

Machine Learning-Driven Service Function Chain Placement and Scaling in MEC-enabled 5G Networks

Tejas Subramanya, Davit Harutyunyan and Roberto Riggio
FBK CREATE-NET

Via Alla Cascata 56/C, 38123, Trento, Italy;
Email: {t.subramanya,d.harutyunyan,rriggio}@fbk.eu

Abstract—5G mobile network technology promises to deliver unprecedented ultra-low latency and high data rate, paving the way for many novel applications and services. *Network Function Virtualization* (NFV) and *Multi-access Edge Computing* (MEC) are two of the technologies that are expected to play a pivotal role in 5G to achieve ambitious Quality of Service requirements of such applications. While NFV provides flexibility by enabling network functions to be dynamically deployed and inter-connected to realize Service Function Chains (SFC), MEC brings the computing capability to the edges of the mobile network thus reducing latency and alleviating the transport network load. However, adequate mechanisms are needed to meet the dynamically changing network service demands, to optimally utilize the network resources while, at the same time, making sure that E2E latency requirement of services is always satisfied.

In this work, we first propose a neural-network machine learning model that can perform auto-scaling by predicting the required number of virtual network function instances based on the traffic demand, using the traffic traces collected over a real-operator experimental network. We then employ Integer Linear Programming techniques to formulate and solve a joint user association and SFC placement problem, where each SFC represents a service requested by a user with E2E latency and data rate requirements. Finally, we propose a heuristic to address the scalability concern of the ILP model.

Index Terms—Auto-scaling, Service Function Chain Placement, Machine learning, Multi-access Edge Computing

I. INTRODUCTION

The 5th generation of mobile networks is expected to support high data rates, extremely low-latency, high reliability, the capability to extend access to distributed computation and storage facilities in addition to connectivity and bandwidth [1]. These characteristics of the 5G networks open the door for many novel Ultra-reliable low-latency (URLLC) applications such as augmented/virtual reality and autonomous driving, whose ambitious QoS requirements cannot be satisfied by the preprocessors of the 5G networks. Therefore, the 5G architecture needs to incorporate new technologies such as Multi-access Edge Computing (MEC) [2] and Network Function Virtualization (NFV) [3], to meet the insatiable data rate and low-latency requirements of the applications mentioned above [4].

The basic idea of MEC is to bring computing capabilities and applications closer to the end-users, from cloud data centers to the edges of the cellular network, therefore, reducing the delay experienced by the users and alleviating the transport network load. Consequently, the ETSI ISG MEC group proposes three possible MEC deployment options, collocated with the gNodeB (*gnb.mec*) or collocated with an aggregation

point (*ap.mec*) or collocated with the 5GC (*5gc.mec*) [5]. It is essential to mention that the integration of MEC into the 5G network requires a User Plane Function (UPF), which is a core network component in the 5G technology responsible for decapsulating the GTP header from user-plane traffic, routing plain IP packets to/from the MEC applications and re-encapsulating the GTP header back, to be collocated with all the MEC nodes in order to reap the benefits of the MEC system [5]. The closer the MEC nodes are towards the end-users, the scarcer their computational resources become.

NFV, on the other hand, decouples network functions (e.g., UPF) or MEC applications from their dedicated proprietary hardware and deploys them as virtualized software entities on commodity servers [6]. In our work, we use the generic term VxFs to refer to a collection of UPF virtualized network functions (VNFs) and virtualized MEC application functions (VMAFs). VxFs require a specific computational capacity (e.g., CPU) to be spawned/instantiated and can be chained together forming a Service Function Chain (SFC) that represents a specific service, with a guaranteed latency and data rate requirements, that can be requested by User Equipments (UEs). As shown in Fig.1, there are multiple locations (e.g., *gnb.mec*, *ap.mec*, *5gc.mec*) for instantiating VxFs, which can be shared by many UEs. In the described network scenario, given the UEs with their SFC demands, the natural question that arises is how to associate UEs, place their SFCs, allocate sufficient VxF instances and resources to make sure that the UEs SFC requirements are satisfied while the network resources are used efficiently? Moreover, the mobility patterns of UEs result in non-uniform traffic distribution within the mobile network [7]. Consequently, the number of VxF instances required to manage load variations and to meet performance guarantees is expected to fluctuate frequently. Towards this end, auto-scaling of VxFs in addition to distributed SFC placement is thought to be an essential requirement for successful orchestration.

Most of the existing literature address either the problem of VxF autoscaling [8] or the placement of SFC in distributed MEC nodes [9]. In this paper, we first advocate that distributed VxFs of SFC has to be proactively scaled in synergy with varying network traffic dynamics to avoid service disruption. Based on those scaling decisions, the VxFs need to be dynamically placed in distributed MEC nodes, to minimize end-to-end latency and to reduce latency violations. To the best of our knowledge, we are the first to address the combined challenges in VxF auto-scaling and placement of

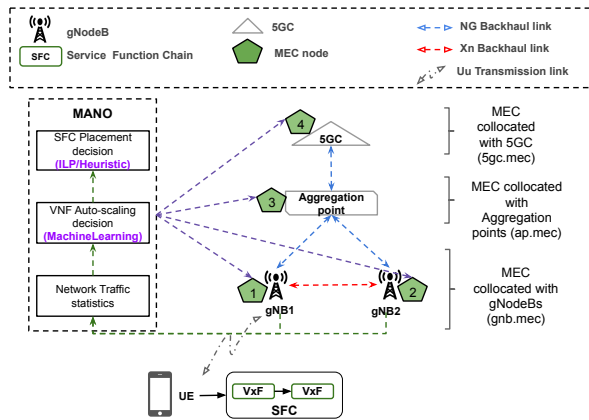


Fig. 1: An example of a distributed MEC-NFV System Architecture.

SFCs within a distributed MEC-NFV environment, based on the real-operator mobile network traces. This work has three main contributions (as depicted in Fig. 1):

(i) Leveraging a real-operator experimental 5G network dataset, we model and implement a neural-network-based Multi-layer Perceptron (MLP) classifier, which *estimates the required number of UPF* instances as a function of the base station network traffic they should process. Furthermore, we compare the performance of MLP classifier with other classification algorithms in Machine Learning (ML). The output from step 1 is fed as an input to the Integer Linear Programming (ILP) problem in step 2.

(ii) Employing ILP technique to formulate and solve a *'joint UE association and SFC placement problem'*, where each SFC is composed of a number of VxVs, with specific latency and data rate demands as requested by UEs positioned in diverse areas of the 5G mobile network. We also develop a complete end-to-end latency model for 5G mobile networks.

(iii) Proposing a heuristic algorithm with the same objective as that of ILP to address the scalability problem of ILP.

Out of the above three contributions, the first one is from our previous work [10], while the second and third contributions are entirely new that proposes an ILP and heuristic models for the joint UE association and SFC placement problem rather than just VxV placement problem like in our previous work. Moreover, we consider 5G mobile networks, user associations to gNodeBs and their SFC requests, and comprehensive end-to-end latency calculations in our models.

The rest of this paper is organized as follows. Section II describes the related work. Section III describes the proposed MLP classifier model and evaluates its performance. In Section IV, we model our latency-optimal SFC placement problem while Section V formulates the problem using ILP and also proposes a heuristic algorithm. In Section VI, we perform several experiments to evaluate our proposed SFC placement solutions. Finally, we conclude the paper in Section VII.

II. STATE OF THE ART

ETSI NFV Industry Specification Group defines network service as a composition of one or more VNFs that are chained together. Each VNF requires a specific amount of

resource to process the traffic flowing through it. To deploy a network service, the operator needs to find the right placement of VNFs complying with various resource constraints and service latency agreements. Once the hosts are selected and the VNFs deployed, resource requirements for the VNFs may vary due to traffic fluctuations. To meet these demands, a resource allocation algorithm is needed that can automatically allocate/release resources to a VNF (vertical scaling) or add/remove one or more VNF instances (horizontal scaling).

A. Virtual Network Function auto-scaling.

Previous works on VNF auto-scaling can be divided into two categories: reactive mode and proactive mode.

In reactive mode, threshold levels can be either statically pre-defined or dynamically updated. In [11], [12], and [13], the authors propose scalability mechanisms based on static thresholds. They define two threshold levels ($scalein_{thr}$ and $scaleout_{thr}$) to determine if the load reduces below or exceeds above the respective limits and accordingly trigger the scaling process. However, such techniques may result in an oscillating behaviour affecting the overall system performance. On the other hand, [14] and [15] propose mechanisms such as queuing theory and reinforcement learning, which allows the scaling policy to be improved based on dynamic or adaptive thresholds. Although it performs better than static approaches, it remains a reactive solution with similar weaknesses.

In proactive mode, forecasting techniques (e.g., machine learning) are applied to allow the systems to automatically learn and to anticipate future needs, based on which scalability decisions are taken. For example, the authors in [16] propose a solution to forecast CPU usage based on a historical dataset using time series model. Other authors such as Mijumbi et al. [17] and Mestres et al. [18] addresses the problem of managing VNF resource fluctuations by predicting resource requirements using ML techniques and thereby enhancing the performance of the resource allocation algorithm.

In contrast to these works which targets data centers, our approach investigates the problem of proactive auto-scaling in a distributed MEC-NFV deployment. Moreover, we use real-operator traffic traces to generate training sets required for predicting auto-scaling decisions, unlike other works that are based on simulated datasets.

B. Service Function Chain placement.

There already exists some literature on the SFC placement problem with certain end-to-end latency needs that need to be satisfied [9], [19], and [20]. In [9], the authors present a delay-aware SFC placement problem such that VxVs forming SFCs are placed so as to satisfy end-to-end latency demands while utilizing network resources in an effective manner. A joint VxV placement and CPU allocation problem is studied in [19] and an optimization problem is formulated by employing a queuing-based model to minimize the ratio between the actual and the maximum allowed latency, for all SFC requests. The authors in [20] study the problem of VxV instantiation and migration with a goal of minimizing SFC delays. However, all these studies do not consider UE processing time, gNodeB processing time, and propagation or transmission time over the air interface. Besides, none of

the above studies consider heterogeneous MEC nodes, which increases the search space causing the SFC placement problem to grow cumbersome.

In contrast to the above SFC placement solutions, we consider the joint problem of user association and SFC placement which allows the optimization of end-to-end latency according to user locations, SFC latency and data rate requirements, and computing/networking resource availabilities. Furthermore, our proposed latency model stands out from the existing delay models within the context of 5G mobile network.

III. MACHINE LEARNING-DRIVEN PROACTIVE 'UPF' AUTO-SCALING

In this section, we create an ML classifier model that can identify and exploit hidden patterns in network traffic load instances to predict UPF scaling decisions ahead of time. In particular, we illustrate on the different steps involved in creating our model and eventually evaluate it based on several performance metrics [21].

A. Problem Description

We investigate how to map traffic load statistics X to VNF scaling decisions Y using supervised learning, which involves learning from a training set of data. The traffic load statistics X include measurements from a real-operator experimental 5G mobile network. The VNF scaling decisions Y refer to the required number of UPFs to process incoming traffic without violating UEs QoS Service Level Agreement (SLA). The details on the composition of X and Y are discussed in Section IV-D.

The X and Y metrics evolve over time, influenced mainly by the mobile network traffic dynamics and the actual number of mobile users. Consequently, the combined evolution of X and Y metrics is modeled as a time series $\{(x_t, y_t)\}$. Our goal is to determine the distribution of scaling decision metric Y constrained on knowing the traffic load metric $x \in X$.

Employing the statistical learning framework, X and Y are modeled as random variables. We assume that each sample (x_t, y_t) in the training set is obtained from the joint probability distribution of (X, Y) . Further, we assume that x_t is multi-dimensional and y_t is one-dimensional (univariate). In this formalism, the inference problem consists of finding a model $F : x \rightarrow P(Y|x)$ for $x \in X$, so as to maximize the likelihood function $L(\{P(y_t|x_t)\})$, which can be attained by minimizing the loss function $E = -\log(L)$ [22].

In this work, a neural-network called Multilayer Perceptron (MLP) is used to estimate the parameters of the model to predict the probability distribution $P(Y|x)$. We select neural-network in our approach for two reasons:

- (i) it has proven its potential in identifying traffic patterns due to its effectiveness in predicting time-series problems, whether periodic or not [23].
- (ii) it can build new customized features through hidden layers and fit nonlinear activation functions when a specific mathematical definition is not available.

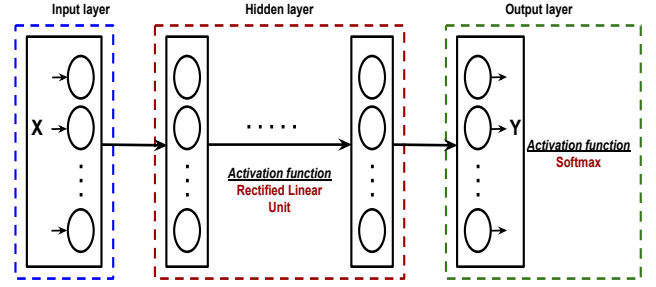


Fig. 2: Structure of the proposed MLP Classifier.

B. Multilayer Perceptron (MLP)

An MLP is a class of feed-forward artificial neural network, consisting of at least three layers of nodes (neurons): an input layer, one or more hidden layers, and an output layer, as shown in Fig. ???. These nodes are fully interconnected in the form of a directed graph, starting from the input to the output. All nodes except the input nodes have an associated activation function, which is used to compute the node output based on the weighted inputs from other nodes. Usually, the relu activation function is used for the hidden layer nodes, and softmax activation function is used for the output layer nodes [24]. The output is a vector containing the probabilities that sample $x \in X$ belongs to each class, which is equivalent to a categorical probability distribution. The final result is the class with the highest probability.

With a categorical cross-entropy loss function, the network parameters are chosen to minimize the following:

$$E = - \sum_{l=1}^C b_{x,l} \log(p_{x,l}) \quad (1)$$

Where C is the number of classes, b is the binary indicator (0 or 1) whether class label l is the correct classification for input x , and p is the predicted probability that input x belongs to class l . Here, a separate loss is calculated for each class label per input, and the result is the sum of all those losses.

An MLP model is trained through a backpropagation mechanism using gradient-descent as an optimization algorithm, where the weights between the nodes are adjusted iteratively for minimizing the error function.

C. Modeling MLP in Keras

Keras is an open-source neural-network Python library capable of running on top of Theano [25] or TensorFlow [26]. It is characterized by a clean, uniform, and streamlined high-level API, allowing users to rapidly define, train, and evaluate neural network models [27].

In Keras, the structure of the neural network model can be defined in a modular way, as a sequence of standalone and fully configurable modules, which can be readily plugged together. Keras offers several predefined neural layers such as a dense layer, a recurrent layer, and a convolutional layer. A wide range of activation functions is also available including relu, sigmoid, softmax, tanh, to name a few. Similarly, many predefined loss functions (e.g., mean squared error,

cross entropy) and regularization schemes (e.g., dropout) are supported. Also, since Keras performs backpropagation automatically, users do not need to implement it. Moreover, numerous approaches are available to partition the dataset into training, validation, and test sets.

To implement an MLP in Keras, we construct a sequential model with a number of predefined dense layers and their corresponding activation functions. We then configure the learning process of the model by choosing an optimizer, a loss function (equation 1), and a list of metrics to be reported. Lastly, the model is trained with an objective to minimize the loss function and then evaluated.

D. Collecting Data and Feature Engineering

The different steps that we followed in creating our MLP model are as follows:

1) *Data Collection*: The dataset utilized in this work is generated from a real-operator by monitoring the mobile network traffic load on 6 base stations, with each base station having 10 cells, for a period of 8 consecutive days. The traces in the dataset are in the form of a time series $\{(x_t, y_t)\}$ and we interpret this time series as a set of samples $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$. The traces are collected on an hourly timescale.

2) *Feature Extraction and Class definition*: We now describe the input feature sets $X_{default}$ and $X_{constructed}$, which when combined is referred as X , as well as the output classes/predictions Y .

The $X_{default}$ feature set includes 8 numeric features that are already available in the dataset as described in Table I. In addition to these default features, we construct 9 numeric features ($X_{constructed}$) from the basic dataset, as shown in Table II, using a process called *feature transformation*. These constructed features contain information or patterns on how the traffic load evolves over time, therefore assisting in proactive VNF scaling decisions.

Furthermore, it is necessary to define the desired output classes of the proposed MLP classifier model. Here, class refers to the number of UPF instances required per cell at time t , such that the auto-scaling decision allocates enough resources to meet the traffic demand until next auto-scaling decision at time $t + 1$, which is defined in equation ??,

$$No. \text{ of VNFs } (Y) = \min(vnf_{max}, \max(\frac{\lambda(t)}{\gamma}, \frac{\lambda(t+1)}{\gamma})) \quad (2)$$

Where $\lambda(t)$ and $\lambda(t + 1)$ are the traffic load in a cell at time t and $t + 1$, respectively, γ is the maximum traffic load a single UPF can handle, and vnf_{max} is the maximum number of UPFs per cell that can be hosted at the adjacent MEC node. Since there is a maximum limit on the hosting of VNFs, we model this as a classification problem rather than a regression problem. Note, we perform auto-scaling decisions once every hour, since our traffic traces are collected on hourly time intervals. However, our model is generic enough to handle lower interval granularities.

Default features ($X_{default}$)
1. gNodeB ID.
2. Date.
3. Time-stamp t .
4. Average number of users between t and $t - 1$ in each cell.
5. Maximum number of users between t and $t - 1$ in each cell.
6. Average downlink user throughput in each cell.
7. Average uplink user throughput in each cell.
8. Traffic load measured in each cell at time t , given by $\lambda(t)$.

TABLE I: Default set of features available in the dataset.

Constructed features ($X_{constructed}$)
9. Traffic load measured in each cell at time $t - 2$, given by $\lambda(t - 2)$.
10. Traffic load measured in each cell at time $t - 1$, given by $\lambda(t - 1)$.
11. Traffic load measured in each cell at time $t + 1$, given by $\lambda(t + 1)$.
12. Traffic load measured in each cell at time $t + 2$, given by $\lambda(t + 2)$.
13. Change in traffic load in each cell from time $t - 2$ to $t - 1$.
14. Change in traffic load in each cell from time $t - 1$ to t .
15. Change in traffic load in each cell from time t to $t + 1$.
16. Change in traffic load in each cell from time $t + 1$ to $t + 2$.
17. Weekday or weekend.

TABLE II: Constructed set of features from the dataset.

3) *Feature Subset Selection*: This process eliminates the redundant features from $X_{default}$ and $X_{constructed}$ feature sets to reduce the dimensionality of the data and also to reduce computational overhead. Therefore, to understand the impact of different features on our ML classifier model, we use *Principal Component Analysis (PCA)* and *Recursive Feature Elimination estimator*. Based on the ranking of these features, we use only 12 features (eliminating 2, 4, 5, 6 and 7 from Table I) that provides the best accuracy for our model.

4) *Dataset Decomposition*: Once data is collected and features extracted, the dataset is decomposed into training and test datasets. We use a rule-of-thumb decomposition conforming to 75%/25% between the training and test datasets, respectively. During the training phase, the MLP classifier model learns the relationship between the features and the classes.

E. Classification using neural networks

Finding the parameters of a neural-network model means searching for the best hyper-parameters of the MLP that can make the best predictions on the input. We applied *grid search* and *baby-sitting* as search strategies to perform an extensive search on the space of hyper-parameters to find the most accurate neural-network classifier. This process included finding the number of hidden layers and nodes, the batch size, the regularization parameter, the learning rate of the optimizer, and the number of epochs. We encountered the process of finding hyper-parameters time-consuming and hard, which assures that this topic still requires significant research.

We eventually found the architecture of the neural network that performs best on our traffic load traces and is described as follows. The structure includes one input layer with 12 nodes, three hidden layers with 12, 24 and 12 nodes, respectively, and an output layer with 10 nodes. The regularization parameter used is 0.01, the optimizer is based on stochastic gradient

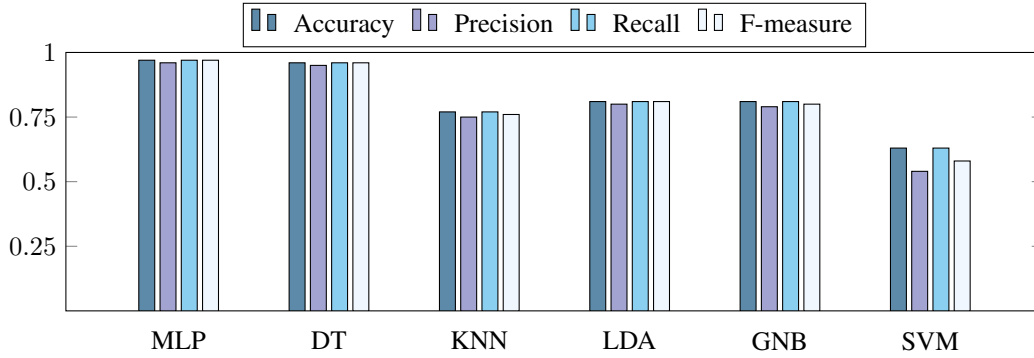


Fig. 3: Performance comparison of different classification algorithms for VNF auto-scaling.

approach with a constant learning rate of 0.001, the batch size is fixed to 100, and the number of epochs equals 300.

F. MLP model evaluation

We consider that MEC nodes in proximity to the gNodeBs are capable of hosting UPFs on their NFV infrastructure. We assume the link bandwidth capacity to be 20 Gbps and each VNF can process a maximum of 200 Mbps traffic without QoS degradation. We consider horizontal VNF auto-scaling with each MEC node capable of hosting 100 (20Gbps/200Mbps) VNFs and $vnf_{max} = 10$, i.e., a maximum of 10 VNFs can be hosted per cell. These assumptions are derived based on the evaluations performed by authors in [28]. If traffic load increases, additional VNF instances are deployed to meet QoS requirements, whereas if traffic load decreases, VNF instances are removed to save operational expenses.

Once the MLP model is created as discussed before, a test dataset is used to assess the performance of the model in predicting outcomes. The test outcomes can be classified into four groups: True Positive (TP) and True Negative (TN) are when the model correctly predicts actual positive and negative instances, respectively. Whereas, False Positive (FP) and False Negative (FN) are when the model makes incorrect predictions for negative and positive actual instances, respectively. Therefore, we consider four performance metrics to evaluate our MLP model: accuracy, precision, recall, and f-measure, as given by equations 5, 6, 7 and 8, respectively.

$$Accuracy = \frac{1}{|C|} \sum_{i=1}^{|C|} \frac{TP_i + TN_i}{TP_i + TN_i + FP_i + FN_i} \quad (3)$$

$$Precision = \frac{1}{|C|} \sum_{i=1}^{|C|} \frac{TP_i}{TP_i + FP_i} \quad (4)$$

$$Recall = \frac{1}{|C|} \sum_{i=1}^{|C|} \frac{TP_i}{TP_i + FN_i} \quad (5)$$

$$F_{measure} = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (6)$$

where C is the number of classes in the MLP model.

Accuracy is the most intuitive performance measure that gives the proportion of true predictions among the total

Class	1	2	3	4	5	6	7	8	9	10
1	1014	1	0	0	0	0	0	0	0	0
2	7	505	8	0	0	0	0	0	0	0
3	0	2	499	0	0	0	0	0	0	0
4	0	0	5	295	2	0	0	0	0	0
5	0	0	0	2	220	8	0	0	0	0
6	0	0	0	0	2	118	5	0	0	0
7	0	0	0	0	0	0	79	0	0	0
8	0	0	0	0	0	0	5	51	1	0
9	0	0	0	0	0	0	0	1	35	0
10	0	0	0	0	0	0	0	0	0	15

TABLE III: Confusion matrix for the proposed MLP classifier model.

number of predictions observed. However, accuracy is an excellent measure only if the datasets are entirely symmetric, i.e., false positives and false negatives are almost the same. Therefore, other performance metrics need to be considered when evaluating a model. *Precision* is a measure of correctly predicted positive observations to the total predicted positive observations. It is a good measure to determine when the cost of FP is high. In the case of VNF auto-scaling, a high number of FPs results in over-provisioning of resources leading to increased operational costs. On the other hand, *Recall* is a measure that calculates how many of the actual positives are captured in our model by labeling it as positive. It is a good measure to determine when the cost of FN is high. In the case of VNF auto-scaling, a high number of FNs results in under-provisioning of resources leading to QoS degradation. Finally, *F-measure* is the weighted average of precision and recall, and it is used when there is an uneven class distribution.

Fig. 5 compares the performance of the proposed MLP classifier model implemented in Keras with other classification algorithms implemented in Scikit-learn [29] (because Keras supports only neural-networks), such as Decision Tree (DT), K-Nearest Neighbour (KNN), Linear Discriminant Analysis (LDA), Naive Bayes (NB), and Support Vector Machine (SVM). We use 6 days of data (i.e., 6 days * 6 gNodeBs * 10 cells * 24 hours = 8640 samples) for training and two days of data (i.e., 2880 samples) for testing. The MLP model outperforms other models in all measures, with 97% accuracy, 96% precision, 97% recall, and 97% f-measure. The closest to MLPs performance was DT with 96% accuracy, 95% precision, 96% recall, and 96% f-measure.

Table III reports the confusion matrix for the proposed

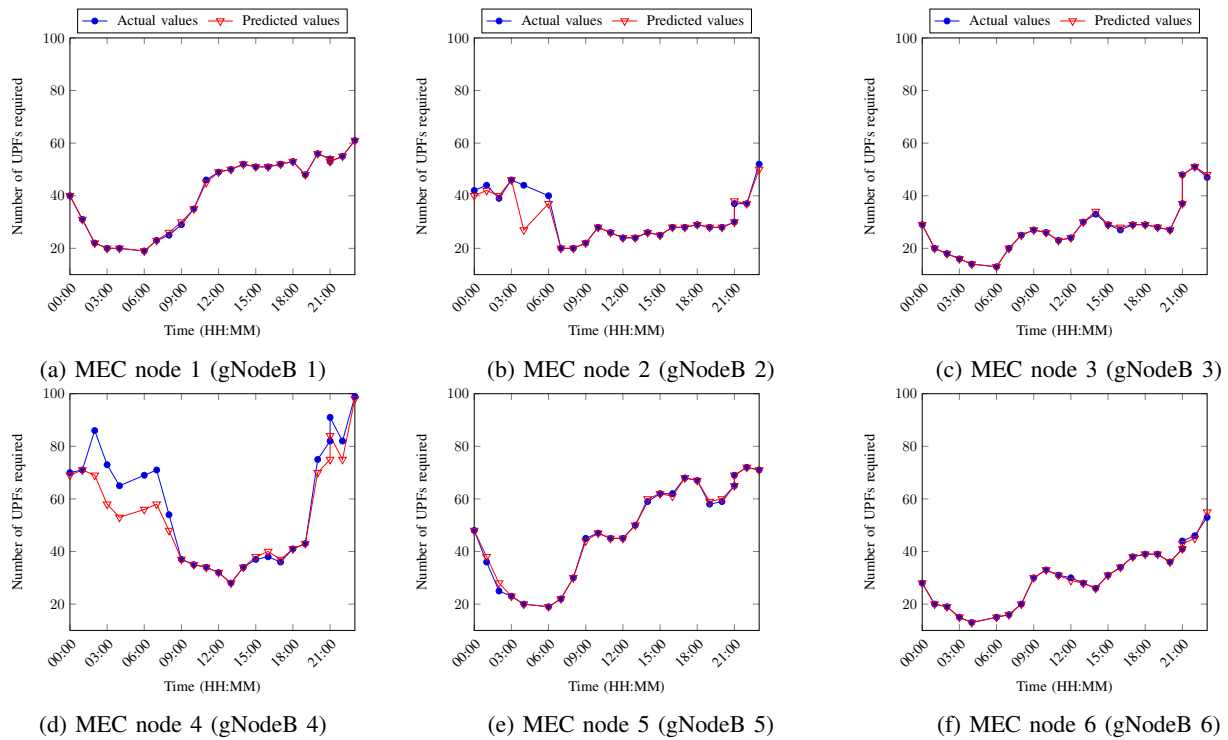


Fig. 4: Prediction results on the number of UPFs required at each MEC node based on the proposed MLP model.

Time	MAE (CI, 95%)	Time	MAE (CI, 95%)	Time	MAE (CI, 95%)
00:00	0.5 (-0.16 to 1.16)	08:00	0.33 (-0.07 to 0.74)	16:00	0.16 (-0.16 to 0.49)
01:00	0.5 (-0.16 to 1.16)	09:00	0 (0)	17:00	0 (0)
02:00	3.5 (-1.87 to 8.87)	10:00	0.16 (-0.16 to 0.49)	18:00	0.16 (-0.16 to 0.49)
03:00	2.5 (-2.4 to 7.4)	11:00	0.16 (-0.16 to 0.49)	19:00	1 (-0.6 to 2.6)
04:00	3.16 (-0.95 to 7.29)	12:00	0 (0)	20:00	1.16 (-1.12 to 3.45)
05:00	2.66 (-1.49 to 6.82)	13:00	0.33 (-0.07 to 0.74)	21:00	1.5 (-0.69 to 3.69)
06:00	2.16 (-2.08 to 6.41)	14:00	0.16 (-0.16 to 0.49)	22:00	1.33 (-0.91 to 3.57)
07:00	1.16 (-0.75 to 3.08)	15:00	0.66 (0.01 to 1.32)	23:00	1 (0.28 to 1.71)

TABLE IV: Mean Absolute Error with 95% Confidence Interval for the proposed MLP model.

MLP classifier model concerning the test data samples. For example, if we observe the 2nd row, on 7 instances, class 2 is misclassified as class 1, and on 8 instances, class 2 is misclassified as class 3.

Fig. 6 shows the prediction results of VNF auto-scaling (for 1 a full day) using MLP classifier model, where we display the prediction performance on all six MEC nodes, aggregated over all 10 cells for each gNodeB. In the figure, the blue line represents the actual output generated from the dataset, and the red line means the predicted VNF scaling decisions. As we can observe, the MLP classifier model introduced in this study can accurately follow the pattern of actual data, which point out the strong predicting capability of the model.

Table ?? presents the mean absolute error (MAE) between the actual and predicted values of VNF scaling decisions, calculated over all six MEC nodes for every hour during the entire day. We also calculate the confidence interval (CI) that determines the 95% likelihood on the range of classification errors that is expected from our model. As shown in the table, the model can have an upper limit of 8.87 MAE (at time

02:00) and a lower limit of -2.08 MAE (at time 06:00) in predicting VNF auto-scaling decisions. *It is worth mentioning that the predicted UPF auto-scaling decisions in our MLP model are used as an input to evaluate the SFC placement model presented in Section IV. However, in doing so, we assume that UEs can be associated and served by any of the cells of a candidate gNodeB, to simplify the problem.*

IV. LATENCY-OPTIMAL SFC PLACEMENT PROBLEM DESCRIPTION AND NETWORK MODEL

In this section, we first define the latency-optimal SFC placement problem and then describe the 5G mobile network model, SFC request model, and UE association, scheduling and delay model employed in formulating the ILP problem.

A. Problem Statement

Consider a 5G mobile network, composed of *six* gNBs (a non-split option of gNB is chosen where RRU, DU and CU are integrated together [30]), *two* aggregation points, and *one* 5GC, as depicted in Fig. 8. A set of three gNBs are interconnected to each other through *Xn*-interfaces. Using

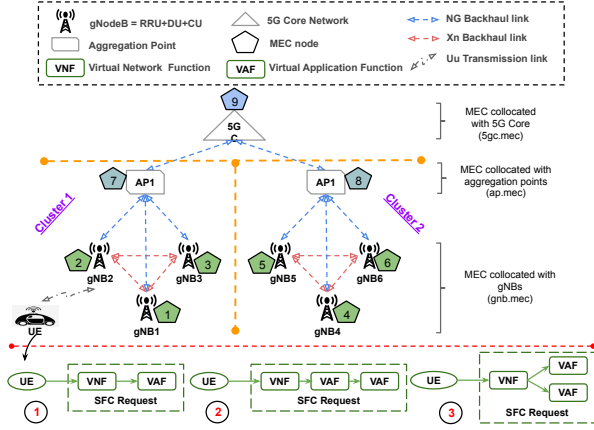


Fig. 5: Substrate network topology and SFC requests

NG-interfaces, the six gNBs are served by two aggregation points, and the 5GC serves both of these aggregation points. For simplicity, in the rest of this paper, we consider gNB1, gNB2, gNB3, and AP1 belong to cluster 1 while gNB4, gNB5, gNB6, and AP2 belong to cluster 2, as represented in Fig 8. Each element in our network topology is equipped with a resource-constrained (e.g., CPU) MEC node that is capable of hosting SFCs composed of one or several VxFs (e.g., VMs, containers). We consider three feasible options for physically deploying MEC nodes in 5G networks, as defined by ETSI [5], i.e., MEC collocated with gNB, MEC collocated with aggregation point, and MEC collocated with 5GC. Furthermore, in the considered hierarchical network topology, we assume that the closer is the MEC node to UE, the less is its computational capacity (e.g., MEC1, MEC2, and MEC3 are identical nodes with least capacity, MEC7 has the medium capacity, and MEC9 has the highest capacity).

Suppose the UE (e.g., autonomous car) is associated with gNB2 and requests for an SFC with an end-to-end latency (i.e., real-time, near real-time or non-real-time) and data rate requirements. The SFC requests considered in our work can be either of the three types as depicted in Fig. 8. Depending on the selected cost function to be minimized, the network provider can choose to place the VxFs of the SFC requested by the UE on either the host node (i.e., MEC2) or any neighboring nodes (i.e., MEC1, MEC3) or distant nodes (MEC7, MEC9) or cluster 2 nodes (i.e., MEC4, MEC5, MEC6, MEC8) by allocating sufficient network resources (e.g., CPU, backhaul bandwidth), efficiently, while also making sure that the end-to-end latency and data rate requirements of the requested SFC is always satisfied. In the first case, no additional delay is introduced in the backhaul since the MEC node collocated with the host gNB is the one hosting the VxFs. Conversely, in the other three cases, backhaul delay is introduced to map the virtual link onto a backhaul path, connecting the host gNB with a neighboring MEC node or a distant MEC node or a MEC node from a different cluster that is hosting VxFs. Formally, the problem of latency-optimal SFC placement is stated as follows:

Given: a small 5G mobile network with gNBs, aggregation

points, 5GC, MEC nodes, the scheduling capabilities of gNBs (e.g., PRBs, TTI duration, subcarrier spacing), the computational capacity of each MEC node, the transport network topology with the capacity of each backhaul link, the number of UEs and their requested SFCs with an end-to-end latency and data rate requirements.

Find: 'where' to allocate resources to VxFs and 'which' network paths to use.

Objective: minimize average end-to-end latency for UEs to access their SFCs in the mobile network.

B. 5G Mobile Network Model

The mobile network infrastructure is modeled as an undirected graph $G_{net} = (N_{net}, E_{net})$, where $N_{net} = N_{gnb.mec} \cup N_{ap.mec} \cup N_{5gc.mec}$ is the union of the set of $|N_{gnb.mec}|$ gNBs collocated with the MEC node, $|N_{ap.mec}|$ aggregation points collocated with the MEC node, and $|N_{5gc.mec}|$ 5GCs collocated with the MEC node and E_{net} is the set of backhaul links such that an edge $e^{mn} \in E_{net}$ only if a connection exists between $m, n \in N_{net}$. Each network node $m \in N_{net}$ is attributed with a weight $w_{cpu}^{net}(m)$, representing its CPU capacity, under the assumption that one VxF requires one CPU unit to be instantiated. Additionally, each network node $m \in N_{gnb.mec}$ is also associated with a weight $w_{prb}^{gnb.mec}(m)$ representing the number of Physical Resource Blocks (PRBs) available at each timeslot that can be scheduled to UEs for transmitting data packets. Furthermore, each edge $e^{mn} \in E_{net}$ is associated with a weight $w_{bw}^{net}(e^{mn})$ representing its bandwidth capacity (in Gbps). Finally, each network node $m \in N_{net}$ is associated with a geographical location $loc(m)$ (in terms of (x, y) coordinates) and each network node $m \in N_{gnb.mec}$ is associated with a coverage area $cov(m)$. Table VII summarizes all the parameters used in the mobile network model.

Notation	Definition
G_{net}	Graph of the mobile network.
N_{net}	Set of all network nodes in G_{net} .
$N_{gnb.mec}$	Set of gNBs collocated with the MEC node in G_{net} .
$N_{ap.mec}$	Set of aggregation points collocated with the MEC node in G_{net} .
$N_{5gc.mec}$	Set of 5GCs collocated with the MEC node in G_{net} .
E_{net}	Set of all backhaul links in G_{net} .
$w_{cpu}^{net}(m)$	Computing capacity of the network node $m \in N_{net}$.
$w_{prb}^{gnb.mec}(m)$	PRBs available for each timeslot at gNB $m \in N_{gnb.mec}$.
$w_{bw}^{net}(e^{mn})$	Bandwidth capacity of the backhaul link $e^{mn} \in E_{net}$.
$loc(m)$	Geographical location of the network node $m \in N_{net}$.
$cov(m)$	Coverage area of the gNodeB $m \in N_{gnb.mec}$.

TABLE V: Parameters in the mobile network model.

C. Service Function Chain Request Model

Let $G_{req} = (N_{req}, E_{req})$ be a directed graph modeling the SFC requests, where $N_{req} = N_{ue} \cup N_{sfc}$ is the union of the set of $|N_{ue}|$ UEs and $|N_{sfc}|$ the set of SFCs requested from the UEs and E_{req} is the set of virtual links between the UEs and their requested SFCs. Each SFC $s \in N_{sfc}$ is composed of a UPF (i.e., for encapsulation and decapsulation

of GPRS Tunnelling Protocol for the user plane (GTP-U) of an UE requesting MEC services [5]) and one or more VMAFs from a set of N_{vnfs} . Each SFC $s \in N_{sfc}$ is characterized by a maximum acceptable end-to-end latency (e.g., real-time, near real-time, non real-time) represented by $D_{E2E,max}(u, s)$ and a minimum guaranteed data rate denoted by $Thr_{req}(u, s)$ that needs to be satisfied. Each UE $u \in N_{ue}$ is associated with a location $loc(u)$ (in terms of (x, y) coordinates). Table VIII summarizes the parameters used in the SFC request model.

Notation	Definition
G_{req}	Graph of the SFC association request.
N_{req}	Set of all UEs and their SFC requests in G_{req} .
N_{ue}	Set of UEs in G_{req} .
N_{sfc}	Set of all the SFCs in G_{req} .
N_{vnfs}	Set of all the VxFs available to compose an SFC.
E_{req}	Set of all virtual links in G_{req} .
$D_{E2E,max}(u, s)$	Maximum acceptable end-to-end latency for a UE $u \in N_{ue}$ on its requested service $s \in N_{sfc}$.
$Thr_{req}(u, s)$	Requested data rate for a UE $u \in N_{ue}$ on the requested service $s \in N_{sfc}$.
$loc(u)$	Geographical location of the UE $u \in N_{ue}$.

TABLE VI: Parameters in the SFC request model.

D. UE Association, Scheduling and Delay Model

In contrast to 4G technology where the goal is only to enhance the throughput of Mobile Broadband (MBB) services, 5G is expected to support low-latency applications with end-to-end latency constraints of 1 – 10ms and error rates of 10^{-3} to 10^{-5} (e.g., connected cars). For cellular communications, two types of latencies are defined in 3GPP: control-plane latency (C-plane) and user-plane latency (U-plane). The C-plane latency is the transition time for the UE to switch from idle mode to connected mode including the establishment of the user plane while U-plane latency is the one-way delay required to transmit a data packet from the UE to the mobile network (uplink) or vice-versa (downlink) [31].

In this work, we consider only the U-plane latency for computing end-to-end delay (D_{E2E}), since it is the major contributor that is hindering the support of URLLC applications. D_{E2E} is computed from the time UEs start transmitting packets in uplink until the time they start being processed in MEC nodes. For a scheduled UE, we assume we have 3 different communication delays contributing to D_{E2E} :

(i) **Radio delay** ($D_{radio}^{ue,gnb}$) is the sum of UE processing delay (t_{ue}^{proc}), over-the-air transmission delay (t_{TTI}), gNodeB processing delay (t_{gnb}^{proc}), scheduler queuing delay (t_q), and HARQ retransmission delay, which is given by equation 13,

$$D_{radio}^{ue,air,gnb} = t_{ue}^{proc} + t_{TTI} + t_{gnb}^{proc} + t_q + 2 \cdot n_{harq} (t_{ue}^{proc} + t_{TTI} + t_{gnb}^{proc} + t_q) \quad (7)$$

where n_{harq} is the number of HARQ retransmissions required to achieve a BLER target of 10^{-3} to 10^{-5} . Similar to 3GPP, we adopt Orthogonal Frequency Division Multiplexing (OFDM) scheme. To satisfy the latency requirements of URLLC, 3GPP also proposes new frame structures with

shorter TTI durations and multiple subcarrier spacings. Scaling up the base subcarrier spacing of $15kHz$ by 2^μ (e.g., $30kHz$, $60kHz$, and $120kHz$), the TTI duration of $1ms$ is scaled down by 2^μ (e.g., $0.5ms$, $0.25ms$, and $0.125ms$), where $\mu = \{1, 2, \dots, n\}$, enabling faster transmission and lower processing time [32]. In our model, we adopt a TTI duration of $0.25ms$ resulting in a subcarrier spacing of $60kHz$ for all URLLC UEs. The resulting t_{ue}^{proc} and t_{gnb}^{proc} processing delays are 3 OFDM symbols and 1 TTI, respectively, as measured in [33]. The scheduler queuing delay (t_q), as represented in equation 14, is the sum of offset time (t_{offset}) i.e., the waiting time (~ 0 to 1 TTI) once the packet is ready for transmission until the beginning of the next TTI and the packet congestion time (t_{pktcon}) i.e., if the scheduler does not have enough PRBs to schedule a requested SFC packet in one TTI, the SFC packets may remain in the gNodeB buffer for longer duration.

$$t_q = t_{offset} + t_{pktcon} \quad (8)$$

To determine t_{pktcon} , we first need to determine the number of PRBs ($N_{prb}(u, s, m)$) required for the SFC $s \in N_{sfc}$ of an UE $u \in N_{ue}$ to be assigned by its associated gNodeB $m \in N_{gnb.mec}$. Given the data rate demand of the SFC s , the number of PRBs ($N_{prb}(u, s, m)$) is computed according to equation 15 as given in 3GPP [34]:

$$N_{prb}(u, s, m) = \frac{Thr_{req}(u, s) * T_s^\mu}{12 * 10^{-6} * N_{cc} * N_{mimo} * N_{mod} * sf * R * (1 - oh)} \quad (9)$$

where, N_{cc} is the number of aggregated component carriers, N_{mimo} is the number of MIMO layers, N_{mod} is the modulation order (e.g., 2 for QPSK, 4 for 16QAM, 6 for 64QAM, 8 for 256QAM), sf is the scaling factor, R is the code rate, oh is the overhead for control channels, and $T_s^\mu = 10^{-3} / (14 * 2^\mu)$ is the average OFDM symbol duration in a subframe for numerology μ ($\mu = 2$ in our case) assuming normal cyclic prefix. Except for N_{mod} and R , which are determined as per the below three steps, all other parameters are predefined according to the radio access capabilities:

SINR measurement: The UE $u \in N_{ue}$ measures the SINR value for a reference signal coming from its associated gNodeB $m \in N_{gnb.mec}$ using equation 16,

$$sindr(u, m) = \frac{\frac{F_m}{|loc(u) - loc(m)|^\alpha}}{\sum_{m' \neq m} \frac{F_{m'}}{|loc(u) - loc(m')|^\alpha} + N} \quad (10)$$

where $|loc(u) - loc(m)|$ is the distance between the UE u and its associated gNodeB m , $|loc(u) - loc(m')|$ is the distance between the UE u and the neighbouring gNodeBs of m (m'), α is a path loss exponent between 2 and 6, F_m and $F_{m'}$ are fading random variables of some distribution, and N is a constant noise term [35].

CQI report: The UE $u \in N_{ue}$ maps the SINR value measured in step 1 to a CQI index from the mapping table in [36], which is expected to be reported to the scheduler of its associated gNodeB $m \in N_{gnb.mec}$. It is to be noted that these mappings are not defined in 3GPP but are vendor specific.

CQI to MCS mapping: The scheduler is now expected to map the reported CQI index to an MCS index and determine

the best combination of modulation order (N_{mod}) and code rate (R) to be used from the mapping table in [37], resulting in a BLER target of 10^{-5} .

Therefore, we have all the necessary parameters in equation 15 to determine the number of PRBs that must be assigned for an SFC $s \in N_{sfc}$ requested by the UE $u \in N_{ue}$ to meet its data rate requirements. If the number of PRBs that needs to be assigned are not available in a particular TTI, they will be assigned during the next TTI and so forth, adding up to the total t_q latency (i.e., $t_{pktcon} = \text{no. of TTIs to schedule SFC packets} * \text{TTI duration}$).

(ii) **Backhaul delay** ($D_{bh}^{X_n, NG}$) is the sum of X_n propagation delay ($t_{X_n}^{prop}$), X_n transmission delay ($t_{X_n}^{tx}$), NG propagation delay (t_{NG}^{prop}), and NG transmission delay (t_{NG}^{tx}), given by equation 17,

$$D_{bh}^{X_n, NG} = t_{X_n}^{prop} + t_{X_n}^{tx} + t_{NG}^{prop} + t_{NG}^{tx} \quad (11)$$

where for a link $e^{mn} \in E_{net}$, $t_{X_n}^{prop}$ refers to the propagation time required to transmit SFC packets from node $m \in N_{gnb.mec}$ to node $n \in N_{gnb.mec}$ while t_{NG}^{prop} refers to the propagation time required to transmit SFC packets from node $m \in N_{gnb.mec}$ to node $n \in N_{ap.mec} | N_{5gc.mec}$. Similarly, $t_{X_n}^{tx}$ and t_{NG}^{tx} refers to the transmission time required to transfer SFC packets from node $m \in N_{gnb.mec}$ and node $m \in N_{ap.mec} | N_{5gc.mec}$, to the outgoing link e^{mn} , respectively.

(iii) **SFC processing delay** (D_{mec}^{sfc}) is the time required for all VxFs in an SFC $s \in N_{sfc}$ to apply a specific network operation on the arriving packets.

Therefore, D_{E2E} is computed according to equation 18.

$$D_{E2E} = D_{radio}^{ue, air, gnb} + D_{bh}^{X_n, NG} + D_{mec}^{sfc} \quad (12)$$

It is to be noted that the same delay model can be used for both downlink and uplink direction.

V. PROBLEM FORMULATION

Once, a batch of UE associations and its SFC requests arrive at the substrate network, it is either approved and embedded onto the network or it is denied. The embedding process includes both node and link mapping and is generally referred to as virtual network embedding problem which is proven to be NP-hard [38]. In the node mapping stage, each virtual node (i.e., UEs, VxFs in the SFCs requested by UEs) is mapped to a substrate node (i.e., gNodeBs, MEC nodes) while in the link mapping stage, each virtual link (i.e., the link between the UE and its requested SFC) is mapped to a single substrate path (i.e. the path between the gNodeB hosting the UE and MEC nodes hosting the VxFs in the SFC). In both stages, the constraints imposed on substrate nodes and substrate links must be satisfied.

A. Integer Linear Programming

The proposed joint UE association and SFC placement problem is formulated employing ILP techniques. Before starting the actual problem formulation, for each UE $u \in N_{ue}$, we first determine the set of candidate gNodeBs ($\bar{N}_{gnb.mec}(u)$) using equation 19,

$$\bar{N}_{gnb.mec}(u) = \{m \in N_{gnb.mec} | (|loc(u) - loc(m)|) \leq cov(m)\} \quad (13)$$

Then, we find neighboring gNodeBs for each gNodeB $m \in N_{gnb.mec}$ and neighboring MEC nodes for each MEC node $m \in N_{net}$ using equations 20 and 21, respectively.

$$nbr_gnbs(m) = \{m' \in N_{gnb.mec} | e^{m, m'} \in E_{net}\} \quad (14)$$

$$nbr_nodes(m) = \{m' \in N_{gnb.mec}, m'' \in N_{ap.mec}, m''' \in N_{5gc.mec} | e^{m, m'}, e^{m, m''}, e^{m'', m'''} \in E_{net}\} \quad (15)$$

Next, we find the candidate MEC nodes that can host VxFs of SFC requested by UEs. For each VxF v of SFC s from the UE u , the set of candidate MEC Nodes ($\bar{N}_{net}(u, v, s)$) can be defined according to equation 22.

$$\bar{N}_{net}(u, v, s) = \{m \in \bar{N}_{gnb.mec}(u), m' \in nbr_gnbs(m), m'' \in N_{ap.mec}, m''' \in N_{5gc.mec} | e^{m, m'}, e^{m, m''}, e^{m'', m'''} \in E_{net}\} \quad (16)$$

Thus, in our ILP model, either the UEs candidate gNodeB MEC node, or the gNodeB MEC node connected to the candidate gNodeB MEC node, or the aggregation point MEC node connected to the candidate gNodeB MEC node, or the 5GC MEC node connected to the aggregation point MEC node serving the candidate gNodeB MEC node can host UEs SFC.

Now, we formulate the SFC placement problem with three binary decision variables, χ_m^u , $\Upsilon_m^{u, v, s}$, and $\Psi_{m, n}^{u, s}$, as represented in Table IX.

Notation	Definition
χ_m^u	To show if $u \in N_{ue}$ is associated to gNodeB $m \in N_{gnb.mec}$.
$\Upsilon_m^{u, v, s}$	To show if $v \in N_{vnfs}$ of $s \in N_{sfc}$ from $u \in N_{ue}$ is assigned to $m \in N_{net}$.
$\Psi_{m, n}^{u, s}$	To show if virtual link between $u \in N_{ue}$ and $s \in N_{sfc}$ is assigned to substrate link between $m \in N_{net}$ and $n \in nbr_nodes(m)$.

TABLE VII: Binary decision variables.

The objective function of the ILP, given in equation 23, is to minimize the overall end-to-end latency from all users to their respective SFCs.

$$ILP : \min[\sum_{u \in N_{ue}} \sum_{m \in N_{gnb.mec}} \chi_m^u * D_{radio}^{ue, air, gnb}(u, m) + \sum_{u \in N_{ue}} \sum_{v \in N_{vnfs}} \sum_{s \in N_{sfc}} \sum_{m \in N_{net}} \Upsilon_m^{u, v, s} * D_{mec}^{sfc}(u, v, s, m) + \sum_{u \in N_{ue}} \sum_{s \in N_{sfc}} \sum_{m \in N_{net}} \sum_{n \in nbr_nodes(m)} \Psi_{m, n}^{u, s} * D_{bh}^{X_n, NG}(u, s, m, n)] \quad (17)$$

In equation 23, $D_{radio}^{ue, air, gnb}(u, m)$ depends on the number of UEs that need to be scheduled in a given time slot by

gNodeB m . For each UE u , we first find the $sinr(u, m)$ from equation 16 and then calculate the needed $N_{prb}(u, s, m)$ from equation 15 to transmit packets of SFC s with a particular size at a requested data rate ($Thr_{req}(s)$). If the total required PRBs exceed the maximum available PRBs in the gNodeB, some UEs are scheduled in the next time slot, thus increasing the Radio delay for those UEs. Since the backhaul links, Xn and NG , in the mobile network, $D_{bh}^{Xn,NG}(u, s, m, n)$ depends on the the number of UEs sharing the same backhaul link.

We will now describe all node and link constraints imposed in our problem formulation. Constraint (24) ensures that each UE is associated to only one gNodeB from its candidate set.

$$\sum_{m \in \bar{N}_{gnb.mec}(u)} \chi_m^u = 1, \forall u \in N_{ue} \quad (18)$$

Constraint (25) guarantees that each VxF of SFC requested from each UE is hosted by only one substrate MEC node from its candidate set.

$$\sum_{m \in \bar{N}_{net}(u,s)} \Upsilon_m^{u,v,s} = 1, \forall u \in N_{ue} \forall s \in N_{sfc}^u \forall v \in N_{vnfs}^s \quad (19)$$

Constraint (26) guarantees that each VxF is at most shared by vn_{fmax}^{shared} number of UEs.

$$\sum_{v \in N_{vnfs}^s} \Upsilon_m^{u,v,s} \leq vn_{fmax}^{shared}, \forall u \in N_{ue} \forall s \in N_{sfc}^u \quad (20)$$

$$\forall m \in \bar{N}_{net}(u, s)$$

Constraint (27) ensures that the amount of CPU resources allocated to VxFs of SFCs adheres to the available CPU capabilities on the substrate node.

$$\sum_{u \in N_{ue}} \sum_{s \in N_{sfc}^u} \sum_{v \in N_{vnfs}^s} \Upsilon_m^{u,v,s} \leq w_{cpu}^{net}(m), \forall m \in N_{net} \quad (21)$$

Constraint (28) makes sure that in each time slot gNodeBs can associate UEs only if they have enough PRBs to meet the data rate demand of the requested SFC by the UE.

$$\sum_{u \in N_{ue}} \sum_{s \in N_{sfc}^u} N_{prb}(u, s, m) * \chi_m^u \leq w_{prb}^{gnb.mec}(m), \quad (22)$$

$$\forall m \in N_{gnb.mec}$$

Flow constraint (29) enforces for each virtual link between UE $u \in N_{ue}$ and its SFC $s \in N_{sfc}$ there exists a continuous path established between the gNodeB to which the UE is associated and the MEC node hosting the Vxfs of SFC s .

$$\sum_{n \in nbr_nodes(m)} (\Psi_{n,m}^{u,s} - \Psi_{m,n}^{u,s}) = \Upsilon_m^{u,s} - \chi_m^u, \quad (23)$$

$$\forall m \in N_{net}, \forall e^{u,s} \in E_{req}$$

Constraint (30) makes sure that virtual links are mapped onto the backhaul substrate links in the mobile network, if

and only if it has enough bandwidth capacity to meet the link demand of virtual links.

$$\sum_{u \in N_{ue}} \sum_{s \in N_{sfc}} Thr_{req}(u, s) (\Psi_{n,m}^{u,s} + \Psi_{m,n}^{u,s}) \leq w_{bw}^{net}(e^{nm}), \quad (24)$$

$$\forall m \in N_{net}, \forall n \in nbr_nodes(m), n < m$$

Constraint (31) ensures that the end-to-end latency from the UEs to its associated SFCs does not exceed the maximum acceptable latency as requested by the UEs.

$$\sum_{m \in N_{gnb.mec}} \chi_m^u * D_{radio}^{ue,air,gnb}(u, m) + \sum_{m \in N_{net}} \sum_{v \in N_{vnfs}^s} \Upsilon_m^{u,v,s} * D_{mec}^{sfc}(u, v, s, m) + \sum_{m \in N_{net}} \sum_{n \in nbr_nodes(m)} \Psi_{m,n}^{u,s} * D_{bh}^{Xn,NG}(u, s, m, n) \leq D_{E2E,max}(u, s), \forall u \in N_{ue}, \forall s \in N_{sfc}^u \quad (25)$$

B. Heuristic

The above ILP formulation took 44 hours to associate 300 UEs including their latency-sensitive SFC requests composed of a number of Vxfs on a mobile network comprised of *six* gNodeBs, *two* aggregation points, and *one* 5G core. The ILP was solved using ILOG CPLEX solver on an Intel Core i7 laptop with 3GHz CPU and 16 GB RAM. To address the issue of scalability in ILP, we propose a heuristic algorithm, as seen in Algorithm 1, that can solve the above association/mapping problem in a couple of seconds. Similar to the ILP-based algorithm, the objective of our heuristic algorithm is to minimize the overall end-to-end latency from all UEs to their requested SFCs.

In the first step (lines 1–10), the algorithm loops through all the UEs to determine a set of candidate gNodeBs considering the location of the UE, location of the gNodeB, and the coverage area of the gNodeB, and then creates a list of $cand_gnb(u)$ for each UE. Next, each UE is mapped to the gNodeB, among the $cand_gnb(u)$, that measures the best signal quality (i.e., SINR) and also has sufficient PRBs to host the UE.

In the second step (lines 11 – 20), the algorithm finds the candidate MEC nodes for each VxF of SFC requests received from all UEs, which is nothing but the union of the MEC node collocated with the host gNodeB of UE (determined from step 1) and all other MEC nodes connected directly or indirectly to the host gNodeB of UE through backhaul links. Another list of $cand_mec(u, s, v)$ is created for each VxF of the SFC received from all UEs.

In the third step (lines 21 – 38), the algorithm begins mapping all Vxfs of SFC requests with real-time latency requirements considering $cand_mec(u, s)$ for each SFC, starting from $gnb.mec$ nodes. Once they run out of computing resources, the algorithm moves on to $ap.mec$ nodes, and finally on to $5gc.mec$ nodes, if and only if the computed end-to-end latency ($D_{E2E}(u, s, m)$) is less than the maximum acceptable end-to-end latency for that SFC ($D_{E2E,max}$). Moreover, if an instance of VxF is already mapped to the

Algorithm 1 Heuristic

Require: G_{net} , G_{req} , and SFC latency budget $[N_{sfc}(rt), N_{sfc}(near_rt), N_{sfc}(non_rt)]$.

Ensure: User association and latency-optimal SFC placement.

Step 1. Find candidate gNodeBs for each UE and perform UE association.

```
1: for  $u$  in  $N_{ue}$  do
2:    $cand\_gnb(u) \leftarrow 0$ 
3:    $map\_gnb(u) \leftarrow 0$ 
4:   for  $m$  in  $N_{gnb.mec}$  do
5:     if  $|loc(u) - loc(m)| \leq cov(m)$  then
6:        $cand\_gnb(u) \leftarrow m$ 
7:     end if
8:   end for
9:    $map\_gnb(u) \leftarrow m$  from the list of  $cand\_gnb(u)$  with
    $max(sinr(u, m))$  and enough PRBs available.
```

Step 2. Find candidate MEC nodes for VxFs of each SFC from each UE.

```
11: for  $u$  in  $N_{ue}$  do
12:   for  $s$  in  $N_{sfc}$  do
13:     for  $v$  in  $N_{sfc}(u)$  do
14:        $cand\_mec(u, s, v) \leftarrow 0$ 
15:       for  $m$  in  $neighbours(map\_gnb(u))$  do
16:          $cand\_mec(u, s, v) \leftarrow m$ 
17:       end for
18:     end for
19:   end for
```

Step 3. Perform SFC placement for each UE.

```
21: for  $u$  in  $N_{ue}$  do
22:   for  $s$  in  $N_{sfc}(rt)$  do /* real-time SFCs. */
23:     for  $v$  in  $N_{sfc}(u)$  do
24:        $map\_mec(u, s, v) \leftarrow 0$ 
25:       for  $m$  in  $cand\_mec(u, s, v)$  do
26:          $compute(D_{E2E}(u, s, m))$ 
27:         if  $D_{E2E}(u, s, m) \leq D_{E2E,max}$  then
28:           if  $inst(v)$  not in  $m$  or  $neighbours(m)$  then
29:              $map\_mec(u, s, v) \leftarrow m$ 
30:           end if
31:         end if
32:          $allocate\_continuous\_path(u, s, m)$ 
33:          $update\_node\_and\_link\_resources()$ 
34:       end for
35:     end for
36:   end for
37: end for
38: Repeat Step 3 for  $s$  in  $N_{sfc}(near\_rt)$  and  $N_{sfc}(non\_rt)$ .
```

candidate MEC Node the UE shares this VxF to realize its SFC instead of instantiating a new VxF. The heuristic then uses the shortest path algorithm to map the virtual link between the UE and its requested SFC onto the substrate link between the gNodeB that the UE is associated to and the MEC node that the SFC is being hosted on. The VxFs of a single SFC might be mapped on different MEC nodes, and therefore further caution is exercised during link mapping. The node and link computational resources are updated after each mapping. The same process is repeated for other SFC requests with near-real-time and non-real-time latency requirements until all SFC requests are mapped.

VI. ILP AND HEURISTIC EVALUATION

The performance of the latency-optimal SFC placement ILP model is evaluated based on the simulations implemented in

Python. We then compare it to the implemented heuristic algorithms performance. Real-operator network topology and realistic latency values are used when modeling the simulation environment to produce realistic simulation results, which can better illustrate the benefits of placing SFCs composed of VxFs at the network edges closer to the end-user.

A. Simulation Environment

A small experimental 5G network composed of 9 network nodes is considered in our simulation, as depicted in Fig. 8. A set of 3 gNodeBs are connected to each other through 20 Gbps Xn backhaul links, while each of the three gNodeBs is connected to the aggregation point using 20 Gbps NG backhaul links which in turn is connected to the 5G core network using 50 Gbps backhaul links. The number of aggregated component carriers is set to 4, and each carrier has a bandwidth capacity of 20 MHz. We assume that the gNodeBs support 4×4 MIMO configuration. We then introduce MEC nodes at each of these 9 network nodes capable of hosting a limited number of VxFs. The MEC nodes collocated with gNodeBs each have 50 CPUs, the MEC nodes collocated with aggregation points each have 100 CPUs, and the MEC node collocated with 5GC has 500 CPUs.

Our simulations are carried out for two scenarios. In the first scenario, we consider that SFC requests arrive in batches of 30 UEs (equally divided among real-time, near-real-time, and non-real-time) with each batch corresponding to 1 timeslot. In every timeslot, the ILP considers the SFC requests received in previous batches and associates all the UEs and their SFC requests onto the mobile network, considering the latency and data rate requirements of each SFC. We consider 10 batches of SFC requests corresponding to 300 UEs. In the second scenario, we consider that SFC requests arrive according to the *predicted number of UPF instances from our MLP neural-network model*, as illustrated in Section III. The performance of ILP is compared with our heuristic algorithm in both scenarios. Additionally, in both scenarios, we assume that each UPF instance corresponds to one UE and each SFC is composed of 2 or 3 VxFs (1 UPF and 1 or 2 VMAF as depicted in Fig. 8) with 1 CPU required to instantiate every VxF. Furthermore, each UPF is shared with 5 UEs while each VMAF is shared among 2 UEs. Since the UEs considered in our model are URLLC UEs, we assume three categories of user-to-SFC one-way delay requirements, i.e., $1ms$, $2ms$, and $5ms$. We assume that each UE is transmitting short packets of size $15Kb$ every TTI and requests a minimum data rate of $200Mbps$. In Section IV, we discussed on how we calculate $D_{Radio}^{ue,air,gnb}$. We calculate $D_{bh}^{Xn,NG}$ by dividing the total packet size generated from all the UEs that are using the same backhaul link with the bandwidth capacity of the link [39]. Finally, we calculate D_{mec}^{sfc} by dividing the packet size that the VxFs of the SFC should process by the CPU speed. We consider a CPU speed of $3GHz$ with 64 bit processor.

B. Simulation Results

CPU Utilization: The CPU utilization is computed by dividing the number of CPUs utilized in $gnb.mec$ nodes or $ap.mec$ nodes or $5gc.mec$ nodes once the VxFs are mapped

to the total number of CPUs available in all *gnb.mec* nodes or all *ap.mec* nodes or all *5gc.mec* nodes, respectively.

Fig. 9 illustrates the CPU utilization of MEC nodes with respect to the number of UEs for simulations carried out in scenario 1. We observe that up to ≈ 90 UEs, the ILP places most of the VxFs on *gnb.mec* nodes because of its proximity to UEs, irrespective of the SFC latency requirements, while some non-real-time VxFs that are shared by UEs associated with cluster 1 (gNB1, gNB2 or gNB3) and cluster 2 (gNB4, gNB5 or gNB6) are placed on *5gc.mec* nodes. Only after *gnb.mec* nodes are depleted with their CPU resources (after 90 UEs), the ILP starts moving VxFs with near-real-time and non-real-time latency requirements initially placed in *gnb.mec* nodes to *ap.mec* nodes and starts placing new SFCs with real-time latency requirements on *gnb.mec* nodes. Similarly, when CPU resources of *ap.mec* nodes are depleted (≈ 180 UEs) the ILP starts moving VxFs with non-real-time latency requirements initially placed in *gnb.mec* or *ap.mec* nodes to *5gc.mec* nodes. On the other hand, heuristic algorithm follows a similar pattern to that of ILP, but instead of placing non-real-time VxFs that are shared by UEs associated to cluster 1 and cluster 2 on *5gc.mec* nodes, those VxFs are initially placed on *gnb.mec* nodes, then on *ap.mec* nodes and finally on *5gc.mec* nodes (in this order depending on the MEC nodes resource availability). This is evident from Fig. 9, where CPU utilization of *gnb.mec* nodes is always higher in heuristic compared to that of ILP. However, this results in the increase of overall latency for heuristic due to the users taking a long path to access their SFC services (e.g., if a user is associated to gNB2 and its VxFs are placed in *gNB6.mec*, the path mapping could be gNB2 \rightarrow AP1 \rightarrow 5GC \rightarrow AP2 \rightarrow gNB6). Consequently, up to 180 UEs, the CPU utilization of *ap.mec* and *5gc.mec* nodes are lower in heuristic compared to ILP.

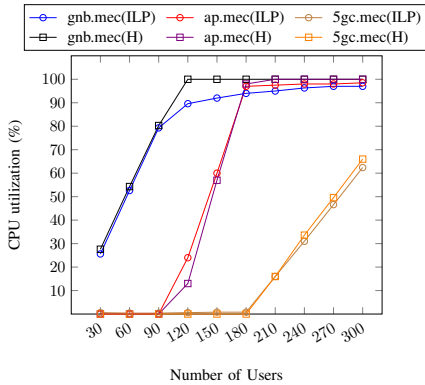


Fig. 6: CPU utilization of MEC nodes (Scenario 1).

Fig. 10 illustrates the CPU utilization of MEC nodes with respect to time over one full day for the network considered in the MLP neural-network model (scenario 2). We observe that the CPU utilization of *gnb.mec* nodes are always full, the *ap.mec* nodes are most of the time full except from 2:00 to 8:00 and *5gc.mec* nodes have low utilization during early morning (2:00 to 8:00) due to the low number of UEs being active and the utilization gradually increases during the day peaking late in the night ($\approx 22:00$). The heuristic follows a similar pattern to that of ILP, but as discussed earlier VxFs

shared by UEs belonging to different clusters are initially placed on *gnb.mec* nodes rather than on *5gc.mec* nodes like in ILP. Therefore, CPU utilization for *gnb.mec* nodes are always higher in heuristic compared to that of ILP, with a tradeoff being the increase in overall latency.

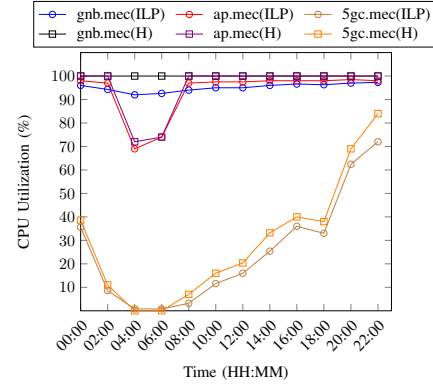


Fig. 7: CPU utilization of MEC nodes (Scenario 2).

Link utilization: Link utilization is calculated by dividing the usage of either X_n or NG backhaul links by UEs for utilizing SFCs in the MEC nodes to the total available capacity of the respective links.

Fig. 11 and Fig. 12 illustrates, respectively, the X_n link utilization and the NG link utilization as a function of the number of UEs for experiments carried out in scenario 1. In Fig. 11, we observe that in ILP, irrespective of the number of UEs, the X_n link utilization remains almost the same ($< 3\%$), which is attributed to the fact that ILP principally places the VxFs of UEs on the *gnb.mec* that is currently serving the corresponding UE over the air interface in order to minimize the end-to-end latency. However, we observe that heuristic algorithm places VxFs of some UEs on *gnb.mec* nodes that are currently not serving the corresponding UE over the air interface, which leads to the usage of X_n links. After a certain point (≈ 90 UEs in Fig. 11), the X_n link utilization remains almost the same for heuristic since the capacity of all *gnb.mec* nodes are depleted, and VxFs are placed on *ap.mec* or *5gc.mec* nodes there on. In Fig. 12, we can observe that both in ILP and heuristic NG links are least utilized up to ≈ 90 UEs since most VxFs of SFCs, irrespective of their latency demands, are always placed on *gnode.mec* nodes until then. Once *gnode.mec* nodes are out of CPU resources, the VxFs of SFCs are moved to *ap.mec* nodes and later to *5gc.mec* nodes considering the latency requirements of SFCs, resulting in the significant usage of NG backhaul links. However, the reason for higher NG link utilization in heuristic is attributed to the fact that some UEs take longer routes, from cluster 1 to cluster 2 or viceversa, in order to access their SFC which is not the case in ILP.

Fig. 13 and Fig. 14 illustrates, respectively, the X_n link utilization and the NG link utilization with respect to time over one full day based on the network considered in the MLP neural-network model (Scenario 2). Since the number of UEs is always more than 90, we observe that both ILP and heuristic algorithm places VxFs of some UEs on *gnb.mec* nodes that

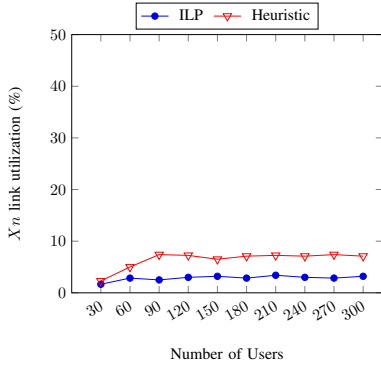


Fig. 8: Xn-link (gNB-to-gNB) utilization (Scenario 1).

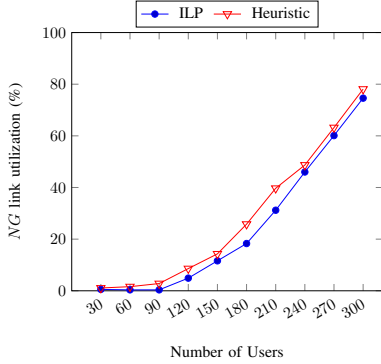


Fig. 9: NG-link (gNB-to-AP-to-5GC) utilization (Scenario 1).

are currently not serving the corresponding UE over the air interface which leads to the usage of X_n links as already explained in Scenario 1. Likewise, NG link utilization for both ILP and heuristic is lowest during early morning (2:00 to 8:00) due to the low number of UEs being active and the utilization gradually increases during the day peaking late in the night ($\approx 22:00$). However, both X_n and NG link utilizations in heuristic are higher compared to ILP because of the long path the UEs take to access SFC like we discussed before.

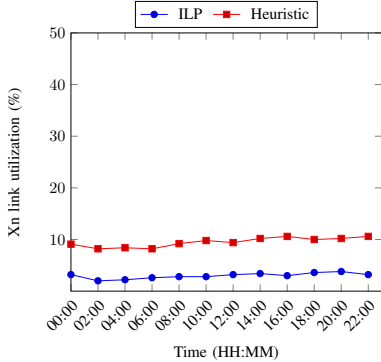


Fig. 10: Xn-link (gNB-to-gNB) utilization (Scenario 2).

Average end-to-end latency: Fig. 15 compares the average user-to-sfc end-to-end delay for ILP and heuristic for Scenario 2 experiments. Like we already discussed, UEs belonging to different clusters share some VxFs. The ILP produces an

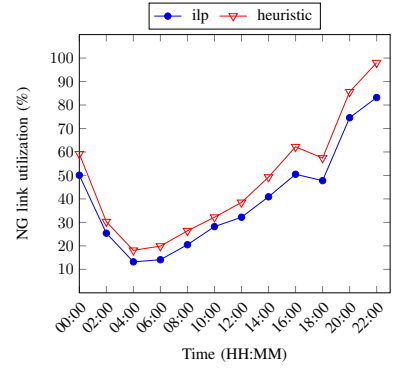


Fig. 11: NG-link (gNB-to-AP-to-5GC) utilization (Scenario 2).

optimal solution by placing such VxFs at $5gc.mec$ nodes to minimize the overall user-to-sfc delay, but the heuristic initially places such VxFs on $gnb.mec$ nodes, and therefore some UEs take the longer path (e.g., from cluster 1 to cluster 2) to access their SFC resulting in increased latency. Therefore, ILP performs better than heuristic in terms of average end-to-end delay between, as seen in Fig 15.

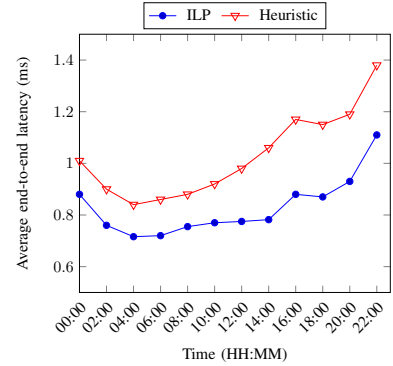


Fig. 12: Average D_{E2E} based on the predicted number of UPFs from the MLP classifier model (Scenario 2).

Execution time: The above ILP formulation took 44 hours to associate 300 UEs including their latency-sensitive SFC requests composed of a number of VxFs on a mobile network comprised of six gNodeBs, two aggregation points, and one 5G core. Therefore, we proposed a heuristic algorithm that performs a comparable association and mapping in a couple of seconds except with sub-optimal outcomes. Both ILP and heuristic were solved using CLOG IPLEX solver on an Intel Core i7 laptop with 3GHz CPU and 16 GB RAM.

VII. CONCLUSIONS

The first part of the paper aims at applying machine learning techniques to optimize network management operations. Towards this end, we proposed a neural-network model to facilitate proactive auto-scaling of VNFs, based on the traffic traces obtained from the commercial MNO. We evaluated the proposed model for its effectiveness in accurately predicting the amount of UPF instances required as a function of the network traffic it should process. Moreover, we compared the performance of the neural-network model with five other

classification learning methods, and the performance metrics show that the neural-network model can achieve higher prediction accuracy (97%) compared to other methods.

In the second part of the paper, we solve a joint UE association and SFC placement problem aiming to minimize the overall user-to-sfc end-to-end latency. We have seen that the ILP improves QoS of all UEs by initially placing their SFCs in MEC nodes closer to gNodeBs (*gnb.mec*) and thereby reducing *NG* backhaul link usage. Once the *gnb.mec* node CPU resources are depleted, near-real-time and non-real-time SFCs are moved/placed in MEC nodes closer to aggregation points and 5GC which results in increased usage of *Xn* and *NG* backhaul links. We evaluated the proposed model using simulations based on real-operator network topology and real-world latency values. Our results show that the average end-to-end latency reduces significantly when SFCs are placed at the MEC nodes according to their latency and data rate demands. Furthermore, we propose a heuristic algorithm to address the issue of scalability in ILP, that can solve the above association/mapping problem in seconds rather than hours.

ACKNOWLEDGEMENTS

This work has been performed in the framework of the European Unions Horizon 2020 project 5G-CARMEN co-funded by the EU under grant agreement No 825012. The views expressed are those of the authors and do not necessarily represent the project. The Commission is not liable for any use that may be made of any of the information contained therein.

REFERENCES

- [1] A. Gupta and R. K. Jha, "A survey of 5g network: Architecture and emerging technologies," *IEEE access*, vol. 3, pp. 1206–1232, 2015.
- [2] Q.-V. Pham, F. Fang, V. N. Ha, M. Le, Z. Ding, L. B. Le, and W.-J. Hwang, "A survey of multi-access edge computing in 5g and beyond: Fundamentals, technology integration, and state-of-the-art," *arXiv preprint arXiv:1906.08452*, 2019.
- [3] R. Mijumbi, J. Serrat, J.-L. Gorricho, N. Bouten, F. De Turck, and R. Boutaba, "Network function virtualization: State-of-the-art and research challenges," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 1, pp. 236–262, 2015.
- [4] I. F. Akyildiz, S. Nie, S.-C. Lin, and M. Chandrasekaran, "5g roadmap: 10 key enabling technologies," *Computer Networks*, vol. 106, pp. 17–48, 2016.
- [5] "MEC in 5G networks," ETSI, Whitepaper, Jun 2018.
- [6] "Network Functions Virtualization (NFV)," ETSI, Whitepaper, Feb 2017.
- [7] I. Farris, T. Taleb, H. Flinck, and A. Iera, "Providing ultra-short latency to user-centric 5g applications at the mobile network edge," *Transactions on Emerging Telecommunications Technologies*, vol. 29, no. 4, p. e3169, 2018.
- [8] E. Casalicchio and L. Silvestri, "Mechanisms for sla provisioning in cloud-based service providers," *Computer Networks*, vol. 57, no. 3, pp. 795–810, 2013.
- [9] A. Alleg, T. Ahmed, M. Mosbah, R. Riggio, and R. Boutaba, "Delay-aware vnf placement and chaining based on a flexible resource allocation approach," in *13th International Conference on Network and Service Management (CNSM)*. IEEE, 2017, pp. 1–7.
- [10] T. Subramanya and R. Riggio, "Machine learning-driven scaling and placement of virtual network functions at the network edges," in *5th International Conference on Network Softwarization (NetSoft)*, 2019, 2019.
- [11] S. Dutta, T. Taleb, and A. Ksentini, "Qoe-aware elasticity support in cloud-native 5g systems," in *IEEE International Conference on Communications (ICC)*. IEEE, 2016, pp. 1–6.
- [12] G. A. Carella, M. Pauls, L. Grebe, and T. Magedanz, "An extensible autoscaling engine (ae) for software-based network functions," in *2016 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*. IEEE, 2016, pp. 219–225.
- [13] M. M. Murthy, H. Sanjay, and J. Anand, "Threshold based auto scaling of virtual machines in cloud environment," in *IFIP International Conference on Network and Parallel Computing*. Springer, 2014, pp. 247–256.
- [14] C. H. T. Arteaga, F. Rissoi, and O. M. C. Rendon, "An adaptive scaling mechanism for managing performance variations in network functions virtualization: A case study in an nfv-based epc," in *2017 13th International Conference on Network and Service Management (CNSM)*. IEEE, 2017, pp. 1–7.
- [15] T. Lorigo-Botran, J. Miguel-Alonso, and J. A. Lozano, "A review of auto-scaling techniques for elastic applications in cloud environments," *Journal of grid computing*, vol. 12, no. 4, pp. 559–592, 2014.
- [16] A. Bilal, T. Tarik, A. Vajda, and B. Miloud, "Dynamic cloud resource scheduling in virtualized 5g mobile systems," in *2016 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2016, pp. 1–6.
- [17] R. Mijumbi, S. Hasija, S. Davy, A. Davy, B. Jennings, and R. Boutaba, "Topology-aware prediction of virtual network function resource requirements," *IEEE Transactions on Network and Service Management*, vol. 14, no. 1, pp. 106–120, 2017.
- [18] A. Mestres, A. Rodriguez-Natal, J. Carner, P. Barlet-Ros, E. Alarcón, M. Solé, V. Muntés-Mulero, D. Meyer, S. Barkai, M. J. Hibbett *et al.*, "Knowledge-defined networking," *ACM SIGCOMM Computer Communication Review*, vol. 47, no. 3, pp. 2–10, 2017.
- [19] S. Agarwal, F. Malandrino, C.-F. Chiasserini, and S. De, "Joint vnf placement and cpu allocation in 5g," in *IEEE INFOCOM 2018-IEEE Conference on Computer Communications*. IEEE, 2018, pp. 1943–1951.
- [20] H. Hawilo, M. Jammal, and A. Shami, "Orchestrating network function virtualization platform: Migration or re-instantiation?" in *2017 IEEE 6th International Conference on Cloud Networking (CloudNet)*. IEEE, 2017, pp. 1–6.
- [21] R. Boutaba, M. A. Salahuddin, N. Limam, S. Ayoubi, N. Shahriar, F. Estrada-Solano, and O. M. Caicedo, "A comprehensive survey on machine learning for networking: evolution, applications and research opportunities," *Journal of Internet Services and Applications*, vol. 9, no. 1, p. 16, 2018.
- [22] C. M. Bishop, *Pattern recognition and machine learning*. springer, 2006.
- [23] S. Basterrech, G. Rubino, and V. Snášel, "Sensitivity analysis of echo state networks for forecasting pseudo-periodic time series," in *2015 7th International Conference of Soft Computing and Pattern Recognition (SoCPar)*. IEEE, 2015, pp. 328–333.
- [24] A. K. Jain, J. Mao, and K. Mohiuddin, "Artificial neural networks: A tutorial," *Computer*, no. 3, pp. 31–44, 1996.
- [25] "Theano." [Online]. Available: <http://deeplearning.net/software/theano/>
- [26] "Tensorflow." [Online]. Available: <https://www.tensorflow.org/>
- [27] "Keras 2018." [Online]. Available: <https://keras.io/>
- [28] S. Rahman, T. Ahmed, M. Huynh, M. Tornatore, and B. Mukherjee, "Auto-scaling vnfs using machine learning to improve qos and reduce cost," in *2018 IEEE International Conference on Communications (ICC)*. IEEE, 2018, pp. 1–6.
- [29] "Scikit-learn." [Online]. Available: <https://scikit-learn.org/>
- [30] "Transport network support of IMT-2020/5G," ITU-T, Technical Report, Oct 2018.
- [31] "Feasibility study for Further Advancements for E-UTRA (LTE-Advanced)," 3GPP TR 36.912, Oct 2012.
- [32] "5G Frame Structure," Nomor, Whitepaper, Aug 2017.
- [33] "UP Latency in NR," 3GPP contribution R2-1711550, Oct 2017.
- [34] "5G; NR; User Equipment (UE) radio access capabilities," 3GPP TS 38.306 version 15.3.0 Release 15, Oct 2017.
- [35] F. Massimo and M. Ronald, *Random networks for communication: from statistical physics to information systems*. Cambridge University Press, 2008, vol. 24.
- [36] A. Othman, S. Y. Ameen, and H. Al-Rizzo, "A new channel quality indicator mapping scheme for high mobility applications in lte systems," *Journal of Modeling and Simulation of Antennas and Propagation*, vol. 1, no. 2, pp. 38–43, 2015.
- [37] "Physical layer procedures for data," 3GPP TS 38.214 version 15.3.0 Release 15, Oct 2018.
- [38] A. Schrijver, *Theory of linear and integer programming*. John Wiley & Sons, 1998.
- [39] D. Harutyunyan, S. Nashid, B. Raouf, and R. Riggio, "Latency-aware service function chain placement in 5g mobile networks," in *IEEE Conference on Network Softwarization (NetSoft 2019)*, 2019.