Received XX Month, XXXX; revised XX Month, XXXX; accepted XX Month, XXXX; Date of publication XX Month, XXXX; date of current version 11 January, 2024.

Digital Object Identifier 10.1109/OJCOMS.2024.011100

Minimum-Cost Design of Auto-Scaling Server Farms Providing Reliability Guarantees

Jesus Perez-Valero¹, Pablo Serrano² (Senior Member, IEEE), Jaime Garcia-Reinoso³, Albert Banchs^{2,4} (Senior Member, IEEE), Xavier Costa-Perez^{5,6,7} (Senior Member, IEEE)

¹Department of Information and Communications Engineering, University of Murcia, 30100 Murcia, Spain ²Department of Telematic Engineering, University Carlos III of Madrid, 28911 Leganés, Spain ³Department of Automatic, University of Alcalá, 28871 Alcalá, Spain. ⁴IMDEA Networks Institute, 28911 Leganés, Spain ⁵NEC Laboratories Europe, 69115 Heidelberg, Germany ⁶ICREA, 08010 Barcelona, Spain ⁷i2cat, 08034 Barcelona, Spain

CORRESPONDING AUTHOR: Jesus Perez-Valero (e-mail: jesus.perezvalero@um.es).

ABSTRACT As next-generation mobile networks increasingly rely on virtualized infrastructure to deliver critical services, ensuring both the efficiency and reliability of server farms becomes essential. These infrastructures must meet stringent reliability guarantees to support time-sensitive applications in emerging 5G and beyond networks. In this paper, we address the design of autoscaling server farms –specifically, selecting the most suitable server type and corresponding number of servers–by considering both service requirements and associated operational and infrastructure costs. To this end, we develop an optimization algorithm that combines (i) a queueing-theoretic model to estimate the resources needed to meet reliability constraints, and (ii) a general cost model that captures both capital and operational expenditures. We validate our approach through extensive simulations, comparing it against classical queueing-based methods and exhaustive numerical searches: our proposal reduces costs by 22% as compared against the benchmark, with solutions that are within 3% of numerical searches at 10% of the computational complexity, offering a new scalable and cost-effective methodology for designing reliable server farms.

INDEX TERMS URLLC, B5G, Auto-scaling, NFV, Reliability

I. INTRODUCTION

With the arrival of Network Function Virtualization (NFV) [1], mobile services will be implemented as interconnected virtual network functions (VNFs) hosted by cloud servers [2]. To make an efficient use of the resources, these VNFs need to be scaled up and down based on their load [3]. This scaling does not impose several challenges when dealing with traditional best effort services, and therefore the impact of non-zero activation times, or the fallibility of servers, is negligible. However, when dealing with services with stringent reliability requirements, such as e.g. Ultra Reliable Low Latency Communications (with reliability levels of up to 5 or 6 nines [4]), the impact of these parameters cannot be neglected.

In fact, in our previous work we analyzed the nonnegligible impact of non-zero start-up times and finite server lifetimes on the reliability of services provided by an auto-scaling server farm [5], [6]. We first assumed a fixed scenario with a given server farm and developed an analytical model to optimize its activation and deactivation thresholds [5]. Later, we proposed a configuration mechanism that dynamically adjusts these thresholds to ensure a given reliability guarantee [6] to the observed conditions. In these works we illustrate that, depending on traffic conditions and service requirements, a few carrier-grade servers (which are powerful and reliable but consume more energy) may be less efficient than many consumer-grade servers (which are less powerful but more energy-efficient).

Motivated by the above finding, in this paper we tackle the design problem: given (i) a specific service to be supported in a server farm, characterized by a target reliability guarantee and an arrival and service rate; and (ii) a set of candidate server types, each one with a given set of characteristics, the challenge is to determine the optimal deployment configuration (i.e., number of servers and their type) that ensures that the service is provided with the required guarantee while minimizing the total cost. To compute this cost we assume a generic cost model that assumes some infrastructure ownership, and therefore includes capital expenditure (CapEx) and operating expenditure (OpEx) terms, but it could be used in other scenarios such as e.g., a cloud-based service provision, which requires pre-booking resources with different cost and performance characteristics.

We illustrate the effectiveness of our proposal assuming five different types of servers, which are modeled after existing equipment. According to the results, our proposal reduces costs by 22% as compared against the benchmark based on classical queueing theory, with solutions that are within 3% of exhaustive numerical searches at 10% of the computational complexity. As compared against previous works considering scaling farms for NFV with reliability requirements (duly reviewed in Section VI), the main novelties and contributions of this work as are follows:

- We develop an analytical model to estimate the number of resources that a server farm requires to support a given service with certain reliability requirements.
- We present a cost model to estimate the cost of providing the service, including CapEx and OpEx terms and leveraging the analytical model to estimate the server occupation.
- We formalize the optimal design of a server farm, which given a set of different server types, selects the type and number of servers that minimize the cost while guaranteeing the required reliability.
- We illustrate how to extend the design for the case of heterogeneous scenarios, where different services with different reliability requirements are supported by different type of services.
- We assess the accuracy of the analytical model and validate the server farm design algorithm through extensive simulations, which assume five different types of servers modeled after real-life service types and equipment.

The rest of the paper is organized as follows. Section II outlines the key motivations and challenges addressed in this work. Section III describes the system model, defining the key components and assumptions. Section IV introduces the cost model used throughout the study, followed by an analytical characterization of the system's key performance variables, and concluding with the proposed algorithm for server farm design. Section V evaluates the accuracy of our analytical model and the effectiveness of the proposed approach through extensive

simulations. Section VI provides a comprehensive review of related work on the performance analysis and configuration of auto-scaling server farms. Finally, Section VII summarizes our key findings and conclusions.

II. Motivation and challenge

The growing reliance on virtualized infrastructures to support critical services in next-generation mobile networks introduces significant challenges, particularly in ensuring the reliability of auto-scaling server farms [7], [8]. Unlike traditional best-effort services, where minor delays in resource activation or occasional hardware failures have little impact, mission-critical applications demand stringent reliability guarantees. Services such as industrial automation, autonomous driving, and remote healthcare require failure probabilities as low as one in a million (99.9999%), making even brief disruptions unacceptable. These requirements impose not only an operational challenge but also a design challenge: the number of active resources must be dynamically adapted to the observed traffic, and the type and number of servers must be selected to meet reliability goals without incurring excessive cost.

Regarding the operational challenges, since traffic loads fluctuate dynamically, adaptive resource allocation is needed to prevent service interruptions while minimizing Operational Expenditures (OpEx). Some solutions rely on auto-scaling techniques to adapt to the traffic load [9], [10] but often assume instantaneous activation of pshysical machines (PMs), which is unrealistic in practical deployments, while other solutions rely on analytical tools while taking into account the non-zero boot up times and fallability of servers [5], [11] (we review the related work in Section VI). Since these works do not tackle the design challenge, one possible strategy would be over-provisioning, where additional servers are deployed to ensure redundancy. While this approach improves reliability, it may result in significant Capital Expenditures (CapEx).

In this paper, we develop a cost-aware optimization framework that selects the optimal type and number of servers, minimizing both CapEx and OpEx while meeting stringent reliability constraints. We focus on the worst-case scenario, assuming that the peak traffic load corresponds to the average traffic load as this imposes the most demanding operational constraints, ensuring that our methodology remains applicable to other scenarios.

III. System model

We assume the same system as in [5], [6], i.e., an autoscaling server farm for NFV, comprising a centralized infrastructure manager (IM) and several physical machines (PMs). Initially, we assume a homogeneous server deployment (i.e., server homogeneity), meaning that all PMs are modeled after the same type of server. This assumption simplifies the operational procedures, as maintaining a uniform hardware profile reduces management complexity and operational costs, making heterogeneous deployments more costly [12]. We describe how our framework can be adapted to heterogeneous scenarios in Section E.

In our paper, we refer with "task" to a traffic session, which refers to an instance of a stringent URLLC service such as remote-assisted driving, an industrial slice in a factory environment, or tele-operated driving. Assuming independence among traffic sessions, following the Palm–Khintchine theorem the aggregate arrival process asymptotically behaves like a Poisson process.

Following the above, a deployment can be modeled with a set of parameters that characterizes its performance: each PM can support a maximum of M tasks, which we refer to as server capacity, and has an exponential lifetime¹ and boot up times, with average $1/\beta$ and $1/\alpha$, respectively (the key variables in our paper are summarized in Table 1). Following [14], [15], we assume a load-proportional power consumption model for the PMs, characterized by a fixed term P_{idle} , and the proportional term P_{load} , which can be computed as the difference in tasks between the so-called peak power consumption P_{peak} and the idle term, i.e.,

$$P_{\text{load}} = \frac{P_{\text{peak}} - P_{\text{idle}}}{M}.$$
 (1)

Each server has a monetary cost K and a lifespan R. We denote with τ_i the set of parameters that characterizes a given server type i, i.e., $\tau_i = \{M_i, \beta_i, \alpha_i, P_{\text{idle},i}, P_{\text{load},i}, K, R\}$, and with $\mathcal{T} = \{\tau_1, \tau_2, \ldots\}$ the set of candidate server types.

Tasks arrive to the system following a Poisson process (note that this assumption is relaxed in our performance evaluation) at a rate λ and require an exponential service time of average $1/\mu$. During the initial planning phase, we rely on parameter estimates to dimension the system as optimally as possible, while at runtime, we employ adaptive mechanisms (such as stochastic optimization [11]) to guide operational decisions, including determining the number of active servers. Moreover, system measurements collected during operation can be used to iteratively refine parameter estimates, thereby improving system dimensioning over time.

The IM balances the load across PMs, seamlessly migrates tasks whenever needed (e.g., a machine is about

to fail), and powers PMs on and off as needed. Each PM can be in one of three states: active (serving tasks), booting up (initiating due to a need for more resources), or stopped (either due to a crash or because they are not needed to handle the current traffic load). To power on/off the PMs, the IM implements the following policy. At least one server is kept active at all times to avoid any delay to serve in arriving tasks. For the rest of the servers, a threshold-based policy is followed, where the thresholds to power on or off a server depend on the number of active servers, denoted by m: a new server is activated when the number of tasks in the system reaches s_m , and deactivated when the number of tasks reaches $s_m - 1$. The motivation behind this lack of hysteresis is motivated by the objective of achieving the largest possible energy savings, which is accomplished by keeping the absolute minimum number of servers active at any given moment to meet performance guarantees [5]. These energy savings are gained at the expense of a higher frequency of server state transitions, which might lead to some hardware wear-and-tear.

The farm serves tasks with high reliability requirements, e.g., an industry 4.0 service [4] or autonomous driving services [16]. When a task arrives, the IM selects an active PM with sufficient resources to handle it. If no PMs are available to handle the task, the IM queues it until either a task finishes service or a PM boots up. In any case, this results in a service disruption, which negatively affects the committed reliability of the task. Since PMs are prone to failure, an active PM may crash at any time. If it was processing tasks, we assume that these can be seamlessly (i.e., with negligible latency and energy overheads) migrated to another active PM, where sufficient resources are available –note that this is already supported in Linux environments [17] with open-source technologies such as ACHO $[18]^2$. In case there are not enough active PMs, the affected tasks are placed on hold until sufficient PMs are activated again, which again is considered a service disruption.

We assume that the successful provisioning of the service requires a minimum reliability level. This reliability level is determined by the probability of a task being disrupted, which is denoted as P_f , being below a maximum threshold, denoted as T_f . Although the terms "reliability" and "failure probability" refer to complementary terms, e.g., a reliability of 3 nines (99,9%) corresponds to a failure probability of 10^{-3} , for readability reasons we will use them interchangeably throughout the paper.

¹Additional experiments (not reported here due to space constraints) indicate that our framework also applies to other scenarios. These include: (i) correlated failures, modeled as a twostate Markov chain with distinct average lifetimes in each state [13]; and (ii) non-exponential distributions for boot-up times and lifetimes. Specifically, we tested constant and Weibull-distributed boot-up times, and Weibull-distributed lifetimes across multiple server classes, using parameters chosen to match the same mean values as in the exponential case.

²For instance, assuming a failover mechanism over a reliable wired network, as in [18], and a memory transfer of 64 MB over a 1 Gbps link, each migration requires approx. 500 ms, which is significantly shorter than typical session durations. Based on our simulation results and the power model in [19], this results in an additional power consumption of approx. 0.18 W, which is negligible compared to the idle consumption of a server.

Variable	Description
N	Maximum number of physical servers
M	Capacity of one server
λ	Task arrival rate
$1/\mu$	Average service time of a task
$1/\beta$	Average lifetime of a server
$1/\alpha$	Average boot time of a server
s_m	Activation threshold with m active servers
C	Total cost
C_{Cap}	Cost associated to CapEx
C_{Op}	Cost associated to OpEx
K	Equipment cost
R	Equipment lifespan
N_u	Average number of users in the system
N_a	Average number of active servers
P_{load}	Energy consumption due to resource usage
$P_{\rm idle}$	Energy consumption due to active server count
P_f	Failure probability of the server farm
T_f	Target failure probability

TABLE 1: Key variables used throughout the pap
--

To guarantee the required reliability, we assume that the server farm executes the algorithm presented in [6], which automatically drives the (de)activation thresholds $\{s_m\}$ to an adequate point of operation (note that in [6] we illustrate that the algorithm provides the most energy saving operation, but we did not provide the actual values of these parameters). In this paper, we address the optimal design of the server farm, i.e., given the set of available candidate server types \mathcal{T} , determine the most appropriate server type, which is denoted by τ^* , and the required number of servers, denoted by N, that guarantees the required performance while minimizing the cost.

IV. Analysis and design of the server farm

In this section, we formalize the optimization problem that we address in this paper, which is the choice of the optimal server type and the number of servers to support a given service while minimizing the cost. To this aim, we first introduce our cost model, and then an analytical model to characterize the operation of the server farm, which is used by the algorithm to select the best server type.

A. Cost model

4

Following the usual assumptions in the literature (e.g., [20]), we assume that the total cost of the server farm composed of PMs of type τ is given by the sum of the Capital Expenditure (CapEx) and Operating Expenditure (OpEx).

$$C(\tau) = C_{\text{Cap}}(\tau) + C_{\text{Op}}(\tau) \tag{2}$$

where the CapEx term C_{Cap} is determined by initial investments in infrastructure, while the OpEx term C_{Op} is determined by the cost of running and maintaining the service. More specifically, C_{Cap} is determined by the number of servers required to provide a service N, their cost K, and the equipment lifespan in hours R, according to:

$$C_{\rm Cap} = N \times \left\lceil \frac{K}{R} \right\rceil,\tag{3}$$

where R is the average lifespan for the type of server considered, taking into account the expected load and wear-and-tear effects; the OpEx term is determined by the resource consumption due to the service provisioning. Following our previous work [6], this is given by the amount of resources required to process the tasks, and therefore the OpEx term can be expressed as

$$C_{\rm Op} = (P_{\rm load}N_u + P_{\rm idle}N_a) \times k_{\rm Wh} \tag{4}$$

where N_u denotes the average number of users in the server farm, N_a denotes the average number of active servers, and $k_{\rm Wh}$ represents the monetary cost of energy per hour.³ Note that we consider that the server farm is working continuously in a time span of a year.

Both (3) and (4) depend on several parameters: a subset of them corresponds to numerical figures that are determined for a given server type τ_i , i.e., K, R, P_{load} , and P_{idle} , while the rest of them depend on the operation conditions, i.e., the average number of users N_u , the total number of servers N, and the average number of active servers N_a . We next present an analytical model to compute these.

B. Average number of users N_u and user distribution

To compute the average number of tasks and their distribution, we make the approximation that incoming tasks never wait to be attended. This approximation is motivated by the fact that the number of tasks that will have to wait is very small and hence this approximation will have a very small impact on the resulting task distribution. Assuming that the impact of server failures is negligible, the system behaves as an $M/M/\infty$ system [21], and therefore the probability of having n users in the system, denoted by p_n , follows a Poisson distribution given by the following expression

$$p_n = \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} e^{-\frac{\lambda}{\mu}} \tag{5}$$

where the average number of users is

$$N_u = \frac{\lambda}{\mu}.$$
 (6)

³For simplicity, our cost model assumes a fixed average energy price. This could be extended to incorporate dynamic pricing (e.g., time-of-day tariffs) by adjusting for load and cost variations across time periods.

C. Total number of servers N

To compute the required number of servers, we look at the number of servers that are needed to be able to serve all tasks with a very high probability (provided that all servers are active). In particular, we enforce that the probability that an arriving task does not have to wait, when all servers are active, is well above $1 - T_f$, i.e., the failure due to all servers being busy is much smaller than the failure due to other reasons, which can be as large as T_f . More specifically, we enforce that the probability that a task does not have to wait is equal to $s = 1 - T_f/X$, where X is a sufficiently large value (unless otherwise stated, in the rest of the paper we take X = 10).

Note that following (5) the s-percentile of the total number of users in the system, denoted as $N_u(s)$, is given by

$$N_u(s) = \min\left\{Q \in \mathbb{N} \left| \sum_{n=0}^Q \frac{1}{n!} \left(\frac{\lambda}{\mu}\right)^n e^{-\frac{\lambda}{\mu}} \ge s \right. \right\}$$
(7)

To meet the requirements stated above, we need to ensure that the system can host up to $N_u(s)$. This means that, if the size of a server in tasks is given by M, the number of servers in the system needs to be dimensioned as follows:

$$N = \left| \frac{N_u(s)}{M} \right| \tag{8}$$

D. Average number of active servers N_a

Given the user distribution probability $\{p_n\}$ provided by (5), and the set of activation thresholds $\{s_m\}$, which is computed below, the average number of active servers N_a is given by

$$N_a = 1 + \sum_{m=2}^{N} m \sum_{n=s_m}^{s_{m+1}-1} p_n.$$
 (9)

where the first term accounts for the fact that there is always at least one active server active, and the second term computes the weighted average of a number of servers m and the probability of having that number of servers active (there are m active servers if the number of users is between s_m and $s_{m+1} - 1$).

E. Computation of the activation thresholds

As described in Section III, the activation thresholds $\{s_m\}$ are automatically configured using a control theoretical mechanism [6] that guarantees that the failure probability P_f is below the target T_f while minimizing the energy consumption. In this section, we provide a theoretical analysis to estimate the value of these thresholds.

To perform the analysis, again we neglect the impact of the finite server lifetime on failures, i.e., we assume that all service failures are due users arriving to the system with not enough resources to immediately start service. This assumption reflects the server farm's operation, where the auto-scaling scheme enables re-routing tasks from a failing server to other available resources; our simulations confirm such direct failures have a negligible impact on overall system reliability. We consider a scenario with m active servers, and assume that the activation threshold for an additional server is set to k. Following the above, a failure will occur right after the system has k users (and triggers the activation of a new server) if the total number of users exceeds the current total capacity of the farm $(m \times M)$ before another server has been fully activated, which requires a time $T_{\rm on}$. We denote this conditional probability upon arrival as $P_f(k, T_{on})$, which corresponds to the probability of reaching a state with more than $m \times M$ users in less than $T_{\rm on}$, starting from a state with k users. By denoting with $P_{i,j}(t)$ the probability of reaching state j from i in less than t, the probability $P_f(k, T_{on})$ can be computed as

$$P_f(k, T_{\rm on}) = \sum_{j=mM+1}^{NM} P_{k,j}(T_{\rm on})$$
(10)

We exemplify the above formulation with the toy example depicted in Fig. 1. The figure illustrates a scenario with m = 2 active servers, where each server has a capacity M = 5. Assuming that the activation threshold for the third server is k = 9 and a total of N = 10 servers (i.e., a maximum of 50 simultaneous tasks), the conditional probability of a failure when the system is in state k = 9 is given by

$$P_f(9, T_{\rm on}) = \sum_{j=11}^{50} P_{9,j}(T_{\rm on})$$
(11)

To compute $P_{i,j}(t)$, we assume that the transition matrix of the system **Q** is the same as the one from a classical M/M/c queue [22] with c = N.

Based on this, assuming an initial probability distribution vector $\mathbf{P}(0)$, the probability distribution vector after t is given by [23]

$$\mathbf{P}(t) = \mathbf{P}(0)e^{tQ},\tag{12}$$

and therefore $P_{i,j}(t)$ can be computed by substituting $\mathbf{P}(0)$ with a distribution vector with 1 in the *i*-th position.

To determine the activation thresholds, we assume that the impact of each threshold is independent of the others. We also assume that all tasks that arrive while server m is booting are assumed to have arrived when the number of users in the system was exactly equal to its activation threshold k (rather than during states with fewer users). Under these conditions, Eq. (10) can be used to compute the failure probability associated with server mwhen threshold k is applied. Under these assumptions, the actual failure probability P_f is upper-bounded by the conditional failure probabilities associated with the activation of each server. As a result, if we ensure that each of these conditional probabilities remains below T_f ,



FIGURE 1: Failure probability with activation threshold k = 9 for m = 2 servers with capacity M = 5.

so will be the actual failure probability. Following this reasoning, we estimate the optimal activation thresholds s_m^\ast as

$$s_m^* = \max k \in [(m-2)M + 1, \dots, (m-1)M]$$
 (13a)

s. t.
$$P_f(k, T_{\text{on}}) < T_f$$
 (13b)

where we select the largest activation threshold out of those fulfilling the reliability requirement to minimize the activation of servers and therefore the energy consumption.

F. Optimal design of the server farm

Following the above, the optimization problem can be formalized as follows. Let \mathcal{T} denote the set of all possible server types, $C(\tau)$ denote the cost of using server type τ to support the service, and $P_f(\tau)$ the corresponding failure probability. The optimization problem is to find the optimal server type τ^* defined as follows:

$$\tau^* = \min_{\tau \in \mathcal{T}} \quad C(\tau) \tag{14a}$$

s. t.
$$P_f(\tau) \le T_f$$
 (14b)

Since server types have no relation with each other, nor between the parameters that characterize their performance, there is no alternative to performing an exhaustive search on all server types. We summarize the operation of the mechanism to design the server farm in Algorithm 1, which is described next.

The algorithm begins by computing the s-th percentile of the number of users (line 1). It then iterates over each server type (line 2), computing: the number of required servers (line 3), the activation thresholds (line 4), used to estimate the average number of active servers (line 5), and finally the total cost per server type (line 6). The computational complexity is dominated by the loop over all server types, which requires $|\mathcal{T}|$ iterations. For each server type, the algorithm computes N thresholds, each involving the solution of a system of N equations, resulting in an overall complexity of $O(|\mathcal{T}|N^4)$. Assuming that the s-th percentile grows linearly with the system load λ/μ , the complexity becomes $O(|\mathcal{T}|(\lambda/\mu)^4)$. Despite Algorithm 1 Server Farm Design

Input: Set $\tau_i = \{M_i, \beta_i, \alpha_i, P_{\text{idle},i}, P_{\text{load},i}, K, R\}$, target failure T_f , service rate μ , arrival rate λ

Output: Optimal server type τ^*

- 1: $N_u(s) \leftarrow \text{Using } (7) \quad \triangleright \text{ Compute } s \text{-percentile of } N_u$ 2: for $\tau \in \mathcal{T}$ do
- 3: $N \leftarrow \text{Using } (8) \triangleright \text{Compute number of servers}$ 4: $\{s_m^*\} \leftarrow \text{Using } (13b) \text{ and } (13a) \triangleright \text{Compute}$ activation thresholds
- 5: $N_a \leftarrow \text{Using } (9) \triangleright \text{Compute avg. number of active servers}$
- 6: $C(\tau) \leftarrow \text{Using (2), (3) and (4)} \triangleright \text{Compute total}$ cost
- 7: end for
- 8: $\tau^* = \min_{\tau \in \mathcal{T}} C(\tau)$ s. t. $P_f(\tau) \leq T_f$ \triangleright Select optimal deployment

this theoretical increase with the load, the results below (Fig. 6) show that the computational time remains practically constant across the tested load values, indicating that the algorithm scales efficiently in practice.

V. Performance evaluation

In this section, we first assess the accuracy of the analytical model, and then validate the proposed algorithm to design the server farm. Results from the analytical model are obtained using MATLAB Release 2023a, while simulation results are obtained using a discrete event simulation written in C++. This simulator was also used in our previous works [5], [6], [24]. We perform as many replications as required until the confidence intervals are below 1% of the average (not shown for clarity). Note that the simulation does not relax certain assumptions of the analytical model such as the server infallibility. All computations are performed on a server equipped with an Intel Core i7 CPU with 4 cores, 8 threads, operating at a base frequency of 1.30GHz, and supported by 16 GB of RAM.

	Carrier	Enterprise	Consumer	Rack	Blade
Capacity M (tasks)	16	4	2	12	32
Boot-up $1/\alpha$ (min)	3	8	18	2	1
Lifetime $1/\beta$ (days)	32	16	8	7	4
P_{peak} (W)	270	78	7.6	70	241
P _{load} (W)	7.5	18.22	1.5	3.33	22.84
P_{idle} (W)	150	5.1	4.6	30	20
Cost K (\in)	2000	500	100	1000	300
Lifespan R (years)	10	6	2	8	4

TABLE 2: Deployments considered in the performanceevaluation.

We assume that boot up times and lifetimes are exponential random variables too, with an average that depends on the server type τ . We focus on the following set of target failure probabilities $T_f = \{10^{-3}, 10^{-4}, 10^{-5}\}$, which correspond to a reliability between three nines (99.9%) and five nines (99.99%).

A. Server types

In our experiments, we consider five different server types, ranging from cost-effective consumer-grade machines to high-performance blade servers. These configurations were previously defined in [6], [24], and for completeness, we summarize their performance parameters in Table 2. The selected server types cover a broad range of characteristics, ensuring that our analysis and design are validated under diverse conditions. However, our analysis is not tied to these specific configurations, i.e., our model can be applied to assess the performance of alternative parameterizations and support the design of other server farm architectures. Lastly, we emphasize that this list is not an exhaustive catalog of possible server types.

B. Model validation

We first confirm the validity of the analytical model to estimate the required performance figures to compute the cost for a given configuration τ , namely, the average number of users N_u , the number of servers N, and the average number of active servers N_a , which in turn depends on the activation thresholds $\{s_m\}$. To this aim, we next compare the results obtained using the analytical model with those obtained using simulations. To obtain these, we use the following methodology. For a given value of λ and μ , we initially set the number of servers to $N = \lambda/\mu$ (i.e., the load in Erlangs), and run the simulations using the configuration algorithm in [6] that aims at minimizing resource consumption while ensuring that the failure probability P_f is below the target value T_f . If P_f is above T_f , we increase the number of servers N by one and repeat the process, until P_f is below T_f . Throughout or model validation, we consider three different inter-arrival distributions, each represented with a distinct symbol in the figures:



FIGURE 2: Average number of users vs. λ using analysis (lines) and simulation (symbols).

- Exponential distribution, with rate $\lambda = \{0.2, 0.4, \dots, 4.0\}$ tasks/min.
- Pareto distribution, with shape parameter $\alpha = 2$ and scale parameter $x_m = \{2.5, 1.25, \dots, 0.125\}$ min/tasks.
- Weibull distribution, with shape parameter k = 2 and scale parameter $\theta = \{5.6, 1.12, \dots, 0.28\} \min/tasks.$

1) Average number of users and distribution of users

We start our model validation by comparing the analytical results of Section B with those obtained via simulations. We first compare the average number of users obtained using (6) with those computed using simulations, for all considered server types and the different arrival rates considered.

We present the results in Fig. 2, using points for the simulation results and lines for the analytical values. The figure confirms the good accuracy of the analytical model, since the results practically overlap.

We also compare the distribution of users using (5) with those computed using simulations, for all the server types and three selected values of λ . We depict the probability mass function of the user distribution in Fig. 3, using bins for simulation results and black lines for the analytical values obtained using (5). These results also confirm the accuracy of the distribution of users of the analytical model, as the differences between the theoretical and experimental distributions are very small.



FIGURE 3: User distribution for different λ values: analysis (line) and simulation (bins and symbols).

2) Total number of servers

Here we assess the validity of our analysis to dimension the server farm. To this aim, we compare the total number of servers N required to guarantee $P_f < T_f$ using the methodology described for the simulations with the values obtained via (8). We perform the comparison for the same values of λ as before and all server types using $T_f = 10^{-4}$, and depict the corresponding results in Fig. 4, using lines for the analytical results and points for those obtained using simulations.



FIGURE 4: Maximum number of servers vs. λ using analysis (lines) and simulation (symbols).

Like in the previous case, the results confirm the accuracy of the model, as the results practically overlap for all considered values of λ and server type, with an average error of 4.6% and the maximum error being below 10%. As expected, the number of required servers grows with the inverse of the server capacity M, with the consumer-grade server type requiring the maximum

8

λ (tasks/min)		s_2	s_3	s_4	s_5	s_6	s_7	s_8	s_9
0.2	Sim.	22	34	46	58	70	82	94	106
	Ana.	21	33	45	57	70	82	94	106
	Δs	1	1	1	1	0	0	0	0
0.4	Sim.	21	33	45	57	69	82	94	106
	Ana.	20	32	44	56	68	80	93	105
	Δs	1	1	1	1	1	2	1	1
1.0	Sim.	16	28	40	52	64	76	88	100
	Ana.	16	28	41	53	65	77	90	102
	Δs	0	0	-1	-1	-1	-1	-2	-2
2.0	Sim.	14	26	38	50	62	74	86	98
	Ana.	13	25	37	49	61	73	85	98
	Δs	1	1	1	1	1	1	1	0
3.0	Sim.	13	25	37	49	61	73	85	97
	Ana.	13	25	37	49	61	73	85	97
	Δs	0	0	0	0	0	0	0	0
4.0	Sim.	12	24	36	48	60	72	84	96
	Ana.	12	24	36	48	60	72	83	96
	Δs	0	0	0	0	0	0	1	0

TABLE 3: Activation thresholds for the *rack* server deployment and different values of λ .

number of servers, and the blade type the minimum. We find that the analytical figures are always above the ones obtained using simulations, and therefore the total number of servers according to our design never falls below the required number of servers. Finally, we conducted the same experiment for the rest of server types and T_f values, obtaining an average error of 4.8% and a maximum error below 12%.

3) Activation thresholds

Here we validate the analysis presented in Section E to estimate the values of the activation thresholds $\{s_m\}$ configured by the algorithm. To this aim, we first focus on the configurations using the rack server type and set $T_f = 10^{-3}$ as the maximum failure probability, and compare the first eight values of the vector $\{s_m\}$ for different values of the arrival rate λ as in the previous sections. The simulation (Sim.) and analytical (Ana.) results are presented in Table 3, as well as the difference (Δs) .

The table illustrates that as the load grows, the activation thresholds decrease, since servers need to be activated with more anticipation to accommodate the incoming tasks. It also shows a good match between the results predicted by the model and those obtained using simulations, with a mean absolute difference of 0.6 tasks and a maximum difference of 2 tasks. We repeated the same experiment for the other types of servers and the



FIGURE 5: Average number of active servers vs. λ using analysis (lines) and simulation (symbols).

considered T_f values, with the mean absolute difference being 0.7 tasks and the maximum absolute difference being 2 tasks. These results again confirm the accuracy of the analytical model to estimate the configuration of the server farm.

4) Average number of active servers

Finally, we assess the accuracy of (9) to estimate the average number of active servers, using the same methodology as before. As in the previous cases, we first assume $T_f = 10^{-4}$ and different values of the traffic load. We represent in Fig. 5 with lines the results from the analytical model and with points those using simulations. The figure confirms the accuracy of the model, with an average absolute error of 0.45 servers and a maximum absolute error of 0.62 servers. Like in the case of Fig. 4, the larger the capacity M of the server type, the smaller the average number of active servers.

Based on the above results, we confirm the accuracy of the analytical model to predict the performance of a server farm for different traffic loads, reliability requirements, and server types. We next assess the performance of the algorithm proposed to design the server farm.

C. Design of the server farm

Following the validation of the analytical model, we next assess the performance of the algorithm presented in Section F to design a server farm. To this aim, we assume the same set of traffic rate values λ and target failure probability levels $T_f = \{10^{-3}, 10^{-4}, 10^{-5}\}$ considered before. To provide an adequate context, we compare the results from our algorithm against two benchmarks:

- An exhaustive search in the configuration space.
- A benchmark based on the Erlang-C [25], which determines for each type of server the number of resources required to ensure that the probability of blocking meets a specific reliability target (i.e., $P_B = 1 R$), and selects the minimum.

For each considered scenario, we compute:

- For each of the five server types considered, the minimum cost (C_s) obtained via simulations using a search on the total number of servers.
- The optimal server type τ^* and corresponding cost (C_a) according to our algorithm.
- The difference between the minimum cost obtained with the numerical search and the one obtained with our algorithm (ΔC).
- The cost of the benchmark design $C_{\rm b}$ based on the Erlang-C, and the corresponding difference vs. $C_{\rm a}$, denoted as ΔB .

We provide the resulting figures for the λ and T_f values considered in Table 4. For each scenario, we highlight in gray the minimum cost obtained using simulations, and in bold font the resulting τ^* whenever the best type of server according to our method corresponds to the one that minimizes costs using simulations.

There are several results that can be derived from Table 4. First, both for simulations and analysis, the minimum cost increases as the load increases, since more resources are needed to accommodate the incoming traffic. Second, there is no optimal server configuration for all scenarios, as the best server type alternates between the five considered types depending on the load and reliability considered. Third, we note that our configuration algorithm provides the optimal server type in 16 out of the 18 scenarios considered (i.e., 88.8% of the scenarios), and that for those two scenarios the relative error in terms of cost is smaller than 2%. The average cost difference between the configuration provided by the algorithm and the simulations is 3%, which confirms the effectiveness of our proposal to design a server farm. Fourth, the comparison of the cost between the analysis (C_a) and benchmark (C_b) highlights the advantages of using our proposed method over a predefined benchmark configuration. We note that the absolute cost for the benchmark is notably higher than the one by the analysis for all combinations. The benchmark consistently results in a higher number of servers due to its more pessimistic assumptions. The relative cost difference (ΔB) demonstrates that relying on a fixed server configuration can lead to significant cost inefficiencies, with the benchmark costing on average 22% more than the optimized analysis-based approach. In some cases, such as $\lambda=0.2$ and $T_f = 10^{-4}$, the cost increase reaches 54%, reinforcing the importance of dynamically selecting the optimal server type rather than adhering to a static deployment

		Simulation Cost (C_s) Analysis						Benchmark			
λ	T_f	Carrier	Enterprise	Consumer	Rack	Blade	τ^*	C_{a}	ΔC	$C_{\rm b}$	ΔB
0.2	10^{-3}	495.98	256.26	665.26	855.40	260.03	Blade	260.62	1.70%	388.35	49.01%
	10^{-4}	498.17	268.04	668.52	858.25	261.73	Enterprise	265.13	1.30%	408.30	54.00%
	10^{-5}	502.53	275.78	671.83	861.07	265.21	Blade	269.69	1.66%	474.01	75.76%
0.4 1	10^{-3}	824.32	712.03	1006.47	1569.52	706.90	Blade	716.08	1.28%	776.71	8.47%
	10^{-4}	828.42	763.08	1110.23	1573.58	710.14	Enterprise	766.43	0.44%	840.32	9.64%
	10^{-5}	835.07	770.93	1125.36	1580.05	770.55	Blade	787.81	2.24%	1084.30	37.63%
1.0	10^{-3}	1634.76	1557.30	1360.30	1860.12	1631.93	Consumer	1365.19	0.36%	1516.15	11.06%
	10^{-4}	1840.56	1729.00	1747.32	1865.57	1731.43	Enterprise	1747.15	1.05%	1860.23	6.47%
	10^{-5}	2055.09	2028.34	2022.25	2056.09	2031.67	Consumer	2189.73	6.50%	2430.62	11.00%
2.0	10^{-3}	2995.68	3131.26	2804.93	2812.14	3121.38	Consumer	3008.56	7.26%	3118.77	3.66%
	10^{-4}	3010.08	3150.12	2899.54	2892.20	3135.09	Rack	3051.56	5.51%	3343.76	9.58%
	10^{-5}	3140.14	3175.16	3126.13	3198.45	3150.17	Consumer	3219.60	2.99%	3684.36	14.44%
3.0	10^{-3}	3861.44	3891.09	3917.21	3901.75	4611.36	Carrier	4100.84	6.20%	4888.86	19.22%
	10^{-4}	3975.12	4005.11	3990.17	3947.93	4625.20	Rack	4148.09	5.07%	5120.98	23.45%
	10^{-5}	4210.21	4188.63	4147.31	4167.34	4660.33	Consumer	4333.52	4.49%	5876.56	35.61%
4.0	10^{-3}	5321.68	5561.98	5432.71	5412.82	5871.76	Carrier	5587.76	5.01%	6006.73	7.50%
	10^{-4}	5640.19	5580.29	5559.92	5512.15	5890.42	Rack	5643.89	2.39%	6345.41	12.43%
	10^{-5}	5980.31	5821.87	5820.19	5847.63	5930.54	Consumer	5899.34	1.36%	6640.85	12.57%

TABLE 4: Minimum cost desgin using simulations (C_s) , our analysis (C_a) , and the benchmark (C_b) .

strategy. Finally, the table also highlights the importance of an adequate selection of the best server type, in addition to its optimal configuration, since there are substantial differences in terms of cost between optimal server deployments with different types. For instance, for $\lambda = 0.2$ and $T_f = 10^{-3}$ (first row of the table), there is a factor of $3.33 \times$ between the minimum cost using the Enterprise type of server and the one using the Rack type of server, while the average difference is $1.64 \times$.

D. Computational time

Finally, we compare the time to determine the optimal design using the method presented in Section IV with a exhaustive search using simulations. To this aim, we compute the total execution time required for each approach. We assume the same scenarios as before with different values of λ and the reliability levels $T_f = \{10^{-3}, 10^{-4}, 10^{-5}\}$. We illustrate the results in Fig. 6, using a logarithmic scale on the y axis.

According to the figure, the proposed method results in significantly shorter execution times for all values of λ and T_f . Furthermore, these times are practically constant, while simulation times increase with the traffic load, and with the inverse of T_f . Based on these results, we conclude that the analysis developed in Section IV offers a cost-effective solution, particularly well suited for scenarios demanding robust, scalable, and resourceefficient methodologies.



FIGURE 6: Execution time vs. lambda required by simulation (dashed lines) and analysis (continuous lines) for different values of T_f .

E. Heterogeneous scenarios

As discussed in Section III, we initially consider homogeneous scenarios, where a single type of service is provided by a single type of server. In this section, we relax this assumption to demonstrate how the proposed framework can be extended to design heterogeneous scenarios. Specifically, we consider a case with two types of services, labeled 1 and 2. Both have the same service rate μ , but differ in their reliability requirements, T_1 and T_2 , as well as their arrival rates, λ_1 and λ_2 , respectively.

One way to design a heterogeneous server farm using our framework is to independently determine the optimal server type for each type of traffic. We denote these as τ_i^* for $i \in 1, 2$. Assuming that the Infrastructure Manager (IM) redirects each type of task to the corresponding server type, the total cost of this heterogeneous design is given by:

$$C_{\rm het} = C(\tau_1^*) + C(\tau_2^*) \tag{15}$$

As a benchmark, we assume a homogeneous design to support the total traffic $\lambda_1 + \lambda_2$ and the most stringent reliability requirement, $\min(T_1, T_2)$. The resulting cost of this design is denoted as C_{hom} . Table 5 presents the resulting values of C_{het} and C_{hom} for different values of $\lambda_1, \lambda_2, T_1$, and T_2 . It also reports the relative difference between the homogeneous and the heterogeneous design, ΔB . Note that for some scenarios, the cost C_{hom} is reused from Table 4.

λ_1	T_1	λ_2	T_2	$C_{\rm het}$	$C_{ m hom}$	ΔB
0.3	10^{-3}	0.7	10^{-4}	1831.92	1729.00	-5.95%
0.3	10^{-4}	0.7	10^{-3}	1714.34	1729.00	0.85%
0.3	10^{-4}	0.7	10^{-5}	2032.31	2022.25	-0.50%
1.5	10^{-3}	1.5	10^{-4}	3892.55	3947.93	1.40%
1.5	10^{-4}	1.5	10^{-3}	3892.55	3947.93	1.40%
1.5	10^{-4}	1.5	10^{-5}	4422.17	4147.31	-6.63%
1.5	10^{-3}	2.5	10^{-4}	5218.46	5512.15	5.33%
1.5	10^{-4}	2.5	10^{-3}	5159.25	5512.15	6.40%
1.5	10^{-4}	2.5	10^{-5}	5509.08	5820.19	5.35%

TABLE 5: Heterogeneous scenarios.

The results confirm that our proposal can be extended to heterogeneous scenarios, as most configurations yield additional cost savings –up to 6.4% in some cases. However, these gains remain moderate due to two main factors. First, traffic is isolated across server types, which prevents potential multiplexing gains. Second, our homogeneous design already performs well, leaving limited room for further improvement except in specific scenarios. As discussed in Section VII, our ongoing work focuses on developing novel analytical models to more effectively address the design of heterogeneous server farms.

VI. Related work

A. Resource Allocation for NFV

The research community has shown significant interest in dynamically managing cloud resources to minimize consumption. For example, in [26], the authors analyze the impact of various static algorithms for activating and deactivating resources, as well as reallocating tasks within a data center, with a focus on reducing energy consumption and minimizing service violations. In a subsequent study [27], they suggest adapting thresholds based on estimated conditions. Finally, control theory has historically been leveraged for energy-efficient resource allocation in cloud computing systems, as outlined in [28]. For example, [29] applies control theory principles for load balancing and CPU frequency selection, which differ from our previous work [6] where control theory is applied to drive the system to the desired reliability levels while minimizing energy consumption. However, the techniques proposed therein tackle distinct challenges compared to those addressed in our paper, and none specifically account for both the waiting queue time of the tasks and the fact that servers has non-zero boot up times.

B. Analysis of reliability in NFV

The study by [30] meticulously examines the reliability of a carrier-grade server system, employing a fault tree model at a high level that intricately links various lowerlevel Markov models. These models account for the inherent fallibility of hardware components such as CPUs and memory modules. A similar methodology is pursued in [31], where the authors examine the reliability of both virtualized and non-virtualized systems comprising two hosts. Furthermore, [32] delves into a related system, conducting a sensitivity analysis to pinpoint parameters that significantly impact reliability. A closer examination akin to our research is presented in [33], where a Markov chain is used to model a server farm, factoring in setup delays concerning response time and power consumption. Similarly, [3] explores the analysis within the realm of 5G/6G networks, using thresholds to manage instance power and performance evaluation in terms of power consumption and waiting time. However, none of these works have proposed a theoretical model to asses the design of an optimal server farm, targeting a desired level of reliability while minimizing infrastructure costs.

C. NFV and Reliability

In prior works [5], [6], [24], we have addressed a system similar to the one studied in this paper. In [5], we characterized service reliability and derived an optimal configuration of the server farm to support a required reliability while minimizing the resource consumption. In [24], we studied the trade-offs of a server farms in terms of reliability and power consumption based on a static configuration. Our analyses of [5], [24] are static and require knowledge about the system load, while other proposal rely on stochastic optimization [11] to find the best trade-off between resource consumption and average waiting time. In [6], we introduced a control theory algorithm that dynamically adapted the configuration to reach the desired point of operation. However, all these papers assumed a fixed server farm where machines are characterized in terms of a number of parameters (e.g. energy consumption, capacity, lifetime). In contrast to this analysis and configuration problems, in this paper we address the design problem: given a given a set of parameters defining a service and a list of candidate server types, that could be used to deploy the server farm, select the most adequate server type to support the service.

VII. Conclusions and future work

Designing auto-scaling server farms requires balancing reliability and cost, complicated by activation delays, finite lifespans, and failure risks. Traditional methods often ignore these factors or over-provision, raising costs. In this paper, we present a framework that combines queueing theory, cost modeling, and server selection to identify minimum-cost deployments. Our approach considers boot-up times and failure probabilities, ensuring reliable service at minimal cost. The proposed methodology is particularly well suited for latency- and reliabilitysensitive applications, such as industrial automation, autonomous vehicle coordination, and mission-critical IoT deployments in 5G and beyond networks.

As part of our ongoing and future work, we are exploring several key directions to enhance the applicability and realism of our framework. First, we aim to improve the design of heterogeneous scenarios. While the current approach assigns a dedicated and independently optimized server type to each service class, we plan to develop more sophisticated strategies that enable resource sharing and exploit multiplexing gains across services. Second, we intend to incorporate network latency into the optimization process by accounting for the performance characteristics of the RAN [34], the underlying network topology, and the VNF chaining architectures. This will involve probabilistic modeling of end-to-end latency to better support URLLC requirements. Finally, we are extending our implementation over Linux-based architectures [17] integrating a load balancing mechanism with ACHO's [18] seamless migration capabilities, enabling more efficient and resilient deployment of virtualized services.

ACKNOWLEDGMENT

This publication is part of the project **6GINSPIRE** PID2022-137329OB-C42. funded bv MCIN/AEI/10.13039/501100011033/. This work has been partly funded by NEC Laboratories Europe Student Research Fellowship program of 2021. This work is also partially supported by the Spanish Ministry of Economic Affairs and Digital Transformation and

the European Union-NextGenerationEU through the UNICO 5G I+D SORUS project. The work of Jaime Garcia-Reinoso has been partly funded by the Spanish Ministry for Science and Innovation through the ADMINISTER (TED2021-131301B-I00) project.

REFERENCES

- P. Rost, C. Mannweiler, D. S. Michalopoulos, C. Sartori, V. Sciancalepore, N. Sastry, O. Holland, S. Tayade, B. Han, D. Bega, D. Aziz, and H. Bakker, "Network slicing to enable scalability and flexibility in 5g mobile networks," <u>IEEE</u> Communications Magazine, vol. 55, no. 5, pp. 72–79, 2017.
- [2] M. Gramaglia, P. Serrano, A. Banchs, G. Garcia-Aviles, A. Garcia-Saavedra, and R. Perez, "The case for serverless mobile networking," in <u>2020 IFIP Networking Conference</u> (Networking), 2020, pp. 779–784.
- [3] Y. Ren, T. Phung-Duc, J. Chen, and Z. Yu, "Dynamic Auto Scaling Algorithm (DASA) for 5G Mobile Networks," in Proceedings of the IEEE Global Communications Conference (GLOBECOM 2016), Washington DC, USA, Dec. 2016.
- [4] Č. Stefanović, "Industry 4.0 from 5g perspective: Use-cases, requirements, challenges and approaches," in <u>2018 11th CMI</u> <u>International Conference: Prospects and Challenges Towards</u> <u>Developing a Digital Economy within the EU</u>. IEEE, 2018, pp. 44–48.
- [5] J. Ortin, P. Serrano, J. Garcia-Reinoso, and A. Banchs, "Analysis of scaling policies for nfv providing 5g/6g reliability levels with fallible servers," IEEE Transactions on Network and Service Management, vol. 19, no. 2, pp. 1287–1305, 2022.
- [6] J. Perez-Valero, A. Banchs, P. Serrano, J. Ortín, J. Garcia-Reinoso, and X. Costa-Pérez, "Energy-aware adaptive scaling of server farms for nfv with reliability requirements," <u>IEEE Transactions on Mobile Computing</u>, vol. 23, no. 5, pp. 4273– 4284, 2024.
- [7] J. Liu, Z. Jiang, N. Kato, O. Akashi, and A. Takahara, "Reliability evaluation for nfv deployment of future mobile broadband networks," <u>IEEE Wireless Communications</u>, vol. 23, no. 3, pp. 90–96, 2016.
- [8] L. Qu, C. Assi, K. Shaban, and M. J. Khabbaz, "A reliabilityaware network service chain provisioning with delay guarantees in nfv-enabled enterprise datacenter networks," IEEE <u>Transactions on Network and Service Management</u>, vol. 14, no. 3, pp. 554–568, 2017.
- [9] Y. Ren, T. Phung-Duc, J.-C. Chen, and Z.-W. Yu, "Dynamic auto scaling algorithm (dasa) for 5g mobile networks," in 2016 <u>IEEE Global Communications Conference (GLOBECOM)</u>, 2016, pp. 1–6.
- [10] A. Gandhi, M. Harchol-Balter, and I. Adan, "Server farms with setup costs," <u>Performance Evaluation</u>, vol. 67, no. 11, pp. 1123–1138, 2010, performance 2010. [Online]. Available: https://www.sciencedirect.com/science/ article/pii/S0166531610000957
- [11] A. Song, W. Wang, and J. Luo, "Stochastic modeling of dynamic power management policies in server farms with setup times and server failures," <u>International Journal of Communication Systems</u>, vol. 27, no. 4, pp. 680–703, 2014. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/ 10.1002/dac.2761
- [12] J. Burge, P. Ranganathan, and J. L. Wiener, "Cost-aware scheduling for heterogeneous enterprise machines (cash'em)," in <u>2007 IEEE International Conference on Cluster Computing</u>, 2007, pp. 481–487.
- [13] K. Goseva-Popstojanova and K. Trivedi, "Failure correlation in software reliability models," in <u>Proceedings 10th</u> <u>International Symposium on Software Reliability Engineering</u> (Cat. No.PR00443), 1999, pp. 232–241.
- [14] A. Vasan, A. Sivasubramaniam, V. Shimpi, T. Sivabalan, and R. Subbiah, "Worth their watts? - an empirical study of datacenter servers," in <u>HPCA - 16 2010 The Sixteenth</u> <u>International Symposium on High-Performance Computer</u> Architecture, 2010, pp. 1–10.

- [15] G. Chen, W. He, J. Liu, S. Nath, L. Rigas, L. Xiao, and F. Zhao, "Energy-aware server provisioning and load dispatching for connection-intensive internet services," in <u>Proceedings</u> of the 5th USENIX Symposium on Networked Systems Design and Implementation, ser. NSDI'08. USA: USENIX Association, 2008, p. 337–350.
- [16] 5G Americas, ""Vehicular Connectivity: C-V2X and 5G," Technical White Paper, 2021.
- [17] S. K. Singh, <u>Linux Yourself: Concept and Programming</u>, 1st ed. Chapman and Hall/CRC, 2021. [Online]. Available: https://doi.org/10.1201/9780429446047
- [18] G. Garcia-Aviles, C. Donato, M. Gramaglia, P. Serrano, and A. Banchs, "Acho: A framework for flexible re-orchestration of virtual network functions," <u>Computer Networks</u>, vol. 180, p. 107382, 2020.
- [19] J. A. Aroca, A. Chatzipapas, A. F. Anta, and V. Mancuso, "A measurement-based characterization of the energy consumption in data center servers," <u>IEEE Journal on Selected Areas</u> in Communications, vol. 33, no. 12, pp. 2863–2877, 2015.
- [20] M. Ananth and R. Sharma, "Cloud management using network function virtualization to reduce capex and opex," in 2016 8th International Conference on Computational Intelligence and Communication Networks (CICN). IEEE, 2016, pp. 43–47.
- [21] L. Kleinrock, Theory, Volume 1, Queueing Systems. USA: Wiley-Interscience, 1975.
- [22] G. Bolch, S. Greiner, H. De Meer, and K. S. Trivedi, <u>Queueing</u> networks and <u>Markov</u> chains: modeling and performance evaluation with computer science applications. John Wiley & Sons, 2006.
- [23] K. Rupp, R. Schill, J. Süskind, P. Georg, M. Klever, A. Lösch, L. Grasedyck, T. Wettig, and R. Spang, "Differentiated uniformization: A new method for inferring markov chains on combinatorial state spaces including stochastic epidemic models," <u>Computational Statistics</u>, pp. 1–21, 2024.
- [24] J. Perez-Valero, J. Garcia-Reinoso, A. Banchs, P. Serrano, J. Ortin, and X. Costa-Perez, "Performance trade-offs of auto scaling schemes for nfv with reliability requirements," Computer Communications, vol. 212, pp. 251–261, 2023.
- [25] T. R. Robbins, D. J. Medeiros, and T. P. Harrison, "Does the erlang c model fit in real call centers?" in Proceedings of the 2010 Winter Simulation Conference. IEEE, 2010, pp. 2853– 2864.
- [26] A. Beloglazov and R. Buyya, "Energy efficient allocation of virtual machines in cloud data centers," in <u>2010 10th</u> <u>IEEE/ACM International Conference on Cluster, Cloud and</u> <u>Grid Computing</u>, 2010, pp. 577–578.
- [27] A. Beloglazov and R. Buyya, "Adaptive threshold-based approach for energy-efficient consolidation of virtual machines in cloud data centers," in Proceedings of the 8th International Workshop on Middleware for Grids, Clouds and e-Science, ser. MGC '10. New York, NY, USA: Association for Computing Machinery, 2010.
- [28] A. Hameed, A. Khoshkbarforoushha, R. Ranjan, P. P. Jayaraman, J. Kolodziej, P. Balaji, S. Zeadally, Q. M. Malluhi, N. Tziritas, A. Vishnu, S. U. Khan, and A. Zomaya, "A survey and taxonomy on energy efficient resource allocation techniques for cloud computing systems," <u>Computing</u>, vol. 98, no. 7, pp. 751–774, Jun. 2014.
- [29] Y. Wang, X. Wang, M. Chen, and X. Zhu, "Power-efficient response time guarantees for virtualized enterprise servers," in <u>2008 Real-Time Systems Symposium</u>, 2008, pp. 303–312.
- [30] W. E. Smith, K. S. Trivedi, L. A. Tomek, and J. Ackaret, "Availability analysis of blade server systems," <u>IBM Systems</u> Journal, vol. 47, no. 4, pp. 621–640, 2008.
- [31] D. S. Kim, F. Machida, and K. S. Trivedi, "Availability modeling and analysis of a virtualized system," in 2009 15th IEEE Pacific Rim International Symposium on Dependable Computing, 2009, pp. 365–371.
- [32] R. d. S. Matos, P. R. M. Maciel, F. Machida, D. S. Kim, and K. S. Trivedi, "Sensitivity analysis of server virtualized system availability," <u>IEEE Transactions on Reliability</u>, vol. 61, no. 4, pp. 994–1006, 2012.

- [33] A. Gandhi, M. Harchol-Balter, and I. Adan, "Server farms with setup costs," <u>Performance Evaluation</u>, vol. 67, no. 11, pp. 1123 – 1138, Nov. 2010.
- [34] O. Adamuz-Hinojosa, L. Zanzi, V. Sciancalepore, A. Garcia-Saavedra, and X. Costa-Pérez, "Oranus: Latency-tailored orchestration via stochastic network calculus in 6g o-ran," in <u>IEEE INFOCOM 2024</u> - <u>IEEE Conference on Computer</u> <u>Communications</u>, 2024, pp. 61–70.



Jesus Perez-Valero received the Ph.D. degree from the Universidad Carlos III de Madrid (UC3M) in 2024. He is currently a Postdoctoral Researcher and lecturer with Universidad de Murcia (UMU). He has been involved in several projects funded by the European Commission through the SNS. His main research interests lie in the performance analysis and optimization of communication systems.



Pablo Serrano(M'09, SM'16) is an Associate Professor at the University Carlos III de Madrid. His research interests lie in the analysis of wireless networks and the design of network protocols and systems. He currently serves as Editor for IEEE Open Journal of the Communication Society and IEEE Transactions on Mobile Computing.



Jaime Garcia-Reinoso(M'04) received the Telecommunications Engineering degree in 2000 from the University of Vigo, Spain and the Ph.D. in Telecommunications in 2003 from the University Carlos III of Madrid, Spain. He is an associate professor at University of Alcala, Spain since 2021 and he has published over 60 papers in top magazines and conferences. He has been involved in many international projects on next generation networks, 5G, SDN and NFV.





Albert Banchs has a double affiliation as Professor at the University Carlos III of Madrid and Deputy Director of the IMDEA Networks institute. Prof. Banchs has served in many TPCs and has also served in the editorial board of a number of journals, including IEEE Transactions in Wireless Communications and IEEE/ACM Transactions on Networking. Dr. Banchs has participated in many European projects and industry contracts.

Xavier Costa-Perez is ICREA Research Professor, Scientific Director at the i2cat Research Center and Head of 5G Networks R&D at NEC Laboratories Europe. His research focuses on the transformation of society driven by the interplay of mobile networks and AI. Currently, he is serving as Associate Editor at IEEE Transactions on Mobile Computing, IEEE Transactions on Communications and Elsevier Computer Communications journals.