAN) maximizes the benefits of the operator by accurately tolerating the cut solicitations. Two kinds of long-haul and momentary interests are thought of. The drawn-out credit is determined by the qualities in the network cut demand and the momentary interest is accomplished by memory structure energy destruction in every frame. The streamlining issue is figured out as mixed-integer nonlinear processing (MINP), then to tackle the issue, the researchers utilize a progressively raised estimation (SCA) technique. In (D'Oro et al. 2018), a near-optimal less complexity distribution is recommended the RAN segmenting obstacle as a crowd resolute. The complication desires to reduce cost. The researchers formulate resource allocation the complication is based on utility desolation, delay, and minimization scale restrictions. To sustain the authenticity impulsion, they employ adaptive modulation with encoding. In (Sun et al. 2020), they suggested the smart delivery mechanism through the implementation of multi-agent Q-Learning in order to reduce delivery expense whereas ensuring diverse service quality demands. To compute the cost, the research identify four kinds of delivery cost: 1) Cost of switching kids of service whenever the user equipment (UE) stands inside coverage along the similar base station (BS). 2) Delivery cost at UE stages base station coverage with the similar kind of service. 3) Cost associated against cliet movement and service change species. 4) The cost of implementing a modern network segment for maintaining quality of service delivery to the client. In (X. Wang and Zhang 2019), the researchers explained NwS resource allocation obstacle in 5G C-RAN for optimization Benefit drivers. Includes problem framework the upward stage that handles the planning of the potential agreement bundle task; A lower layer controls the radio unit remotely linkage, authority and sub-channel assignment. To alleviate complexity of the Q-value table, the authors used multiple proxy Q-learning technology. In(Lee et al. 2018), the researchers investigated an effective NS network for heterogeneous multi-tenant downlink cloud RAN (H-CRAN) by regarding each of the little cells also total cell layers. The suggested structure contains: dual stages, a higher stage of admittance regulation management, baseband resource allocation, client association, with less a level to handle the allocation of radio resources among clients. The objective is to optimize the rate of tenants throughput with consideration of quality of

service limitation, forward connectivity using delivery capabilities, tenant priorities, baseband resources, with intervention.

## 3.2. Resource allocation at the Core network

(Ko, Lee, and Pack 2022) propose a priority-based dynamic resource allocation scheme (PDRAS) in 5G NS environments to optimize quality of service (QoS) and distribution efficiency. by considering different priorities of slices and dynamically adjusts allocated resources based on slicing information and formulates a constrained Markov decision process problem to obtain the optimal allocation policy using linear programming. Extensive evaluation results demonstrate that PDRAS with the optimal policy outperforms other schemes in terms of QoS and resource usage efficiency .it addresses the challenge of resource allocation in 5G network slicing environments, ensuring that sufficient resources are allocated to resource-intensive service slices to prevent degradation of QoS . (Chen et al. 2021), the essential objective of this research is to minimize the overall power consumption of the cloud hub, that includes of fixed power consumption as well as Load-dependent dynamic power consumption. The researchers regard resource budgeting, career creation, with flood routing, ensuring E2E latency for every service, whereas the E2E delay includes of the overall NFV delay on the cloud hubs with total communication delay on links., (Ebrahimi et al. 2020), the researchers introduced a modern model for the slice admission and distribution of resources technique in multi-tenant scenario to reduce the expenses of bandwidth and power consumption of each running cloud hubs. In(Reddy, Baumgartner, and Bauschert 2017), to deal with unpredictability in the traffic demand, a new method of optimization is proposed, which depending on the *Γ-robustness* notion is developed. A MILP formulation for the Γ-robust optimization is used. To improve the scalability of the model, a modified MIP-based variable neighborhood search (VNS) heuristic is offered. The authors in (Baumgartner et al. 2017) reveal a new model that uses the idea of light robustness to solve the ambiguity of traffic in NS scalability problems and to gain a better understanding of the balancing that exists between the expense of ensuring adaptability and the level of adaptability achieved.

## Table 1: Summary of related works

| Ref | Objective | Method | The scope of slicing | Evaluation method |
|---|---|---|---|---|
| (Chien et al. 2020) | Improves utilization of resources per slice to meet delay requirements without wasteful allocation. | Heuristic algorithm | E2E | Emulation |
| (Ko, Lee, and Pack 2022) | Enhance the Quality of Service (QoS) for slices simultaneously ensuring that the total resources assigned remain below a specific limit. | CMDP and Convex Optimization (LP) | Core | Numerical analysis |
| (Li, Zhu, and Liu 2020) | Increase the total rate of access accumulation across all slices | Deep-RL | E2E | Simulation |

| (Messaoud et al. 2021) | An Improvement of QoS for slices | Deep federated Q-Learning | RAN | Simulation |
|---|---|---|---|---|
| (Halabian 2019) | Improvement of the whole system utility (focusing on the fairness) | Convex Optimization and Game Theory | C-RAN | Numerical analysis |
| (Sciancalepore et al. 2017) | Ensure the achievements of a Service Level Agreements | Deep-RL | E2E | simulation |
| (Van Huynh et al. 2019) | Optimizes the overall average revenue of the operator of the network | Deep-RL | E2E | Numerical analysis |
| (L. Zhou et al. 2020) | SLA contract rate maximization | Lagrangian dual | RAN | Numerical analysis |
| (Tang, Shim, and Quek 2019) | Revenue maximization | SCA | C-RAN | Numerical analysis |
| (D'Oro et al. 2018) | Cost minimization | Game theory | RAN | Numerical analysis |
| (Sun et al. 2020) | Handover cost minimization | Multi-agent Q-learning | RAN | Numerical analysis |
| (X. Wang and Zhang 2019) | Utility maximization | Multi-agent Q-learning | C-RAN | Numerical analysis |
| (Lee et al. 2018) | Throughput maximization | Greedy & Lagrangian dual | H-CRAN | Numerical analysis |
| (Ebrahimi et al. 2020) | Cost minimization | Heuristic | Core | Numerical analysis |
| (Reddy, Baumgartner, and Bauschert 2017) | Reduction of Cost | Heuristic-based on Γ-robustness | Core | Numerical analysis |
| (Baumgartner et al. 2017) | Reduction of Cost | Light robustness-based on Γ-robustness | Core | Numerical analysis |

## Conclusion

The increasing of network communications with various needs, has made the allocation of resources more complex problem. NFV and SDN are the primary technologies that enabled network softwarization into different slices. These slices have represented as emerged technique for efficient management of

resources. Further enhancements on the process of resource managements are produced by Machine learning approaches. Some of the benefits of these approaches is improvements of performance network by reducing delay, increase allocated bandwidth and capacity and reducing the operational costs by optimized allocation of resources. Additional advantage of incorporating of deep reinforcement learning in 5G+ network to get optimal and adaptive allocation policy without needs to prior knowledge. These enhancements of Machine Learning approaches reflect great contributions to the overall network by ensure its efficiency and stability. So, with increasing growth of 5G communications with these advancements of machine learning techniques will play an effective role in satisfying the diverse requirements for different services of communications.

## References

1) Alfoudi, Ali Saeed Dayem et al. 2019. "An Efficient Resource Management Mechanism for Network Slicing in a LTE Network." *IEEE Access* 7: 89441–57.

2) Asakipaam, Simon Atuah, Jerry John Kponyo, and Kwame Oteng Gyasi. 2023. "Resource Provisioning and Utilization in 5G Network Slicing: A Survey of Recent Advances, Challenges, and Open Issues." *International Journal of Computer Networks and Applications* 10(2): 201–16.

3) Azimi, Yaser, Saleh Yousefi, Hashem Kalbkhani, and Thomas Kunz. 2022. "Applications of Machine Learning in Resource Management for RAN-Slicing in 5G and Beyond Networks: A Survey." *IEEE Access* 10(October): 106581–612. https://ieeexplore.ieee.org/document/9904606/.

4) Bakhshi, Taimur. 2017. "State of the Art and Recent Research Advances in Software Defined Networking." *Wireless Communications and Mobile Computing* 2017: 1–35. https://www.hindawi.com/journals/wcmc/2017/7191647/.

5) Baldi, Pierre. 2012. "Autoencoders, Unsupervised Learning, and Deep Architectures." *ICML Unsupervised and Transfer Learning*: 37–50.

6) Baumgartner, Andreas, Thomas Bauschert, Abdul A. Blzarour, and Varun S. Reddy. 2017. "Network Slice Embedding under Traffic Uncertainties - A Light Robust Approach." *2017 13th International Conference on Network and Service Management, CNSM 2017* 2018-Janua: 1–5.

7) Bega, Dario et al. 2017. "Optimising 5G Infrastructure Markets: The Business of Network Slicing." *Proceedings - IEEE INFOCOM*.

8) ———. 2020. "DeepCog: Optimizing Resource Provisioning in Network Slicing with AI-Based Capacity Forecasting." *IEEE Journal on Selected Areas in Communications* 38(2): 361–76.

9) Bektas, Caner, Stefan Monhof, Fabian Kurtz, and Christian Wietfeld. 2019. "Towards 5G: An Empirical Evaluation of Software-Defined End-to-End Network Slicing." *2018 IEEE Globecom Workshops, GC Wkshps 2018 - Proceedings*: 1–6.

10) Cao, Haotong et al. 2024. "Softwarized Resource Allocation in Digital Twins-Empowered Networks for Future Quantum-Enabled Consumer Applications." *IEEE Transactions on Consumer Electronics* 70(1): 800–810.

11) Chen, Wei Kun et al. 2021. "Optimal Network Slicing for Service-Oriented Networks with Flexible Routing and Guaranteed E2E Latency." *IEEE Transactions on Network and Service Management* 18(4): 4337–52.

12) Chien, Hsu Tung, Ying Dar Lin, Chia Lin Lai, and Chien Ting Wang. 2020. "End-to-End Slicing with Optimized Communication and Computing Resource Allocation in Multi-Tenant 5G Systems." *IEEE Transactions on Vehicular Technology* 69(2): 2079–91.

13) D'Oro, Salvatore, Francesco Restuccia, Tommaso Melodia, and Sergio Palazzo. 2018. "Low-Complexity Distributed Radio Access Network Slicing: Algorithms and Experimental Results." *IEEE/ACM Transactions on Networking* 26(6): 2815–28.

14) Ebrahimi, Sina, Abulfazl Zakeri, Behzad Akbari, and Nader Mokari. 2020. "Joint Resource and Admission Management for Slice-Enabled Networks." *Proceedings of IEEE/IFIP Network Operations and Management Symposium 2020: Management in the Age of Softwarization and Artificial Intelligence, NOMS 2020*.

15) van Engelen, Jesper E., and Holger H. Hoos. 2020. "A Survey on Semi-Supervised Learning." *Machine Learning* 109(2): 373–440. https://doi.org/10.1007/s10994-019-05855-6.

16) Ericsson. 2021. "Network Slicing: Top 10 Use Cases to Target." www.ericsson.com/en/ mobility-report/reports/june-2021.

17) ETSI. 2017. "Network Functions Virtualisation ( NFV ) Release 3 ; Evolution and Ecosystem ; Report on Network Slicing Support with ETSI NFV Architecture Framework." *Etsi Gr Nfv-Eve 012 V3.1.1* 1: 1–35.

18) Gharehgoli, Amir et al. 2023. "AI-Based Resource Allocation in End-to-End Network Slicing under Demand and CSI Uncertainties." *IEEE Transactions on Network and Service Management*: 1.

19) Halabian, Hassan. 2019. "Distributed Resource Allocation Optimization in 5G Virtualized Networks." *IEEE Journal on Selected Areas in Communications* 37(3): 627–42.

20) Han, Bin, Di Feng, and Hans D. Schotten. 2018. "A Markov Model of Slice Admission Control." *IEEE Networking Letters* 1(1): 2–5.

21) Haque, Israat Tanzeena, and Nael Abu-Ghazaleh. 2016. "Wireless Software Defined Networking: A Survey and Taxonomy." *IEEE Communications Surveys and Tutorials* 18(4): 2713–37.

22) Hassine, Nesrine Ben. 2017. "Machine Learning for Network Resource Management." *Http://Www.Theses.Fr* 2.

23) Hernandez-Leal, Pablo, Bilal Kartal, and Matthew E. Taylor. 2019. 33 Autonomous Agents and Multi-Agent Systems *A Survey and Critique of Multiagent Deep Reinforcement Learning*. Springer US. https://doi.org/10.1007/s10458-019-09421-1.

24) Van Huynh, Nguyen, Dinh Thai Hoang, Diep N. Nguyen, and Eryk Dutkiewicz. 2019. "Optimal and Fast Real-Time Resource Slicing with Deep Dueling Neural Networks." *IEEE Journal on Selected Areas in Communications* 37(6): 1455–70.

25) Ji, Hyoungju et al. 2018. "Ultra-Reliable and Low-Latency Communications in 5G Downlink: Physical Layer Aspects." *IEEE Wireless Communications* 25(3): 124–30.

26) Jin, Xin, Li Erran Li, Laurent Vanbever, and Jennifer Rexford. 2013. "SoftCell." In *Proceedings of the Ninth ACM Conference on Emerging Networking Experiments and Technologies*, New York, NY, USA: ACM, 163–74. https://dl.acm.org/doi/10.1145/2535372.2535377.

27) Kamel, Mahmoud I., Long Bao Le, and Andre Girard. 2014. "LTE Wireless Network Virtualization: Dynamic Slicing via Flexible Scheduling." *IEEE Vehicular Technology Conference*: 1–5.

28) Kazemifard, Nasim, and Vahid Shah-Mansouri. 2021. "Minimum Delay Function Placement and Resource Allocation for Open RAN (O-RAN) 5G Networks." *Computer Networks* 188(October 2020): 107809. https://doi.org/10.1016/j.comnet.2021.107809.

29) Khamse-Ashari, Jalal, Gamini Senarath, Irem Bor-Yaliniz, and Halim Yanikomeroglu. 2022. "An Agile and Distributed Mechanism for Inter-Domain Network Slicing in Next Generation Mobile Networks." *IEEE Transactions on Mobile Computing* 21(10): 3486–3501.

30) Ko, Haneul, Jaewook Lee, and Sangheon Pack. 2022. "PDRAS: Priority-Based Dynamic Resource Allocation Scheme in 5G Network Slicing." *Journal of Network and Systems Management* 30(4): 1–20. https://doi.org/10.1007/s10922-022-09681-5.

31) Kreutz, Diego et al. 2008. "OpenFlow: Enabling Innovation in Campus NetworksSoftware-Defined Networking: A Comprehensive Survey." *Proceedings of the IEEE* 103(1): 14–76. http://ccr.sigcomm.org/online/files/p69-v38n2n-mckeown.pdf.

32) Ksentini, Adlen, and Navid Nikaein. 2017. "Toward Enforcing Network Slicing on RAN: Flexibility and Resources Abstraction." *IEEE Communications Magazine* 55(6): 102–8.

33) Lee, Ying Loong, Jonathan Loo, Teong Chee Chuah, and Li Chun Wang. 2018. "Dynamic Network Slicing for Multitenant Heterogeneous Cloud Radio Access Networks." *IEEE Transactions on Wireless Communications* 17(4): 2146–61.

34) Li, Taihui, Xiaorong Zhu, and Xu Liu. 2020. "An End-to-End Network Slicing Algorithm Based on Deep Q-Learning for 5G Network." *IEEE Access* 8: 122229–40.

35) Liu, Yongshuai, Jiaxin DIng, and Xin Liu. 2021. "Resource Allocation Method for Network Slicing Using Constrained Reinforcement Learning." *2021 IFIP Networking Conference, IFIP Networking 2021*: 1–3.

36) Luu, Quang Trung, Sylvaine Kerboeuf, and Michel Kieffer. 2021. "Foresighted Resource Provisioning for Network Slicing." *IEEE International Conference on High Performance Switching and Routing, HPSR* 2021-June: 1–8.

37) Luu, Quang Trung, Sylvaine Kerboeuf, Alexandre Mouradian, and Michel Kieffer. 2020. "A Coverage-Aware Resource Provisioning Method for Network Slicing." *IEEE/ACM Transactions on Networking* 28(6): 2393–

2406.

38) Messaoud, Seifeddine et al. 2021. "Deep Federated Q-Learning-Based Network Slicing for Industrial IoT." *IEEE Transactions on Industrial Informatics* 17(8): 5572–82.

39) Nguyen, Hoa T.T. et al. 2021. "DRL-Based Intelligent Resource Allocation for Diverse QoS in 5G and toward 6G Vehicular Networks: A Comprehensive Survey." *Wireless Communications and Mobile Computing* 2021.

40) Nurcahyani, Ida, and Jeong Woo Lee. 2021. "Role of Machine Learning in Resource Allocation Strategy over Vehicular Networks: A Survey." *Sensors* 21(19): 6542. https://www.mdpi.com/1424-8220/21/19/6542.

41) Nyanteh, Andrews O., Maozhen Li, Maysam F. Abbod, and Hamed Al-Raweshidy. 2021. "CloudSimHypervisor: Modeling and Simulating Network Slicing in Software-Defined Cloud Networks." *IEEE Access* 9: 72484–98.

42) Oladejo, Sunday O., and Olabisi E. Falowo. 2018. "Profit-Aware Resource Allocation for 5G Sliced Networks." *2018 European Conference on Networks and Communications, EuCNC 2018*: 43–47.

43) Ordonez-Lucena, Jose et al. 2017. "Network Slicing for 5G with SDN/NFV: Concepts, Architectures, and Challenges." *IEEE Communications Magazine* 55(5): 80–87. http://ieeexplore.ieee.org/document/7926921/.

44) Oulahyane, Hafsa Ait et al. 2024. "Towards an SDN-Based Dynamic Resource Allocation in 5G Networks." *Procedia Computer Science* 231(2023): 205–11. https://doi.org/10.1016/j.procs.2023.12.194.

45) Pang, Xue, and Peiying Zhang. 2020. "Resource Allocation Strategy of IoT Based on Network Slicing." *2020 IEEE Computing, Communications and IoT Applications, ComComAp 2020*.

46) Reddy, Varun S., Andreas Baumgartner, and Thomas Bauschert. 2017. "Robust Embedding of VNF/Service Chains with Delay Bounds." *2016 IEEE Conference on Network Function Virtualization and Software Defined Networks, NFV-SDN 2016*: 93–99.

47) Van Rossem, Steven et al. 2016. "Deploying Elastic Routing Capability in an SDN/NFV-Enabled Environment." *2015 IEEE Conference on Network Function Virtualization and Software Defined Network, NFV-SDN 2015*: 22–24.

48) Rost, Peter et al. 2017. "Network Slicing to Enable Scalability and Flexibility in 5G Mobile Networks." *IEEE Communications Magazine* 55(5): 72–79.

49) Sánchez, Johanna Andrea Hurtado, Katherine Casilimas, and Oscar Mauricio Caicedo Rendon. 2022. "Deep Reinforcement Learning for Resource Management on Network Slicing: A Survey." *Sensors* 22(8).

50) Sciancalepore, Vincenzo et al. 2017. "Mobile Traffic Forecasting for Maximizing 5G Network Slicing Resource Utilization." *Proceedings - IEEE INFOCOM* (671584).

51) Shen, Shuyi, Ticao Zhang, Shiwen Mao, and Gee Kung Chang. 2021. "DRL-Based Channel and Latency Aware Radio Resource Allocation for 5G Service-Oriented RoF-MmWave RAN." *Journal of Lightwave Technology* 39(18): 5706–14.

52) Sun, Yao et al. 2020. "Efficient Handover Mechanism for Radio Access Network Slicing by Exploiting Distributed Learning." *IEEE Transactions on Network and Service Management* 17(4): 2620–33.

53) Tang, Jianhua, Byonghyo Shim, and Tony Q.S. Quek. 2019. "Service Multiplexing and Revenue Maximization in Sliced C-RAN Incorporated With URLLC and Multicast EMBB." *IEEE Journal on Selected Areas in Communications* 37(4): 881–95.

54) Thirupathi, V., Ch Sandeep, S. Naresh Kumar, and P. Pramod Kumar. 2019. "A Comprehensive Review on Sdn Architecture, Applications and Major Benifits of SDN." *International Journal of Advanced Science and Technology* 28(20): 607–14.

55) Wang, Xiaofei, and Tiankui Zhang. 2019. "Reinforcement Learning Based Resource Allocation for Network Slicing in 5G C-RAN." *2019 Computing, Communications and IoT Applications, ComComAp 2019*: 106–11.

56) Wang, Zhaoying, Yifei Wei, F. Richard Yu, and Zhu Han. 2020. "Utility Optimization for Resource Allocation in Edge Network Slicing Using DRL." *2020 IEEE Global Communications Conference, GLOBECOM 2020 - Proceedings*.

57) Wen, Ruihan et al. 2019. "On Robustness of Network Slicing for Next-Generation Mobile Networks." *IEEE Transactions on Communications* 67(1): 430–44.

58) Wu, Dapeng et al. 2019. "Biologically Inspired Resource Allocation for Network Slices in 5G-Enabled Internet of Things." *IEEE Internet of Things Journal* 6(6): 9266–79.

59) Wu, Zong Xun, Yun Zhe You, Chien Chang Liu, and Li Der Chou. 2024. "Machine Learning Based 5G Network Slicing Management and Classification." *6th International Conference on Artificial Intelligence in Information*

*and Communication, ICAIIC 2024*: 371–75.

60) Xiong, Zehui et al. 2019. "Deep Reinforcement Learning for Mobile 5G and beyond: Fundamentals, Applications, and Challenges." *IEEE Vehicular Technology Magazine* 14(2): 44–52.

61) Ye, Hao, Geoffrey Ye Li, and Biing Hwang Fred Juang. 2019. "Deep Reinforcement Learning Based Resource Allocation for V2V Communications." *IEEE Transactions on Vehicular Technology* 68(4): 3163–73.

62) Zhang, Shunliang, and Dali Zhu. 2020. "Towards Artificial Intelligence Enabled 6G: State of the Art, Challenges, and Opportunities." *Computer Networks* 183(October): 107556. https://doi.org/10.1016/j.comnet.2020.107556.

63) Zhou, Liushan, Tiankui Zhang, Jing Li, and Yutao Zhu. 2020. "Radio Resource Allocation for RAN Slicing in Mobile Networks." *2020 IEEE/CIC International Conference on Communications in China, ICCC 2020* (Iccc): 1280–85.

64) Zhou, Yuan et al. 2020. "Subcarrier Assignment Schemes Based on Q-Learning in Wideband Cognitive Radio Networks." *IEEE Transactions on Vehicular Technology* 69(1): 1168–72.