

# Latency-aware and Survivable Mapping of VNFs in 5G Network Edge Cloud

Prabhu Kaliyammal Thiruvassagam, Abhishek Chakraborty, and C. Siva Ram Murthy  
Indian Institute of Technology Madras, Chennai 600036, India  
prabhut@cse.iitm.ac.in, abhishek2003slg@ieee.org, murthy@iitm.ac.in

**Abstract**—Network Functions Virtualization (NFV) and Multi-access Edge Computing (MEC) play crucial roles in 5G networks for dynamically provisioning diverse communication services with heterogeneous service requirements. In particular, while NFV improves flexibility and scalability by softwarizing physical network functions as Virtual Network Functions (VNFs), MEC enables to provide delay-sensitive/time-critical services by moving computing facilities to the network edge. However, these new paradigms introduce challenges in terms of latency, availability, and resource allocation. In this paper, we first explore MEC cloud facility location selection and then latency-aware placement of VNFs in different selected locations of NFV enabled MEC cloud facilities in order to meet the ultra-low latency requirements of different applications (e.g., Tactile Internet, virtual reality, and mission-critical applications). Furthermore, we also aim to guarantee the survivability of VNFs and an edge server against failures in resource limited MEC cloud facility due to software bugs, configuration faults, etc. To this end, we formulate the problem of latency-aware and survivable mapping of VNFs in different MEC cloud facilities as an Integer Linear Programming (ILP) to minimize the overall service provisioning cost, and show that the problem is NP-hard. Owing to the high computational complexity of solving the ILP, we propose a simulated annealing based heuristic algorithm to obtain near-optimal solution in polynomial time. With extensive simulations, we show the effectiveness of our proposed solution in a real-world network topology, which performs close to the optimal solution.

**Index Terms**—NFV, VNF, MEC, Network latency, Survivability, Closeness centrality, Simulated annealing.

## I. INTRODUCTION

Network Functions Virtualization (NFV) and Multi-access Edge Computing (MEC)<sup>1</sup> have emerged as promising key technology enablers for 5G networks and services. NFV replaces hardware middleboxes as Virtual Network Functions (VNF) that can be run on general purpose hardware, which increases flexibility and reduces capital and operational expenditures [1]. MEC enables network operators to support delay-sensitive services by moving cloud computing facilities from the core to the network edge [2] [3]. The primary reason behind the introduction of MEC is to reduce the network latency and bandwidth consumption, and also to leverage the advantage of faster computing and decision-making at the edge of the access network. Since more data are generated at the edge of the network, processing data at the network

<sup>1</sup>Note that MEC is also known as Mobile Edge Computing. Also note that “MEC” and “MEC cloud facility” are synonymous and we use them interchangeably throughout this paper.

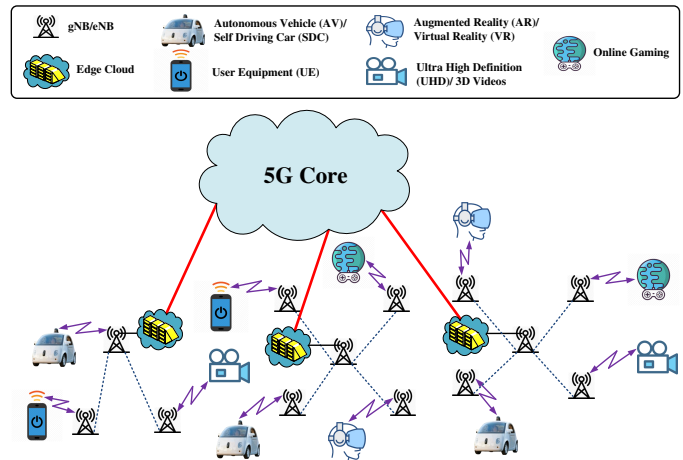


Fig. 1: 5G network architecture based on NFV and MEC, where a few MEC cloud facilities are placed at the selected influential locations.

edge can be efficient solution to accommodate more service requests with extreme service requirements.

Although NFV and MEC based deployment enables to meet stringent and extreme delay requirements of future demands and to reduce capital and operational expenditures, still there are challenges that need to be addressed. In particular, availability, service continuity, and resource allocation are major concerns in MEC enabled 5G networks due to vulnerability in softwarization/cloudification of network functions, sharing of common resources to provide multiple services, and limited resources in MEC cloud facilities. VNFs are subject to failures due to software bugs, configuration faults, etc [4]. Random and unexpected failure of a VNF may disrupt the service abruptly and lead to Service Level Agreement (SLA) violations [5]. Hence, ensuring survivability of network functions is of paramount importance in 5G use case scenarios to guarantee the service continuity and to improve the quality of experience, which is one of the major requirements of 5G systems [6] [7].

Fig. 1 depicts a 5G network architecture that leverages the features of both NFV and MEC technologies. NFV enabled 5G core network is accessed through the Radio Access Networks (RANs) in multihop fashion for providing services. However, NFV based softwarized infrastructure at the core cloud alone may not meet the communication delay requirement to process a service request from a distant source node, as it takes ample amount of time that may not be feasible in the context of latency-critical services such as autonomous driving and virtual reality-based services. Alternatively, edge

cloud facilities can be established at the network edge for provisioning delay-sensitive services. As the service requests can be generated from any corner of a network, placement of a few MEC clouds near to the Base Stations (BSs) of RAN also demands attention [3]. Therefore, MECs should be placed in such a way that the overall network can be covered in a manageable way and it also reduces the capital and operational costs [8]. From Fig. 1, it can be observed that among the set of BS (eNB/gNB) locations, only a few potential BSs (in terms of influences based on the topological structure) are identified as the MEC cloud establishment locations. We note that the identified MEC cloud locations are very close to other remaining BSs in the network and, thus, reduce the response time for accessing a time-critical service. Hence, these MECs can be used to provide services in a timely manner. Furthermore, along with the proper selection of MECs, the required VNFs should also be optimally placed/instantiated on the MECs to provide the requested services efficiently. In this paper, we explore the problem of latency-aware and survivable mapping of VNFs in different selected locations of MEC cloud to provide delay-sensitive services in reliable manner.

The significant contributions of this paper are as follows:

- We first propose an algorithm to select and establish a few MEC cloud facilities in the potential locations of BSs to cover the entire region and meet the delay requirements.
- Then, we explore latency-aware placement of VNFs onto the servers of a few established MEC cloud facilities. In addition, to ensure survivability, we jointly place primary and backup VNFs in different edge servers, which enhances the availability of communication services and provides protection from random and unexpected failures.
- We formulate the problem as an ILP for provisioning services with minimal cost, and show that the problem is an NP-hard problem.
- We use CPLEX solver to find the optimal solution for the problem, and propose a simulated annealing based heuristic algorithm to provide near-optimal solution in polynomial time for large input instances.
- With extensive simulations, we show the effectiveness of our proposed solution in a real-world network topology.

The remaining part of the paper is structured as follows: We review the related work in Section II. We describe the system model and problem definition in Section III. We present MEC selection algorithm in Section IV. We formulate the ILP and present our proposed algorithm for latency-aware and survivable mapping of VNFs onto MEC edge servers in Section V. We evaluate the performance of our proposed algorithms in Section VI. Finally, we conclude the paper with some directions for future work in Section VII.

## II. RELATED WORK

In [4] and [9], latency-aware and reliable VNF chain placement problem is considered, but MEC scenario is not

taken into account. In [10], an efficient VNF chain placement problem is considered in an MEC-NFV environment with the goal of maximizing the resource utilization. In [11], VNF placement and resource allocation problem is considered in NFV/SDN-enabled MEC networks with the goal of minimizing the overall placement and resource cost. In [12], VNF placement problem is considered at the network edges with the goal of minimizing end-to-end latency, and neural-network based model is used to proactively predict the number of VNFs required to process the network traffic. In [13], joint user association and VNF placement problem is considered for providing latency sensitive applications using MEC in 5G networks with the goal of minimizing the service provisioning cost. In [14], QoS-aware VNF placement problem is considered in edge-central cloud architecture with the goal of efficiently allocating resources for provisioning services.

In [15], latency-aware and availability driven VNF placement problem is considered in MEC-NFV environment with the goal of minimizing the cost. The work deals with availability of resources in MEC or core cloud and the latency associated with it. In [16], latency-aware VNF composition problem is considered in 5G edge network with the goal of minimizing the overall latency. In [17], dynamic latency optimal VNF placement problem is considered at the network edge with the goal of minimizing the end-to-end latency. In [18], latency-aware VNF deployment problem is considered at edge for IoT services with the goal of minimizing the end-to-end latency. In [19], latency-aware VNF placement and assignment problem is considered in MEC with the goal of maximizing the number of admitted service requests. In [20] and [21], resilient placement of VNFs in MEC is considered with the goal of minimizing the overall service provisioning cost.

In the literature, MEC cloud network location is assumed to be given and the research is focused only on either latency-aware or resilient VNF placement in MEC. In this work, we first consider selection of MEC cloud facility location and then focus on both latency-aware and survivability aspects together in the placement of VNFs onto the selected/established MEC cloud facility with minimum cost.

## III. SYSTEM MODEL AND PROBLEM DEFINITION

We model the physical network as an undirected graph  $G = (N, E)$ , where  $N$  denotes the set of BSs in the region and  $E$  denotes the set of physical links that interconnect the BSs. The BSs can be interconnected by SDN based backhaul network. A small subset of BSs is chosen to establish MEC cloud facilities. We use the notation  $L$  to denote the set of locations where MEC cloud network facilities are being established, where  $L \subset N$ . At each MEC cloud network facility, a set of limited number of servers  $S$  are used to place VNFs in order to provide service for user requests. We use the symbol  $C_s$  to denote the available resource capacity (e.g., CPU, RAM, and storage space) of each server  $s \in S$ . We consider a set of VNFs, denoted by  $V$ ,

that process data traffic to provide services for user requests. Each VNF  $v \in V$  requires a certain amount of resource to process the packets. The amount of resource required by VNFs is denoted by  $C_v$  and it should be less than the available resource capacity of MEC cloud network servers.

We consider that VNFs are subject to failures due to software bugs, configuration faults, unexpected failures of network functions, and cyber attacks (e.g., denial of service). Abrupt failure of a VNF may disrupt communication services and results in customer dissatisfaction and revenue loss. In order to enhance the reliability of communication services, backups are assigned to VNFs such that they meet SLAs and improve the service continuity.

Multiple users are connected to the network through near by base stations and their service requests come through these base stations. We assume that each user service request  $r \in R$  is represented as  $(v^r, n^r, t^r, d^r)$ , where  $v^r \in V$  denotes service type VNF,  $n^r \in N$  denotes which BS user connects to and requests for service,  $t^r$  denotes the data rate demand, and  $d^r$  denotes the maximum allowed delay/latency.

**Problem Definition:** Given a physical network graph  $G = (N, E)$  and a set of service requests  $r \in R$  with  $(v^r, n^r, t^r, d^r)$ , find a joint efficient mapping of primary and backup VNFs (for ensuring survivability) in different edge servers at MEC cloud facility locations to minimize the overall provisioning cost while meeting the SLAs.

---

**Algorithm 1** MEC cloud facility location selection and establishment

---

**Input:** Graph  $G = (N, E)$  and maximum allowed delay requirement  $D_{max}$  to reach MEC cloud facility  
**Output:** Number of established MEC cloud facilities and its locations

```

1: for  $i = 1 \rightarrow |N|$  do  $\triangleright$  Estimate the CC value of all BSs in  $G$ 
2:   for  $j = 1 \rightarrow |N|$  do
3:      $CC[i] = \frac{1}{\sum_j \text{distance}(i, j)}$   $\triangleright$  distance( $i, j$ ) is the shortest path distance between nodes  $i$  and  $j$ 
4:   end for
5: end for
6:  $S = \text{Sort the nodes (i.e., } N \in G) \text{ in the descending order of the CC values}$ 
7: Current node locations  $L = \{l_1, l_2, \dots, l_{|N|}\}$  based on the CC value of nodes  $\triangleright$  the set of locations for establishing MEC cloud facility
8: Delay =  $\infty$ ,  $i=1$ 
9: while Delay  $\leq D_{max}$  do
10:   Select the location  $l_i$   $\triangleright$  high CC node location is selected
11:   Establish MEC cloud facility at the location  $l_i$ 
12:   Delay = maximum delay from BSs in the network to reach one of the established MEC cloud facilities
13:    $i=i+1$ 
14: end while
15: Return the number of MEC cloud facilities established and their corresponding locations
```

---

#### IV. MEC SELECTION AND ESTABLISHMENT

In this paper, to identify a few potential MEC cloud facilities/locations  $L \subset N$  [3], we select a set of influential BS nodes on the basis of high Closeness Centrality (CC) [22] [23] values. The procedure to select high influential BSs from the network is given in Algorithm 1. Note that a node with a high CC value can be reached, from any distant node in

a network, with a few hops (or by traversing less distance). Therefore, we choose high CC-valued BSs as potential MEC cloud network locations.

Algorithm 1 first identifies the CC value of each BS to identify the influential nodes in the network (lines 1 to 5). Influential nodes are selected based on the connectivity and closeness with respect to all other nodes in the network. Then, a few MEC cloud facilities are established on the locations of high CC-valued BSs to reach MEC cloud facility from nodes in the network within the maximum allowed delay ( $D_{max}$ ) requirement (lines 6 to 14). In this work,  $D_{max}$  is set to 2 ms.

#### V. ILP FORMULATION AND PROPOSED SOLUTION

##### A. ILP Formulation

The objective is to place the required VNFs onto the MEC servers in reliable manner such that the placement strategy minimizes the overall provisioning cost while meeting the SLAs of diverse service requests.

1) **Decision Variables:** We define the following decision variables to formulate our problem of survivable placement of VNFs in MEC.

- $w_l$ : Binary variable that equals 1 if an MEC cloud facility  $l \in L$  is chosen for providing service, and 0 otherwise.
- $x_{ls}$ : Binary variable that equals 1 if a server  $s \in S$  is activated in the MEC cloud facility  $l \in L$  to deploy VNF, and 0 otherwise.
- $y_{lsv}$ : Integer variable that equals  $\mathbb{N}$  if number of instances of VNF  $v \in V$  are deployed on a server  $s \in S$  in the MEC cloud facility  $l \in L$ , and 0 otherwise.
- $y_{lsv^b}$ : Integer variable that equals  $\mathbb{N}$  if number of backup instances of VNF  $v^b \in V$  are deployed on a server  $s \in S$  in the MEC cloud facility  $l \in L$ , and 0 otherwise.
- $z_{lsv}^{nr}$ : Binary variable that equals 1 if a request  $r \in R$  through the base station  $n \in N$  is served by the VNF  $v \in V$  which is placed on the server  $s \in S$  in the the MEC location  $l \in L$ , and 0 otherwise.
- $z_{lsv^b}^{nr}$ : Binary variable that equals 1 if a request  $r \in R$  through the base station  $n \in N$  is served by the backup VNF  $v^b \in V$  which is placed on the server  $s \in S$  in the the MEC location  $l \in L$ , and 0 otherwise.

2) **Objective Function:** The objective is to minimize the cumulative costs of number of physical MEC servers activated, number of VNFs deployed, and amount of traffic being forwarded on each link for provisioning reliable and delay-sensitive communication services.

i) Activation Cost of Physical MEC Server: It includes design, procurement, deployment, and maintenance costs of MEC cloud, where multiple servers are acti-

vated to host VNFs to provide reliable communication services. It can be expressed as follows:

$$SC = c_{sc} \sum_{l \in L} \sum_{s \in S} x_{ls}, \quad (1)$$

where  $c_{sc}$  denotes activation cost of a single server in MEC cloud locations.

ii) Deployment Cost of VNF Instance: It includes the deployment/license cost of both primary and backup VNFs hosted on the physical MEC servers, which can be expressed as follows:

$$VC = c_{vc} \sum_{l \in L} \sum_{s \in S} \sum_{v, v^b \in V} (y_{lsv} + y_{lsv^b}), \quad (2)$$

where  $c_{vc}$  denotes deployment cost of a VNF on any physical MEC server.

iii) Forwarding Cost of Service Traffic: It is the cost for forwarding service request traffic from the base station of user to the MEC cloud facility where the VNF is hosted on the server to provide service, which can be expressed as follows:

$$TC = c_{tc} \sum_{l \in L} \sum_{s \in S} \sum_{v, v^b \in V} \sum_{n \in N} \sum_{r \in R} (z_{lsv}^{nr} + z_{lsv^b}^{nr}) \times t^r, \quad (3)$$

where  $c_{tc}$  denotes traffic forwarding cost for the service request  $r$  and it is calculated per Mbps and  $t^r$  denotes the data rate requirement of the service request.

The objective is to minimize the overall cost of the aforementioned costs, which can be expressed as follows:

$$P: \min (\gamma_1 \times SC + \gamma_2 \times VC + \gamma_3 \times TC), \quad (4)$$

where  $\gamma_1$ ,  $\gamma_2$ , and  $\gamma_3$  are weighing factors to give relative importance to the objective functions.

3) **Capacity Constraints:** The resource requirement of VNFs should be within the limit of resources available in the MEC servers, and the processing capacity requirement of service requests should be within the limit of available processing capacity of VNFs.

i) The total resource requirement of VNFs (both primary and backup) to be placed should not exceed the available resource capacity of the MEC server which hosts VNFs. It can be expressed as follows:

$$\sum_{v, v^b \in V} C_v \times (y_{lsv} + y_{lsv^b}) \leq C_s \times x_{ls}, \forall l \in L, \forall s \in S, \quad (5)$$

where  $C_v$  denotes resource requirement of VNFs and  $C_s$  denotes the available resource capacity of the MEC server.

ii) The total processing capacity requirement of service requests should not exceed the available processing capacity of VNFs (both primary and backup), which can be expressed as follows:

$$\sum_{n \in N} \sum_{r \in R} t^r \times z_{lsv}^{nr} \leq pc_v \times y_{lsv}, \forall l \in L, \forall s \in S, \forall v \in V \quad (6)$$

$$\sum_{n \in N} \sum_{r \in R} t^r \times z_{lsv^b}^{nr} \leq pc_v \times y_{lsv^b}, \forall l \in L, \forall s \in S, \forall v^b \in V \quad (7)$$

where  $pc_v$  denotes the processing capacity of VNFs.

4) **Delay Constraint:** The delay requirement of service requests should be less than or equal to the delay between the service requesting base stations and MEC server locations, which can be expressed as follows:

$$\sum_{l \in L} \sum_{s \in S} \sum_{v \in V} d_{ln} \times z_{lsv}^{nr} \leq d^r, \forall n \in N, \forall r \in R \quad (8)$$

$$\sum_{l \in L} \sum_{s \in S} \sum_{v^b \in V} d_{ln} \times z_{lsv^b}^{nr} \leq d^r, \forall n \in N, \forall r \in R \quad (9)$$

where  $d_{ln}$  denotes the communication delay between the service request carrying base station  $n \in N$  and the MEC server location  $l \in L$ .

5) **Placement Constraint:** Each service request from the user through base station is assigned to two instances of the same VNF type (primary and backup) to ensure survivability, which can be expressed as follows:

$$\sum_{l \in L} \sum_{s \in S} \sum_{v \in V} z_{lsv}^{nr} = 1, \forall n \in N, \forall r \in R \quad (10)$$

$$\sum_{l \in L} \sum_{s \in S} \sum_{v^b \in V} z_{lsv^b}^{nr} = 1, \forall n \in N, \forall r \in R \quad (11)$$

6) **Anti-affinity VNF Mapping Constraint:** The primary and backup VNFs should be placed in different edge servers in order to handle failure of VNFs or an edge server, which can be expressed as follows:

$$\sum_{v, v^b \in V} (z_{lsv}^{nr} + z_{lsv^b}^{nr}) \leq 1, \forall n \in N, \forall r \in R, \forall s \in S, \forall l \in L \quad (12)$$

7) **Other Constraints:**

i) The MEC cloud location is chosen if at least one MEC server is activated in that location to place VNFs, which can be expressed as follows:

$$w_l = 1 \text{ if } \sum_{l \in L} \sum_{s \in S} x_{ls} > 0, \forall l \in L \quad (13)$$

ii) The MEC server is activated if at least one VNF (either primary or backup) is placed on it, which can be expressed as follows:

$$x_{ls} = 1 \text{ if } \sum_{v, v^b \in V} (y_{lsv} + y_{lsv^b}) > 0, \forall l \in L, \forall s \in S \quad (14)$$

iii) VNFs (both primary and backup) are deployed if at least one service request from the user through base station requires the VNF to provide a particular service, which can be expressed as follows:

$$y_{lsv} = 1 \text{ if } \sum_{n \in N} \sum_{r \in R} z_{lsv}^{nr} > 0, \forall l \in L, \forall s \in S, \forall v \in V \quad (15)$$

$$y_{lsv^b} = 1 \text{ if } \sum_{n \in N} \sum_{r \in R} z_{lsv^b}^{nr} > 0, \forall l \in L, \forall s \in S, \forall v^b \in V \quad (16)$$

**Theorem 1.** *Latency-aware and survivable mapping of VNFs in MEC is an NP-hard problem.*

*Proof.* Let A be the problem of latency-aware and survival mapping of VNFs in MEC and B be the Reliable Capacitated Facility Location (RCFL) problem. RCFL problem is an optimization problem and it is NP-hard [24]. In RCFL problem, it is considered that facilities fail with equal probability and the model assigns primary and backup facilities for the demand to enhance the reliability. RCFL problem is defined as follows: the problem is to select facilities from the given set of potential facility locations, where each facility has limited capacity and subject to failure, to provide services to the demands such that the model is robust against failures and minimizes the cost of establishing facilities (primary and backup) and of transportation of goods from the facilities to the demand points. To prove that the problem A is NP-hard, it is sufficient to show that an instance of the problem B can be reduced to an instance of the problem A in polynomial time, i.e.,  $B \leq_P A$  [25].

We can transform an instance of the problem B into an instance of the problem A in the following way: i) consider each facility in the problem B as equivalent to an MEC cloud facility in the problem A, ii) set the capacity of the facility in the problem B to be equal to the capacity of the MEC cloud facility in the problem A, iii) set the cost of activating facility in the problem B is equivalent to the activation cost of servers and deployment cost of VNFs (primary and backup) at MEC cloud in the problem A, and iv) set the transportation cost in the problem B as the traffic forwarding cost in the problem A. The transformation operation can be done in polynomial time of the input size. Hence, the problem B is reducible to the problem A in polynomial time. If A is not NP-hard, then B is also not NP-hard (since B is reducible to A), which is a contradiction. Therefore, it can be concluded that the problem A is also an NP-hard problem.  $\square$

### B. Proposed Heuristic Solution

As latency-aware and survivable mapping of VNFs onto the selected MEC cloud facility locations is an NP-hard problem, we develop a Simulated Annealing (SA) [26] [27] based heuristic algorithm to obtain near-optimal solution in polynomial time for providing services to user requests. Algorithm 2 gives the procedure for latency-aware and survivable mapping of VNFs onto the MEC cloud facilities, and it is based on the concept of SA. SA is a probabilistic method for finding the global minimum cost function and the process may consist of multiple local minima. The SA algorithm is similar to traditional local search algorithms, but SA allows upward moves occasionally with the hope to come out of local minima. Although upward moves lead to increase in cost, it will help to escape from local minima.

SA mathematically mirrors the physical process whereby a solid is slowly cooled to a frozen state of minimum energy. The minimum energy state (or ground state configuration) in statistical mechanics corresponds to the minimum cost function in combinatorial optimization problems, where the cost function plays the role of energy [26].

---

**Algorithm 2** Simulated annealing based latency-aware and survivable mapping of VNFs onto the MEC cloud facilities

---

**Input:**  $G = (N, E)$  and a set of service requests with information  $(v^r, n^r, t^r, d^r) \forall r \in R$   
**Output:** Latency-aware and survivable mapping of VNFs onto the MEC cloud servers with minimum cost

- 1:  $T = T_0$
- 2: currentSol = Generate a current solution randomly
- 3: Evaluate the current solution using the objective function (Equation 4), i.e.,  $c_1 = \text{cost}(\text{currentSol})$
- 4: **while**  $T > T_{min}$  **do**
- 5:   **for**  $i = 1 \rightarrow \text{maxIterations}$  **do**
- 6:     nextSol = Generate a next solution
- 7:     **if** nextSol does not violate any SLAs/constraints **then**
- 8:       Evaluate the next solution generated using the objective function (Equation 4), i.e.,  $c_2 = \text{cost}(\text{nextSol})$
- 9:       **if**  $c_2 \leq c_1$  **then**
- 10:          currentSol = nextSol, i.e.,  $c_1 = c_2$
- 11:       **else**
- 12:           $r = \text{random}(0, 1)$
- 13:           $p = e^{(c_1 - c_2)/T}$
- 14:          **if**  $r < p$  **then**
- 15:            currentSol = nextSol, i.e.,  $c_1 = c_2$
- 16:          **end if**
- 17:       **end if**
- 18:     **end for**
- 19:    **end while**
- 20:     $T = \alpha \times T$
- 21: **end while**
- 22: Return  $c_1$

---

In Algorithm 2, first an initial temperature is set and current feasible solution is generated and evaluated using the objective function (Equation 4). The current solution is generated by sorting the requests based on the latency requirement in ascending order and placing the required VNFs (primary and backup) onto different MEC cloud servers. We follow first fit principle to reuse the activated edge servers and deployed VNFs. Then, different next solutions are explored for maximum number of iterations. At each iteration in the inner loop, the cost of a new next solution is computed using the same objective function (Equation 4). If the cost difference  $(c_1 - c_2)$  is less than or equal to zero, then the solution is accepted directly, and the configuration of the new solution is set as the current solution. In the case that cost difference  $(c_1 - c_2)$  is greater than zero, a new solution is accepted with a certain probability. First, a random number  $r$  is generated that is uniformly distributed between  $(0, 1)$ . Then,  $r$  is compared with the probability value  $p$  that is a function of the temperature and the cost difference of current and new solutions. If  $r < p$ , then the configuration of the new solution is set as the current solution; otherwise the original configuration is retained. At the end of maximum of number of iterations, the temperature value is updated. The process continues till the temperature goes below the

minimum threshold temperature. The proposed heuristic follows the annealing process of cooling the temperature in a controlled manner and allowing bad movement with a certain probability to come out of a local minima.

## VI. PERFORMANCE EVALUATION

In this section, we evaluate the performance of our proposed solution for solving latency-aware and survivable VNF mapping problem in MEC cloud networks.

### A. Simulation Setup

For the evaluation purpose, we use germany50 real-world network topology from SNDlib [28] which is a library of test instances for telecommunications network. The germany50 network topology consists of 50 nodes that are interconnected by 88 links. We assume that each node is a base station node. The selection of MEC cloud facility locations is based on closeness centrality metric from complex network theory as explained in Algorithm 1, hence the chosen MEC cloud network locations are closely situated with respect to all other base stations in the network.

We use CPLEX solver (version 12.8) and Concert Technology in Java to solve the ILP formulation, latency-aware and survivable mapping of VNFs in MEC. The proposed heuristic is implemented using Matlab and we run the simulations multiple number of times and take the average for evaluation. In the SA based heuristic design in Algorithm 2, we set the initial temperature ( $T_0$ ) as 100, the minimum temperature ( $T_{min}$ ) for termination as 0.1, the number of inner loop iterations ( $maxIterations$ ) as 50, and the cooling rate ( $\alpha$ ) as 0.9. We have observed that going beyond 50 iterations do not improve the quality of the solution significantly.

### B. Performance Analysis of MEC Cloud Facility Location Selection

As explained in Algorithm 1 earlier, MEC cloud facilities are chosen based on the centrality metric which depends on the topological structure of the network. We use CC to select the set of potential MEC cloud locations from germany50 network to provide latency-aware services. Fig. 2 compares the performance of CC-based MEC selection with random selection. Average delay is the mean minimum delay to reach one of the MEC facilities from all the nodes in the network and max delay is the highest minimum delay from any node in the network to reach one of the MEC cloud facilities. As it can be seen that as we increase the number of MEC cloud facilities both average and max delays reduce. Since CC-based selection chooses high centrality nodes in the network, it takes less time for other nodes to reach the MEC facilities. Since random selection method chooses the MEC facilities randomly and it may choose distant corner node as MEC cloud facility location, average and max delays are high compared to CC-based selection method. As it can be seen from Fig. 2, the maximum delay is within 2 ms when 5 MEC cloud facilities are established using CC-based selection method for providing services to users.

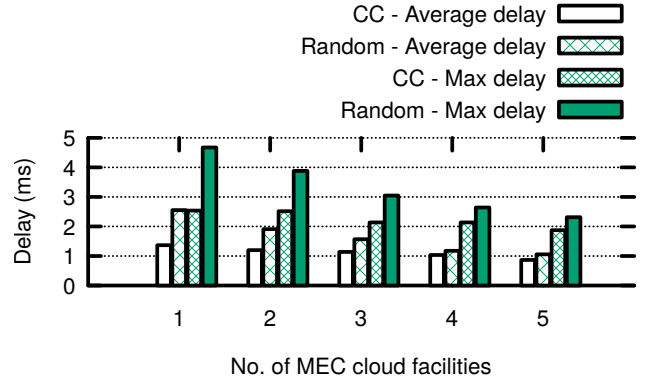


Fig. 2: CC-based selection vs. Random selection.

TABLE I: Simulation parameters [29]

Service Types	Data Rate	Max Allowed Delay
AR/VR	200 Mbps	2 ms
V2X	100 Mbps	3 ms
e-health	100 Mbps	5 ms
8K TV and Gaming	200 Mbps	10 ms

### C. Performance Analysis of Latency-aware and Survivable Mapping of VNFs in MEC

We analyze the performance of our heuristic solution proposed in Section V for solving latency-aware and survivable placement of VNFs in MEC with minimal cost. Table I shows the simulation parameters considered in this work, which are based on the requirements given in [29]. Four service types of VNFs are considered and each type has its corresponding data rate and maximum allowed delay requirements. We consider that each MEC cloud network has 10 MEC servers and each server has the resource capacity of 16 cores, and each VNF requires 4 cores and has the processing capacity of 1 Gbps [30]. Hence, 4 VNFs can be placed in a server and multiple services can share the same VNF. Each VNF service type has different data rate and maximum allowed delay requirement, and each user service request through the base station is randomly associated with one of the four service types with equal probability. The propagation delay between the base stations and MEC cloud network locations is computed based on the distance between them and considered that base stations are interconnected using optical fiber. We assume that processing and forwarding delay of the VNF is  $50 \mu s$  approximately [31]. Our latency-aware VNF placement strategy satisfies the maximum allowed delay requirement by giving priorities to service requests of low delay requirements. For reliable service provisioning, service requests are associated with two different VNFs (active and backup) in different edge servers such that if the active VNF fails unexpectedly in random manner then the backup VNF takes charge to continue providing services without service interruption and disruption. From the above MEC cloud facility selection analysis, 5 MEC cloud facilities are enough to meet the required maximum delay requirement of 2 ms from the nodes in the network to reach one of the MEC

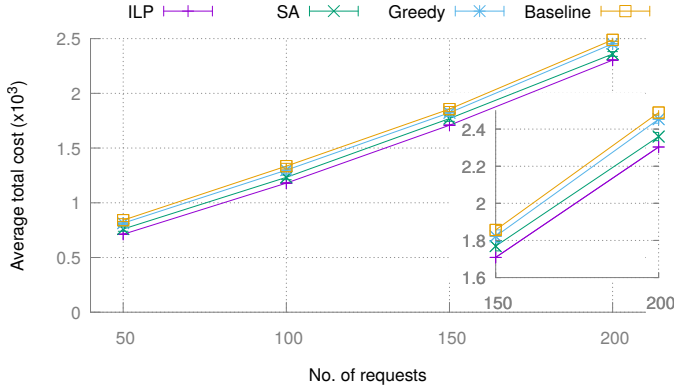


Fig. 3: Comparison of average total cost for provisioning services.

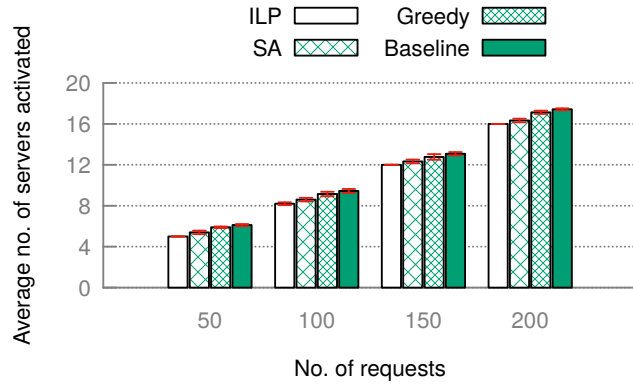


Fig. 4: Comparison of average no. of MEC servers activated.

cloud facilities.

Meeting the extreme requirements as well as effectively reusing the available resources to provide reliable services is a challenging task. The proposed algorithm aims to primarily meet the SLA requirements and reuse the shareable resources as efficiently as possible while ensuring survivability against failures.

We compare the performance of our proposed Simulated Annealing (SA) based heuristic solution against the following:

- ILP: Formulated ILP problem is solved using CPLEX solver and it provides optimal solution.
- Greedy: This approach always places VNF in the nearest MEC cloud among all the possible MEC clouds which meet the SLA requirements. It incurs the minimal delay to provide services.
- Baseline: This approach places VNF on the MEC cloud server that encounters first in the search space and satisfies the SLA requirements of the service request.

For the evaluation purpose, we consider that MEC cloud server activation cost is 100, VNF deployment cost is 10, and traffic forwarding cost is 1 per Mbps with respect to the delay between the service request carrying base station and the MEC cloud that provides service [32]. Fig. 3 shows the comparison of total cost with respect to the number of requests being served. ILP takes the least cost for providing services. Although ILP provides the best solution, it is com-

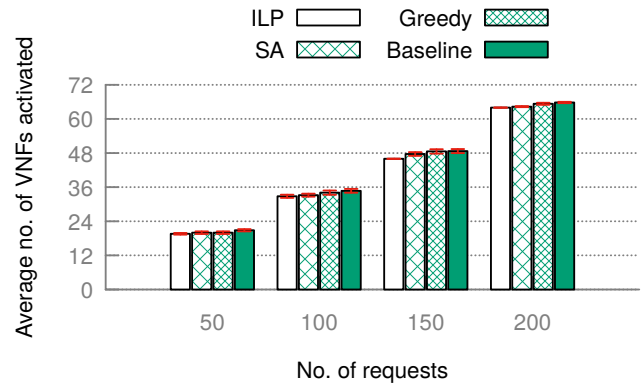


Fig. 5: Comparison of average no. of VNFs activated.

TABLE II: Average running time comparison of different approaches (in seconds)

#Service requests	50	100	150	200
ILP	28.073	110.102	2298.805	21312.82
SA	25.756	47.782	74.164	96.889
Greedy	0.3549	0.4813	0.7592	0.8546
Baseline	0.4832	0.6288	0.8771	0.9981

putationally expensive and impractical for large-scale input instances. Hence, we developed the Simulated Annealing (SA) based heuristic algorithm to provide near-optimal solution in polynomial time. We compare the performance of the proposed heuristic with the ILP (optimum solution). As shown in Fig. 3, our proposed SA based heuristic solution provides near-optimal solution (maximum optimality gap is 3.5%) with lesser cost. Fig. 4 shows the comparison of average number of MEC servers activated to place the required VNFs, and it is clear that our proposed SA based heuristic solution performs close to the ILP and reduces the overall average cost as shown in Fig. 3. Fig. 5 shows comparison of average number of VNFs activated. Although the average number of VNFs activated on different MEC servers are close for different approaches, the number of MEC servers activated to deploy the required VNFs for providing different classes of services to users differ clearly as shown in Fig. 4 and thus influences the overall cost for providing services.

We compare the average running time (in seconds) of ILP with different approaches for solving the latency-aware and survivable VNF placement problem in Table II. Solving ILP using CPLEX provides optimal solution in reasonable amount of time for small input instances. However, the running time to solve ILP increases exponentially as we scale the number of requests as shown in Table II. Owing to the high computational complexity of solving large instance of the ILP problem, we propose a SA based heuristic algorithm. The running time of heuristic algorithms are insignificant (in the order of seconds) compared to the running time of ILP. Since SA based heuristic algorithm explores different neighborhood solutions with respect to the temperature and the number of inner loop iterations, running time is in the order of seconds to provide near-optimal solution. Hence, from Table II it is clear that the proposed algorithm solves

the problem in polynomial time. Since greedy and baseline approaches are executed once, they take less than a second to provide solution.

## VII. CONCLUSION

In this work, we focused on latency-aware and survivable placement of VNFs in 5G network edge cloud. We first proposed an algorithm to select a few MEC cloud facility locations from the set of base stations to establish MEC cloud infrastructure and meet the user requirements of delay-sensitive services. Then, we explored both latency and survivability aspects together by leveraging the features of NFV and MEC cloud based technologies. We formulated the problem as an ILP to minimize the overall service provisioning cost (including both computing and communication resource cost). In order to overcome the high computational complexity of the ILP problem, we proposed a simulated annealing based heuristic algorithm which provides near-optimal and reliable solution to delay-sensitive heterogeneous service requests from users and industry verticals in polynomial time. We evaluated our proposed algorithm in terms of total provisioning cost and running time. With extensive simulations, we showed that our proposed solution performed close to the optimal solution (optimality gap is 3.5%) in real-world network topology.

In this work, we designed an offline algorithm to process the batch of service requests in order to place VNFs such that the SLA requirements are met. As a future work, we plan to design machine learning based online algorithm by considering the fact that the future service requests are not well known in advance. In addition, we would like to explore failure detection and rerouting mechanisms to analyze the actual delay incurred in the recovery process after failure of a network component in NFV/SDN-enabled 5G networks.

## ACKNOWLEDGEMENT

This research work was supported by the Department of Science and Technology (DST), New Delhi, India.

## REFERENCES

- [1] ETSI NFV ISG, "Network Functions Virtualization: An Introduction, Benefits, Enablers, Challenges and Call for Action," Oct. 2012.
- [2] ETSI ISG MEC, "Mobile edge computing: A key technology towards 5G," ETSI, Sep. 2015.
- [3] —, "MEC in 5G networks," ETSI, June 2018.
- [4] P. Kaliyammal Thiruvassagam, V. J. Kotagi, and S. R. Murthy, "A reliability-aware, delay guaranteed, and resource efficient placement of service function chains in softwarized 5G networks," *IEEE Transactions on Cloud Computing*, pp. 1–1, 2020.
- [5] P. K. Thiruvassagam, V. J. Kotagi, and C. S. R. Murthy, "The more the merrier: Enhancing reliability of 5G communication services with guaranteed delay," *IEEE Networking Letters*, vol. 1, no. 2, pp. 52–55, June 2019.
- [6] ETSI GS NFV-REL 001 V1.1.1, "Network functions virtualization: Resiliency requirements," ETSI, Jan. 2015.
- [7] ETSI GS NFV-REL 003 V1.1.2, "Network Functions Virtualization; Reliability; Report on Models and Features for End-to-End Reliability," ETSI, July 2016.
- [8] Y. Li and S. Wang, "An energy-aware edge server placement algorithm in mobile edge computing," in *Proc. IEEE International Conference on Edge Computing (EDGE)*. IEEE, July 2018, pp. 66–73.
- [9] P. K. Thiruvassagam, A. Chakraborty, A. Mathew, and C. S. R. Murthy, "Reliable placement of service function chains and virtual monitoring functions with minimal cost in softwarized 5G networks," *IEEE Transactions on Network and Service Management*, pp. 1–1, 2021.
- [10] M. Wang, B. Cheng, W. Feng, and J. Chen, "An Efficient Service Function Chain Placement Algorithm in a MEC-NFV Environment," in *Proc. IEEE Global Communications Conference (GLOBECOM)*, 2019, pp. 1–6.
- [11] N. Kiran, X. Liu, S. Wang, and C. Yin, "VNF Placement and Resource Allocation in SDN/NFV-Enabled MEC Networks," in *Proc. IEEE Wireless Communications and Networking Conference Workshops (WCNCW)*, 2020, pp. 1–6.
- [12] T. Subramanya and R. Riggio, "Machine Learning-Driven Scaling and Placement of Virtual Network Functions at the Network Edges," in *Proc. IEEE Conference on Network Softwarization (NetSoft)*, 2019, pp. 414–422.
- [13] R. Behraves, E. Coronado, D. Harutyunyan, and R. Riggio, "Joint User Association and VNF Placement for Latency Sensitive Applications in 5G Networks," in *Proc. IEEE International Conference on Cloud Networking (CloudNet)*, 2019, pp. 1–7.
- [14] F. Ben Jemaa, G. Pujolle, and M. Pariente, "QoS-Aware VNF Placement Optimization in Edge-Central Carrier Cloud Architecture," in *Proc. IEEE Global Communications Conference (GLOBECOM)*, 2016, pp. 1–7.
- [15] L. Yala, P. A. Frangoudis, and A. Ksentini, "Latency and Availability Driven VNF Placement in a MEC-NFV Environment," in *Proc. IEEE Global Communications Conference (GLOBECOM)*, 2018, pp. 1–7.
- [16] B. Martini, F. Paganelli, P. Capanera, S. Turchi, and P. Castoldi, "Latency-aware composition of virtual functions in 5G," in *Proc. IEEE Conference on Network Softwarization (NetSoft)*, 2015, pp. 1–6.
- [17] R. Cziva, C. Anagnostopoulos, and D. P. Pezaros, "Dynamic, Latency-Optimal vNF Placement at the Network Edge," in *Proc. IEEE Conference on Computer Communications (INFOCOM)*, 2018, pp. 693–701.
- [18] M. Emu, P. Yan, and S. Choudhury, "Latency Aware VNF Deployment at Edge Devices for IoT Services: An Artificial Neural Network Based Approach," in *Proc. IEEE International Conference on Communications Workshops (ICC Workshops)*, 2020, pp. 1–6.
- [19] D. Harris, J. Naor, and D. Raz, "Latency Aware Placement in Multi-access Edge Computing," in *Proc. IEEE Conference on Network Softwarization and Workshops (NetSoft)*, 2018, pp. 132–140.
- [20] H. D. Chantre and N. L. Saldanha da Fonseca, "The Location Problem for the Provisioning of Protected Slices in NFV-Based MEC Infrastructure," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 7, pp. 1505–1514, 2020.
- [21] P. Zhao and G. Dan, "Resilient placement of virtual process control functions in mobile edge clouds," in *Proc. IFIP Networking Conference (IFIP Networking) and Workshops*, 2017, pp. 1–9.
- [22]
- [23] B. S. Manoj, A. Chakraborty, and R. Singh, *Complex Networks: A Networking and Signal Processing Perspective*. Prentice Hall PTR, New Jersey, USA, Feb. 2018.
- [24] Z.-J. M. Shen, R. L. Zhan, and J. Zhang, "The reliable facility location problem: Formulations, heuristics, and approximation algorithms," *INFORMS Journal on Computing*, vol. 23, no. 3, pp. 470–482, 2011.
- [25] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. MIT Press, USA, July 2009.
- [26] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, "Optimization by simulated annealing," *SCIENCE*, vol. 220, no. 4598, pp. 671–680, 1983.
- [27] Z. Michalewicz and David B. Fogel, *How to Solve It: Modern Heuristics*. Springer-Verlag, Germany, 2004.
- [28] S. Orłowski, M. Pióro, A. Tomaszewski, and R. Wessäly, "SNDlib 1.0 – Survivable network design library," in *Proc. International Network Optimization Conference (INOC)*, Apr. 2007, pp. 1–15.
- [29] ETSI White Paper, "Cloud RAN and MEC: A Perfect Pairing," ETSI, Feb. 2018.
- [30] A. Varasteh, M. De Andrade, C. M. Machuca, L. Wosinska, and W. Kellerer, "Power-aware virtual network function placement and routing using an abstraction technique," in *Proc. IEEE Global Communications Conference (GLOBECOM)*, Dec. 2018, pp. 1–7.
- [31] "Network I/O Latency on VMware vSphere 5," 2020. [Online]. Available: <https://www.vmware.com/content/dam/digitalmarketing/vmware/en/pdf/techpaper/network-io-latency-perf-vsphere5-white-paper.pdf>
- [32] W. Vizarreta, M. Condoluci, C. M. Machuca, T. Mahmoodi, and W. Kellerer, "QoS-driven function placement reducing expenditures in NFV deployments," in *Proc. IEEE International Conference on Communications (ICC)*, 2017, pp. 1–7.