



HAL
open science

Optimized Sampling Strategies to Model the Performance of Virtualized Network Functions

Steven van Rossem, Wouter Tavernier, Didier Colle, Mario Pickavet, Piet Demeester

► **To cite this version:**

Steven van Rossem, Wouter Tavernier, Didier Colle, Mario Pickavet, Piet Demeester. Optimized Sampling Strategies to Model the Performance of Virtualized Network Functions. 2020. hal-02354401v2

HAL Id: hal-02354401

<https://hal.science/hal-02354401v2>

Preprint submitted on 19 Jan 2020 (v2), last revised 21 Jun 2020 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Optimized Sampling Strategies to Model the Performance of Virtualized Network Functions

Steven Van Rossem · Wouter Tavernier · Didier Colle · Mario Pickavet · Piet Demeester

Received: date / Accepted: date

Abstract Modern network services make increasing use of virtualized compute and network resources. This is enabled by the growing availability of softwarized network functions, which take on major roles in the total traffic flow (such as caching, routing or as firewall). To ensure reliable operation of its services, the service provider needs a good understanding of the performance of the deployed softwarized network functions. Ideally, the service performance should be predictable, given a certain input workload and a set of allocated (virtualized) resources (such as vCPUs and bandwidth). This helps to estimate more accurately how much resources are needed to operate the service within its performance specifications. To predict its performance, the function should be profiled in the whole range of possible input workloads and resource configurations. However, this input can span a large space of multiple parameters and many combinations to test, resulting in an expensive and overextended measurement period. To mitigate this, we present a profiling framework and a sampling heuristic to help select both workload and resource configurations to test. Additionally, we compare several machine-learning based methods for the best prediction accuracy, in combination with the sam-

pling heuristic. As a result, we obtain a reduced dataset which can still model the performance of the network functions with adequate accuracy, while requiring less profiling time. Compared to uniform sampling, our tests show that the heuristic achieves the same modeling accuracy with up to five times less samples.

Keywords Sampling Heuristic · Network Function Virtualization · Performance Profiling · Machine Learning · Regression

1 Introduction

In the telecom industry, there is an increasing adoption of cloud-native services and network functions based on Software Defined Networking (SDN) and Network Function Virtualization (NFV) techniques. By virtualizing compute and network resources, a very flexible environment can be created to deploy Virtual Network Functions (VNFs) with an optimal amount of allocated resources, adapted to the realtime incoming workload. The recent rise of 5G enabled services also advocates the use of cloud-native functions, which are deployed

This work has been performed in the framework of the NG-PaaS and 5GTANGO project, funded by the European Commission under the Horizon 2020 and 5G-PPP Phase2 programmes, resp. under Grant Agreement No. 761 557 and 761 493 (<http://ngpaas.eu>) (<https://www.5gtango.eu>). This work is partly funded by UGent BOF/GOA project 'Autonomic Networked Multimedia Systems'.

All authors are at
Ghent University - imec, IDLab
E-mail: {firstname.surname}@ugent.be

The corresponding author is:
Steven Van Rossem E-mail: stevanvanrossem@gmail.com

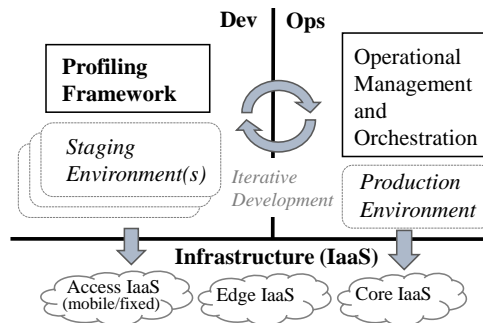


Fig. 1 The profiling framework can use a similar infrastructure compared to operations, as part of a DevOps workflow.

over a virtualized infrastructure [1] [10]. This illustrates the growing need to map the amount of allocated resources and incoming workload to the Key Performance Indicators (KPIs) of the deployed network service, specified in the Service Level Agreement (SLA). To characterize this relation as good as possible, we propose that the virtualized function is profiled on the targeted infrastructure, before being deployed in production. This is illustrated in Fig. 1 as a DevOps inspired workflow. Cloud-native techniques enable a very flexible use of the Infrastructure as a Service (IaaS) for a wide variety of use cases, as envisioned in [22]. This enables the profiling framework (Dev) to do its testing on a representative (but isolated) part of the IaaS, compared to the operational environment (Ops).

In this paper, we propose an optimized sampling procedure to shorten the profiling time as much as possible, without losing sampling accuracy. This optimized profiling procedure greatly benefits a DevOps-inspired deployment cycle: frequent VNF updates can be quickly and representatively validated, before the updated VNF is handed over to the service provider for deployment in production. We hereby consider the VNF as a black-box, with no formal way to deterministically calculate the performance metrics. VNF Profiling is then basically a form of load testing, where representative workloads are emulated and performance is recorded. The data analysis investigated in [23], shows that the output of the profiling tests can be further used to map resource allocations and expected workload to performance levels for the profiled VNFs. The validation and modelling of this relation in the profiling environment, brings additional trust when deploying the profiled VNF in the production environment. The research goal of this article is to optimize the data gathering part of the profiling procedure as much as possible. This is done by investigating whether the total test space of possible workloads and resource allocations can be limited, without losing too much prediction accuracy. We compare several machine-learning based methods for the best performance modeling accuracy, in combination with several sampling heuristics. But first, we present the architecture and implementation of our profiling tool, which we used to gather the profiling data. The profiling tool also contributes to shortening the total profiling time, by parallelizing multiple measurement campaigns over equal hardware nodes.

2 VNF Profiling Framework

The situation of the Profiling Framework in the overall service creation process is given in Fig. 1. Using cloud-native techniques, a light weight profiling platform can

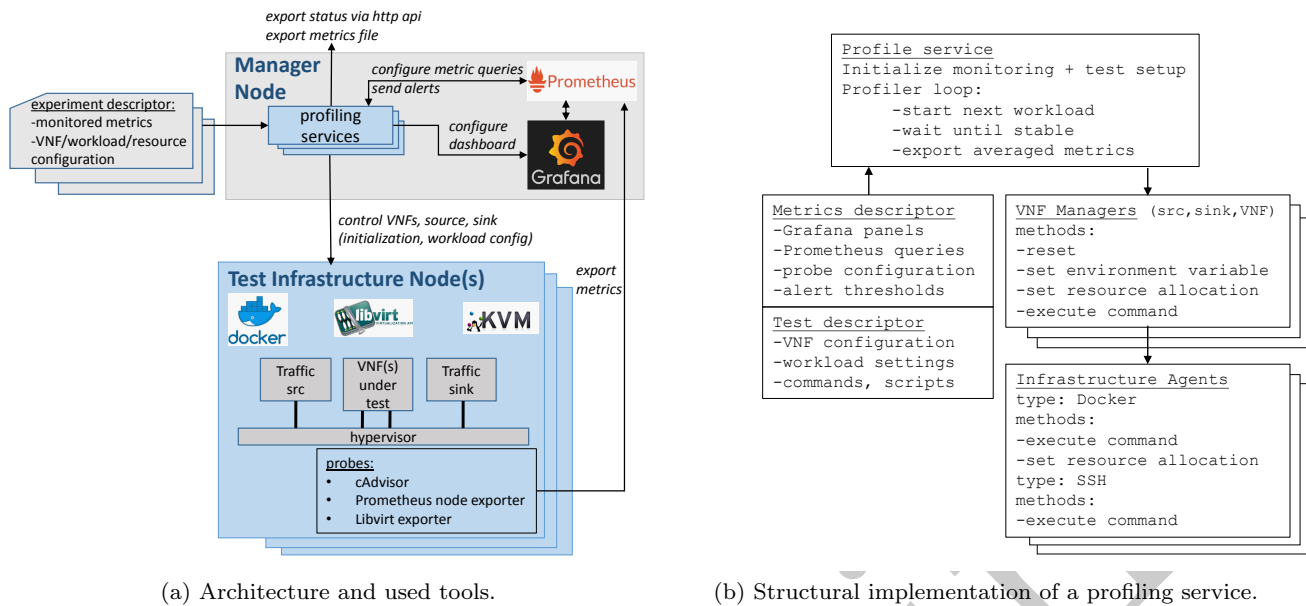
be quickly set up, to test and measure VNFs on available IaaS nodes as a staging environment. To generate the profiled dataset, the framework automates a set of measurements where varying resources are allocated, workload is generated and KPIs are measured. We hereby give a concise overview of this framework. The architecture and implementation builds further upon a previous tool developed in [21]. We re-factored the tool following a Service Oriented Architecture (SOA) approach, where multiple profiling tests can run independently and in parallel, thereby controlling multiple infrastructure nodes. We also adopted several existing functional modules, as seen in the architecture of Fig. 2a. Each test setup consists of a VNF under test, where the test traffic is routed from a traffic source to a traffic sink. The traffic source and sink can be considered as VNFs also, custom built and deployed for our test purpose. The Infrastructure Node should be pre-provisioned. This means that before the profiling tool can operate, the VNFs under test should be pre-deployed on one or more infrastructure nodes. This can be done by a common orchestration framework (e.g. OpenStack, Kubernetes). Also the used probes must be running on each Infrastructure Node: cAdvisor, Prometheus Node Exporter and a probe to export Virtual Machine metrics gathered by KVM and libvirt. Care should be taken that resources (e.g. assigned vCPUs) are well isolated between VNFs under test and other components.

Prometheus is used as monitoring framework and metrics database. For every started profiling test, a *Grafana* dashboard is generated, to visually check the status of the defined metrics being monitored.

The main 'ancillary' services are deployed as Docker containers (Prometheus, Grafana, traffic source/sink, probes). The actual VNF under test can be deployed as Docker container or as VM under KVM.

The Manager Node is where each profiling test is running as a separate instance. Each running profiling test can be considered as a micro-service, from which the test status can be queried through an HTTP API. At the end of the profiling loop, the measurements are stored in .csv format for further analysis (See Section 4). A more detailed implementation diagram is depicted in Fig. 2b. We can distinguish following important functional blocks in the profiling service:

- Every profiling test is defined by a *Test descriptor*. Here we describe which scripts or commands to execute in each VNF, in order to generate the requested workload. See the example Listings 1 and 2 in next subsection.
- A *Metrics descriptor* defines for each VNF in the Test descriptor which metrics should be monitored



(a) Architecture and used tools.

(b) Structural implementation of a profiling service.

Fig. 2 Architecture and implementation of the profiling framework.

and recorded. For each requested metric, the correct Prometheus Query (PromQL) should be given, this is a Prometheus specific syntax to retrieve the metric from the database. Care should be taken that all required metrics are exported by the pre-deployed probes. Prometheus is also configured to send alerts back to profiling service when measurement stability or overload of the traffic src/sink is detected. See the example Listing 3 in next subsection.

- Every VNF in the test setup (source, sink and VNF under test) has its own *VNF Manager* instance, which has methods available to control the VNF state. The most important function is to execute commands which start/stop traffic workloads. The actual commands or scripts are specified in the descriptor files, the VNF Manager only executes the defined command using the correct Infrastructure Agent.
- Every VNF Manager, can have attached *Infrastructure Agents* to address a specific API in the remote Test Node. This agent is used to execute commands inside a container or set container resources via the Docker API of the remote node or execute commands through SSH.

Our Test Infrastructure Nodes are equal compute nodes with 2x 8core Intel E5-2650v2 (2.6GHz) CPU with Ubuntu 18.04. Linux Bridge is used as the hypervisor switch. We do not change default OS options (e.g. we leave hyperthreading enabled). Depending on the virtualization of the VNF (container or Virtual Machine (VM)) we use the configuration options of Docker resp. KVM to isolate the CPU cores between the DUT and the traffic sink/source. The Manager Node is a lighter

machine: 4 core Intel E3-1220 CPU with Ubuntu 18.04. The main bottleneck resource of the Manager Node is the disk space used by Prometheus, to store all the metrics gathered from various running profiling tests.

2.1 Profiling Descriptors

To have a more practical idea of the profiling execution, we give a succinct overview of the format used to describe the profiling tests. The syntax of the test definition is loosely based on the implementation in [15] and includes some adaptations to map it better to our envisioned sampling heuristic.

```

1 agents:
2   docker1:
3     class: DockerClient
4     url: 'tcp://docker.api.url:port'
5     # ...
6   ssh1:
7     class: SshClient
8     host: 'infrastructure.node.url'
9
10 managers:
11   src:
12     class: DockerVnf
13     agent: docker1
14     cpuset_cpus: '8-15'
15   # ...
16   pfsense:
17     class: Vm
18     agent: ssh1
    
```

Listing 1 YAML based test descriptor - interface configuration

A first part of the *Test descriptor* is given above in Listing 1. This part defines a configuration agent for each VNF in the profiling setup. As previously explained, a class instance is made for each specified VNF manager and each manager uses an agent as interface to the underlying VNF. The names of the VNF managers defined here, are referred to in the remainder of the descriptor. The declarations in the descriptor hold test-specific settings. For the agents this includes: api endpoints, credentials, authentication methods, ... For the VNF managers we need specific settings such as: the agent to communicate with the VNF, container or VM uid, resource or operational initialization to configure via the agent, ...

A second part of the *Test descriptor* is then given in Listing 2. This is the most important part from the profiling perspective, as it defines the actual values for each relevant parameter in the test. For each parameter we define:

- A list or range of values to test (Line 4,10,14,25).
- The method of the manager instance to call, in order to practically configure this value into the VNF test (Line 5,11,15,26). In our tests, we use a common technique based on environment variables. The workload settings are stored in environment variables in the VNFs, later when the workload script is called, these variables are read and the configured workload is started.
- A list of fixed initialization commands, which are executed every time a new setting is configured (e.g. to stop/start a workload generating script in the traffic generator) (Line 18,29).

The sampling heuristic takes the defined parameter ranges into account and will iterate through all combinations in an optimized order. For this reason, the total configuration space is categorized into three sections:

- *resource_parameters*
- *workload_parameters*
- *primary_workload_parameter*

The main reason for this categorization is to guide the sampling heuristic through which metric value to sample next, as will be explained in the coming sections.

In Listing 3, we illustrate the structure of the *Metrics descriptor*. This file is translated to the needed configuration directives for the Prometheus monitoring framework to gather the required metrics. The list of all the required metrics is given (Line 1). For each given metric, a template should be defined, which maps to the correct Prometheus Query (Line 9). Deployment or test specific parameters such as id's should be dynamically filled in the template. The metrics descriptor also

defines the probes where Prometheus can get the metric values from (Line 21). The Profile service will query all defined metrics from the Prometheus database, once a configured workload shows stabilized measurements. (Measurement stability is assessed by the implementation described in [23].) The queried metric values are exported to a file and kept for online analysis by our sampling heuristic.

```

1 resource_parameters:
2   - name: pfsense_cpu_limit
3     # allocated cpu in %
4     values: [25,50,75,100,200,300,400,500]
5     function: set_cpu
6     manager: pfsense
7
8 workload_parameters:
9   - name: packetsize
10    values: [64,128,256,512,1024,1500] #
11    Bytes
12    function: set_environment_var
13    manager: src
14   - name: flows
15    values: [1, 2, 10, 100, 1000, 10000]
16    function: set_environment_var
17    manager: src
18  initialization:
19    - manager: src
20      cmd: 'pkill -9 -f start_src.sh'
21    # ...
22
23 primary_workload_parameter:
24   name: packetrate
25   # Values will be chosen in the defined
26   # interval for this parameter
27   range: [0.1,500] #kpps
28   function: set_environment_var
29   manager: src
30
31  initialization:
32    - manager: src
33      cmd: 'bash start_src.sh'
34    # ...

```

Listing 2 YAML based test descriptor - test configuration space

It is beneficial if each measured and exported metric is explicitly mentioned in the *Metrics descriptor* and *Test descriptor*. When exporting the test results after the profiling, there should be a clear link between the metric name in the exported results and where/how this metric is exactly measured. Our approach is that the exact name of the measured metrics can be found back in the two descriptor files. The descriptor files become the reference for what each metric stands for and where/how it is exactly gathered: (i) The test descriptor explains the metrics which define active settings such as the workload generation and resource allocation. (ii) The metrics descriptor explains the metrics gathered passively by all probes.

```

1 metrics:
2   - sink:cpu
3   - src:cpu
4   - pfsense:cpu
5   - sink:packetrate_receive:eth1
6   - pfsense_packetrate_loss
7   # ...
8
9 definitions:
10  docker:
11    cpu:
12      template: 'sum(rate(
13        container_cpu_usage_seconds_total{id="/
14        docker/{{ docker_id }}"})[10s])*100'
15      unit: '%'
16
17    packetrate_receive:
18      template: 'sum(rate(
19        container_network_receive_packets_total{id
20        ="/docker/{{ docker_id }}",interface="{{
21        interface_id }}" }[10s]))'
22      unit: 'pps'
23
24    # ...
25
26 probes:
27  node_exporter:
28    job_name: node_exp_pfsense1
29    scrape_interval: 1s
30    static_configs:
31      - targets:
32        - 'infrastructure.node.url:9100'
33
34  cadvisor:
35    job_name: cAdvisor_pfsense1
36    scrape_interval: 1s
37    static_configs:
38      - targets:
39        - 'infrastructure.node.url:8080'
40
41    # ...

```

Listing 3 YAML based metrics descriptor

The use of the above explained descriptor files, makes it easy to modify and repeat profiling tests. The configuration of monitored metrics, workload and resource parameters is kept very generic to allow a wide applicability in VNF testing. In the remainder of this article we will present measurement results gathered by using the above explained framework and descriptor formats.

3 VNFs Under Test

To choose exemplary VNFs for our profiling tests, we looked at some typical use cases defined in [1]. The adoption of 5G technologies enables new possibilities for the telco industry to diversify their network services to new markets. To enhance the security of these services we look at the deployment of a virtual firewall (pfSense). As a large portion of the traffic over 5G will

be media based, we also look at the deployment of a virtual streaming server (Nginx).

For all the tested VNFs we consider CPU and bandwidth as the most important resource metrics, as we assume these are more likely to become a bottleneck resource than memory. This is also confirmed in [9]. Our measurements also show little to none variation in the memory usage of the VNFs while they are under test. In the next subsection we will discuss each VNF more in detail.

3.1 Firewall - pfSense

We use pfSense¹ as a free and open source firewall solution example, deployed as a VM. We stress the firewall by generating multiple unique parallel flows. Also the packet size is varied. Using the tool *Scapy* we assemble a .pcap file with a stream of packets of varying mac addresses and unique destination IP/port in the packet header. *TcpReplay* is then used to stream the .pcap file at a given packetrate from the traffic source. There is also an iperf stream running, with an iperf server in the traffic sink. This is used to monitor packet loss. For the firewall to function properly, we need to make sure the ARP table of the VNF contains the mac addresses of the generated packets, so the firewall forwards the packets properly to the traffic sink. This is done by arp spoofing the firewall from the traffic sink. To have an idea of the baseline performance of the firewall, we install no specific firewall rules and let the traffic pass.

Generated workload metrics:

- **packetrate:** [0.1-500]kpps. 50 different packetrate values are selectively chosen, spaced evenly along the log scale.
- **packetsize:** [64,128,256,512,1024,1500] bytes
- **flows:** [1,2,10,100,1000,10000] unique parallel flows (with unique IP/port combination in the header).

Resource metrics:

- **CPU allocation:** [0.25, 0.5, 0.75, 1, 2, 3, 4, 5] vCPUs

The bandwidth allocation is not a dedicated setting in this test. It is determined by the workload, since we specify the generated packetrate and packet size up front.

Performance metric:

We choose packet loss (%) as the main KPI to reflect the performance of the firewall.

All combinations of above metrics result in 12000 measurement points. If we need about 30sec per measurement to get a stable reading, the total profiling time reaches up to 100h to measure each combination once.

¹ <https://www.pfsense.org/>

3.2 Streaming Server - Nginx

We set up a live streaming service using Nginx², a well known open source, all-in-one load balancer, web server, content cache and API gateway solution. Nginx is deployed in a Docker container. We configure Nginx to accept incoming live movie streams via the Real-Time Messaging Protocol (RTMP [13]) protocol. The incoming RTMP live stream is then transcoded to a specific video bitrate and resolution (Nginx uses ffmpeg for this purpose). Next, Nginx serves the newly encoded movie chunks live, through the HTTP Live Streaming (HLS [12]) protocol. In our test setup, the traffic source sends 1 - 5 movies in realtime to Nginx over RTMP. On the client side, the traffic sink opens many concurrent sessions to Nginx, to download playing live movies over HLS (We use *Locust.io* to emulate the HLS clients and download the stream requests). This use-case exemplifies the situation where a small number of incoming live movies is temporarily cached in an edge server and then streamed with a certain quality to a large number of clients.

Generated workload metrics:

- **streams:** [10-5000] parallel client HLS streams. 80 different stream values are selectively chosen, spaced evenly along the log scale.
- **movies:** [1,2,3,4,5] number of different source movie streams, input via RTMP.
- **quality:** [1,2,3,4,5] indicator for the quality of the streams (resolution and video bitrate ranging from 1280x720/2500kbps to 426x240/200kbps).

Resource metrics:

The streaming performance is determined by both the available bandwidth and vCPU. It is unpredictable how the balance between cpu time for ffmpeg transcoding and cpu time for serving the movie chunks will be scheduled (as we consider this a black-box VNF). Therefore we have no way to deterministically predict the influence of both the allocated bandwidth and cpu on the KPI. We need to profile the performance with several combinations of allocated vCPU and bandwidth. This also reflects the availability of different flavours to deploy the VNF.

- **flavours:** [(0.5, 0.5), (0.5, 1), (1, 1), (1, 2), (2, 1), (2, 2), (3, 2), (3, 3), (3, 4), (4, 5), (6, 5)] (vCPUs, Gbps). Eleven different flavours to deploy the VNF, defined by their given vCPU and bandwidth allocation and encoded from [0-10].

² <https://www.nginx.com/products/nginx/modules/rtmp-media-streaming/>

Performance metric:

We choose **lag ratio** (%) as the main KPI to reflect the performance of the streaming server. This indicator is a measure for the risk of 'hiccups' or lagging during video playback. It is the ratio of downloaded video playback time over the last period. If the video time is less than the waiting time, the playback buffer will empty and the risk of lagging will increase:

$$\text{lag ratio (\%)} = \max\left(1 - \frac{T_{\text{video}}}{T_{\text{wait}}}, 0\right)$$

We measure the lag ratio in a moving average over 20s (we assume 20s buffer time). If the KPI gets above zero, it means that during the last 20s, the playback buffer was addressed because less than 20s of video stream was downloaded. Increasing KPI values mean more buffer time is continuously needed, resulting in video rebuffering and thus 'lagging'. The HLS protocol will try to keep the lag ratio at zero by varying the size of the served movie chunks and maximizing the bandwidth over all clients.

In order to get a stable measurement, a certain ramp-up time is needed to generate the required number of clients and to let the HLS based streaming stabilize. In our setup this takes up to 100sec per measurement point. To test all above combinations once, takes then over 500h to complete.

We only evaluate the KPIs below 30% packet loss or lag ratio, as we assume that above this threshold the VNF is practically unusable. Therefore there is no need to accurately model the KPI above 30%.

The long profiling times of the above introduced VNFs show the need to optimize both: (i) the parallel execution of measurement runs by the profiling framework (as explained in Section 2) and (ii) the sampling strategy to limit the number of needed sampling points. The latter will be explained next.

4 VNF Data Analysis

As proposed in [23], we have classified the tested VNF metrics under three groups in the previous section:

- **Workload metrics** reflect the configuration of the incoming traffic to be processed by the VNF.
- **Resource metrics** quantify the allocated resources which determine the cost and processing capabilities of the VNF. For our analysis we express this as *resource usage*, which is the averaged used portion (%) of allocated vCPU and bandwidth.
- **Performance metrics** monitor the Key Performance Indicators (KPIs), to assure that the performance of the VNF remains within the SLA.

From the obtained VNF measurements, we want to derive a model which predicts the performance KPI in function of the given workload and resource allocation. From an abstract and generalized viewpoint, the VNF performance model f can be described as:

$$f(wl, res) = perf \quad (1)$$

where:

wl = input workload (e.g. packetrate, filesize, streams)

res = resource allocation (e.g. number of allocated vCPUs, bandwidth or flavour)

$perf$ = KPI metrics (e.g. packet loss, lag ratio)

Figure 3 shows a subset of our measurements: for each VNF a certain workload configuration is executed on varying resource allocations. The measurements in Fig. 3 confirms these trends (which were also described in [23], but on other VNF examples):

- The resource usage is correlated with the rising workload (on the x-axis) until saturation (Fig. 3a and 3b). Either CPU or bandwidth gets saturated first, which explains why the averaged resource usage can saturate below 100%.
- Before resource saturation, the KPI levels remain stable and flat. When resource contention starts, the KPI levels start to vary more rapidly (Fig. 3c and 3d).

We can also distinguish two other interesting facts from the plots:

- For pfSense (Fig. 3c) we can see that the performance does not increase with more than 3 allocated vCPUs. This points to a deployment limitation where it makes no sense to allocate more vCPUs, because it is not exploited by the VNF implementation.
- When it comes to Nginx (Fig. 3d), we see that some resource flavors have overlapping performance curves. This indicates that the workload is bounded by a common resource limit of those flavors, namely bandwidth in this case.

It is challenging to discover the above mentioned phenomena automatically, without visual inspection of the data plots. A KPI prediction can be made by training the model with the obtained profiled datasets. But a common adagio from the Machine Learning domain is that the model will only be as smart as its training data, meaning that we must provide representative training data in all foreseeable situations. This implies that we must also profile in the regions where resource saturation occurs or where resource flavors overlap , or

where the performance is limited by the internal VNF implementation.

In the next subsection we will outline methods to model the performance of the trends shown above. Regarding the accuracy, it is important to note that the breakpoint of the KPI curve is the area of most interest. This is the maximum workload possible by the VNF, just before the performance declines more severely.

4.1 Modeling Methods

We look for an appropriate modeling method to predict the KPI values from a given workload and resource configuration, as explained earlier by Eq. 1. A first idea of the trends to be modelled can be seen in Fig. 3c and 3d. From a pure mathematical perspective, we can consider the KPI values to be a response surface, defined by a multi-variate function where the workload and resource allocation metrics are the input parameters. The total input space is multivariate, since all workload and resource metrics can influence the resulting KPI value. We compare several generic methods from the machine learning domain, which are capable of modeling generic, non-linear and multi-variate response surfaces. The used methods have also shown promising applications in regression modeling, where the amount of training samples is limited. We have used the implementations available in the library Scikit-learn [14]. We also include the Interpolation and Curve Fit method, which have shown promising results in [23]. The investigated modeling methods are:

- **Support Vector Regression (SVR)**: This method has shown promising results in estimating non-linear relationships using limited, sparse datasets. SVR selects samples to form a flexible tube of minimal radius, symmetrically around the estimated function, such that the absolute values of errors less than a certain threshold (ϵ) are ignored both above and below the estimate. Points outside the tube are penalized and not taken into account for the regression. The hyperparameters of this method are C , the penalty parameter, ϵ and the standard RBF kernel. More details can be found in [19]. The use of SVR for modeling VNF performance has also been applied in [16] with limited success. We also include it here for verification on our data sets.
- **Random Forest (RF)**: The basic idea behind this method is to combine multiple decision trees in determining the final output rather than relying on an individually built decision tree. Maximum tree depth is set to 10, and the number trees in the forest is 100.

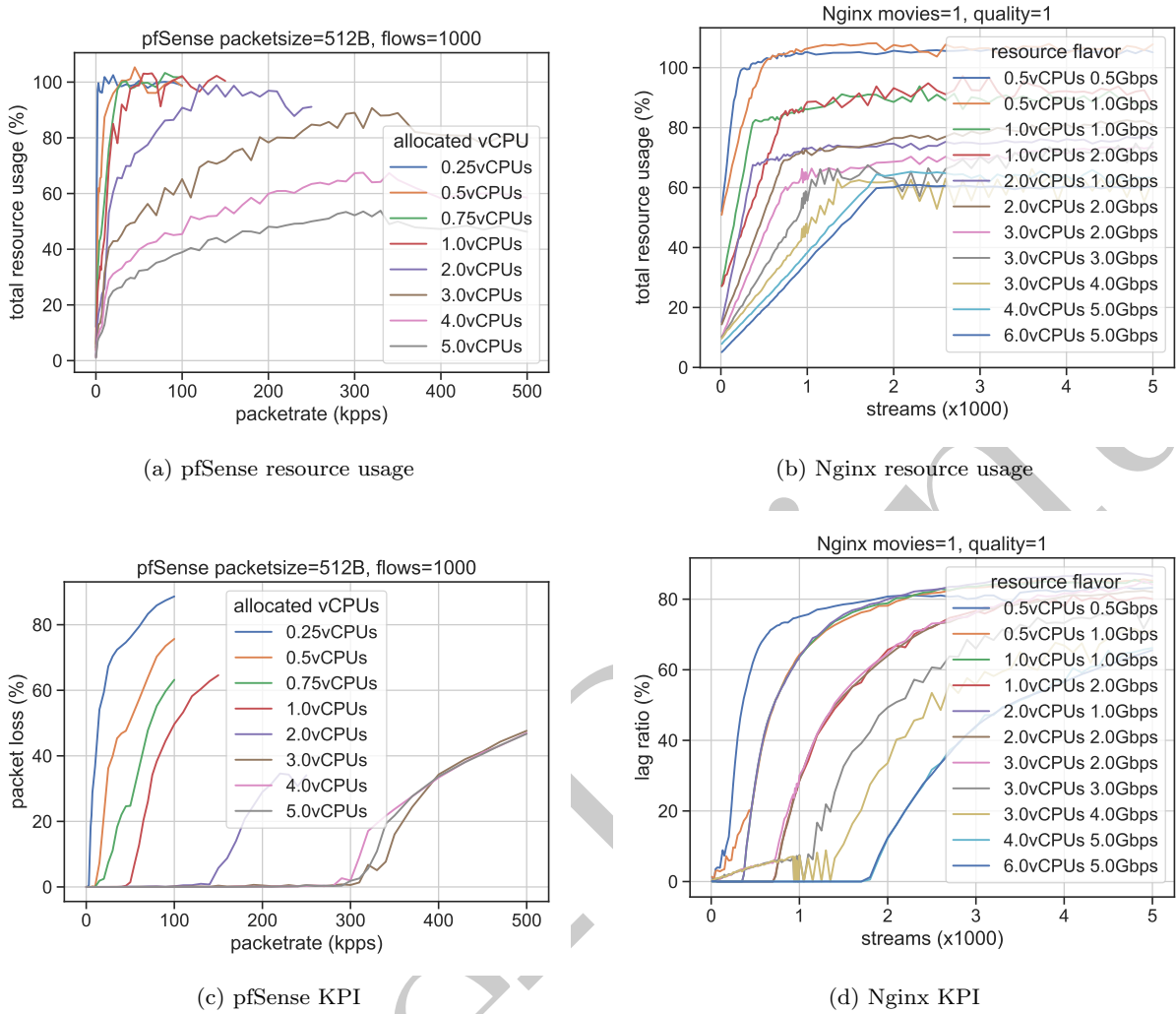


Fig. 3 Subset of measured VNF metrics under different resource allocations.

The use of decision trees for modeling VNF performance has been investigated in [16], [7] and [6] with promising results. We include the RF method here for verification on our data sets.

- **Gaussian Process (GP)**: This method implements a Bayesian approach to (non-)linear regression. A GP defines a prior over functions, which can be converted into a posterior over functions once it has seen some data. The covariance between training samples is a given kernel function. The kernel function we use: $Constant * RBF + WhiteNoise$. This is a generic kernel function used in many GP examples. One main advantage of using GPs, is that the kernel hyperparameters (RBF lengthscale, noise level, constant) can be learnt automatically via evidence maximisation from the training points themselves, no exhaustive search is needed as with other methods. The key idea is that if the training samples are deemed by

the kernel to be similar, then we expect the output of the function around those points to be similar, too. More information on this method is available in [24]. The use of GP for modeling software performance has been proposed in [5]. We also use GP here for verification on our datasets.

- **k-Nearest Neighbors (kNN)**: is a commonly used technique due to its simplicity and often accurate results. In kNN regression, the output value is the average of the values of k nearest neighbors in the input space (weighted by distance). We use the Euclidean distance metric and standardize the input values. For our tests we use $k = 2$.
- **Interpolation method**: Regression is done by interpolating linearly between surrounding samples. The interpolant is constructed by triangulating the input data using Delaunay triangulation, and on each triangle performing linear barycentric interpolation.

This method also works in multiple dimensions, so the total workload and resource input space is taken into account to interpolate an intermediate KPI value. This method has been used for VNF modeling in [23] and we also include the method here for verification on our new datasets.

- **Curve Fit method:** This method has been successfully used for VNF modeling in [23] and we reuse the method here for verification on our new datasets. Based on the analysis in [23], we fit the KPI trends to the analytical functions of this form:

$$\begin{cases} f(x)_{non-saturated} = a + exp[b(x - c)] \\ f(x)_{saturated} = 100 \left(1 - \frac{d}{x-c}\right) \end{cases}, \text{ where } x > d + e \quad (2)$$

where the output is clipped in the range $[0, 100]$:

$$\begin{cases} 0 & , f(x) < 0 \\ 100 & , f(x) > 100 \end{cases}$$

The full justification for the functions in Curve Fit method is explained in [23], we summarize here shortly the fitting procedure. Figure 4 shows how the trends in Eq. 2 are fitted to data subsets of both investigated VNFs. The saturated region starts when the workload shows a higher covariation with the KPI than with the resource usage metric. The parameters a, b, c, d, e in Eq. 2 are fitted to the profiled data points in the respective (non) saturated regions. Note that two possible intersections can exist between $f_{non-saturated}$ and $f_{saturated}$. When the fitted curve is known, all possible intersections are estimated. For each intersection point, the accuracy (RMSE) of the resulting piecewise model is calculated. Finally, the intersection point which yields the best accuracy is chosen and stored into the Curve Fit model. The clipped output represents that the KPI values (representing packet loss and lag ratio) are limited between $[0, 100]$ %. We further use a weighted curve fitting for the saturated part, to prioritize more accuracy at lower KPI values, just after resource saturation. Since this is the region where most accuracy is wanted.

5 VNF Sampling Strategies

In this section, we present our results after investigating different sampling methods. As illustration, the maximum available samples for one example workload are shown in Fig. 4. Each of the tested sampling strategies would yield a different set of sampling points, following the same trend curves. As baseline reference we pick the samples uniformly, then our heuristic tries to focus

on more interesting parts of the trend curves, like the breakpoint between the (non)saturated regions.

We also check the effect of the sampling heuristic on the accuracy of the modeling methods proposed in previous section. We have executed twice all workflows defined in previous Section 3. We thus obtain two datasets per VNF, of which one is used as training set and the other as test set. From the training set, we select data points using our sampling heuristics and use the selected samples to train a prediction model for the KPI values. We then check the Root Mean Squared Error (RMSE) of the trained model on the test dataset. The reported RMSE in the next sections is then an indication for the accuracy of the model, trained from a limited sample set. We only evaluate the RMSE below 30% packet loss or lag ratio, as we assume that above this threshold the VNF is practically unusable. Therefore there is no need to assess the model accuracy above a KPI threshold of 30%.

5.1 Uniform Sampling

As baseline measurement and benchmark for our sampling heuristic, we first test the accuracy of a generic uniform sampling strategy. For each VNF, we pick resource and workload metric values uniformly in the range described above in Section 3. The result is shown in Fig. 5. On the x-axis, the number of measurements indicates an increasing amount of uniformly chosen values per metric. First we measure one value, then two values for each metric and so forth...

As can be seen in Fig. 5, different modeling methods yield varying accuracy. Support Vector Regression (SVR) proves the least accurate. The tuning of hyperparameters (C, ϵ) in the SVR model is a tedious task which we do using an exhaustive search and also seems very sensitive to the number of training samples. We experience long training times (tens of minutes with more than 1000 samples). We do not succeed to reach the same accuracy as the other methods. The best performing methods, which we select for further use, are interpolation and curve fit. This is also in line with previous research done in [23].

Uniform sampling is however not an online sampling method, as the number of samples must be chosen in advance. In the next sections, we will try to find heuristics to select samples in an online way, and try to reach the same accuracy with less training samples.

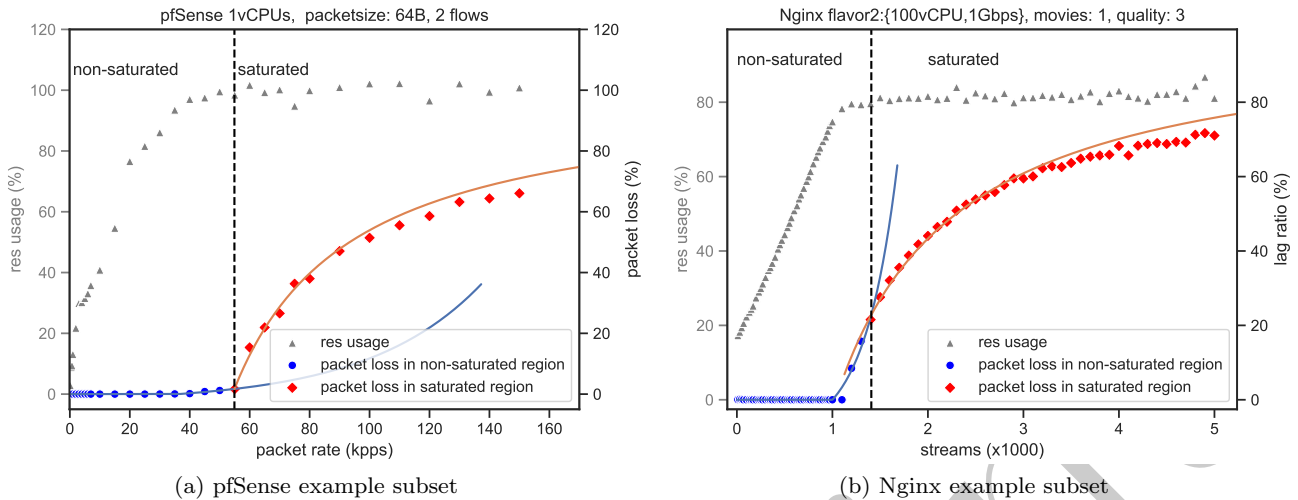


Fig. 4 The Curve fit model exemplified in two example VNF subsets. For each plot, the y-axis on the left is for the saturating resource usage, the one on the right depicts the increasing KPI value. The boundary between (non) saturated regions is where the workload covariance with the KPI becomes larger than with the resource usage.

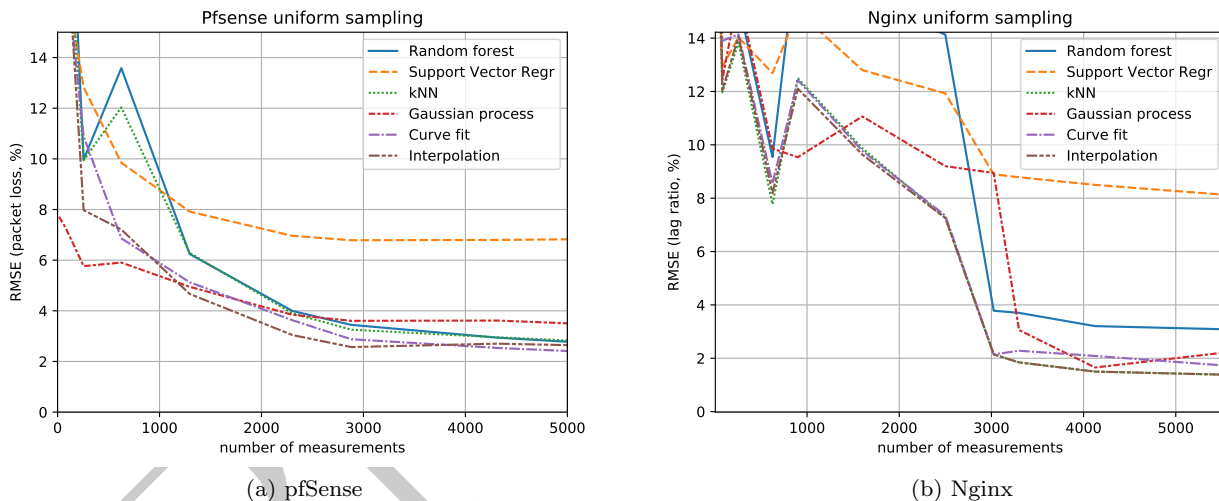


Fig. 5 The accuracy of different modeling methods applied to a uniformly sampled dataset.

5.2 Unsuccessful Strategies

It is worth noting that uniform sampling, exemplified above, is not the least performing method we encountered during our tests. When following the approach outlined in [5], a Gaussian Process (GP) is used to select online which samples to measure next. In fact, GP is a modeling method, but it also produces an estimate of uncertainty in the prediction, a confidence interval for the predicted values as indicated in Fig. 6a. The points where the predicted values have the largest confidence interval, is considered as the most interesting region to sample next. The confidence interval calculated by the

GP is an indication in which region the uncertainty of the predicted value is larger. For our datasets, this selection method is not optimal, since it favors regions where a flat response is measured, and not the transition zone to the steeper part of the trend.

Fig. 6b and 6c show the effect of a GP chosen sample set on different modeling methods, compared to the uniformly sampled benchmark. We see indeed that GP does not optimally select samples in our datasets, as the accuracy is worse compared to the benchmark.

Another commonly used adaptive sampling method is based on surrogate modeling using a gradient based approach [3]. The drawback of this methodology is that

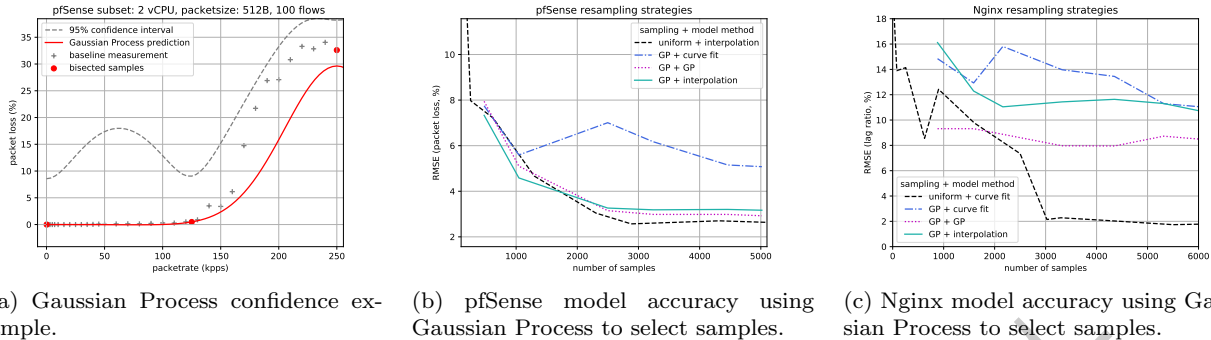


Fig. 6 The Gaussian Process sampling method does not improve the accuracy.

sampling efforts risk to focus too hard on local extrema, and therefore take an unbalanced number of samples in wrong regions. We also can also visualize why some strategies are under-performing. The KPI response surface in a subset of the samples is illustrated in Fig. 7.

- Sample selection using GP will likely select points in the flat (blue) area of Fig. 7, as explained in previous paragraph.
- On the other hand, gradient-based sampling strategies will likely choose sample points in the steep (green/red) regions of Fig. 7.

For our datasets, the most interesting region to focus on, is where the trend break happens, i.e. where the flat region transforms into the steep region. In the next sections we try to find heuristics which focus on this area.

pfSense response surface, subset: 512B packetsize, 100 flows

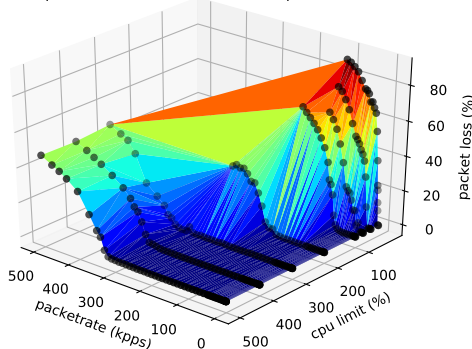


Fig. 7 A subset of the pfSense dataset to illustrate the response surface of the KPI (packet loss), under different vCPU allocations.

5.3 Feature Selection

In our black-box approach, we try to cover the operational boundaries at the start of the profiling procedure. We measure each combination of the minimal and maximal values of the VNF workload and resource metrics mentioned in Section 3. This comes down to a full factorial sampling space with two levels per factor. This gives us a minimal sample set where an initial analysis can indicate which parameters are important or not. We then use this information to determine which metric configurations to profile next.

In this first phase of metrics exploration, it is primarily the intention to select the workload and resource metrics (X) which have the most influence on the performance metrics (Y). Several mathematical methods are generally used to capture the relation between two sets of metrics ($Y = A.X + B$):

- **Multiple Linear Regression (MLR)**: achieves maximum correlation between X and Y , using the well known Ordinary Least Squares (OLS) method.
- **Principle Component Regression (PCR)**: captures the maximum variance in only X , using the well known Principle Components Analysis (PCA) method. Then it uses OLS to predict Y from the main components in X , calculated by PCA.
- **Partial Least Squares (PLS)**: tries to do both by maximizing the covariance between X and Y . The power of this method lies in the fact that influential factors in X can be determined using relatively few samples compared to the other methods. PLS remains robust when the influence on multiple response metrics needs to be assessed (there are multiple columns in Y), or when there is multicollinearity between input variables (collinear columns in X) [4].

For our purpose, we need a method which can estimate the relation between X and Y , from a small

initial dataset. Also it is not our main goal in this initial phase to derive a (linear) model. We only need have an idea of the main factors in X that seem to have the most effect on Y . We therefore select PLS (also sometimes called Projection to Latent Structures) [4] as most generic method for this use case. The results are shown in Fig. 8, where the height of the bars indicates how much each workload or resource metric influences the KPI variation. The bar heights can be thought of as coefficients of a regression model. But from this small dataset, the only reasonable conclusion for now is that a higher bar indicates a more influential factor for the KPI.

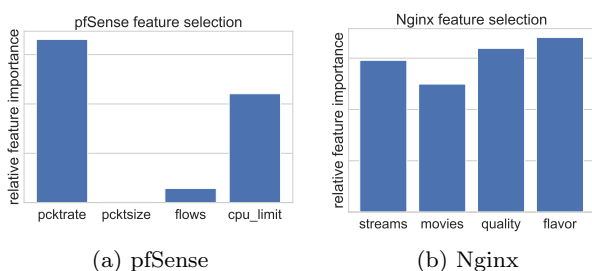


Fig. 8 Feature selection based on the initial sample set of the profiled VNFs.

- For pfSense in Fig. 8a, we clearly see that workload metrics *packetsize* and *flows* have little to no influence the KPI (packet loss).
- For Nginx in Fig. 8b, we observe that all parameters have a reasonable contribution to the variation of the KPI (lag ratio). This means that we should not be selective in the input metric space here.

We must note that the analysis in the previous subsection takes only the variation between the edge values into account. The sample size is too limited to assess if there are any local extrema in between the edge values. In the next sections we will further investigate where to pick next profiling points.

5.4 Primary Workload Metric Selection (wl_p)

We further limit the sampling space, by focusing on selected metrics. For each VNF we can prioritize one specific workload metric, which is more likely to vary in real-world traffic. This is the x-axis metric used on the plots in Fig. 3: *packet rate* (pfSense) and *streams* (Nginx). With the VNF deployed in production, we assume that in the most realistic scenarios both resource allocation and workload configuration are relatively stable,

while the respective x-axis metrics in Fig. 3 are most likely to vary. We call these the *primary workload metric* wl_p . We focus our profiling efforts in such a way, that the VNF model can predict more fine grained at which level of wl_p the performance is outside SLA bounds (using Eq. 1). The other workload and resource allocation metrics are sampled more coarse grained. We favor the respective wl_p metrics during our profiling measurements by picking more samples in their specified range, in order to estimate more accurately the performance breakpoint. For pfSense the $wl_p = \text{packetrate}$, for Nginx: $wl_p = \text{streams}$. Intuitively, this corresponds to how we expect a VNF to be generally used: with resource usage mainly driven by increasing packet or request rate.

The coarse grained values, specified for workload and resource allocation metrics apart from wl_p , allow the inclusion of 'expert knowledge' into the profiling procedure. Next to the edge values, a service developer or tester should include values which are considered interesting to monitor or benchmark. The resource allocation metrics are also likely to be limited by availability. The sampling heuristic will then iterate through the defined workload and resource allocation values in an optimized way.

Algorithm 1: Online sampling of workload metric wl_p

Data: wl_p boundaries, ϵ , max
Result: S_{wl_p} = sampled wl_p values

- 1 $S_{wl_p,old} \leftarrow$ the two wl_p boundary values;
- 2 $M_{old} \leftarrow$ train model with $S_{wl_p,old}$;
- 3 **while** $len(S_{wl_p}) < max$ **do**
- 4 $s \leftarrow$ bisect a new wl_p value (Algorithm 2);
- 5 $S_{wl_p,new} = S_{wl_p,old}.append(s)$;
- 6 $RMSE_{old} \leftarrow RMSE(M_{old}, S_{wl_p,new})$;
- 7 $M_{new} \leftarrow$ train model with $S_{wl_p,new}$;
- 8 $RMSE_{new} \leftarrow RMSE(M_{new}, S_{wl_p,new})$;
- 9 $\Delta RMSE = |RMSE_{old} - RMSE_{new}|$;
- 10 **if** $\Delta RMSE < \epsilon$ **then**
- 11 | stop while loop;
- 12 **end**
- 13 $S_{wl_p,old} = S_{wl_p,new}$;
- 14 $M_{old} = M_{new}$;
- 15 **end**
- 16 **return** $S_{wl_p,new}$;

Algorithm 1 describes how we select online which values of wl_p to sample next. Also a stop criterion is included, to assess when to stop sampling:

- Line 6: M_{old} is trained *without* the latest added sample. The $RMSE_{old}$ is calculated on how M_{old} predicts all samples, including the new sample.

- Line 8: M_{new} is trained *with* the latest added sample. The $RMSE_{new}$ is calculated on how M_{new} predicts all samples, including the new sample.
- Line 10: $\Delta RMSE$ is a measure of how much the latest sample influences the model accuracy. If the accuracy delta is small enough (threshold ϵ is used), we assume that the model will not improve further by adding extra samples, and we can stop sampling this particular workload configuration. Throughout our tests we have used $\epsilon = 0.5$.

In Algorithm 2, we bisect the most interesting workload level to sample next. The goal is to select samples to estimate the point where the KPI trend breaks from flat into a steeper curve, as shown in Fig. 4. Hence, we try to find points around the breakpoint, where the KPI starts to rise.

Algorithm 2: Bisect next value of workload metric wl_p

Data: already sampled wl_p values, KPI_{max}
Result: s = newly sampled wl_p value

- 1 $X \leftarrow$ already sampled x values (workload, ascending);
- 2 $Y \leftarrow$ already sampled y values (KPI) ;
- 3 $S \leftarrow$ empty set ;
- 4 **foreach** *sample* in (X, Y) **do**
- 5 $\Delta x = |X[i] - X[i + 1]|$;
- 6 $\Delta y = |Y[i] - Y[i + 1]|$;
- 7 $S.append(i, \Delta x, \Delta y)$
- 8 **end**
- 9 sort S descending by $\Delta x, \Delta y$;
- 10 take i in S with largest $\Delta x, \Delta y$ and $y[i] < KPI_{max}$;
- 11 $x_{new} \leftarrow \frac{X[i] + X[i + 1]}{2}$;
- 12 $y_{new} \leftarrow$ measure KPI in x_{new} ;
- 13 $s \leftarrow (x_{new}, y_{new})$;
- 14 return s ;

- Line 4-8: The already sampled values are grouped per interval. In each interval we calculate $\Delta x, \Delta y$ as a measure how wide and how steep this interval is.
- Line 10: The most interesting interval has both: a large KPI change (Δy) *and* a wide gap between measured workloads (Δx).
- Line 11: The next sample is bisected in the selected interval.

In Fig. 9 we investigate which modeling method works best to determine $\Delta RMSE$ (line 2 and 7 in Algorithm 1). For this test we decrease ϵ from 2 to 0.01. Then for each value of ϵ :

- We filter each total VNF dataset for each unique combination of profiled workload and resource metrics (except wl_p). For each of these sample subsets, we use Algorithm 1 to gather samples for wl_p .

- In each subset, we compare three methods: curve fit, interpolation and GP, to calculate the $\Delta RMSE$ as stop criterion for the sampling. Each of these three methods thus yields a different number of samples. Note that this is only a one-dimensional model, used here only to calculate $\Delta RMSE$ on training subsets (cf. wl_p on the x-axis and KPI on the y-axis as in Fig. 6a), as wl_p varies while other workload metrics and resource allocation remain fixed.
- When the sampling is completed (threshold ϵ is reached), we combine the sampled subsets to train both a curve fit (Fig. 9a and 9c) and an interpolation model (Fig. 9b and 9d) on the total input space and check their accuracy.

The GP method emerges here, rather surprising, as the most accurate method in all four tests. Apparently, the GP needs the least amount of samples to capture accurately the KPI trend. Intuitively, this seems to indicate that curve fit and interpolation methods tend to overfit here and require more samples until $\Delta RMSE$ stabilizes. The GP method seems to generalize better in this local one-dimensional model and converges faster. Note that we do not use GP to select the samples here (this would be a bad practice as explained in Section 5.2). Instead we use our own bisection method (Algorithm 2) to select samples and only use GP to calculate the $\Delta RMSE$. Another remark, as shown in earlier measurements, is that the GP method also works less well when the input space is increased with more workload and resource metrics. In this multi-dimensional space, the accuracy of the GP method decreases (we will detail this later in Section 5.6). In the next section we will apply our previous learnings to implement the on-line sampling heuristic.

5.5 Sampling Heuristic Overview

We mentioned earlier that the KPI value can be considered as a response surface of a multi-variate input space defined by the workload and resource metrics. Several generic techniques are available to sample response surfaces, referred to as *Design of Experiments* (DoE). An important learning from this research domain is that edge values are the best starting points to characterize a black-box response surface. Based on this strategy, we have also constructed our sampling heuristic so far:

1. We defined the important workload, resource metrics and their (edge) values.
2. As an initial measurement, each combination of the different edge values is measured (a 2-level, full factorial measurement). A feature selection analysis on

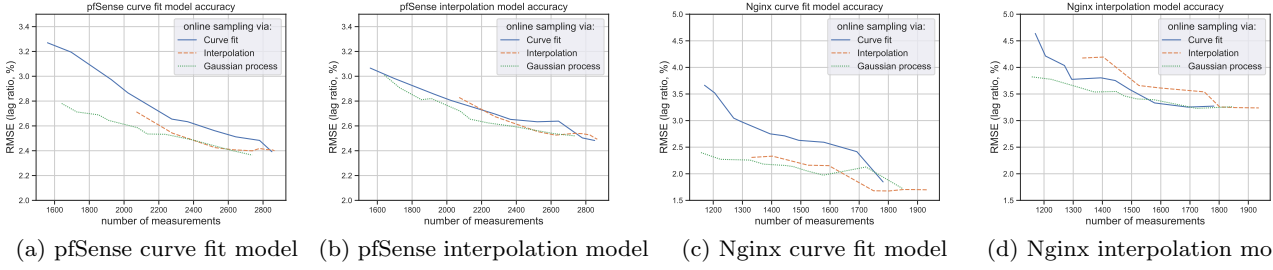


Fig. 9 When sampling the workload metric wl_p online, a Gaussian Process (GP) works best to assess when enough samples are taken. The GP produces the best model for the KPI trend, under varying wl_p values, and with the least amount of samples.

this initial sample set reveals which metrics contribute most to the KPI variation and are thus more interesting to sample.

- One specific workload metric, wl_p , is bisected more fine grained, while the other metrics are iterated through fixed pre-defined values.

We want to select in an online way, which next configuration to measure. This means we need to decide from the previous samples, which next workload and resource configuration will reveal thee most information to model the KPI trends. We use following procedure to achieve this: (i) the total sample space is divided into configurations (all possible combinations of workload and resource settings), and (ii) in each configuration, a number of values of wl_p is sampled through bisection. Practically, for our measured VNFs, this means:

- **pfSense**: Each combination of these metrics settings is a configuration: [CPU allocation, flows, packetsize]. In each configuration, we bisect new values for $wl_p = packetrate$, between the given boundary values.
- **Nginx**: Each combination of these metrics is a configuration: [flavor, quality, movies]. In each configuration, we bisect new values for $wl_p = streams$, between the given boundary values.

Per VNF we have thus defined a configuration space. We use Algorithm 3 to select which configuration in this space is most interesting to profile next. In this algorithm, the feature F_n with values $f \in F$, refers to a configuration metric of the VNFs as defined above. We select the next configuration metric value based on previously measured values. The next value is located between the two previously measured values which had the most influence on the mean KPI value (characterized by $\Delta m_{f,i}$). Or in other words, we bisect the new value between the points where the average KPI value showed the highest change. The value returned is then used as the next configuration setting where new values of wl_p are sampled using Algorithm 1.

To calculate the RMSE, we have validated our sampling heuristic using an extra test dataset. However,

Algorithm 3: Online configuration selection

Data: F = list of possible values of feature F_n

Result: f_{new} = new value of F_n to sample

- $F_s \leftarrow$ already sampled values of F_n ;
 - $S \leftarrow$ empty set ;
 - foreach** value f in F_s **do**
 - $M \leftarrow$ find all KPI measurements where $f == F_n$
 - ;
 - $m_f \leftarrow mean(M)$;
 - $S.append(f, m_f)$;
 - end**
 - sort S by f ;
 - for each interval in S calculate:
 - $\Delta m_{f,i} = |m_f[i] - m_f[i+1]|$;
 - find interval in S where $\Delta m_{f,i}$ is the largest;
 - $f_{new} \leftarrow$ middle value in F between $f[i]$ and $f[i+1]$;
 - while** $f_{new} \in F_s$ or $f_{new} == empty$ **do**
 - take f_{new} in interval with next largest $\Delta m_{f,i}$;
 - end**
 - return f_{new} ;
-

when using the online sampling mechanism on a new VNF, we likely have no test samples available to asses the accuracy of the model being profiled. To alleviate this, we investigate if we can use the $\Delta RMSE$ of the obtained samples so far, as a stop criterion for the profiling procedure, similar to Algorithm 1. In Fig. 10, we calculate the $\Delta RMSE$ when a new workload or resource metric value is added to the profiled measurements. This is an indication if the last obtained samples contributed much to total accuracy of the model so far. We see indeed on the plots that $\Delta RMSE$ stabilizes around the same number of samples where the RMSE stabilizes in Fig. 5.5. This indicates that $\Delta RMSE$ could indeed be used as a stop criterion, to assess when enough samples are profiled and more samples will not improve the accuracy much more.

Figure 11 summarizes the workflow of our presented sampling heuristic. After feature selection, an initial subset of the configuration space is sampled, as depicted by the black dots in Fig. 11b. The grey dots are added by Algorithm 3. At any point in the sampling proce-

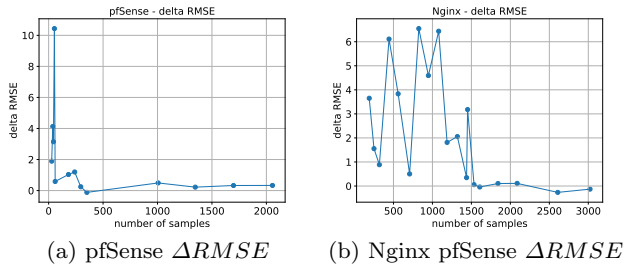


Fig. 10 Variation in the $\Delta RMSE$, to be used as stop criterion for the sampling.

ture, more samples can be added for any of the metrics using the loop in the last two blocks of Fig. 11a. This can be done to increase the accuracy of the prediction model, if there is still budget available for additional profiling time.

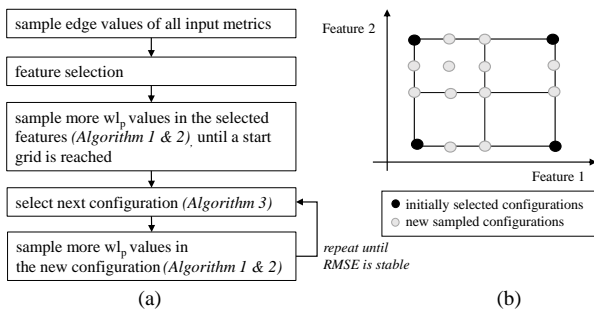


Fig. 11 (a) The sampling workflow with (b) simplified representation of the sampled configuration space.

We show the result of the sampling heuristic in Fig. 12. Two selected modeling methods (interpolation and curve fit) are used to indicate the accuracy of the trained model on the test dataset. The heuristic is used to select training samples from the training set. the test dataset is afterwards used to calculate the model accuracy (RMSE). It can be clearly seen that the used sampling heuristic has a positive effect on the needed sample size, compared to the uniform sampling strategy we put forward as benchmark. This beneficial for the total profiling time needed to model the KPI trends of these VNFs. This result also further extends the findings in [23], showing that the curve fit method is also well performing at smaller datasets, with optimally chosen training samples.

In Table 1 and 2 the subsequent rows indicate which metrics have been measured and what the resulting RMSE is (using the best method seen in Fig. 12). Near the bottom of the tables, more workload and resource configurations have been measured, resulting in a lower RMSE and thus better model accuracy.

While sampling the pfSense dataset (Table 1), we can see that most improvement comes after sampling only the first feature (cpu allocation). Sampling the next workload features (flows, packetsize) brings subsequently less improvement. This is aligned with the feature selection done in the beginning of the profiling procedure. As indicated in Fig. 8, cpu allocation is indeed the most important metric in this dataset.

For the Nginx dataset (Table 2), there were no features which jumped out during the feature selection phase. Each metric (flavor, quality and movies) is considered to have significant effect on the KPI. We therefore pick the metrics values round-robin across all input metrics.

It can be noted that also other, more practical, reasons can influence the order in which metrics are selected. For example, it might be more beneficial to limit many resource allocation changes in the sampling procedure. Often, a VNF restart is required to make use of the newly allocated resources such as CPU cores. This reboot (and possible state reconfiguration) might induce a lot of delay in the test procedure. To mitigate this, it can be opted to change the resource allocation as very last feature after all other workload metric values have been sampled.

5.6 Modeling Method Impact

In Fig. 13 the same sampled dataset by our heuristic is used to train different modeling methods. It can be clearly seen that the sample set also has an effect on the used model method. During our first tests using uniform sampling, we already found large variations in accuracy between several modeling methods (as seen in Fig. 5). Using uniform sampling, the interpolation method seemed to provide the best accuracy. Earlier, in Section 5.4, we selected GP as the best method to model a limited one-dimensional sample subset. This does not hold any more when using GP on the total, multi-dimensional input space.

We selected interpolation and curve fit as most promising methods to be used with our sampling heuristic. In order to use the curve fit method, some expert knowledge or experience is required to define a priori the trends, which need to be fitted to the KPI measurements (as explained in Section 4.1). If such an analytic function cannot be specified, interpolation seems the next best choice.

In Fig. 12 and 13, the curve fit method shows improved accuracy over the interpolation method. The curve fit method has two major advantages:

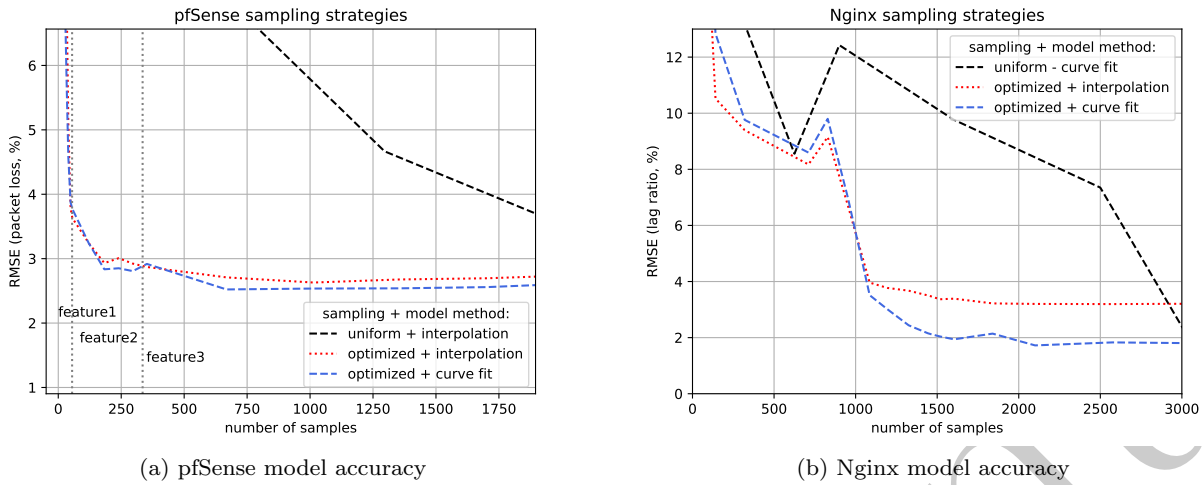


Fig. 12 Optimized sampling strategies compared to uniform sampling. Less samples are needed to achieve a stable and high model accuracy.

	cpu alloc. (%)	flows	packetsize (B)	packetrates (kpps) (wl_p)	# samples	RMSE(packet loss)
Edges only	[0.25, 5]	[1, 10000]	[64, 1500]	[0.1, 500]	16	(feature selection)
Feature1	[0.25, ..., 5]	10	256	bisect max 10 values (via algorithm 2)	56	3.77
Feature2	[0.25, ..., 5]	[1, ..., 10k]	256		352	2.92
Feature3	[0.25, ..., 5]	[1, ..., 10k]	[64, ..., 1500]		2054	2.61

Table 1 Configuration picking order for pfSense. Features are sampled in order of significance in Fig. 8. Each feature is completely sampled before the next one.

	flavor	quality	movies	streams (x1000) (wl_p)	# samples	RMSE (lag ratio)
Edges only	[0, 11]	[1, 5]	[1, 5]	[0.01, 5]	16	(feature selection)
Feature1/2/3	[0, 11]	[1, 5]	[1, 5]	bisect max 10 values (via algorithm 2)	40	25.0
Feature1/2/3	[0, 5, 11]	[1, 3, 5]	[1, 3, 5]		141	12.8
Feature1/2/3	[0, 2, 5, 11]	[1, 3, 4, 5]	[1, 3, 4, 5]		320	9.8
Continue to bisect a new value per feature, alternately, until all feature values are sampled...						
Then also allow more wl_p samples.						
Feature1/2/3	[0, ..., 11]	[1, ..., 5]	[1, ..., 5]	max 10 values	1443	2.2
Feature1/2/3	[0, ..., 11]	[1, ..., 5]	[1, ..., 5]	max 50 values	3025	1.8

Table 2 Configuration picking order for Nginx. Features are sampled in order of significance in Fig. 8. New feature values are bisected round-robin over the resource and workload metrics.

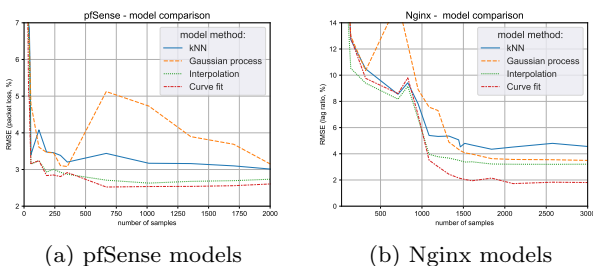


Fig. 13 Comparison of the prediction accuracy of different methods trained with the same sampled dataset.

1. The analytical functions which are fitted have a guaranteed monotonicity. This makes it easier to invert the model and predict a recommended resource allocation from a given KPI and workload target, as no local extrema are present in the modeled KPI. This is explained more elaborately in [23].
2. The curve fit method succeeds to model more efficiently the KPI trend break, happening when resources get saturated at increasing wl_p (see also Fig. 4), using less training samples. It provides better accuracy at low KPI values, where resource saturation is starting but not yet completely dominant.

6 Related Work

A generic architecture for profiling frameworks has been described earlier in [17, 11, 16, 2, 7], where also its relevance to DevOps related workflows is highlighted. This previous work has been extended here with more insights for using a Service Oriented Architecture and analysis of both sampling and modeling methods. Also we exemplify VNF profiling in a more elaborate space of both workload and resource metrics.

The modular architecture of our profiling framework (see Section 2) allows an easy deployment on cloud-based infrastructures, which further enhances the applicability of this framework in cloud-native platform environments as described in [22]. This cloud-native nature allows the profiling framework to be used in light-weight test environments, as well as in production-grade staging environments, with realistic workloads. The use-cases for profiling, in the context of a cloud-native orchestration platform, have been outlined in [10]. It is explained how a profiling workflow can be supported by an NVF based orchestration platform, by leveraging existing Monitoring and DevOps related processes.

Several sampling strategies for VNF chain profiling have been investigated in [16], but no method is found which significantly beats a generic uniform sampling strategy. Decision tree based models are put forward as a promising solution, however our tests show that random forest is one of the less performing methods in our use-case of profiling a single VNF. The profiling tests in [16] are also limited to varying only one resource metric (allocated vCPU), under only one fixed workload.

The VNF profiling platform in [7] exhaustively executes a list of input benchmark tests, without trying to optimize the test time. Decision trees are used as a classification technique to divide tested VNFs into pre-defined categories for resource demand. The authors prefer decision trees over other algorithms because they are very interpretable and a system administrator can clearly understand what the final model is doing. This methodology however, does not try to predict the KPI values, but tries to map existing resource flavors to required performance categories. A similar concept is implemented in [18], where the VNF performance is mapped into discrete hardware flavors. While the VNF is monitored during operation, deviations compared to earlier profiled data can be detected. This can trigger the migration of the VNF to a different node to avoid resource saturation. The KPIs are however only based on resource utilization, with no consideration for SLA based performance metrics. In the context of anomaly detection, a VNF profile is used in [11] as baseline performance model. Any performance

deviation compared to the VNF profile is the trigger for further troubleshooting analysis. It is clear that these use cases could benefit from our sampling heuristic to generate the required VNF profiles for later reference.

The profiling method in [6] focuses on the deployment space of big data applications, with up to seven configuration dimensions. In this deployment space, areas are clustered in which the performance metric can be approximated using a linear model. The total model is then a piecewise combination of different linear models and the space partitioning into different linear regions is done by a decision tree. More samples are adaptively chosen in the areas where a linear model has bad accuracy. When trying to model non-linear functions, a deteriorating performance is reported. Moreover, new sample points are drawn randomly in a uniform manner, without exploiting any expert knowledge. Our sampling heuristic tries to mitigate this by using a curve fit model with piecewise (non-linear) functions and online sample selection using bisection.

The authors in [5] recommend the use of Gaussian Processes (GP) as a flexible method to both model performance and choose the next configurations to sample. Our tests show however that the trends witnessed in VNF profiling are not easily captured by GPs. Also when we let the GP model determine the next samples, the method remains sub-optimal compared to generic uniform sampling (see Section 5.2). Also in [20] it is experimentally tested that uniform sampling and Gaussian Processes are not the best methods to sample the large parameter space of computing programs. The used method is however aimed at finding a global optimum, which is not applicable to our use case, which tries to model the complete response surface.

In [8], it is also experimentally observed that a VNF performance trend break occurs when resources are getting saturated (as in our Fig. 4). This is exploited by estimating at which point the increasing workload is no longer correlated with the saturating resource usage. This way, a prediction model can be made, mapping workload and resource usage to profiled VNF performance. During profiling, this approach increases the workload in fixed steps until saturation occurs. We argue that our approach, based on bisection, provides a faster way to model the maximum performance, needing less samples.

Our research complements and extends this related work by investigating sampling strategies for both the workload and resource allocation space, specifically to model VNF performance. The sampling heuristic described in our paper can be adopted to generate profiled training data for more VNF related models and use cases. For example, to gather benchmark datasets

for anomaly detection or capacity planning. Optimized sampling could also complement the profiling workflow in [23], where a trend is profiled to extrapolate the performance of the VNFs on increasing and untested resource allocations.

7 Summary and Conclusion

We have outlined several optimizations to develop and execute a VNF profiling workflow:

- A micro-service oriented architecture to implement a VNF profiling procedure. The benefits of this architecture are that multiple profiling tests can run in parallel, allowing scaled up and faster generation of measurement data. Additionally, a quick, light weight test set up and metric definition is pursued. Also increased robustness can be expected, since the profiling tests run isolated from each other and a failure in one test will not affect other running tests.
- An optimized sampling heuristic to select online which workload and resource allocation to measure next in the profiling procedure. This results in a smaller sample set and thus less profiling time needed to model accurate KPI predictions. Compared to generic uniform sampling, our tests show that up to 5x less samples are needed for the same accuracy.
- We have analyzed the effect of our proposed sampling heuristic on several modeling methods. Careful use is advised for some typical methods used in machine learning such as Gaussian Processes, Random Forest, Support Vector Machines. Other methods can be trained more accurately with the obtained sample sets, as illustrated by our proposed interpolation and curve fit method.

Our presented profiling framework also offers good integration possibilities into cloud-native platforms, this allows to profile VNFs in representative environments, under realistic workloads. We have selected a set of sampling and modeling methods which we believe to be generic enough, to be applied to a wide range of VNF functionalities. By fitting the measured VNF trends to analytic functions, we add expert knowledge into the VNF performance model. As a result, we see clear improvement of the curve fit method compared to generic parameterized methods typically used in machine learning approaches. This also confirms earlier research results presented in [23].

Our proposed sampling heuristic needs a set of inputs to start from, namely a specified range for workload and resource allocation metrics. One primary workload metric wl_p is bisected in a fine grained way, the

other ones will iterate through a list of specified values which the sampling heuristic will pick online in an optimized way. By using feature selection and bisection, the heuristic will select from previous measurements which workload and resource configuration to profile next. This is implemented by our presented Algorithms: 1, 2 and 3. A stop criterium is included to automatically decide when enough samples are gathered to achieve good accuracy. The presented heuristic is well suited to assist in VNF profiling tasks which can be completed faster due to a reduced testing time. We hope this work can inspire ongoing research in the optimization of VNF profiling.

Acknowledgements This work has been performed in the framework of the NGPaaS (<http://ngpaas.eu>) and 5GTANGO (<https://www.5gtango.eu>) project, funded by the European Commission under the Horizon 2020 and 5G-PPP Phase2 programmes, resp. under Grant Agreement No. 761 557 and 761 493. This work is partly funded by UGent BOF/GOA project 'Autonomic Networked Multimedia Systems'.

Conflict of interest

On behalf of all authors, the corresponding author states that there is no conflict of interest.

References

1. 5G-PPP Software Network Working Group (2019) Cloud-native and verticals' services. Tech. rep., 5G-PPP Software Network Working Group, URL https://5g-ppp.eu/wp-content/uploads/2019/09/5GPPP-Software-Network-WG-White-Paper-2019_FINAL.pdf, accessed: 2019-10-01
2. Blenk A, Basta A, Henkel L, Zerwas J, Kellerer W, Schmid S (2018) perfbench: A tool for predictability analysis in multi-tenant software-defined networks. In: Proceedings of the ACM SIGCOMM 2018 Conference on Posters and Demos, ACM, pp 66–68
3. Crombecq K, Gorissen D, Deschrijver D, Dhaene T (2011) A novel hybrid sequential design strategy for global surrogate modeling of computer experiments. *SIAM Journal on Scientific Computing* 33(4):1948–1974
4. Dunn K (2019) Process improvement using data. URL <https://learnche.org/pid>, accessed 2019-10-01
5. Duplyakin D, Brown J, Ricci R (2016) Active learning in performance analysis. In: 2016 IEEE Inter-

- national Conference on Cluster Computing (CLUSTER), pp 182–191
6. Giannakopoulos I, Tsoumakos D, Koziris N (2017) A decision tree based approach towards adaptive modeling of big data applications. In: 2017 IEEE International Conference on Big Data (Big Data), IEEE, pp 163–172
 7. Gokan Khan M, Bastani S, Taheri J, Kassler A, Deng S (2018) Nfv-inspector: A systematic approach to profile and analyze virtual network functions. In: 2018 IEEE 7th International Conference on Cloud Networking (CloudNet), pp 1–7, DOI 10.1109/CloudNet.2018.8549333
 8. Iglesias JO, et al. (2017) Orca: an orchestration automata for configuring vnfs. In: Proceedings of the 18th ACM/IFIP/USENIX Middleware Conference, ACM, pp 81–94
 9. Kulkarni SG, et al. (2017) Nfvnic: Dynamic backpressure and scheduling for nfv service chains. In: Conference of the ACM Special Interest Group on Data Communication, ACM, pp 71–84
 10. Mimidis-Kentis A, Soler J, Veitch P, et al. (2019) The next generation platform as a service: Composition and deployment of platforms and services. *Future Internet* 11(5):119
 11. Nam J, Seo J, Shin S (2018) Probius: Automated approach for vnf and service chain analysis in software-defined nfv. In: Proceedings of the Symposium on SDN Research, ACM, p 14
 12. Pantos R, May W (2017) Http live streaming, rfc 8216. Tech. rep., URL <https://tools.ietf.org/html/rfc8216>, accessed 2019-10-01
 13. Parmar H, Thornburgh M (2012) Real-time messaging protocol (rtmp) specification. *Adobe specifications, December*
 14. Pedregosa F, Varoquaux G, et al. (2011) Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12:2825–2830
 15. Peuster M, Karl H (2017) Profile your chains, not functions: Automated network service profiling in devops environments. In: Network Function Virtualization and Software Defined Networks (NFV-SDN), IEEE
 16. Peuster M, Karl H (2018) Understand your chains and keep your deadlines: Introducing time-constrained profiling for nfv. In: 2018 14th International Conference on Network and Service Management (CNSM), pp 240–246
 17. Rosa RV, Bertoldo C, et al. (2017) Take your vnf to the gym: A testing framework for automated nfv performance benchmarking. *IEEE Communications Magazine* 55(9):110–117
 18. Sciancalepore V, Yousaf FZ, et al. (2018) z-torch: An automated nfv orchestration and monitoring solution. *IEEE Transactions on Network and Service Management* 15(4):1292–1306, DOI 10.1109/TNSM.2018.2867827
 19. Smola AJ, Schölkopf B (2004) A tutorial on support vector regression. *Statistics and computing* 14(3):199–222
 20. Thiagarajan JJ, Jain N, et al. (2018) Bootstrapping parameter space exploration for fast tuning. In: Proceedings of the 2018 International Conference on Supercomputing, ACM, pp 385–395
 21. Van Rossem S, Tavernier W, Peuster M, Colle D, Pickavet M, Demeester P (2016) Monitoring and debugging using an sdk for nfv-powered telecom applications. In: IEEE NFV-SDN2016, the IEEE Conference on Network Function Virtualization and Software Defined Networks, pp 1–5
 22. Van Rossem S, Sayadi B, et al. (2018) A vision for the next generation platform-as-a-service. In: 2018 IEEE 5G World Forum (5GWF), IEEE, pp 14–19
 23. Van Rossem S, Tavernier W, Colle D, Pickavet M, Demeester P (2019) Profile-based resource allocation for virtualized network functions. *IEEE Transactions on Network and Service Management* pp 1–1, DOI 10.1109/TNSM.2019.2943779
 24. Williams CK, Rasmussen CE (2006) *Gaussian processes for machine learning*, vol 2. MIT press Cambridge, MA, URL <http://www.gaussianprocess.org/gpml/chapters/>, accessed 2019-10-01