

Energy-Efficient Deployment of Stateful FaaS Vertical Applications on Edge Data Networks

Claudio Cicconetti
IIT, CNR – Pisa, Italy
c.cicconetti@iit.cnr.it

Raffaele Bruno
IIT, CNR – Pisa, Italy
r.bruno@iit.cnr.it

Andrea Passarella
IIT, CNR – Pisa, Italy
a.passarella@iit.cnr.it

Abstract—5G and beyond support the deployment of vertical applications, which is particularly appealing in combination with network slicing and edge computing to create a logically isolated environment for executing customer services. Even if serverless computing has gained significant interest as a cloud-native technology its adoption at the edge is lagging, especially because of the need to support stateful tasks, which are commonplace in, e.g., cognitive services, but not fully amenable to being deployed on limited and decentralized computing infrastructures. In this work, we study the emerging paradigm of stateful Function as a Service (FaaS) with lightweight task abstractions in WebAssembly. Specifically, we assess the implications of deploying inter-dependent tasks with an internal state on edge computing resources using a stateless vs. stateful approach and then derive a mathematical model to estimate the energy consumption of a workload with given characteristics, considering the power used for both processing and communication. The model is used in extensive simulations to determine the impact of key factors and assess the energy trade-offs of stateless vs. stateful.

Index Terms—Stateful FaaS, Local Data Networks, Beyond 5G, Vertical Applications, Serverless Computing

I. INTRODUCTION

Since long edge computing has outgrown the role of being a surrogate of the cloud for offloading specific services it was originally assigned (e.g., [1]). Nowadays, it is a thriving reality with a transformative effect on business and society, whose adoption is constantly spreading across more and more domains [2], especially thanks to the opportunities offered by cognitive Artificial Intelligence (AI) services at the edge, which can benefit immensely from the data being consumed close to where they are generated and from the use of decentralized computing resources with embedded Graphics Processing Units (GPUs)/Tensor Processing Units (TPUs) [3].

Indeed, edge computing is expected to remain a key player in the future evolution of communication technologies [4] and has attracted the interest of many standardization bodies, including 3GPP. In particular, a reference architecture for supporting edge computing within a cellular network is standardized in [5], enabling Edge Computing Service Providers (ECSPs) to deploy EDNs, namely 5G network infrastructures that contain edge servers (called EAS in 3GPP jargon) and management entities, which can offer edge computing services to end users. A typical deployment scenario for an EDN is shown Figure 1, together with the interactions between the 5G transport functions and the edge computing entities. At

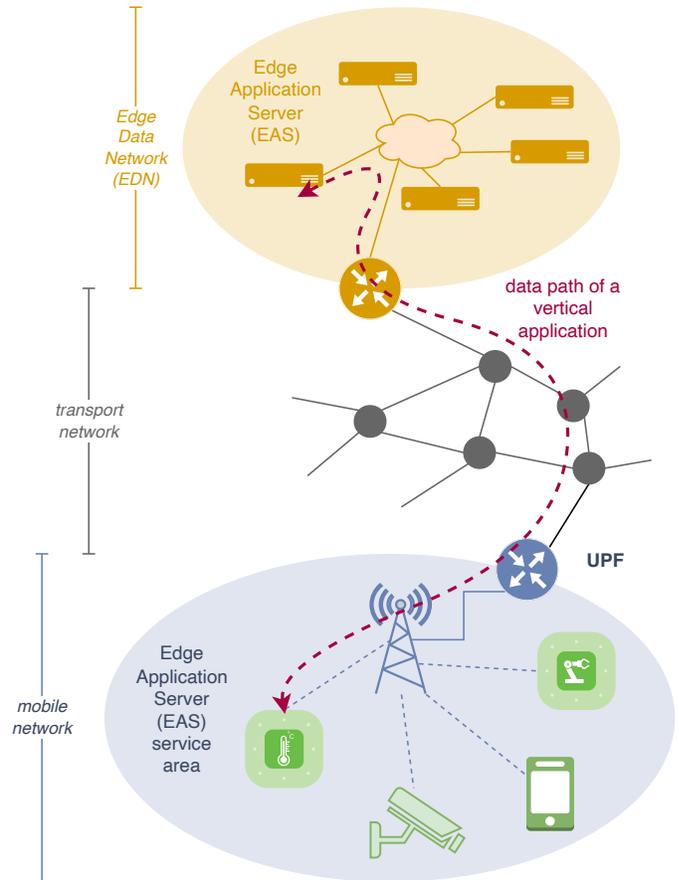


Fig. 1. Simplified beyond 5G architecture of an Edge Data Network (EDN) mobile users with Edge Application Server (EAS) services.

the bottom layer of this reference architecture, there is the mobile network with Internet of Things (IoT) users, possibly consisting of a network slice [6], that is a virtual network dedicated to a customer in logical isolation, and providing specific Service Level Agreements (SLAs), while sharing the physical Radio Access Network (RAN) and 5G Core (5GC) resources [7]. Virtual applications can be assigned by the customer to a given EAS service area of the mobile network, whose traffic is handled by one User Plane Function (UPF) (there can be more than one, not shown in the figure) towards a transport network until it reaches the EDN, at the top of the

diagram. The EDN is a pool of computing, networking, and storage resources contributing to the EAS.

From the point of view of software architectures, cloud-native approaches based on micro-services and containers, instead of monolithic applications, have become very popular, especially relying on serverless computing. The latter enables full automation and flexible billing [8], building on a functional programming model, called Function as a Service (FaaS), which allows a developer to compose elementary and reusable *functions*, whose instances are stateless and can be shared among different users. Serverless computing has been adopted successfully in many application domains in the cloud [9], but it is not yet fully developed at the edge [10], because of two main reasons. First, real applications often require functions to maintain a state, of some kind, associated with a given user or session. In the cloud, this is implemented through storage services, which are readily available in data centers, but incur much more overhead in a decentralized edge computing environment. Second, efficient scheduling and resource management are much more challenging at the edge due to the scarcity of resources and the resulting contention between concurrent applications, which hampers the promises of full automation and infinite scalability of serverless computing.

These observations lead to the following research questions:

- Q1** *Deploying stateful FaaS through stateless runners accessing the state on an external service is the state-of-the-art deployment option. Is there an alternative solution to obtain the benefits of edge computing and the flexibility of FaaS, despite the need to maintain a state?*
- Q2** *Improving the efficiency of Information and Communication Technologies (ICT) is one of the key goals for global sustainable development [12]. What are the implications of different deployment/run-time options to support stateful FaaS on the energy consumption of the EDN infrastructure?*

We address both questions through the paper contributions:

- A1** We consider an alternative paradigm that realizes *stateful FaaS* (also known as “actor model” in the literature, e.g., [11]) exploiting recent lightweight virtualization technologies and allowing efficient multiplexing of many instances of isolated user-space applications. This is explored in Section II below, where we provide a background on the deployment of inter-dependent stateful tasks through stateless vs. stateful FaaS and report the results obtained in a testbed with embedded edge nodes on the execution of concurrent tasks in WebAssembly.
- A2** We formulate in Section III a novel mathematical model for the energy consumption of EDNs using the two stateless vs. stateful FaaS approaches above, which is then used to carry out a performance evaluation study under realistic trace-inspired workload conditions and analyze the impact on energy consumption of the main system factors (Section IV).

The paper continues with a review of the relevant state-

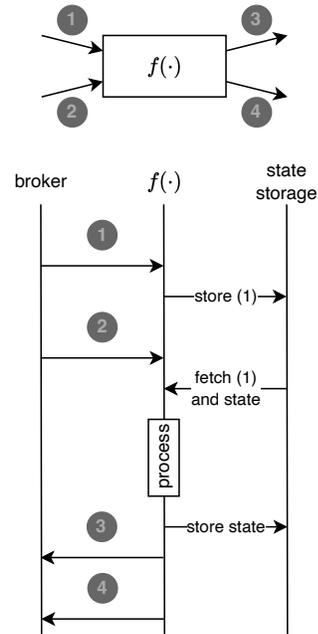


Fig. 2. Example of how to realize stateful processing with stateless FaaS.

of-the-art in Section V, and then comes to the summary in Section VI.

II. STATELESS VS. STATEFUL FAAS: BACKGROUND AND MOTIVATION

Serverless computing and the FaaS programming model are popular in the cloud [9] and they have attracted significant interest also at the edge [13]. As already mentioned in Section I, with FaaS an application is made of a sequence of stateless function calls, which can be arranged in chains (i.e., $f_1 \rightarrow f_2 \rightarrow \dots \rightarrow f_N$) or more complex structures, like Directed Acyclic Graphs (DAGs) [14]. However, realistic applications typically *do need* function execution to be associated with some state, especially for edge applications, such as AI and real-time analytics [15].

A straightforward solution to this problem, which we call **stateless FaaS**, is to maintain the state on an external storage system to be accessed on demand by the functions as part of their execution, as explained, e.g., in [16]. Such a deployment option is illustrated in the example in Figure 2, where function $f(\cdot)$ requires input from two dependencies (1 and 2) and has two outputs (3 and 4). When the function receives input 1, it is kept temporarily in the state storage. Once input 2 is received, full processing can occur combining the latter with the previous input 1 and the state, to produce the final outputs 3 and 4, after updating the state on the storage.

In common serverless computing platforms, function invocation happens through an HTTP command issued on a web server running in a container. Due to the lack of state, the same container can serve multiple users/sessions seamlessly, and the orchestration platform can easily perform autoscaling of such runners, i.e., decreasing or increasing the number of instances per function to match the instantaneous demand.

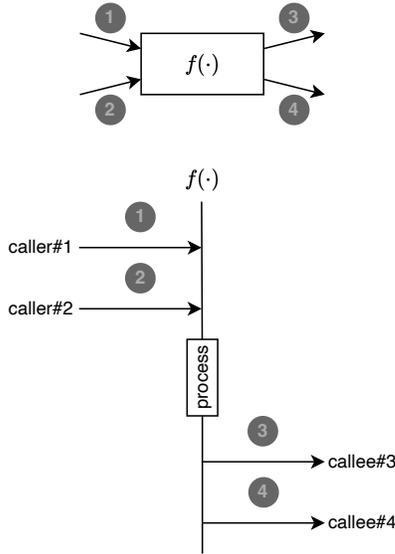


Fig. 3. Example of stateful FaaS.

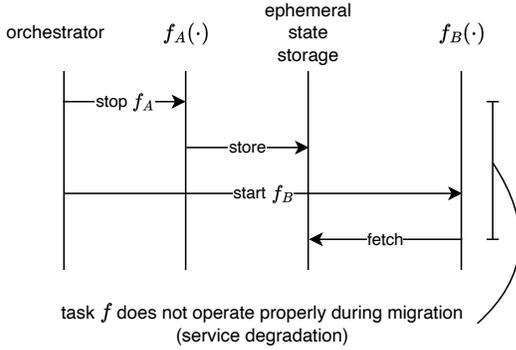


Fig. 4. Migration of a stateful FaaS runner from node A to node B .

An alternative to this strategy is dedicating each user/session to a runner, thus realizing what we call **stateful FaaS**. As illustrated in the example in Figure 3, with this model there is no need to fetch/update the state or store temporary input from previous function calls. In principle, the stateful FaaS model has two inconveniences. First, the number of runners may be much higher than that with stateless FaaS, because the former cannot exploit statistical multiplexing of multiple users/sessions like the latter. We study the practical impact of this issue with a specific virtualization technology in Section II-A. Second, if a runner is migrated from one node to another for any reason, e.g., system resource optimization, its internal state must be moved to the target host. We show an example in Figure 4, where the orchestrator migrates a runner for the function $f(\cdot)$ from node A to node B . First, when stopping $f(\cdot)$ on node A the state is stored temporarily on an external system, which is then queried by the new instance of function $f(\cdot)$ on node B upon creation. With this solution, there would be a period during which the task performed by $f(\cdot)$ is not available. More sophisticated protocols can be

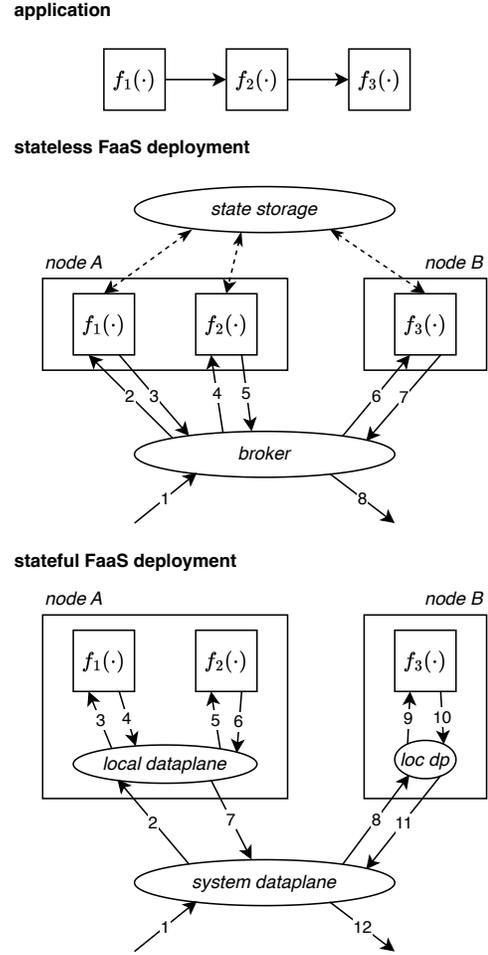


Fig. 5. Deployment of a three-function chain (top) on two processing nodes through stateless FaaS (middle) and stateful FaaS (bottom).

devised [17], but, in any case, they would incur additional complexity or overhead, which is not needed with stateless FaaS. The impact of state migration on energy consumption is captured by the mathematical model defined in Section III and evaluated in Section IV.

We now illustrate deployment with stateless vs. stateful FaaS with the help of the example in Figure 5, with a three-function chain application running on two nodes A and B . In the example, we have one runner per function: node A hosts functions f_1 and f_2 and node B hosts function f_3 . With stateless FaaS, an intermediate layer is needed to dispatch function invocations to one of the matching runners: this is represented by a logical component called *broker*, borrowing the terminology from [18], which is an early study on the realization of distributed computing in pervasive systems. As can be seen, network traffic is generated at each function call for state access, on the state storage, and for invoking the next runner through the broker. On the other hand, with stateful FaaS, we need logical components to mesh together the runners, which can be within a node or at a system level.

Network access for accessing the state is unnecessary because the state is embedded within the runner. Furthermore, when a runner invokes another on the same node no network access is needed, too.

A. Stateful FaaS benchmarking

We now discuss the potential issue of stateful FaaS requiring one runner per user/session, which can be substantially mitigated with lightweight abstractions. A candidate technology for this is WebAssembly, which is new but whose adoption is increasing very fast. With WebAssembly the runner is a user-space application, whose isolation is provided by the runtime environment at a small fraction of the cost of traditional virtualization means, such as Docker containers [19].

In the remainder of this section, we report experimental results obtained by running an application with four stateful functions in a chain, as illustrated in Figure 6. The implementation was done with EDGELESS¹, an open-source framework to develop and orchestrate stateful FaaS applications in the edge-cloud continuum, developed within a project funded by the European Union with the same name. The state of each function is a square matrix of real values and the processing consists of multiplying it by the input argument and providing the result as output. The size N of the input/output vector, and hence the matrix, is 500 and all the states and inputs are initialized with random values in $[0, 1]$. The chain is triggered by a client every time it receives the final output of the chain execution on the previous round, thus each app generates a rate of invocation as high as allowed by the underlying platform/hardware. We ran the experiments on two types of embedded devices: a Raspberry Pi 3 and an NVIDIA Jetson NX. Each experiment lasted 10 seconds after an initial warm-up period (of variable size) and was repeated 10 times. The error bars report low (0.025) and high (0.975) quantiles of measured data.

In Figure 7 we show the invocation rate with an increasing number of apps. The NVIDIA Jetson NX curve lies above the Raspberry Pi 3 because it has a more powerful CPU. With both hardware, the curve increases from 1 app to 2 apps, because the former is not sufficient to saturate the CPUs on the host node. However, the invocation rate afterward remains stable, even with 30 apps, corresponding to 120 concurrent runners (see Figure 6) in the same embedded device, with no noticeable drop in aggregated performance. This confirms that stateful FaaS runners with WebAssembly can scale with negligible overhead until high load levels, relative to the node capabilities. This assumption is used in the next section.

III. SYSTEM MODEL

We now define a mathematical model to estimate the energy consumed in a time horizon T for executing the applications that enter/leave the system during that period. The model is intended to be used to evaluate high-level deployment strategies and run-time orchestration policies and, as such,

it is not intended to provide quantitatively accurate results, but rather qualitative guidelines to drive algorithm design and high-level resource provisioning.

We assume the workload is made of applications (apps for short) that enter and leave the system dynamically at given times t_a^\downarrow and t_a^\uparrow , for app a . An app a consists of some functions (or tasks) arranged in a directed dependency graph $G_a(V_a, E_a)$. Each vertex $v \in V_a$ is a task that depends on its predecessors (incoming edges) and produces output towards its successors (outgoing edges). The amount of data exchanged when task u calls its successor task v is d_{uv} , in bits. Without loss of generality, to have a more compact notation, we assume that the invocation rate is common for all the tasks within app a and equal to λ_a . Task v has an internal state of size s_{av} , in bits, and a processing request equal to r_{av} , in fractions of CPU. An example of a dependency graph is illustrated in Figure 8.

In the following, we consider the system as dynamic, characterized by a series of discrete events happening at time $t_k \in \{t_1, \dots, t_N\}$, where t_N is the end of the period of interest and the other events correspond to an application entering or leaving the system. Between two consecutive events the power consumption remains stable (in a statistical sense) and we can characterize its average value through two step-wise functions, which are constant from time t_k until the next event t_{k+1} : $\alpha(t_k)$ is number of edge nodes used at time t_k to serve the active applications, where each node has a processing capacity C , in fractions of CPU; $\beta_a(t_k)$ is the average network traffic consumed by application a in the unit of time. We assume that the power consumption of an edge node is binary: if it is used, i.e., it serves at least one stateless FaaS or hosts at least one stateful FaaS runner, then it consumes a peak power; otherwise, if it is unused, it does not consume power at all. Such an assumption does not capture the features offered by state-of-the-art power management technologies, which allow nodes to be in multiple states (active, idle, sleep, etc.) and to scale the frequency of individual CPU cores, but contributes to keeping the model simpler and compact. In our future work, we will expand to more sophisticated models.

Regardless of the deployment strategy, we can then define the total energy consumed in the system as follows:

$$E = \sum_{k=1}^{N-1} \left[P_N \cdot \alpha(t_k) + E_B \sum_{a \in A} \beta_a(t_k) \mathbb{I}(t_a^\downarrow \leq t_k \leq t_a^\uparrow) \right] (t_{k+1} - t_k), \quad (1)$$

where P_N is the power consumption of an edge node and E_B is the per-bit network transfer energy, and $\mathbb{I}(\cdot) \in \{0, 1\}$ is an indicator function equal to 1 if and only if the condition is true, which in Equation (1) means that the app a is active at time t_k . In this work, we focus on energy consumption assuming that there are no constraints on the availability of processing and network resources. In other words, we assume that the system can accommodate all the incoming requests, hence no

¹<https://github.com/edgeless-project/edgeless>

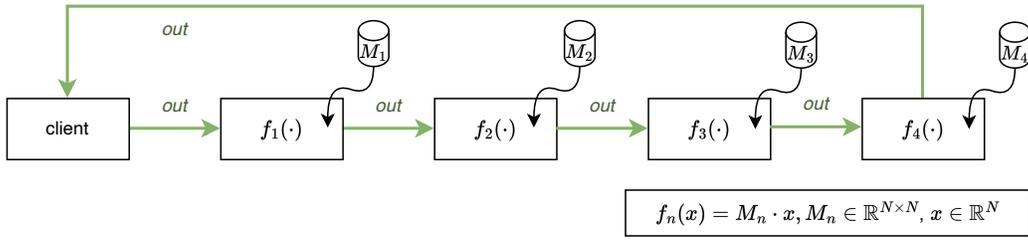


Fig. 6. WebAssembly stateful FaaS: app consisting of four functions in a chain.

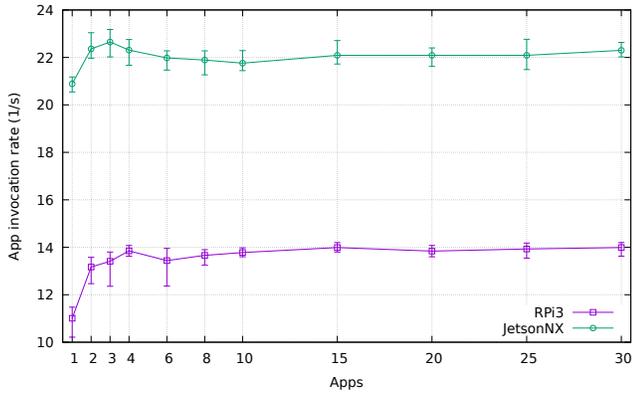


Fig. 7. WebAssembly stateful FaaS: app invocation rate vs. number of apps.

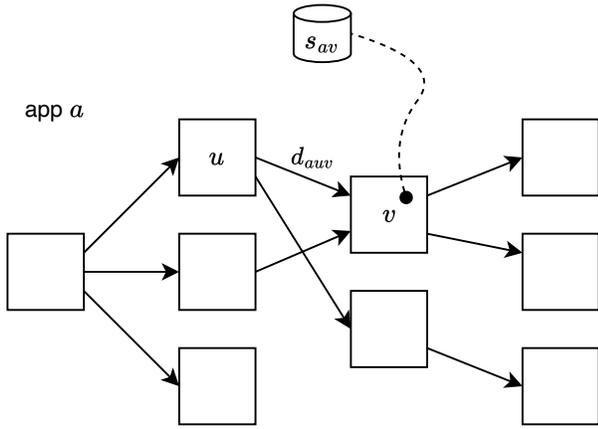


Fig. 8. Application model. An app a consists of functions arranged in a graph. If function u calls function v then an edge exists, and its weight d_{auv} is the amount of data exchanged. Each function v has a state of size s_{av} .

admission control is needed. The notation used in the paper is summarized in Table I.

In the following we define $\alpha(t_k)$ and $\beta_a(t_k)$ separately for stateless FaaS (Section III-A) and stateful FaaS (Section III-B), based on the considerations in Section II above.

A. Stateless FaaS

For stateless FaaS we adopt a simple model that captures well its distinguishing features. Specifically, we assume that the number of active nodes needed at time t_k is the minimum possible, i.e.:

TABLE I
NOTATION USED IN THE PAPER. THE LAST TWO ROWS ARE USED ONLY WITH STATEFUL FAAS.

Parameter	Description	Unit
$G_a(V_a, E_a)$	Task graph of app a . V_a is the set of tasks, E_a represents the invocation dependencies	
λ_a	Invocation rate of app a	s^{-1}
r_{av}	Processing request of task v at app a	CPU
s_{av}	State size of task v at app a	b
d_{auv}	Invocation data size from task u to v at app a	b
t_a^\downarrow	Arrival time of app a	s
t_a^\uparrow	Leaving time of app a	s
A	Set of all the applications in the period of interest	
$A(t_k)$	Set of applications active at time t_k	
$\alpha(t_k)$	Number of edge nodes active at time t_k	
$\beta_a(t_k)$	Traffic rate of app a at time t_k	b/s
P_N	Power consumption of a node	W
E_B	Per-bit network transfer energy	J/b
E	Total energy consumed in the period of interest	J
C	Processing capacity of a node	CPU
t_k	Time instant of the k -th event	s
Δ	Defragmentation interval	s
$x_{av}(t_k)$	Mapping function indicating the index of the node to which task v of app a is allocated at time t_k	

$$\alpha(t_k) = \left\lceil \frac{1}{C} \sum_{a \in A(t_k)} \sum_{v \in V_a} r_{av} \right\rceil, \quad (2)$$

where $A(t_k)$ is set of applications active at time t_k . The inner summation in Equation (2) is the total processing request of app a , which is then summed over all the applications and, finally, divided by the edge node capacity C . This implicitly assumes that no edge effects exist in horizontal scalability and the broker layer can distribute the load appropriately among the multiple task instances. On the other hand, the traffic rate of app a at time t_k is given by:

$$\beta_a(t_k) = \lambda_a \left[\sum_{v \in V_a} s_{av} + \sum_{(u,v) \in E_a} d_{uav} \right], \quad (3)$$

which is the sum of the traffic generated for the state access (first term) and function invocation between each node and its successors (second term), in the unit of time, as given by the invocation rate λ_a .

B. Stateful FaaS

The model with stateful FaaS is more complicated because it depends on how tasks are assigned to edge nodes for three reasons. First, function invocation only consumes network resources if the two tasks are not assigned to the same edge (see Figure 5). Second, since a stateful FaaS runner cannot be split/recombined, assigning the active tasks to available nodes to minimize the number of nodes used is akin to the bin-packing problem, which is known to be NP-complete [20]. Finally, as active apps leave the system, *fragmentation* occurs (a term inspired by the similar effect in the memory management process of operating systems), i.e., edge nodes are only partially allocated: this is sub-optimal for energy consumption. To solve this problem, we foresee a defragmentation process to happen periodically, with the period equal to Δ , which is a system configuration parameter: during defragmentation, the active apps are rearranged to reduce the number of edge nodes needed, thus saving energy in the future. However, this process *consumes* energy because the state of some runners may have to be migrated from one node to another (see Figure 4). In Section IV we will study the trade-off in choosing the value of Δ .

Now we introduce a last bit of notation: let $x_{av}(t_k)$ be a variable that indicates what edge node (using an arbitrary indexing scheme) hosts the runner for the task v of app a at time t_k . In time intervals where the app a is inactive, i.e., before it enters or after it leaves the system, the variable is undefined. The values of $x_{av}(t_k)$ must be determined through two orchestration decision-making algorithms: i) when an app enters the system, the algorithm chooses where to deploy each of its tasks, by either selecting edge nodes already active (hosting other tasks) with sufficient residual capacity or activating new edge nodes; ii) upon defragmentation, the tasks of active applications can be migrated to other edge nodes to reduce the total number of the active ones. Determining an optimal policy for either of these decision processes has the same complexity as finding an optimal allocation for a bin-packing problem, as already mentioned. To keep the focus of this work compact, we defer the study of those problems to future work and we propose to use the following simple heuristic based on the *best-fit* policy, which is widely employed and has bounded performance [21].

Stateful|best-fit algorithm:

- When an app enters the system, for each task we select the active node that hosts one of the predecessor tasks, if any (to save network traffic for function invocation). Otherwise, we select the active node that leaves the smallest residual capacity, if any, breaking ties arbitrarily. Otherwise, we deploy the task on an inactive node.
- Upon defragmentation, we apply the above algorithm policy to all the active apps, in arbitrary order.

We then derive the number of active nodes at time t_k as:

$$\alpha(t_k) = \left| \left\{ x_{av}(t_k), \forall a \in A(t_k), \forall v \in V_a \right\} \right|, \quad (4)$$

where $|\cdot|$ indicates the cardinality of the corresponding set, and the traffic rate of app a at time t_k is:

$$\beta_a(t_k) = \frac{1}{t_{k+1} - t_k} \sum_{v \in V_a} s_{av} \cdot \mathbb{I}(x_{av}(t_k) \neq x_{av}(t_{k-1})) + \lambda_a \sum_{(u,v) \in E_a} d_{auv} \cdot \mathbb{I}(x_{au}(t_k) \neq x_{av}(t_k)), \quad (5)$$

where the first addend takes into account the state migration if the task was moved since the previous time event (by design, this can happen only during the defragmentation procedure) and the second addend considers the network traffic for function invocation, only if the task u and its successor v do not belong to the same node (see Figure 5, again).

IV. PERFORMANCE EVALUATION

In this section, we evaluate the performance, in terms of energy consumption, of the stateless vs. stateful approaches, as modeled in Section III, indicated as *stateless|min-nodes* and *stateful|best-fit*, respectively. For reference purposes, we also include two alternatives: *stateless|max-balancing*, as implied by the name, refers to a stateless FaaS system that seeks to maximize load balancing [22]; *stateful|random* is a variation of the stateful policy in Section III-B, where there is no periodic defragmentation and the tasks of incoming apps are assigned to edge nodes at random, respecting the maximum capacity C , and a new node is made active only if there is none with sufficient residual capacity. For full reproducibility of results, the source code of the simulator and the scripts and artifacts are available publicly as open-source on GitHub².

Methodology and assumptions The workload is created following the model in [23], which is inspired by real traces made available by Alibaba and broadly used in the literature, tuned as follows: the arrival and lifetime of apps follow a Poisson distribution, with average 1 s and 60 s, respectively; both the state size s_{av} and the data invocation size d_{auv} are derived from the memory requirements produced by [23], by applying multiplicative factors called S (state) and D (data invocation), where D is always set to 100, which corresponds to the range [2, 303] kB, and S is expressed through the ratio S/D , which is 100 by default, in which case S would be in the range [0.2, 30.3] MB. The invocation rate is $\lambda_a = 5/s$ and the capacity of a node C is set to 1000, which is sufficient to host any single task, whose requested capacity is drawn from an empiric distribution with a maximum value of 800. The edge node power consumption was set to 100 W, which is typical for a small device such as an Intel NUC; estimating the network consumption is much more complicated because it depends not only on the devices but also on the overall networking infrastructure: based on the results from a recent study [24], we have experimented with different values in the range [0.05, 5] $\mu\text{W}/\text{b}/\text{s}$. Each experiment lasted 1 day of simulated time and was repeated 1000 times; the plots show

²<https://github.com/cciconetti/stateful-faaS-sim> (experiment 001)

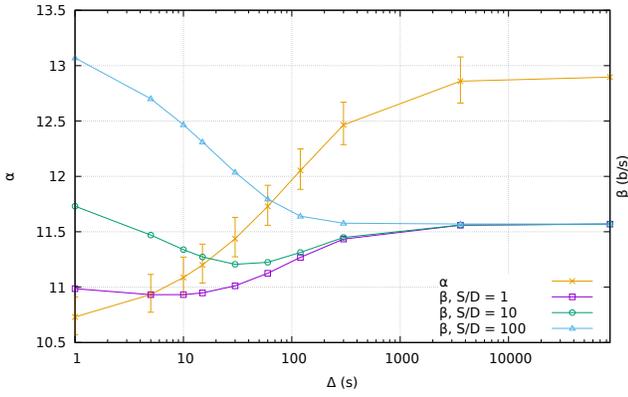


Fig. 9. Simulations: α and β vs. defragmentation period Δ .

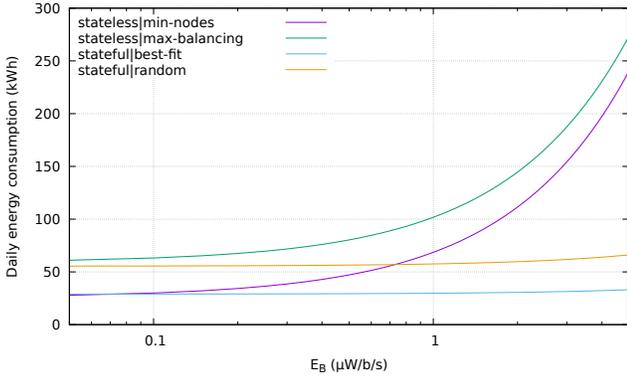


Fig. 10. Simulations: energy consumption vs. E_B .

the average value across the repetitions with a symbol and the low (0.025) and high (0.975) quantiles as error bars. All the values above are to be considered unless specified otherwise.

Defragmentation period In Figure 9 we show α and β with different combinations of Δ and the S/D ratio, only with stateful|best-fit. β is affected significantly by both S/D and Δ : when the state is heavier ($S/D = 100$), the network traffic is very high with small values of Δ (note the log scale on the y -axis) because frequent migrations are expensive. This effect is much less prominent with $S/D = 10$ and $S/D = 1$, because of the smaller state sizes compared to the invocation data sizes. With increasing Δ , all the curves initially decrease and, then, increase again until they converge to the same value (as the defragmentation becomes more sporadic, the state size becomes less important). The minima of the curves depend on the specific value of S/D . The number of active nodes α is independent of S/D and always increases with Δ . **Key message:** *The choice of Δ incurs a trade-off in the energy consumption of computation vs. network.* Devising an algorithm to set Δ at run-time is a possible spin-off research activity. In the following, we set the value to 120 s, i.e., twice the average app lifetime, which appears as a reasonable trade-off between network vs. processing consumption.

Energy per bit-rate In Figure 10 we show the energy

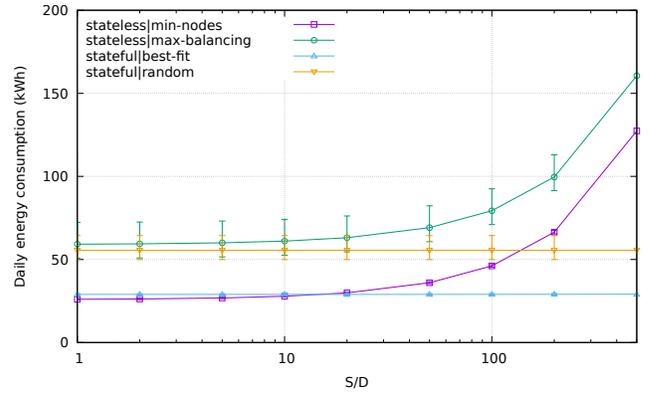


Fig. 11. Simulations: energy consumption vs. S/D , $E_B = 0.05 \mu\text{W/b/s}$.

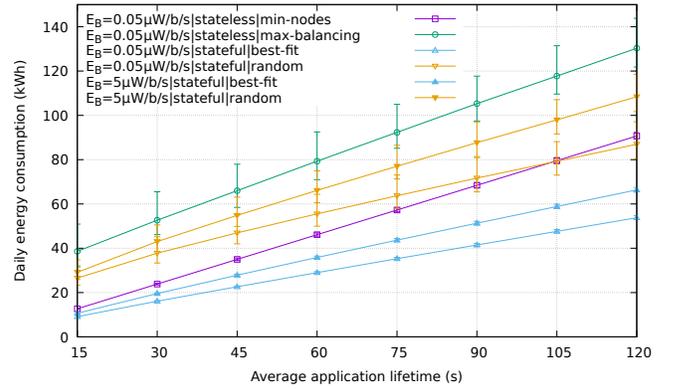


Fig. 12. Simulations: energy consumption vs. average application lifetime.

consumption with increasing E_B while keeping $P_N = 100 \text{ W}$. **Key message:** *The energy consumption increase with a higher per-bit-rate cost is higher with a stateless deployment, especially in the max-balancing flavor, and is very modest with a stateful deployment.* In the latter case, we can see that the best-fit policy reduces energy consumption by about 2 compared to random, for all values of E_B . In the following, we only consider the two extremes of the E_B range.

State size The impact of the state size, compared to the data invocation size, is exposed in Figure 11, with low per-bit-rate energy cost, i.e., $E_B = 0.05 \mu\text{W/b/s}$. A stateless deployment, with a min-nodes policy, is the best option only for $S/D \leq 10$ and only by a small margin compared to stateful|best-fit. On the other hand, as S/D increases significantly above 10, stateless deployment becomes significantly more energy-hungry, due to the cost of accessing the state upon each function invocation. With $S/D > 100$, stateless is outperformed even by stateful|random. The max-balancing policy follows the same trend as min-nodes and is always above the latter, though the gap reduces slightly as S/D increases. **Key message:** *From an energy consumption perspective, stateful deployments are almost insensitive to the size of the applications' states.*

System load In Figure 12 we report the measurements obtained with min/max E_B values for stateful policies (with

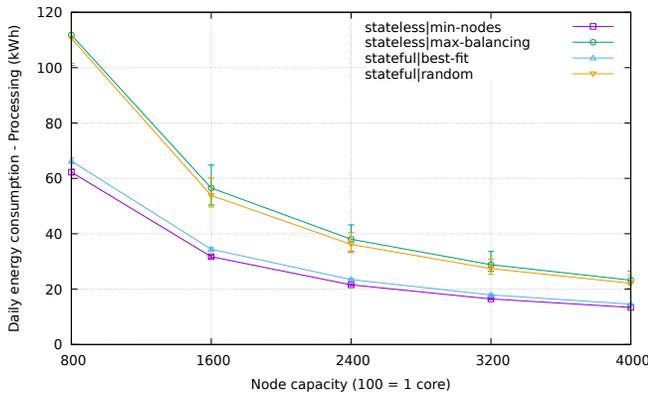


Fig. 13. Simulations: energy consumption vs. node capacity.

stateless, the values with maximum E_B are well above the plot y -axis range) when increasing the application lifetime from 15 s to 120 s. As expected, all the curves increase with the load. **Key message:** *Both stateful|best-fit curves lie at the bottom and gain an increasing margin compared to all the others as the load increases.* The second-best option is stateless|min-nodes (only with minimum E_B), while the stateless|max-balancing performs worst.

Node capacity Finally, in Figure 13 we show the energy consumption (only due to processing) with increasing C from 800 to 4000. All the curves decrease because with increasing C the number of nodes required decreases, as well, while we keep the power consumption per node P_N constant. It is interesting to note that the curves are almost overlapping in pairs. At the bottom (less energy consumed) we find the two systems described in Section III: in fact, they both aim at reducing the edge computing infrastructure energy consumption. Stateless has a slight gain compared to stateful, but it is more than compensated by a lower energy efficiency from the network traffic perspective. At the top (more energy consumed), the two comparison systems show similar performance, which can be explained by the fact that they both try to spread as much as possible the load among the active nodes: stateless|max-balancing does this explicitly, stateful|random implicitly. **Key message:** *A stateful deployment, with a best-fit allocation strategy, can be as efficient as a stateless one despite the fragmentation issue.*

V. RELATED WORK

A. Edge computing support in 5G systems

The Multi-access Cloud Computing (MEC) concept was defined by ETSI [25] to offer a standardized framework to include cloud-computing capabilities to edge networking elements, such as base stations, access points, switches, and routers. Several papers have investigated system architecture and technology enablers to build MEC-enabled shared pools of computing resources at the network edge, especially within cellular networks [26]. More recently, the 3GPP has extended the reference design of MEC systems to support the deployment of edge computing applications within 5G and post-5G

networks [27], by defining the underlying layer that facilitates communication between application clients (ACs) running on the user terminal and application servers (EAS) hosted by an EDN. Furthermore, the 3GPP technical specifications delineate the capabilities for edge server discovery, configuration, and management, encompassing service continuity and the exposure of networking properties. In addition to standardization activities, research studies have investigated possible reference architectures for EDN. For instance, the problem of optimal placement of edge servers in a 5G network is addressed in [28] to maximize the total revenues or to minimize the deployment costs throughout a planning horizon [29]. In [30], the authors argue that an EDN can improve the utilization of computing resources by cooperating with a regional cloud through a wholesale and buyback pricing scheme.

B. Application deployment

The authors in the recent survey [31] identified resource allocation and scheduling as an open research challenge to unlock the full potential of IoT applications at the edge. In [15], the authors argue that most AI applications are stateful and propose an Integer Linear Programming (ILP) to minimize the deployment cost, defined as a combination of processing, transmission, and storage, which then solve via a heuristic algorithm with provable approximation ratio, also putting forward an online learning version under uncertain data volumes and network delays. Even if we start from the same observation, our work proceeds in a different direction of estimating the energy consumption, which is very important in EDNs, under the assumption of a simpler orchestration policy that allows us to focus more specifically on the comparison between stateless and stateful frameworks. On the other hand, in [32] the authors put under the spotlight the trade-off between energy consumption and latency, which is addressed through a distributed algorithm that optimizes at the same time the task offloading decisions, local CPU frequency, and approximate computing accuracy, and is proven to be effective in a small-scale testbed and with numerical experiments. A similar problem has been studied in [33], though the decision-making algorithm uses deep reinforcement learning. However, these papers assume a microservice, not serverless, architecture: this makes the contribution orthogonal to ours, albeit possibly relevant for further investigation. Serverless is instead the specific subject of [34], where the authors design a FaaS platform for the edge-cloud continuum, called *Serverledge*, supporting vertical/horizontal computation offloading. The solution proposed is a possible implementation of the stateless deployment option discussed in this paper, without state management.

VI. CONCLUSIONS

In this paper, we have studied two alternative deployment options to realize stateful FaaS in EDNs. The first relies on stateful FaaS runners interconnected via a service mesh of brokers to dispatch function invocations, where the applications' state resides on an external service accessed as on demand.

The second one exploits WebAssembly to assign one runner to each application instance, state included, thus directly realizing stateful FaaS. Results from a testbed with small devices showed no noticeable drop in performance with several runners. We have then defined models to estimate the energy consumption of processing and network traffic for the two deployment options, with known statistical characterization of the applications (arrival process, DAG task dependencies, and state/invoke sizes). For stateful FaaS, we have proposed a simple, yet effective, heuristic for allocating tasks to nodes upon application arrival and periodically to reduce fragmentation of edge node resources. Extensive simulations have shown that stateful FaaS is more efficient than stateless FaaS, except for minimal applications' state or negligible network energy consumption compared to processing.

ACKNOWLEDGMENT

The work of A. Passarella was funded by the European Union under the Italian National Recovery and Resilience Plan (NRRP) of NextGenerationEU, partnership on "Telecommunications of the Future" (PE00000001 – program "RESTART"). The work was partly funded by the European Union under the projects EDGELESS (GA no. 101092950), 6Green (GA 101096925), and GreenDIGIT (GA 101131207).

REFERENCES

- [1] H. Chang, A. Hari, S. Mukherjee, and T. V. Lakshman, "Bringing the cloud to the edge," *IEEE INFOCOM*, pp. 346–351, 2014.
- [2] R. Ramalingam and T. Tung, "Leading with Edge Computing – How to reinvent with data and AI," Accenture, Tech. Rep., 2023.
- [3] T. Otto, "Edge AI: How AI is sparking the adoption of edge computing," STL Partners, Tech. Rep., Oct. 2023.
- [4] C.-X. Wang, X. You, X. Gao, X. Zhu, Z. Li, C. Zhang, H. Wang, Y. Huang, Y. Chen, H. Haas, J. S. Thompson, E. G. Larsson, M. D. Renzo, W. Tong, P. Zhu, X. Shen, H. V. Poor, and L. Hanzo, "On the Road to 6G: Visions, Requirements, Key Technologies, and Testbeds," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 2, pp. 905–974, 2023.
- [5] "3GPP TS 28.538 Edge Computing Management (ECM) (Release 18)," Dec. 2023.
- [6] X. Foukas, G. Patounas, A. Elmokashfi, and M. K. Marina, "Network Slicing in 5G: Survey and Challenges," *IEEE Communications Magazine*, vol. 55, no. 5, pp. 94–100, May 2017.
- [7] E. Zeydan, J. Manges-Bafalluy, J. Baranda, M. Requena, and Y. Turk, "Service Based Virtual RAN Architecture for Next Generation Cellular Systems," *IEEE Access*, vol. 10, pp. 9455–9470, 2022.
- [8] S. Kounev, N. Herbst, C. L. Abad, A. Iosup, I. Foster, P. Shenoy, O. Rana, and A. A. Chien, "Serverless Computing: What It Is, and What It Is Not?" *Communications of the ACM*, vol. 66, no. 9, pp. 80–92, 2023.
- [9] S. Eismann, J. Scheuner, E. Van Eyk, M. Schwinger, J. Grohmann, N. Herbst, C. Abad, and A. Iosup, "The State of Serverless Applications: Collection, Characterization, and Community Consensus," *IEEE Trans. on Software Engineering*, 2021.
- [10] P. Raith, S. Nastic, and S. Dustdar, "Serverless Edge Computing—Where We Are and What Lies Ahead," *IEEE Internet Computing*, vol. 27, no. 3, pp. 50–64, May 2023.
- [11] J. Spenger, P. Carbone, and P. Haller, "A Survey of Actor-Like Programming Models for Serverless Computing," in *Active Object Languages: Current Research Trends*, F. De Boer, F. Damiani, R. Hähnle, E. Broch Johnsen, and E. Kamburjan, Eds. Cham: Springer Nature Switzerland, 2024, vol. 14360, pp. 123–146.
- [12] C. Freitag, M. Berners-Lee, K. Widdicks, B. Knowles, G. S. Blair, and A. Friday, "The real climate and transformative impact of ICT: A critique of estimates, trends, and regulations," *Patterns*, vol. 2, no. 9, p. 100340, Sep. 2021.
- [13] S. Risco, C. Alarcón, S. Langarita, M. Caballer, and G. Moltó, "Rescheduling serverless workloads across the cloud-to-edge continuum," *Future Generation Computer Systems*, vol. 153, pp. 457–466, Apr. 2024.
- [14] A. Mahgoub, E. B. Yi, K. Shankar, E. Minocha, S. Elnikety, S. Bagchi, and S. Chaterji, "Wisefuse: Workload Characterization and DAG Transformation for Serverless Workflows," *ACM POMACS*, vol. 6, no. 2, 2022.
- [15] Z. Xu, L. Zhou, W. Liang, Q. Xia, W. Xu, W. Ren, H. Ren, and P. Zhou, "Stateful Serverless Application Placement in MEC with Function and State Dependencies," *IEEE Trans. on Computers*, pp. 1–14, 2023.
- [16] R. Xie, Q. Tang, S. Qiao, H. Zhu, F. Richard Yu, and T. Huang, "When Serverless Computing Meets Edge Computing: Architecture, Challenges, and Open Issues," *IEEE Wireless Communications*, pp. 1–8, 2021.
- [17] C. Cicconetti, M. Conti, and A. Passarella, "FaaS execution models for edge applications," *Pervasive and Mobile Computing*, vol. 86, 2022.
- [18] D. Schafer, J. Edinger, J. M. Paluska, S. VanSyckel, and C. Becker, "Tasklets: "Better than Best-Effort" Computing," in *IEEE ICCCN*, 2016.
- [19] P. Gackstatter, P. A. Frangoudis, and S. Dustdar, "Pushing Serverless to the Edge with WebAssembly Runtimes," in *IEEE CCGrid*, 2022, pp. 140–149.
- [20] M. R. Garey and D. S. Johnson, *Computers and intractability: a guide to the theory of NP-completeness*, 27th ed., ser. A series of books in the mathematical sciences. New York [u.a]: Freeman, 2009.
- [21] S. Martello and P. Toth, "A Bound and Bound algorithm for the zero-one multiple knapsack problem," *Discrete Applied Mathematics*, vol. 3, no. 4, pp. 275–288, Nov. 1981.
- [22] C. Cicconetti, M. Conti, and A. Passarella, "A Decentralized Framework for Serverless Edge Computing in the Internet of Things," *IEEE Trans. on Network and Service Management*, pp. 1–1, 2020.
- [23] H. Tian, Y. Zheng, and W. Wang, "Characterizing and Synthesizing Task Dependencies of Data-Parallel Jobs in Alibaba Cloud," *ACM SoCC*, pp. 139–151, 2019.
- [24] E. Ahvar, A.-C. Orgerie, and A. Lebre, "Estimating Energy Consumption of Cloud, Fog, and Edge Computing Infrastructures," *IEEE Trans. on Sustainable Computing*, vol. 7, no. 2, pp. 277–288, Apr. 2022.
- [25] ETSI ISG MEC, "Multi-access Edge Computing (MEC); Framework and Reference Architecture," European Telecommunications Standards Institute (ETSI), GS MEC 003 V3.1.1, mar. 2022.
- [26] H. Liu, F. Eldarrat, H. Alqahtani, A. Reznik, X. de Foy, and Y. Zhang, "Mobile edge cloud system: Architectures, challenges, and approaches," *IEEE Systems Journal*, vol. 12, no. 3, pp. 2495–2508, 2018.
- [27] 3GPP, "Multi-access Edge Computing (MEC); Framework and Reference Architecture," 3rd Generation Partnership Project (3GPP), TS 23.558 V19.0.0, dec. 2023.
- [28] Y. Zhang, W. Wang, J. Ren, J. Huang, S. He, and Y. Zhang, "Efficient revenue-based mec server deployment and management in mobile edge-cloud computing," *IEEE/ACM Trans. on Networking*, vol. 31, no. 4, pp. 1449–1462, 2023.
- [29] W. Ren, Y. Sun, H. Luo, and M. Guizani, "A demand-driven incremental deployment strategy for edge computing in iot network," *IEEE Trans. on Network Science and Engineering*, vol. 9, no. 2, pp. 416–430, 2022.
- [30] Y. Zhang, X. Lan, J. Ren, and L. Cai, "Efficient computing resource sharing for mobile edge-cloud computing networks," *IEEE/ACM Trans. on Networking*, vol. 28, no. 3, pp. 1227–1240, 2020.
- [31] A. Hazra, A. Kalita, and M. Gurusamy, "Meeting the Requirements of Internet of Things: The Promise of Edge Computing," *IEEE Internet of Things Journal*, pp. 1–1, 2023.
- [32] A. Younis, S. Maheshwari, and D. Pompili, "Energy-Latency Computation Offloading and Approximate Computing in Mobile-Edge Computing Networks," *IEEE Trans. on Network and Service Management*, pp. 1–1, 2024.
- [33] Z. Aghapour, S. Sharifian, and H. Taheri, "Task offloading and resource allocation algorithm based on deep reinforcement learning for distributed AI execution tasks in IoT edge computing environments," *Computer Networks*, vol. 223, p. 109577, Mar. 2023.
- [34] G. R. Russo, T. Mannucci, V. Cardellini, and F. L. Presti, "Serverledge: Decentralized Function-as-a-Service for the Edge-Cloud Continuum," in *IEEE PerCom*, 2023, pp. 131–140.