

Stringer: Balancing Latency and Resource Usage in Service Function Chain Provisioning

Freddy C. Chua, Julie Ward,
Ying Zhang, Puneet Sharma, Bernardo A. Huberman
Hewlett Packard Labs, Hewlett Packard Enterprise,
Palo Alto, CA 94304, USA

June 9, 2016

Abstract

Network Functions Virtualization (NFV) enables telecommunications infrastructure providers to replace special-purpose networking equipment with commodity servers running virtualized network functions (VNFs). A provider utilizing NFV faces the Service Function Chain (SFC) provisioning problem of assigning VNF instances to nodes in the physical infrastructure (e.g. datacenters), and routing Service Function Chains (sequences of functions required by customers, a.k.a. SFCs) in the physical network. The provider must balance competing goals of performance and resource usage. We present an approach to SFC provisioning, consisting of three elements. The first element is a fast and scalable round-robin heuristic. The second element is a Mixed Integer Programming (MIP) based approach. The third element is a queueing-theoretic model to estimate the average latency associated with any SFC provisioning solution. Our SFC provisioning system, called *Stringer*, allows providers to balance the conflicting goals of minimizing infrastructure resources and end-to-end latency for meeting their respective SLAs.

1 Introduction

Telecommunications providers are making a strong push towards Network Functions Virtualization (NFV) of their infrastructure to reduce both CAPEX and OPEX while maintaining high carrier-grade service levels. The savings in CAPEX and OPEX come from being able to dynamically assign Virtualized Network Functions (VNFs) to various standard servers in their infrastructure to meet varying workload demands. Similar to Cloud Service resource allocation, such dynamic VNF placement can be automated to optimize various goals of a Telco Operator. Adoption of NFV by Telcos allows dynamic fine-grained Service Function Chaining (SFC) where various service functions chains can be strung together with deployed VNFs using SDN-enabled dynamic route control. Similarly Enterprises are adopting NFV for deployment of network services in their infrastructure.

As a popular use case of NFV in the Telco, SFC is usually deployed in the Telco's datacenters in their PoPs or central offices. The prosperity of SFCs highly depend on its performance. We provide an approach to SFC provisioning within a datacenter. SFC provisioning comprises determining how VNFs are placed on nodes in the datacenter and how VNF instances are assigned to SFCs. The placement and assignment affects the traffic routing from the SFC through the datacenter's network. Our SFC provisioning system, called *Stringer*, allows operators to balance the conflicting goals of minimizing infrastructure resources and end-to-end latency for meeting their respective SLAs.

In any SFC, the relative location of the VNFs will affect the end-to-end latency incurred by the packets traversing the particular SFC. A poor placement will cause the flow to traverse the same path-segments back and forth inside the network, increasing the network delay and consuming more bandwidth.

The SFC provisioning problem entails choosing where to place instances of VNFs on servers in a NFV infrastructure to accommodate the traffic for a given set of SFC requests. Each service chain is a sequence of VNFs that processes a stream of network packets flowing it at a certain rate of network packets flowing through a sequence of VNFs at a certain rate. Network traffic for a given service chain must visit the chain's sequence of VNFs in

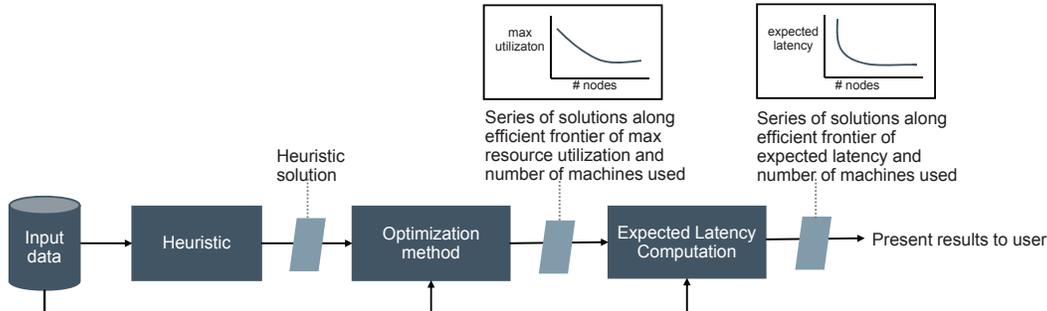


Figure 1: Flowchart of System: The input data, consisting of the service chains, their requirements and the data center network topology, are given to the heuristic. The basic solution from the heuristic provides an initial solution for the mixed integer program. Expected latency is then computed for each of the solutions from the optimizer. The user can then choose from a menu of solutions differing in expected latency and number of servers used.

the specified order. For example, a service chain may require packets to follow the VNF sequence: load balancer, network address translator, and firewall. In the SFC provisioning problem, one must place (possibly multiple) instances of each VNF on servers, and choose the route(s) for each service chain, in such a way that the network can accommodate the traffic for as many service chains according to their priorities. Service chains may share VNF instances. Moreover, the traffic for a given service chain may be split among multiple paths in the network when multiple instances of a specific VNF are used.

Our work differs from prior work in VNF placement in several important ways. One key difference is in the placement objective. Operators have multiple competing goals to consider when placing VNFs. A service provider may want to use as few servers as possible in order to minimize operating costs and leave open servers for future needs ([1, 2]). At the same time, the operator must ensure low end-to-end network latency for his customers. These objectives are in direct conflict. While some prior work proposes multiple alternative objectives ([3, 4, 5]), ours is the first, to our knowledge, that provides a flexible way to trade-off these competing goals in SFC provisioning. Moreover, unlike [5], our optimization model employs only linear constraints to model maximum utilization.

Another important difference is in the way packet delays are modeled. Most prior approaches [1, 2, 3, 4] model network latency with a known fixed delay when packets pass through VNFs, nodes and edges. They do not consider that latency depends on network traffic: packets traveling through congested network resources face much longer queueing delays than at uncongested ones. Expected latency depends on VNF placement and routing decisions, and the implied utilization of network resources, in a complex and non-linear way, which explains why prior work models latency in a simplified, utilization-independent way. Figure 3, which shows the non-linear relationship between expected latency and utilization, highlights what is lost in this simplified approach. A single congested server or switch can dramatically increase latency for all service chains using that resource. If congestion is not explicitly modeled, such effects are ignored.

Some prior work (e.g., [6]) decomposes VNF placement into two separate problems: first determining the number of instances of each VNF and routing among instances, and then placing VNF instances. Steering[7] assumes the number of instances of each VNF is given. In contrast, our approach considers both problems simultaneously. Moreover, our approach considers the utilization of servers whereas [6] considers only switch traffic. In this work, we focus on the chaining of inline services, e.g. firewall, load balancer, IDS. These VNFs operate on their own, with little dependencies across VNFs. Those VNFs with complex inter-dependences, e.g. EPC in cellular core network [8] are not the focus of this paper.

2 *Stringer: Our SFC Placement System*

Stringer provides the ability to the operators to select their operating point for trading-off resource usage and end-to-end SFC latency. Figure 1 shows the architecture and flowchart of *Stringer* system.

There are three main contributions of this work. The first is a scalable heuristic that seeks to minimize the

maximum utilization over all nodes. The second is our MIP-based placement approach that competing objectives of minimizing congestion-induced latency and minimizing the number of servers used. It minimizes a weighted combination of two metrics: (1) the number of servers used to host VNF instances, and (2) the maximum utilization over network resources, which we use as a proxy for latency.

The optimization method generates multiple SFC provisioning solutions for different relative weightings of the two objectives, thereby generating solutions along the efficient frontier of number of servers and maximum resource utilization. The MIP and heuristic each have advantages: The MIP provides optimal benchmarks, but does not scale to very large size networks. The heuristic is fast and scalable, can be used to provide an initial solution that speeds up the MIP solution process, and generates solutions that are close to the efficient frontier of latency and node usage.

Our third contribution is a method to evaluate a SFC provisioning solution. Evaluating the performance of a placement strategy on a real world testbed of large size is not likely in practice. Therefore, we propose a queueing-theoretic model of the network which allows us to simulate the average expected latency associated with any given SFC provisioning solution under mild assumptions on the network traffic. Our model differs from standard M/M/1 queueing models in that it accounts for the fact that network elements have finite buffers and packets are dropped when they arrive to full buffers. Our expression for average expected latency reflects the possibility of packets being dropped and re-sent.

Combined, these three elements create our *Stringer* system that generates a set of SFC provisioning solutions varying in resource usage and performance. When presented with an array of solutions reflecting different tradeoffs between competing objectives, the operator can then make an informed choice about how to place VNFs and route the SFCs accordingly.

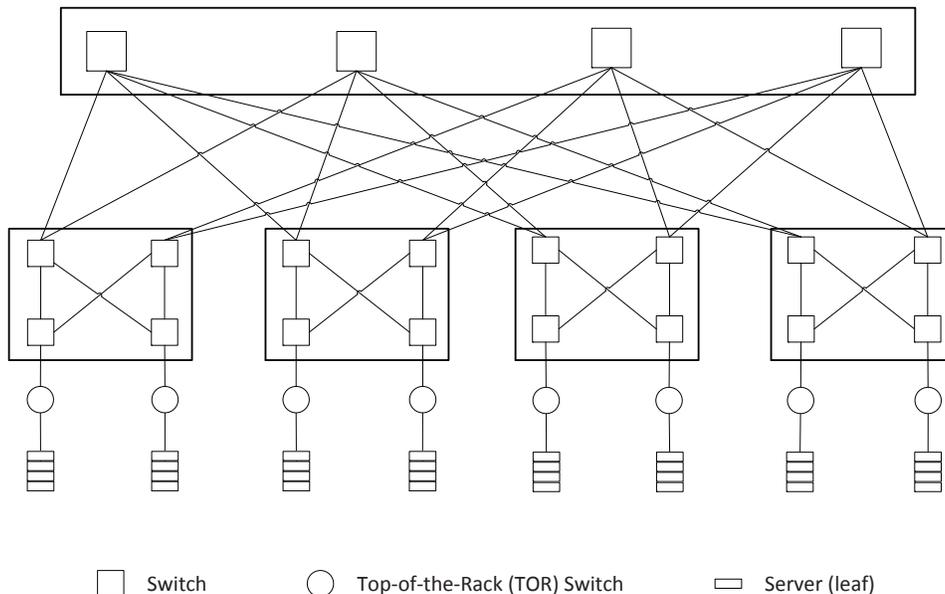


Figure 2: Tree Topology: The most common network topology used in datacenters is the FAT-tree topology where servers communicate with one another through a hierarchical arrangement of network switches as shown in the figure. The servers are connected to the TOR (top-of-the-rack) switches. Our model abstracts the underlying FAT tree topology by aggregating each level and cluster of switches into a single large virtual switch. This switch aggregation simplifies the model because each pair of servers has a unique path connecting them in the modified network.

3 Preliminaries for Stringer

The inputs to Stringer fall into three categories: the physical network topology, the virtualized network functions, and the service chains.

Physical network The physical network topology is a bi-directional graph with the property that each ordered pair of nodes has a unique acyclic directed path to each other. The underlying structure is a tree, consisting of switches (including a root switch r), and servers which are leaves in the tree. Let N denote the set of nodes (switches and servers) and $L \subset N$ be the set of servers. An example network is shown in Figure 2. Let μ_n be the processing rate, in packets per second, associated with any node $n \in N$.

Virtualized Network Functions Let V denote a set of VNF types. Instances of these VNF types must to be assigned to servers in the physical network in order to accommodate service chains. Multiple instances of a given VNF type v may be assigned. We assume that a server in the network can accommodate at most one virtual network function instance, although that constraint can easily be relaxed.

Service Chains Let C denote the set of service chains to be mapped to the network. Service chain $c \in C$ comprises a (possibly repeating) sequence of VNF types. The service chain c is a Poisson process with arrival rate of λ_c packets per second. Traffic for service chain c enters the physical network through the root node, visits each function according to the chain's function sequence, and then departs the network from the root node. Let $\Lambda = \sum_c \lambda_c$ be the sum of arrival rates of all service chains.

4 Expected Latency Evaluation

We show how to compute the expected latency of any packet entering the system, assuming that VNF placement and service chain routing has already been determined. The expected latency of a packet entering the network depends on the service chain with which the packet is associated. Let $E(T_c)$ represent the expected latency of packets in a given service chain $c \in C$. The expected latency $E(\mathcal{J})$ of a randomly selected arriving packet is equal to the sum over all service chains $c \in C$ of the probability (λ_c/Λ) that the packet is associated with chain c times $E(T_c)$:

$$E(\mathcal{J}) = \sum_{c \in C} \frac{\lambda_c}{\Lambda} E(T_c) \quad (1)$$

$$E(T_c) := E(T_{1 \rightarrow n}) \quad (2)$$

where $E(T_{1 \rightarrow n})$ represents the expected latency for a packet to visit the sequence of nodes as $\{1, 2, \dots, n\}$ in N_c , for $n = 1, 2, \dots, |N_c|$.

The model to estimate $E(T_c)$ has two key considerations: 1) the latency τ_n at each node $n \in N_c$, which is independent of the latency at other nodes but is dependent on all service chains' traffic through node n and 2) the probability that a packet may drop at any node n , which would require a resend of the packet from the source up to n . The retransmission of packets is due to the Transmission Control Protocol (TCP). TCP ensures that all packets will arrive at the destination. If any packet is dropped during transmission, TCP will resend the packet from the source until they reach the destination. The expected latency computation must factor in a packet's expected queuing delay at each node as well as extra time incurred due to resent packets.

$$E(T_{1 \rightarrow 1}) = \tau_1 \quad (3)$$

$$E(T_{1 \rightarrow n}) = \tau_n + E(R_n)E(T_{1 \rightarrow n-1}) \text{ for } n = 2, \dots, |N_c| \quad (4)$$

The recursion in Equation 4 is due to the TCP protocol which resends dropped packets when buffers are full and $E(R_n)$ is the expected number of retries.

$$\tau_n = \frac{\rho_n - [1 + K_n(1 - \rho_n)]\rho_n^{K_n+1}}{\lambda_n(1 - \rho_n)(1 - \rho_n^{K_n})} \quad (5)$$

$$\rho_n = \frac{\lambda_n}{\mu_n} \quad (6)$$

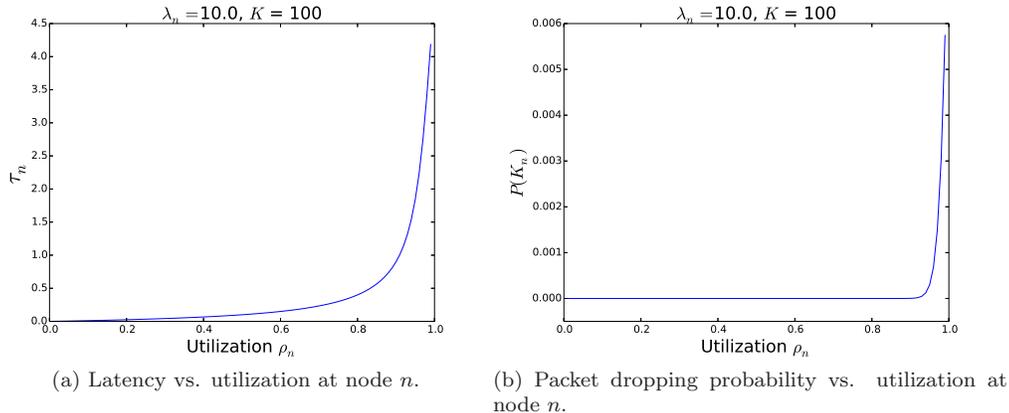


Figure 3: Effects of node utilization. From the two figures shown here, one can see that the node utilization ρ_n has a non-linear effect on both the latency and the probability of dropped packets at each node n . Beyond a certain threshold, the value of latency and probability grows exponentially. Prior to the placement of VNFs, it is hard to predict where the threshold is. So the observation of these charts motivates our optimization strategy to minimize ρ_n as much as possible across every node $n \in N$ in the network.

where λ_n is the incoming rate of packets to the node n , μ_n is the rate which the node n processes packets and K_n is the buffer capacity at node n . Equation 5 is a standard formula for the M/M/1/K queueing model. Figure 3a shows τ_n vs ρ_n .

$$E(R_n) = \frac{1}{1 - P(K_n)} \quad (7)$$

$$P(K_n) = \frac{1 - \rho_n}{1 - \rho_n^{K_n+1}} \rho_n^{K_n} \quad (8)$$

$P(K_n)$ gives the probability that the buffer at node n is full when a packet arrives at node n . Appendix A provides more details on the derivations of these equations and their significance. Figure 3b shows $P(K_n)$ vs ρ_n .

5 Round Robin Heuristic

We propose a heuristic that seeks feasible placements for the incoming service chains while making an effort to minimize the overall latency. The basic principle of this heuristic is to distribute the network traffic among different top-of-the-rack (TOR) switches as much as possible, to reduce the maximum node utilization over all nodes. The heuristic achieves this in two ways: first by distributing service chains across TORs and then by limiting the utilization of each machine.

Throughout the execution, we maintain an upper limit on machine utilization. The heuristic considers each service chain in succession, placing the chain in a TOR which is different from the TOR of the previous service chain. Each VNF of the service chain is placed on a machine that can accommodate its traffic in the chosen TOR. If a machine has hit its utilization limit, we allow multiple machines that host the same type of VNF to fulfill the service chain. If the machines within a TOR cannot handle all the VNFs for a service chain, then the remaining VNFs are placed on the next TOR. If all machines in the datacenter have been used, the upper limit is adjusted upwards to accommodate more traffic and fulfill more service chains.

6 Optimization Method

The MIP-based optimization method produces a set of solutions to the SFC provisioning problem, each representing a different tradeoff between network performance and resource usage. Here we provide an overview of the formulation; details are presented in [9].

The decisions variables include: binary variables indicating whether an instance of a particular VNF is hosted on a particular server; continuous variables representing the fraction of a service chain’s traffic that passes through a particular server and between any pair of servers; continuous variables representing the total bandwidth entering any node in the network; and lastly, the maximum utilization over all nodes in the network. The constraints ensure that flow for each service chain is conserved at each node and that the solution does not use more than the available network resources.

The objective is to minimize a weighted combination of the fraction of servers used to host VNFs and the maximum utilization over all nodes in the network. A weighting parameter $\beta \in [0, 1]$ is used to set the relative priority of these two objectives. When $\beta = 0$, the objective reduces to minimizing the maximum utilization over all nodes in the network, thus distributing the traffic as uniformly as possible in order to reduce the highest utilization over all nodes. When $\beta = 1$, the objective becomes minimizing the total number of nodes used to host VNFs. A placement which minimizes the number of VNFs tends to concentrate traffic in part of the network, leaving other network resources unused. Solving the MIP over a range of $\beta \in [0, 1]$ yields a set of solutions along the efficient frontier of maximum node utilization and number of servers, each representing a different tradeoff between performance and server usage. For each new value of β , the preceding solution can be used as a starting point for the MIP run, speeding its execution. The heuristic solution is used as a starting solution for the first MIP.

Our formulation has two novel features compared to prior MIP approaches to SFC provisioning. One is the use of a weighted objective function to generate alternative solutions trading off performance and resource cost. A second novel feature is a method of modeling maximum utilization using only linear constraints. Node utilization is the ratio of a node’s incoming bandwidth to its capacity. Both bandwidth and capacity are functions of decision variables (capacity at a node depends on the type of VNF assigned to it), and thus utilization is naturally nonlinear in decision variables. Nonlinear constraints make MIP models significantly less tractable. We employ a novel approach to linearize the maximum utilization by including constraints for each possible VNF type assigned to a node, and using penalties to activate only the applicable constraints. Details can be found in [9].

Table 1: Summary of Results

# Servers	# SFCs	Avg Random Time(s)	Avg Heuristic Time(s)	Avg MIP Time(s)	Avg Random Success (%)	Avg Heuristic Success (%)	Avg MIP Success (%)	Avg Optimality Gap	Avg Latency Gap
8	5	0.35	0.49	4.63	70%	94%	90%	9%	0%
16	10	0.35	0.51	299.63	73%	96%	100%	17%	8%
32	20	0.36	0.49	602.42	69%	100%	100%	18%	7%
64	40	0.32	0.45		75%	100%			
128	80	0.37	0.52		74%	100%			
256	160	0.39	0.54		74%	100%			
512	320	0.42	0.76		74%	100%			
1024	640	0.64	1.15		76%	100%			
2048	1280	1.41	5.02		75%	100%			
4096	2560	5.30	4.59		75%	100%			

7 Numerical Results

7.1 Efficient Frontier

We present an example of the efficient frontier that is generated by the MIP along with results from the round robin heuristic. This example corresponds to the network in Figure 2, in which each cluster of switches is aggregated into a single switch. There are 10 service chains to be deployed, each with up to 4 VNFs.

The MIP generates a range of SFC provisioning solutions using between 10 and 32 servers for this example. Properties of these solutions are shown in Figure 4. Figure 4a shows MIP solutions along the efficient frontier of server usage vs. maximum node utilization.

Unlike the MIP, the heuristic generally produces only one solution. However, by applying it to successively smaller subnetworks, we can generate multiple solutions that tradeoff latency and server usage, just as the MIP does. Three such heuristic solutions are also shown in Figure 4a, corresponding to three different versions of the original network: (1) the full network shown in Figure 2, (2) a subnetwork in which one aggregation switch and its descendants are removed, and (3) a subnetwork in which two aggregation switches and descendants are removed.

The heuristic solutions are not far from the efficient frontier, indicating that it achieves low maximum utilization relative to the number of servers it uses to host VNFs. The chart in Figure 4b shows the same set of solutions, in this case highlighting the tradeoff between expected latency and server usage. Note that the heuristic compares even more favorably to the MIP solutions, in that its solutions lie very close to the MIP solution curve. Figure 4c shows directly how latency varies with maximum node utilization in the MIP and heuristic solutions. In particular, it shows how expected latency of the MIP solutions increases with maximum utilization, and grows steeply as maximum utilization approaches 100%, as in Figure 3a. These properties support the choice of maximum node utilization as a good proxy objective for expected latency.

7.2 Comparison of Solution Quality

Table 1 summarizes the timing and results for the round-robin heuristic and the MIP for 100 randomly generated test problems. We also share results for a random placement approach as a baseline.

We randomly generated 10 problems for each of 10 topologies that vary in number of servers, from 8 through 4096, and number of service chains to be routed. The 10 problems for each topology vary in the specific VNFs and volume of traffic required for each service chain. For each problem, we ran the heuristic, and then ran the MIP with the objective of minimizing the maximum node utilization subject to the constraint that it uses no more servers than the heuristic used for the same problem. This approach allows us to compare the maximum utilization of the heuristic with that of the MIP for a fixed number of servers. We also ran random placement for each problem.

The average time required by the heuristic is at most 5 seconds for all topologies, even for problems with over 4,000 servers and 2,000 service chains. For the MIP, we set a time limit of 300 seconds for problems with 16 or fewer servers and 600 seconds for 32-server problems; it is too computationally intensive to run for larger topologies.

Three aspects of solution quality are shown in the table. A first measure is the percentage of service chains deployed. The heuristic may not deploy all service chains, if it runs out of servers. Its success rate, shown in column 5 of Table 1, depends on the network capacity, service chain VNF requirements and the service chain traffic requirements. For the MIP we require that all service chains are routed and thus if a solution is found for a given test problem, the MIP service chain deployment rate is 100%, and otherwise 0%. Random placement never successfully deployed all service chains; its average success rate was at most 75%.

A second measure of quality is in the optimality gap, available only for the set of problems for which the MIP was run. This is the percentage difference between the maximum node utilization achieved by the heuristic and that achieved by the MIP, for the same number of servers used. We present the optimality gap for problems which the heuristic deploys all service chains. The average optimality gap is 18% or less. The worst case optimality gap over all test problems (not shown in the table) is 34%. Since random placement never deployed all service chains, its solution is not comparable to other approaches that did; we do not include its optimality gap in the table.

A third quality measure is in the latency gap shown for problems on which the MIP and heuristic deployed all service chains. The latency gap is the percentage difference between the latency of the heuristic solution and the MIP solution as a percentage of MIP latency. The average latency gap ranges from 0-8%, highlighting the effectiveness of the heuristic and the choice of minimizing maximum utilization as a proxy objective for minimizing latency.

8 Related Work

VM placement: The VM placement has been studied extensively in the cloud computing literature [10, 11, 12, 13, 14, 15]. These work develop heuristics to assign VMs of a tenant close to each other to reduce the overall bandwidth

consumption. Different from existing work, our work focus on the placement of VNFs, in the context of service chaining.

Service chaining: Simple [16] proposes a SDN framework to route traffic through a flexible set of service chains while balancing the load across Network Functions. FlowTags [17] can support dynamic service chaining. Our work is complimentary to these service chain implementation mechanisms. While these work focus on the techniques to realize flexible routing, we provide algorithms that decide on the routes.

NFV deployment: Lukovszki et al. [18] presents an approximate algorithm and an integer programming exact solution. [19] also present a MIP model for NFV placement that minimizes the maximum utilization over all links and switches. However, these models considers only the utilization of links and of servers, excluding the switches' utilization. [8, 7] measured bottlenecks and maximize throughput in inter-datacenter networks while we focus on intra-datacenter networks.

9 Conclusions

This work addresses the problem of choosing the physical locations of virtual network functions required for service chains, and routing service chain traffic which we term as SFC provisioning problem. We offer a system Stringer consisting of a scalable placement heuristic, an optimization-based approach for generating a series of alternative placement solutions reflecting different tradeoffs between performance and resource usage, and a queueing-theory-based method for estimating average latency per packet under a given VNF placement and routing solution.

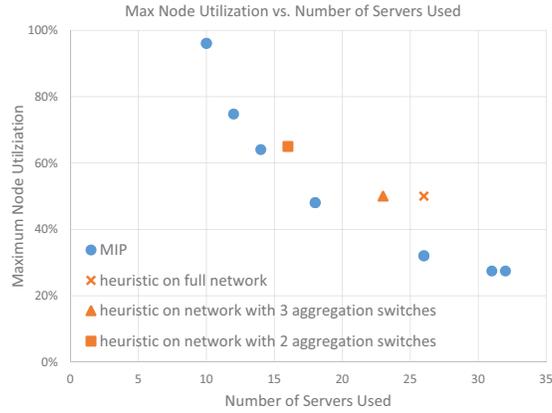
Our experiments comparing the performance of MIP and heuristic show that the heuristic is significantly faster and has an average optimality gap of at most 18%. The heuristic can also be used to generate an efficient frontier of solutions, by running it on a succession of subnetworks.

Stringer has several potential extensions. One way to improve scalability of the optimization is to apply a hierarchical approach, in which we first assign service chains to subnetworks associated with aggregation switches, and then solve the SFC provisioning problem within each subnetwork. We are also planning to collect data from real networks to validate our queueing theoretic latency estimation.

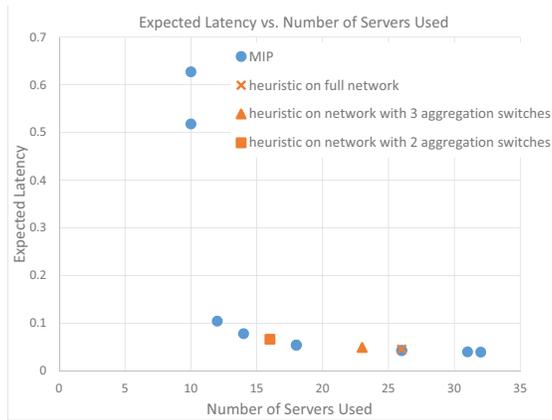
References

- [1] H. Moens and F. De Turck, “Vnf-p: A model for efficient placement of virtualized network functions,” ser. CNSM 2014, Nov 2014, pp. 418–423.
- [2] M. C. Luizelli, L. R. Bays, L. S. Buriol, M. P. Barcellos, and L. P. Gasparry, “Piecing together the nfv provisioning puzzle: Efficient placement and chaining of virtual network functions,” ser. IM 2015, May 2015, pp. 98–106.
- [3] M. F. Bari, S. R. Chowdhury, R. Ahmed, and R. Boutaba, “On orchestrating virtual network functions,” ser. CNSM 2015, M. Tortonesi, J. Schönwälder, E. R. M. Madeira, C. Schmitt, and J. Serrat, Eds. IEEE Computer Society, 2015, pp. 50–56. [Online]. Available: <http://dx.doi.org/10.1109/CNSM.2015.7367338>
- [4] S. Mehraghdam, M. Keller, and H. Karl, “Specifying and placing chains of virtual network functions,” *CoRR*, vol. abs/1406.1058, 2014. [Online]. Available: <http://arxiv.org/abs/1406.1058>
- [5] V. Sekar, N. Egi, S. Ratnasamy, M. K. Reiter, and G. Shi, “Design and implementation of a consolidated middlebox architecture,” in *Presented as part of the 9th USENIX Symposium on Networked Systems Design and Implementation (NSDI 12)*, 2012, pp. 323–336.
- [6] S. Palkar, C. Lan, S. Han, K. Jang, A. Panda, S. Ratnasamy, L. Rizzo, and S. Shenker, “E2: a framework for nfv applications,” in *Proceedings of the 25th Symposium on Operating Systems Principles*. ACM, 2015, pp. 121–136.
- [7] Y. Zhang, N. Beheshti, L. Beliveau, G. Lefebvre, R. Manghirmalani, R. Mishra, R. Patneyt, M. Shirazipour, R. Subrahmaniam, C. Truchan *et al.*, “Steering: A software-defined networking for inline service chaining,” in *Network Protocols (ICNP), 2013 21st IEEE International Conference on*. IEEE, 2013, pp. 1–10.

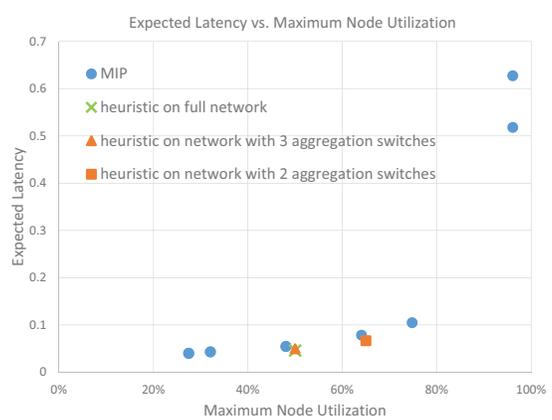
- [8] A. S. Rajan, S. Gobriel, C. Maciocco, K. B. Ramia, S. Kapury, A. Singhy, J. Ermanz, V. Gopalakrishnan, and R. Janaz, “Understanding the bottlenecks in virtualizing cellular core network functions,” in *Local and Metropolitan Area Networks (LANMAN), 2015 IEEE International Workshop on*. IEEE, 2015, pp. 1–6.
- [9] F. C. Chua, J. Ward, Y. Zhang, P. Sharma, and B. A. Huberman, “Stringer: Balancing latency and resource usage in service function chain provisioning,” *Hewlett Packard Labs Tech Report HPE-2016-31*, 2016.
- [10] J. Lee, Y. Turner, M. Lee, L. Popa, S. Banerjee, J.-M. Kang, and P. Sharma, “Application-driven bandwidth guarantees in datacenters,” ser. SIGCOMM ’14. New York, NY, USA: ACM, 2014, pp. 467–478. [Online]. Available: <http://doi.acm.org/10.1145/2619239.2626326>
- [11] H. Ballani, P. Costa, T. Karagiannis, and A. Rowstron, “Towards predictable datacenter networks,” ser. SIGCOMM 2011, 2011.
- [12] V. Jeyakumar, M. Alizadeh, D. Mazières, B. Prabhakar, C. Kim, and A. Greenberg, “Eyeq: Practical network performance isolation at the edge,” ser. NSDI 2013, 2013, pp. 297–312.
- [13] S. H. M. C. M. C. Joe Wenjie Jiang, Tian Lan, “Joint vm placement and routing for data center traffic engineering,” in *INFOCOM*, 2012.
- [14] L. Popa, G. Kumar, M. Chowdhury, A. Krishnamurthy, S. Ratnasamy, and I. Stoica, “Faircloud: Sharing the network in cloud computing,” *SIGCOMM C.C.R.*, vol. 42, no. 4, pp. 187–198, August 2012.
- [15] D. Xie, N. Ding, Y. C. Hu, and R. Kompella, “The only constant is change: Incorporating time-varying network reservations in data centers,” ser. SIGCOMM ’12, 2012, pp. 199–210.
- [16] Z. A. Qazi, C.-C. Tu, L. Chiang, R. Miao, V. Sekar, and M. Yu, “Simple-fying middlebox policy enforcement using sdn,” ser. SIGCOMM ’13, 2013, pp. 27–38.
- [17] S. K. Fayazbakhsh, L. Chiang, V. Sekar, M. Yu, and J. C. Mogul, “Enforcing network-wide policies in the presence of dynamic middlebox actions using flowtags,” ser. NSDI’14, 2014, pp. 533–546.
- [18] T. Lukovszki, M. Rost, and S. Schmid, “It’s a match!: Near-optimal and incremental middlebox deployment,” *SIGCOMM C.C.R.*, vol. 46, no. 1, pp. 30–36, January 2016.
- [19] A. Mohammadkhan, S. Ghapani, G. Liu, W. Zhang, K. Ramakrishnan, and T. Wood, “Virtual function placement and traffic steering in flexible and dynamic software defined networks,” ser. LANMAN 2015. IEEE, 2015, pp. 1–6.



(a) Maximum node utilization vs. number of servers used



(b) Expected latency vs. number of servers used



(c) Expected latency vs. maximum node utilization

Figure 4: Example Solution Metrics

A Additional Details for Expected Latency Derivation

Recall from Section 4 that we showed the steps to calculate $E(\mathcal{T})$.

$$E(\mathcal{T}) = \sum_{c \in \mathcal{C}} \frac{\lambda_c}{\Lambda} E(T_c) \quad (9)$$

The question now is how to estimate $E(T_c)$, the average amount of time each packet takes to go through the set of services required by service chain c . This latency depends on the placement of the service chain's VNFs in the network, the nodes along the paths between successive VNFs, and the amount of traffic (for all service chains) through each of those nodes. Because the VNF placement and chain routing has already been determined, we know which nodes (switches and servers) N_c that packets of c will flow through. Denote the sequence of nodes in N_c as $\{1, 2, \dots, n, \dots, |N_c|\}$ where n indicates the n th node that the packets will pass through, 1 as the source of the packets and $|N_c|$ as the final destination of the packets.

The model to estimate $E(T_c)$ has two key considerations: 1) the latency at each node $n \in N_c$, which is independent of the latency at other nodes but is dependent on all service chains' traffic through node n and 2) the probability that a packet may drop at any node n , which would require a resend of the packet from the source up to n .

Queueing Theoretic Latency Estimation We model each node n as a finite capacity, single server queue where packets are processed one at a time while other packets wait in the queue of size $K_n - 1$ for their turn to be processed on a First-In-First-Out (FIFO) policy. We assume that the packets arrive at node n according to a Poisson process with rate λ_n , a rate which reflects the traffic from all service chains routed through node n .

Unlike the case of M/M/1 queues, we model each node with a finite queue size, and packets can drop from the node if the queue is full. Such packet drops result in an outgoing rate that is less than the incoming rate. Although it is possible to derive an approximation for the outgoing rate for an acyclic network, typical networks in data centers have cyclic dependencies between the outgoing and incoming rate, which makes it hard to estimate it correctly analytically. An example of a network cyclic dependency is traffic which flow through a switch and would later flow back to the same switch after processing at the leaves beneath the switch itself. To simplify the model, we assume that the outgoing rate from node n is equal to the incoming rate.

The time to process each packet at the node n (excluding queueing time) follows an exponential distribution μ_n , and it is assumed that $\lambda_n < \mu_n$. These assumptions allow us to use the well-known formula in the M/M/1/K queueing literature to estimate the expected latency τ_n of a packet at node n , including both queueing time and service time at n .

$$\tau_n = \frac{\rho_n - [1 + K_n(1 - \rho_n)]\rho_n^{K_n+1}}{\lambda_n(1 - \rho_n)(1 - \rho_n^{K_n})} \quad (10)$$

$$\rho_n = \frac{\lambda_n}{\mu_n} \quad (11)$$

Packet Loss at Each Node n Since we model each node n with a finite capacity queue of length $K_n - 1$, packets that arrive to find the queue full will be discarded. The probability of packets dropping in this manner is equal to the probability that there are K_n packets in the system (one packet being processed and $K_n - 1$ packets in the queue). For an M/M/1/K queue, the probability of having K_n packets in the system is,

$$P(K_n) = \frac{1 - \rho_n}{1 - \rho_n^{K_n+1}} \rho_n^{K_n} \quad (12)$$

In software applications that use the Transmission Control Protocol (TCP) for transferring network packets, TCP ensures that all packets will arrive at the destination. If any packet is dropped during transmission, TCP will resend the packet from the source until they reach the destination. The expected latency computation must factor in a packet's expected queueing delay at each node as well as extra time incurred due to resent packets. Let $E(T_{1 \rightarrow n})$ represent the expected latency for a packet to visit the sequence of nodes as $\{1, 2, \dots, n\}$ in N_c , for $n = 1, 2, \dots, |N_c|$.

Thus, $E(T_c) = E(T_{1 \rightarrow |N_c|})$. We define a recursive formula for the latency as follows:

$$E(T_{1 \rightarrow 1}) = \tau_1 \tag{13}$$

$$E(T_{1 \rightarrow n}) = \tau_n + E(R_n)E(T_{1 \rightarrow n-1}) \text{ for } n = 2, \dots, |N_c| \tag{14}$$

where $E(R_n)$ is the expected number of resends required to transmit the packet from node 1 to node n . To compute $E(R_n)$, note that $P(R_n = m)$ is the probability of the packet dropping $m - 1$ times at node n and succeeding on the m th time. Thus $E(R_n)$ is derived as follows:

$$P(R_n = m) = P(K_n)^{m-1} [1 - P(K_n)] \tag{15}$$

$$E(R_n) = \sum_{m=1}^{\infty} m \cdot P(R_n = m) \tag{16}$$

$$= \frac{1}{1 - P(K_n)} \tag{17}$$

Using Equations 9 to 17, we can evaluate the expected latency $E(\mathcal{J})$.

While the derivation in this section allows us to evaluate the expected latency of each service chain given a particular SFC provisioning solution, it does not lend itself to optimizing for latency when making placement decisions, since the placement decisions (and implied congestion) affect latency in a complex and nonlinear way.

However, the derivation offers insight into the importance of node utilization in expected latency. Consider, for example, the relationship between the expected latency τ_n at node n and the utilization ρ_n at node n , illustrated in Figure 3a for arrival rate $\lambda_n = 10$ and queue capacity $K_n = 100$. Latency grows steeply as utilization approaches 100%. Moreover, the relationship between packet dropping probability $P(K_n)$ and node utilization ρ_n reveals the importance of utilization ρ_n in preventing packets from being dropped. Figure 3b illustrates how the packet dropping probability $P(K_n)$ grows abruptly with as node utilization approaches 100% under the same assumptions on λ_n and K_n .

These objectives suggest a simple but powerful objective to use in SFC provisioning. By making placement decisions to minimize the maximum node utilization in the physical network, we can both avoid packet loss and reduce latency. The SFC provisioning methods described below pursue the goal of minimizing the maximum node utilization.

B Additional Details for Mixed Integer Program

This appendix describes the constraints of the MIP formulation introduced in Section 6.

Model Parameters:

- N : the set of all nodes in the network (servers and switches).
- $L \subset N$: the set of servers, which are leaves in the tree network.
- $r \in N$: the root node.
- $P_{n,m}$: the set of nodes in the unique acyclic path from node n to m , including the destination m but excluding the origin n .
- μ_n : the processing rate, in packets per second, associated with switch $n \in N \setminus L$.
- S : the set of different server types.
- $s_l \in S$: the machine type associated with server $l \in L$.
- V : the set of VNF types. Instances of these VNF types must to be assigned to servers in the physical network in order to accommodate service chains.
- γ_v^s : the processing rate, in packets per second, of VNF type $v \in V$ when assigned to server type $s \in S$.

- C : the set of service chains to be mapped to the network.
- q_c : the length of the sequence of VNFs in service chain c .
- $\alpha_{i,v}^c$: a binary parameter indicating whether the i th service in chain c is of type v .
- λ_c : arrival rate, in packets per second, for chain c .
- M : a large positive scalar. For example, any $M > \max\{1, \max_{s,v}\{\gamma_v^s\}\}$ is suitable.
- $\beta \in [0, 1]$: a parameter representing the relative weight between two metrics, number of servers used and maximum utilization, in the objective function.

Decision Variables: The decision variables describe the assignment of VNF instances to leaf nodes, the mapping of each service chain to one or more paths in the network, the volume of flow for each chain along each of its paths, the rate of traffic into each node, and performance metrics associated with the solution.

- $x_{v,l} \in \{0, 1\}$ indicates whether an instance of VNF type v is placed on leaf l .
- $y_{i,l}^c \in [0, 1]$ is the fraction of traffic for the i th function in service chain c that is served by leaf node l .
- $z_{i,k,l}^c \in [0, 1]$ is the fraction of traffic going from the i th to $(i+1)$ st function in service chain c that travels from leaf node k to leaf node l .
- $b_k \geq 0$ is the total traffic rate in packets per second into node $k \in N$.
- ρ is the maximum node utilization over all nodes in the network.

Constraints The MIP constraints ensure that flow for each service chain is conserved at each node, that the solution does not use more than the available network resources, and that the maximum utilization metric is measured.

$$\sum_{v \in V} x_{v,l} \leq 1, \quad l \in L \quad (18)$$

$$y_{i,l}^c \leq \sum_v \alpha_{i,v}^c x_{v,l}, \quad c \in C, i \leq q_c, l \in L \quad (19)$$

$$\sum_{l \in L} y_{i,l}^c = 1, \quad c \in C, i \leq q_c \quad (20)$$

$$z_{i,k,l}^c \leq y_{i,k}^c, \quad c \in C, i < q_c, k, l \in L \quad (21)$$

$$z_{i,k,l}^c \leq y_{i+1,l}^c, \quad c \in C, i < q_c, k, l \in L \quad (22)$$

$$\sum_{k,l \in L} z_{i,k,l}^c = 1, \quad c \in C, i < q_c \quad (23)$$

$$y_{1,k}^c + \sum_{\substack{m \in L \\ i < q_c}} z_{i,m,k}^c = \sum_{\substack{m \in L \\ i < q_c}} z_{i,k,m}^c + y_{q_c,k}^c, \quad c \in C, k \in L \quad (24)$$

$$b_k = \sum_{c \in C} \lambda_c \left(1 + \sum_{\substack{l \in L: \\ k \in P_{r,l}}} y_{1,l}^c + \sum_{\substack{m \in L: \\ k \in P_{m,r}}} y_{q_c,m}^c + \sum_{\substack{i < q_c \\ l,m \in L: \\ k \in P_{l,m}}} z_{i,l,m}^c \right), \quad k \in N \setminus r \quad (25)$$

$$b_r = \sum_{c \in C} \lambda_c \left(1 + \sum_{m \in L} y_{q_c,m}^c + \sum_{\substack{i < q_c \\ l,m \in L: \\ r \in P_{l,m}}} z_{i,l,m}^c \right) \quad (26)$$

$$b_n \leq \mu_n, \quad n \in N \setminus L \quad (27)$$

$$b_l \leq \sum_{v \in V} \gamma_v^{sl} x_{v,l}, \quad l \in L \quad (28)$$

$$\rho \geq \frac{b_n}{\mu_n}, \quad n \in N \setminus L \quad (29)$$

$$\rho \geq \frac{b_l}{\gamma_v^{sl}} - M(1 - x_{v,l}), \quad l \in L, v \in V \quad (30)$$

The constraint 18 ensures that each server $l \in L$ can have at most one VNF type assigned to it. Constraint 19 enforces that the i th function in service chain c can only be placed on a server hosting the VNF type associated with the i th function. Constraint 20 requires that then the total traffic for its i th function must be placed.

We need inequalities 21 and 22 to ensure that $z_{i,k,l}^c$ does not exceed $y_{i,k}^c$ or $y_{i+1,l}^c$ for each chain c , for each function index $i < q_c$, and each physical server pair $k, l \in L$. Constraint 23 implies that the total required traffic rate from the i th function to the $(i+1)$ st function in service chain c must be allocated. Flow conservation constraint 24 requires that the traffic for service chain c into server k (the left hand side) must equal the traffic service chain c exiting k . Constraint 25 defines the total traffic rate b_k into each non-root node $k \in N \setminus \{r\}$. For a given service chain c , the first term ($\lambda_c \sum_{l \in L: k \in P_{r,l}} f_l^c$) captures the traffic into switch k coming from the root to any server l hosting the first function in the chain, the second term ($\lambda_c \sum_{m \in L: k \in P_{m,r}} h_m^c$) captures traffic into switch k heading toward the root from any server m hosting the last function in the chain, and the remaining term captures traffic between any pair of servers l and m hosting consecutive functions in the chain for which their path passes through switch k . Constraint 26 defines the total traffic rate b_r into the root node r . The first term captures the traffic into the root r coming from outside the network (λ_c), the second describes the traffic into r from any server m hosting the last function in the chain (term $\lambda_c \sum_{m \in L} y_{q_c,m}^c$), and the final term captures traffic between any pair of servers l and m hosting consecutive functions in the chain for which their path passes through the root.

Constraints 27 and 28 enforce that the traffic rate into a switch or a server must not exceed the available processing rate. In the case of inequality 28, the server's processing rate is governed by the VNF assigned to it. Constraints 29 and 30 help define the maximum utilization ρ over network resources: ρ must be at least as great as the utilization at any switch $n \in N \setminus L$, and at least as great as the utilization at any server $l \in L$. Because the processing rate of a server l depends on the VNF v assigned to it, we must have a separate constraint of type 30 for $l \in L$ and $v \in V$. If VNF type v is assigned to server l , then $M(1 - x_{v,l}) = 0$ and 30 requires that $\rho \geq b_l / \gamma_v^{sl}$, where γ_v^{sl} is the processing rate of VNF v if assigned to the server l . If v is not assigned to server l , then the right hand side of 30 is negative, and so imposes no restriction on ρ .

Note that while no constraint forces ρ to equal the maximum utilization over all network nodes, the objective function will drive the value of ρ down to the smallest value satisfying the constraints 30 and 29, thus ensuring that it equals the true maximum utilization over all nodes in the network.

Model Objectives: The objective is to minimize a weighted combination of the number of nodes utilized and the maximum utilization over all nodes in the network.

$$w = (1 - \beta)\rho + \beta \frac{1}{|L|} \sum_{v \in V, l \in L} x_{v,l} \quad (31)$$

When $\beta = 0$, the objective reduces to minimizing the maximum utilization over all nodes in the network. This choice of objective has the effect of distributing the traffic as uniformly as possible in order to reduce the highest utilization over all nodes. If instead $\beta = 1$, the objective becomes minimizing the total number of nodes used to host VNFs. A placement which minimizes the number of VNFs tends to concentrate traffic in part of the network, leaving other network resources unused. Solving the MIP over a range of $\beta \in [0, 1]$ yields a set of solutions that represent different tradeoffs between performance and server usage.

Extensions There are several possible extensions to the MIP model. One such extension is handling the case that only a subset of the service chains can be deployed. It may happen that not all service chains can be accommodated by the network. In that case, we still want to produce a solution that deploys a subset of service chains. We assume that there is a priority order among service chains. Let π_c denote the priority weight of service chain c , where higher priority weight corresponds to higher priority. We introduce a new binary decision variable $d_c \in \{0, 1\}$ for each service chain $c \in C$ indicating c is deployed in the physical network. For constraints 20 and 23, we change the right hand side to d_c . We also introduce a new constraint that ensures service chains are deployed according to the given priority:

$$d_c \geq d_{c'} \quad (32)$$

for all service chains $c, c' \in C$ for which $\pi_c > \pi_{c'}$.

In this extension, the primary objective is to deploy all service chains if possible, and if not, to deploy as many service chains as possible according to the given priority. To that end, our solution procedure would change slightly. We would first solve the MIP (including the new decision variables and constraints) with the objective of maximizing $w' = \sum_c d_c$, the number of service chains deployed. We then refine the solution by re-solving with the objective in Equation 31 while fixing the d_c variables to the values obtained in the first solution.

Other extensions that can be easily accommodated include:

- Imposing constraints on link bandwidth.
- Limiting the length of the path(s) travelled by a service chain.
- Including edge utilization when computing maximum utilization ρ .
- Allowing multiple VNFs to be hosted on each server.
- Requiring that service chain traffic flows are not split across multiple paths. (This extension requires continuous variables to become binary.)
- Enforcing redundancy by prohibiting select pairs of service chains from sharing subnetworks. For example, we could restrict a pair of service chains from using servers under a common TOR switch or aggregation switch.
- Deploying additional VNFs and service chains while keeping existing deployments fixed.