

SplitPlace: AI Augmented Splitting and Placement of Large-Scale Neural Networks in Mobile Edge Environments

Shreshth Tuli, Giuliano Casale and Nicholas R. Jennings

Abstract—In recent years, deep learning models have become ubiquitous in industry and academia alike. Deep neural networks can solve some of the most complex pattern-recognition problems today, but come with the price of massive compute and memory requirements. This makes the problem of deploying such large-scale neural networks challenging in resource-constrained mobile edge computing platforms, specifically in mission-critical domains like surveillance and healthcare. To solve this, a promising solution is to split resource-hungry neural networks into lightweight disjoint smaller components for pipelined distributed processing. At present, there are two main approaches to do this: semantic and layer-wise splitting. The former partitions a neural network into parallel disjoint models that produce a part of the result, whereas the latter partitions into sequential models that produce intermediate results. However, there is no intelligent algorithm that decides which splitting strategy to use and places such modular splits to edge nodes for optimal performance. To combat this, this work proposes a novel AI-driven online policy, SplitPlace, that uses Multi-Armed-Bandits to intelligently decide between layer and semantic splitting strategies based on the input task's service deadline demands. SplitPlace places such neural network split fragments on mobile edge devices using decision-aware reinforcement learning for efficient and scalable computing. Moreover, SplitPlace fine-tunes its placement engine to adapt to volatile environments. Our experiments on physical mobile-edge environments with real-world workloads show that SplitPlace can significantly improve the state-of-the-art in terms of average response time, deadline violation rate, inference accuracy, and total reward by up to 46, 69, 3 and 12 percent respectively.

Index Terms—Mobile Edge Computing, Neural Network Splitting, Container Orchestration, Artificial Intelligence, QoS Optimization.



1 INTRODUCTION

Modern Deep Neural Networks (DNN) are becoming the backbone of many industrial tasks and activities [2]. As the computational capabilities of devices have improved, new deep learning models have been proposed to provide improved performance [3], [4]. Moreover, many recent DNN models have been incorporated with mobile edge computing to give low latency services with improved accuracies compared to shallow networks, particularly in complex tasks like image segmentation, high frame-rate gaming and traffic surveillance [5]. The performance of such neural models reflects directly on the reliability of application domains like self-driving cars, healthcare and manufacturing [6], [2]. However, to provide high accuracy, such neural models are becoming increasingly demanding in terms of data and compute power, resulting in many challenging problems. To accommodate these increasing demands, such massive models are often hosted as web services deployed on the public cloud [7], [8].

Challenges. Recently, application demands have shifted from either high-accuracy or low-latency to both of these together, termed as HALL (high-accuracy and low-latency) service delivery [2]. Given the prevalence and demand of

DNN inference, serving them on a public cloud with tight bounds of latency, throughput and cost is becoming increasingly challenging [9]. In this regard, recent paradigms like mobile edge computing seem promising. Such approaches allow a robust and low-latency deployment of Internet of Things (IoT) applications close to the edge of the network. Specifically, to solve the problem of providing HALL services, recent work proposes to integrate large-scale deep learning models with modern frameworks like edge computing [9], [10], [11]. However, even the most recent approaches either provide a low Service Level Agreement (SLA) violation mode or a high-accuracy mode [10], [9] and struggle to provide the benefits of both modes at the same time.

Another challenge of using edge computing is that mobile edge devices face severe limitations in terms of computational and memory resources as they rely on low power energy sources like batteries, solar or other energy scavenging methods [12], [13]. This is not only because of the requirement of low cost, but also the need for mobility in such nodes [5]. In such systems, it is possible to handle the processing limitations of massive DNN models by effective preemption and prolonged job execution. However, memory bottlenecks are much harder to solve [14]. In a distributed edge environment, storage spaces are typically mapped to network-attached-storage (NAS) media. Thus, prior work that runs inference on a pre-trained DNN without memory-aware optimizations leads to high network bandwidth overheads due to frequent overflow of memory

A preliminary version of this work was presented at the Student Research Competition in ACM SIGMETRICS Conference 2021 [1].

- S. Tuli, G. Casale and N. R. Jennings are with the Department of Computing, Imperial College London, United Kingdom.
- N. R. Jennings is also with Loughborough University, United Kingdom. E-mails: {s.tuli20, g.casale}@imperial.ac.uk, n.r.jennings@lboro.ac.uk.

Manuscript received —; revised —.

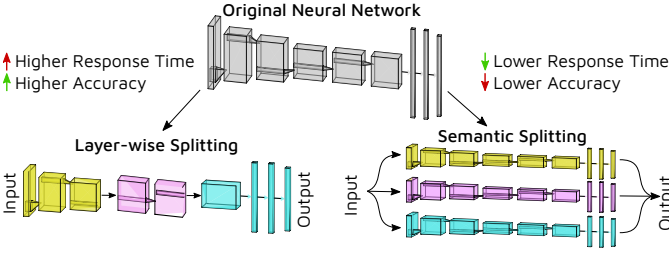


Figure 1: Overview of layer and semantic splitting strategies

and the use of virtual memory (swap space) on NAS media, making high fidelity inference using DNNs hard [15], [10]. To deploy an upgraded AI model, tech giants like Amazon, Netflix and Google usually consider completely revamping their infrastructures and upgrading their devices, raising many sustainability concerns [2]. This has made the integration of massive neural network models with such devices a challenging and expensive ordeal.

Solution. A promising solution for this problem is the development of strategies that can accommodate large-scale DNNs within legacy infrastructures. However, many prior efforts in this regard [16], [17], [18] have not yet tackled the challenge of providing a holistic strategy for not only distributed learning, but also inference in such memory-constrained environments. Recently, research ideas have been proposed like Cloud-AI, Edge-AI and Federated learning that aim to solve the problem of running enormous deep learning models on constrained edge devices by splitting them into modular fragments [17], [18]. However, in Cloud-AI where AI systems are deployed on cloud machines, the high communication latency leads to high average response times, making it unsuitable for latency-critical applications like healthcare, gaming and augmented reality [6], [19], [20]. Instead, Edge-AI provides low-latency service delivery, thanks to edge devices being in the same Local Area Network (LAN), where the input data from multiple edge nodes are combined to a single fixed broker node for processing. Edge-AI based methods aim at scheduling deep neural networks for providing predictable inference [21], [22]. However, due to the centralized collection of data, these solutions typically suffer from high bandwidth overheads and poor service quality [17]. Federated learning depends on data distribution over multiple nodes where the model training and inference are performed in a decentralized fashion. However, this paradigm assumes that neural models with data batches can be accommodated in the system memory. This is seldom the case for common edge devices like Arduinos or Raspberry Pis [23].

Other recent works offer lower precision models that can fit within the limited memory of such devices by using methods like Model Compression or Model Pruning [24], [9], [25]. However, compressed and low-precision models lose inference accuracy, making them unsuitable for accuracy-sensitive applications like security and intrusion detection [26]. Recently, split neural network models have been proposed. They show that using semantic or layer-wise splitting, a large deep neural network can be fragmented into multiple smaller networks for dividing network parameters onto multiple nodes [27], [28], [16], [29]. The former

partitions a neural network into parallel disjoint models that produce a part of the result. The latter partitions a neural network into sequential models that generate intermediate results. We illustrate the accuracy and response time trade-offs through sample test cases in Section 2. Our experiments show that using layer and semantic splitting gives higher inference accuracies than previously proposed model compression techniques (see Section 6). However, no appropriate scheduling policies exist that can intelligently place such modular neural fragments on a distributed infrastructure to optimize both accuracy and SLA together. The placement of such split models is non-trivial considering the diverse and complex dynamism of task distribution, model usage frequencies and geographical placement of mobile edge devices [30].

Research Contributions. This work proposes a novel neural splitting and placement policy, *SplitPlace*, for enhanced distributed neural network inference at the edge. *SplitPlace* leverages a mobile edge computing platform to achieve low latency services. It allows modular neural models to be integrated for best result accuracies that could only be provided by cloud deployments. *SplitPlace* is the *first* splitting policy that dynamically decides between semantic and layer-wise splits to optimize both inference accuracy and the SLA violation rate. This decision is taken for each incoming task and remains unmodified until the execution of all split fragments of that task are complete. The idea behind the proposed splitting policy is to decide for each incoming task whether to use the semantic or layer-wise splitting strategy based on its SLA demands. Due to their quick adaptability, *SplitPlace* uses Multi-Armed-Bandits to model the decision strategy for each application type by checking if the SLA deadline is higher or lower than an estimate of the response time for a layer split decision [31]. Further, *SplitPlace* optimizes the placement decision of the modular neural network fragments using a split decision aware surrogate model. Compared to a preliminary extended abstract of this work [1], this paper provides a substantially expanded exposition of the working of MABs in *SplitPlace*. We also present techniques to dynamically adapt to non-stationary workloads and mobile environments. We present a gradient-based optimization approach for task placement decision conditioned on split decisions. Experiments on real-world application workloads on a physical edge testbed show that the *SplitPlace* approach outperforms the baseline approaches by reducing the SLA violation rate and improving the average inference accuracy.

Outline. The rest of the paper presents a brief background with motivation and related work in Section 2. Section 3 presents the system model assumptions and formulates the problem. Sections 4 and 5 give the model details of the proposed *SplitPlace* approach. We then validate and show the efficacy of the placement policy in Section 6. Finally, Section 7 concludes the work and proposes future directions. Additional experimental results are given in the Appendix A in the supplementary text.

2 BACKGROUND AND RELATED WORK

As discussed in Section 1, there is a need for frameworks that can exploit the low latency of edge nodes and also high

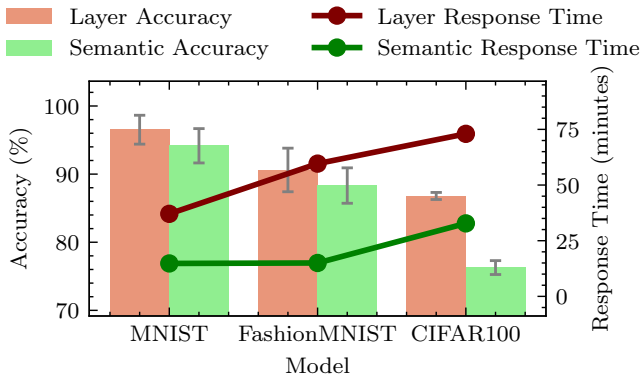


Figure 2: Comparison of Layer and Semantic Splits.

inference performance of DNNs to provide HALL services. However, complete neural models with the input batch can seldom be accommodated in the random-access-memory (RAM) of edge devices. Thus, ideas like model compression or splitting are required to make inference over large-scale neural networks plausible in such environments. Frameworks that aim at achieving this must maintain a careful balance between accuracy requirements and response times for different user tasks. For such a framework, real-time analysis of the incoming tasks is required for quick decision making of task placement. This requires robust algorithms to seamlessly integrate different paradigms and meet the user’s service level agreements.

Semantic and Layer Splitting. In this work, we leverage the only two available splitting schemes for neural networks: layer and semantic splitting [32], [16]. An overview of these two strategies is shown in Figure 1. Semantic splitting divides the network weights into a hierarchy of multiple groups that use a different set of features (different colored models in Figure 1). Here, the neural network is split based on the data semantics, producing a tree structured model that has no connection among branches of the tree, allowing parallelization of input analysis [16]. Due to limited information sharing among the neural network fragments, the semantic splitting scheme gives lower accuracy in general. Semantic splitting requires a separate training procedure where publicly available pre-trained models cannot be used. This is because a pre-trained standard neural network can be split layer wise without affecting output semantics. For semantic splitting we would need to first split the neural network based on data semantics and re-train the model. However, semantic splitting provides parallel task processing and hence lower inference times, more suitable for mission-critical tasks like healthcare and surveillance. Layer-wise splitting divides the network into groups of layers for sequential processing of the task input, shown as different colored models in Figure 1.

Layer splitting is easier to deploy as pre-trained models can be just divided into multiple layer groups and distributed to different mobile edge nodes. However, layer splits require a semi-processed input to be forwarded to the subsequent edge node with the final processed output to be sent to the user, thus increasing the overall execution time. Moreover, layer-wise splitting gives higher accuracy

compared to semantic splitting. Comparison of accuracies and average response times for the two strategies is shown in Figure 2. The figure shows results for 10 edge worker nodes using popular image classification datasets: MNIST, FashionMNIST and CIFAR100 [33], [34], [35] averaged over ResNet50-V2, MobileNetV2 and InceptionV3 neural models [9]. As is apparent from the figure, layer splits provide higher accuracy and response time, whereas semantic splits provide lower values for both. SplitPlace leverages this contrast in traits to trade-off between inference accuracy and response time based on SLA requirements of the input tasks. Despite the considerable drop in inference accuracy when using semantic splitting scheme, it is still used in the proposed SplitPlace approach as it is better than model compression or early-exit strategies for quick inference. This is acceptable in many industrial applications [36], [25] where latency and service level agreements are more important performance metrics than high-fidelity result delivery. In this work, we consider a system with both SLA violation rates and inference accuracy as optimization objectives. This makes the combination of layer and semantic splitting a promising choice for such use cases.

2.1 Related Work

We now analyze the prior work in more detail. We divide our literature review into three major sections based on the strategy used to allow DNN inference on resource-constrained mobile-edge devices: model compression, layer splitting and semantic splitting. Moreover, we compare prior work based on whether they are suitable for edge-only setups (*i.e.*, without leveraging cloud nodes), consider heterogeneous and mobile nodes and work in settings with adaptive Quality of Service (QoS). See Table 1 for an overview.

Model Compression: Efficient compression of DNN models has been a long studied problem in the literature [43]. Several works have been proposed that aim at the structural pruning of neural network parameters without significantly impacting the model’s performance. These use approaches like tensor decomposition, network sparsification and data quantization [43]. Such pruning and model compression approaches have also been used by the systems research community to allow inference of massive neural models on devices with limited resources [44]. Recently, architectures like BottleNet and Bottlenet++ have been proposed [37], [38] to enable DNN inference on mobile cloud environments and reduce data transmission times. BottleNet++ compresses the intermediate layer outputs before sending them to the cloud layer. It uses a model re-training approach to prevent the inference being adversely impacted by the lossy compression of data. Further, BottleNet++ classifies workloads in terms of compute, memory and bandwidth bound categories and applies an appropriate model compression strategy. Other works propose to efficiently prune network channels in convolution neural models using reinforcement learning [39], [45]. Other efforts aim to prune the weights of the neural models to minimize their memory footprint [46]. Such methods aim at improving the accuracy per model size as a metric in contrast to the result delivery time as in BottleNet++. However, model compression does

Table 1: Comparison of related works with different parameters (✓ means that the corresponding feature is present).

Work	Edge Only	Mobility	Heterogeneous Environment	Adaptive QoS	Layer Split	Semantic Split	Model Compression	Optimization Parameters		
								Accuracy	SLA	Reward
[37], [38]							✓		✓	
[39]	✓		✓				✓	✓		
[27], [40], [36], [41]		✓			✓				✓	
[9]	✓		✓	✓	✓			✓	✓	✓
[11], [42]	✓			✓	✓					
[16], [28], [36], [25]	✓		✓	✓		✓		✓	✓	
[32]	✓		✓	✓	✓		✓	✓	✓	
SplitPlace	✓	✓	✓	✓	✓	✓		✓	✓	✓

not leverage multiple compute nodes and has poor inference accuracy in general compared to semantic split execution (discussed in Section 6). Thus, SplitPlace does not use the model compression technique.

Layer Splitting: Many other efforts aim at improving the inference time or accuracy by efficient splitting of the DNN models. Some methods aim to split the networks layer-wise or vertically, *viz.*, that the different fragments correspond to separate layer groups and hence impose the constraint of sequential execution. Most work in this category aims at segregating these network splits into distinct devices based on their computational performance [27], [40], [36], [42], [41]. In heterogeneous edge-cloud environments, it is fairly straightforward to split the network into two or three fragments each being deployed in a mobile device, edge node or a cloud server. Based on the SLA, such methods provide early-exits if the turnaround time is expected to be more than the SLA deadline. This requires a part of the inference being run at each layer of the network architecture instead of traditionally executing it on the cloud server. Other recent methods aim at exploiting the resource heterogeneity in the same network layer by splitting and placing DNNs based on user demands and edge worker capabilities [9]. Such methods can not only split DNNs, but also choose from different architectural choices to reach the maximum accuracy while agreeing to the latency constraints. Other works aim at accelerating the model run-times by appropriate scheduling of a variety of DNN models on edge-clusters [11]. The state-of-the-art method, *Gillis* uses a hybrid model, wherein it employs either model-compression or layer-splitting based on the application SLA demands [32]. The decision is taken using a reinforcement-learning model which continuously adapts in dynamic scenarios. As the model partitioning is also performed dynamically, the *Gillis* model cannot work with semantic splitting strategy that requires a new model to be trained for each partitioning scheme. It is a serverless based model serving system that automatically partitions a large model across multiple serverless functions for faster inference and reduced memory footprint per function. The *Gillis* method employs two model partitioning algorithms that respectively achieve latency optimal serving and cost-optimal serving with service-level agreement compliance. However, this method cannot jointly optimize both latency and SLA. Moreover, it does not consider the mobility of devices or users and hence is ineffective in efficiently managing large DNNs in mobile edge computing environments.

Semantic Splitting: The early efforts of semantic split-

ting only split the neural network at the input layer to allow model parallelization and size reduction [29]. Some methods divide the data batch itself across multiple nodes addressing computational contention problems but not memory limitations of fitting neural networks in the RAM [47]. Other methods use progressive slicing mechanisms to partition neural models into multiple components to fit in heterogeneous devices [25]. Such methods produce the complete output from each split or fragment of the DNN, adversely impacting the scalability of such methods to high-dimensional output spaces such as image segmentation applications [16], [28]. Moreover, these methods require significant cross-communication among network splits, significantly increasing the communication overheads. Recently, more intelligent approaches have been developed which hierarchically split neural networks such that each fragment produces a part of the output using an intelligently chosen sub-part of the input [16]. Such schemes use the “semantic” information of the data to create the corresponding links between input and output sub-parts being given to each DNN fragment, hence the name *semantic splitting*. Such splitting schemes require minimal to no interaction among network fragments eliminating the communication overheads and increased latency due to stragglers. As semantic splitting can provide results quickly, albeit with reduced accuracy, SplitPlace uses it for tasks with tight deadlines (Section 4).

3 SYSTEM MODEL AND PROBLEM FORMULATION

In this work, we assume a scenario with a fixed number of multiple heterogeneous edge nodes in a broker-worker fashion, which is a typical case in mobile-edge environments [48], [49], [38], [27]. Here, the broker node takes all resource management related decisions, such as neural network splitting and task placement. The processing of such tasks is carried out by the worker nodes. Examples of broker nodes include personal laptops, small-scale servers and low-end workstations [48]. Example of common worker nodes in edge environments include Raspberry Pis, Arduino and similar System-on-Chip (SoC) computers [2]. All tasks are received from an IoT layer that includes sensors and actuators to collect data from the users and send it to the edge broker via the gateway devices. Akin to typical edge configurations [50], the edge broker then decides which splitting strategy to use and schedules these fragments to various edge nodes based on deployment constraints like sequential execution in a layer-decision. The data to be processed comes from the IoT sensors/actuators, which with

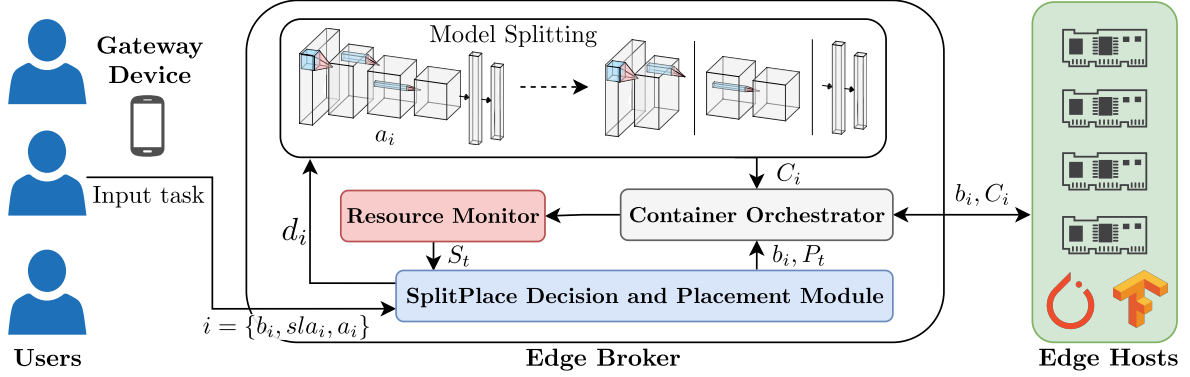


Figure 3: SplitPlace System Model

the decision of which split fragment to use is forwarded by the broker to each worker node. Some worker nodes are assumed to be mobile, whereas others are considered to be fixed in terms of their geographical location. In our formulation, we consider mobility only in terms of the variations in terms of the network channels and do not consider the worker nodes or users crossing different networks. We assume that the CPU, RAM, Bandwidth and Disk capacities of all nodes are known in advance, and similarly the broker can sample the resource consumption for each task in the environment at any time (see *Resource Monitor* in Fig. 3). As we describe later, the broker periodically measure utilizations of CPU, RAM, Bandwidth and Disk for each task in the system. The broker is trusted with this information such that it can make informed resource management decisions to optimize QoS. Moreover, we consider that tasks include a batch of inputs that need to be processed by a DNN model. Further, for each task, a service level deadline is defined at the time the task is sent to the edge environment. We give an overview of the SplitPlace system model in Figure 3. We decompose the problem into deciding an optimal splitting strategy and a fragment placement for each application (motivation in Appendix A.6 and more details in Section 1).

Workload Model. We consider a bounded discrete time control problem where we divide the timeline into equal duration intervals, with the t -th interval denoted as I_t . Here, $t \in \{0, \dots, \Gamma\}$, where $\Gamma + 1$ is the number of intervals in an execution. We assume a fixed number of worker machines in the edge layer and denote them as \mathcal{H} . We also consider that new tasks created at the interval I_t are denoted as \mathcal{N}_t , with all active tasks being denoted as \mathcal{T}_t (and $\mathcal{N}_t \subseteq \mathcal{T}_t$). Each task $i \in \mathcal{T}_t$ consists of a batch input b_i , SLA deadline sla_i and a DNN application a_i . The set of all possible DNN applications is denoted by \mathcal{A} . For each new task $i \in \mathcal{N}_t$, the edge broker takes a decision d^i , such that $d^i \in \{L, S\}$, with L denoting layer-wise splitting and S denoting semantic split strategy. The collection of all split decisions for active tasks in interval I_t is denoted as $\mathcal{D}_t = \{d^i\}_{i \in \mathcal{N}_t}$. Based on the decision d^i for task i , this task is realized as an execution workflow in the form of containers C^i . Similar to a VM, a container is a package of virtualized software that contains all of the necessary elements to run in any environment. The set of all containers active in the interval I_t is denoted as $\mathcal{C}_t = \cup_{i \in \mathcal{T}_t} C^i$. The set of all utilization metrics of CPU, RAM, Network Bandwidth and Disk for all containers and

workers at the start of the interval I_t defines the state of the system, denoted as S_t . A summary of the symbols is given in Table 2.

3.1 Split Nets Placement

We partition the problem into two sub-problems of deciding the optimal splitting strategy for input tasks and that of placement of active containers in edge workers (see Figure 4). Considering the previously described system model, at the start of each interval I_t , the SplitPlace model takes the split decision d^i for all $i \in \mathcal{N}_t$. Moreover, it also takes a placement decision for all active containers \mathcal{C}_t , denoted as an adjacency matrix $P_t : \mathcal{C}_t \times \mathcal{H}$. This is realized as a container allocation for new tasks and migration for active tasks in the system.

The main idea behind the layer-wise split design is first to divide neural networks into multiple independent splits, classify these splits in preliminary, intermediate and final neural network layers and distribute them across different nodes based on the node capabilities and network hierarchy. This exploits the fact that communication across edge nodes in the LAN with few hop distances is very fast and has low latency and jitter [51]. Moreover, techniques like knowledge distillation can be further utilized to enhance the accuracy of the results obtained by passing the input through these different classifiers. However, knowledge distillation needs to be applied at the training stage, before generating the neural network splits. As there are many inputs in our assumed large-scale deployment, the execution can be performed in a pipelined fashion to further improve throughput over and above the low response time of the nodes at the edge of the network. For the semantic split, we divide the network weights into a set or a hierarchy of multiple groups that use disjoint sets of features. This is done by making assignment decisions of network parameters to edge devices at deployment time. This produces a tree-structured network that involves no connection between branched sub-trees of semantically disparate class groups. Each sub-group is then allocated to an edge node. The input is either broadcasted from the broker or forwarded in a ring-topology to all nodes with the network split corresponding to the input task. We use standard layer [32] and semantic splitting [16] methods as discussed in Section 2.

We now outline the working of the proposed distributed deep learning architecture for edge computing environ-

Table 2: Table of Main Notation

Notation	Description
I_t	t -th scheduling interval
\mathcal{T}_t	Active tasks in I_t
\mathcal{N}_t	New tasks received at the start of I_t
\mathcal{H}	Set of workers in the edge layer
$i = \{b_i, sla_i, a_i\}$	Task as a collection input batch, SLA deadline and application type
$d^i \in \{L, S\}$	Splitting decision for input task i
C^i	Container realization of task i based on decision d^i
$C_t = \cup_{i \in \mathcal{T}_t} C^i$	Set of all active containers in the interval I_t
S_t	State of the system at the start of I_t
$P_t : C_t \times \mathcal{H}$	Placement decision of C_t to \mathcal{H} as an adjacency matrix
O_t	Objective score for interval I_t

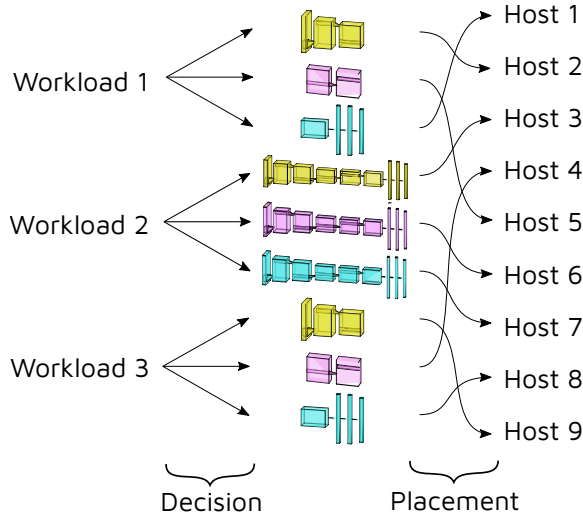


Figure 4: SplitPlace decision and placement problems.

ments. Figure 3 shows a schematic view of its working. As shown, there is a shared repository of neural network parameters which is distributed by the broker to multiple edge nodes. The layer and semantic splits are realized as Docker container images that are shared by the broker to the worker nodes at the start of each execution trace. The placement of tasks is realized as spinning up a Docker container using the corresponding image on the worker. As the process of sharing container images is a one-time event, transferring all semantic and layer split fragments does not impose a high overhead on network bandwidth at runtime. This sharing of containers for each splitting strategy and dataset type are transferred to the worker nodes is performed at the start of the run. At run-time, only the decision of which split fragment to be used is communicated to the worker nodes, which executes a container from the corresponding image. The placement of task on each worker is based on the resource availability, computation required to be performed in each section and the capabilities of the nodes (obtained by the *Resource Monitor*). For intensive computations with large storage requirements (Gated Recurrent Units or LSTMs) or splits with high dimension size of input/output (typically the final layers), the splits are sent to high-resource edge workers. The management of allocation and migration of

neural network splits is done by the *Container Orchestrator*. Other attention based sub-layer extensions can be deployed in either edge or cloud node based on application requirements, node constraints and user demands. Based on the described model assumptions, we now formulate the problem of taking splitting and placement decisions to optimize the QoS parameters. Implementation specific details on how the results of layer-splits are forwarded and outputs of semantic splits combined across edge nodes are given in Section 5.

3.2 Problem Formulation

The aim of the model is to optimize an objective score O_t (to be maximized), which quantifies the QoS parameters of the interval I_t , such as accuracy, SLA violation rate, energy consumption and average response time. This is typically in the form of a convex combination of energy consumption, response time, SLO violation rates, etc. [50], [52]. The constraints in this formulation include the following. Firstly, the container decomposition for a new task $i \in \mathcal{N}_t$ should be based on d^i . Secondly, containers corresponding to the layer-split decisions $\{C^i | d^i = L\}$ should be scheduled as per the linear chain of precedence constraints. This means that a container later in the neural inference pipeline should be scheduled only after the complete execution of the previous containers in the pipeline. This is because the output of an initial layer in an inference pipeline of a neural network is required before we can schedule a latter layer in the pipeline. Thirdly, the placement matrix $P_t : C_t \times \mathcal{H}$ should adhere to the allocation constraints, *i.e.*, it should not allocate/migrate a container to a worker where the worker does not have sufficient resources available to accommodate the container. Thus, the problem can be formulated as

$$\begin{aligned}
 & \underset{P_t, \mathcal{D}_t}{\text{maximize}} && \sum_t^T O_t \\
 & \text{subject to} && \forall t, \forall i \in \mathcal{N}_t, C^i \text{ containers created} \\
 & && \text{based on splitting decision } d^i, \\
 & && \forall t, P_t \text{ is feasible,} \\
 & && \forall d^i = L, C^i \text{ follow precedence chain.}
 \end{aligned} \tag{1}$$

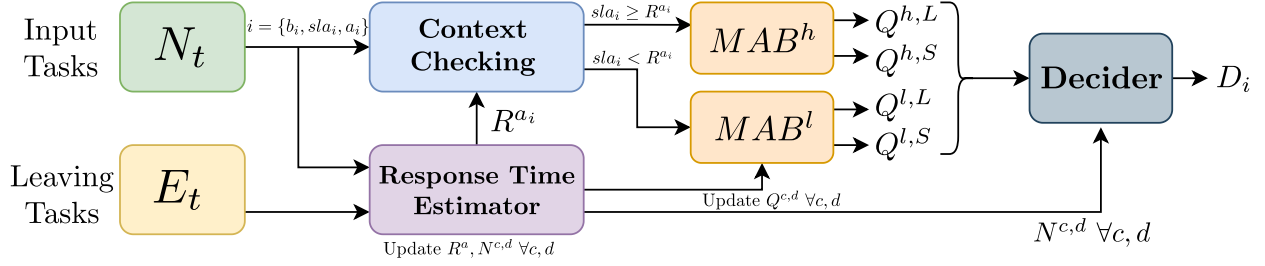


Figure 5: MAB decision workflow.

4 SPLITPLACE POLICY

We now describe the SplitPlace decision and placement policy. For the first sub-problem of deciding optimal splitting strategy, we employ a Multi-Armed Bandit model to dynamically enforce the decision using external reward signals¹. Our solution for the second sub-problem of split placement uses a reinforcement-learning based approach that specifically utilizes a surrogate model to optimize the placement decision (agnostic to the specific implementation)². This two-stage approach is suboptimal since the response time of the splitting decision depends on the placement decision. In case of large variation in terms of the computational resources, it is worth exploring joint optimization of both decisions. However, in our large-scale edge settings, this segregation helps us to make the problem tractable as we describe next.

The motivation behind this segregation is two-fold. First, having a single reinforcement-learning (RL) model that takes both splitting and placement decisions makes the state-space explode exponentially, causing memory bottlenecks in resource-constrained edge devices [55]. Having a simple RL model does not allow it to scale well with several devices in modern IoT settings (see Section 6 with 50 edge devices and the Gillis RL baseline). One of the solutions that we explore in this work is to simplify this complex problem by decomposing it into split decision making and task placement. Second, the response time of an application depends primarily on the splitting choice, layer or semantic, making it a crucial factor for SLA deadline based decision making. To minimize the SLA violation rates we only use the response time based context for our Multi-Armed bandit model. Other parameters like CPU or RAM utilization have high variability in a volatile setting and are not ideal choices for deciding which splitting strategy to opt. Instead, the inference accuracy is another key factor in taking this decision. Thus, SLA violation and inference accuracy are apt objectives for the first sub-problem. Further, the energy consumption and average response time largely depend

1. Compared to other methods like A/B testing and Hill Climbing search [53], Multi-Armed Bandits allow quick convergence in scenarios when different cases need to be modelled separately, which is the case in our setup. Thus, we use Multi-Armed Bandits for deciding the optimal splitting strategy for an input task.

2. In contrast to Monte Carlo or Evolutionary methods, Reinforcement learning allows placement to be goal-directed, *i.e.*, aims at optimizing QoS using it as a signal, and allows the model to adapt to changing environments [54]. Hence, we use a RL model, specifically using a surrogate model due to its scalability, to decide the optimal task placement of network splits.

on the task placement, making them an ideal objective for optimization in the task placement sub-problem.

4.1 Multi-Armed Bandit Decision Module

Multi-Armed Bandit, in short MAB, is a policy formulation where a state-less agent is expected to take one of many decisions with each decision leading to a different reward. The objective of such an agent is to maximize the expected long-term reward [31]. However, in our case, the most important factor to consider when making a decision of whether to use layer or semantic splits for a task is its SLA deadline.

4.1.1 Estimating Response Time of Layer-Splits

The idea behind the proposed SplitPlace approach is to maintain MABs for two different contexts: 1) when SLA is greater than the estimate of the response time for a layer decision, 2) when SLA is less than this estimate. The motivation behind these two contexts is that in case of the SLA deadline being lower than the execution time of layer split, a "layer" decision would be more likely to violate the SLA as result delivery would be after the deadline. However, the exact time it takes to completely execute all containers corresponding to the layer split decision is a priori unknown. Thus, for every application type, we maintain estimates of the response time, *i.e.*, the total time it takes to execute all containers corresponding to this decision.

Let us denote the tasks leaving the system at the end of I_t as E_t . Now, for each task $i \in E_t$, we denote response time and inference performance using r_i and p_i . We denote the layer response time estimate for application $a \in \mathcal{A}$ as R^a . To quickly adapt to non-stationary scenarios, for instance due to the mobility of edge nodes in the system, we update our estimates using new data-points as exponential moving averages using the multiplier $\phi \in [0, 1]$ for the most recent response time observation. Moving averages presents a low computational cost and consequently low latency compared to more sophisticated smoothing functions.

$$R^a \leftarrow \phi \cdot r_i + (1 - \phi) \cdot R^a, \forall i \in E_t \wedge d^i = L, \forall a \in \mathcal{A}. \quad (2)$$

Compared to simple moving average, the above equation gives higher weights to the latest response times, allowing the model to quickly respond to recent changes in environment and workload characteristics.

4.1.2 Context based MAB Model

Now, for any input task $i \in \mathcal{N}_t$, we divide it into two cases: $sla_i \geq R^{a_i}$ and $sla_i < R^{a_i}$. Considering that the response

time of a semantic-split decision would likely be lower than the layer-split decision, in the first case both decisions would most likely not lead to an SLA violation (*high SLA* setting). However, in the second case, a layer-split decision would likely lead to an SLA violation but not the semantic-split decision (*low SLA* setting). To tackle the problem for these different contexts, we maintain two independent MAB models denoted as MAB^h and MAB^l . The former represents a MAB model for the high-SLA setting and the latter for the low-SLA setting.

For each context and decision $d \in \{L, S\}$, we define reward metrics as

$$O^{h,d} = \frac{\sum_{i \in E_t} (\mathbb{1}(r_i \leq sla_i) + p_i) \cdot \mathbb{1}(sla_i \geq R^{a_i} \wedge d^i = d)}{2 \cdot \sum_{i \in E_t} \mathbb{1}(sla_i \geq R^{a_i} \wedge d^i = d)}, \quad (3)$$

$$O^{l,d} = \frac{\sum_{i \in E_t} (\mathbb{1}(r_i \leq sla_i) + p_i) \cdot \mathbb{1}(sla_i < R^{a_i} \wedge d^i = d)}{2 \cdot \sum_{i \in E_t} \mathbb{1}(sla_i < R^{a_i} \wedge d^i = d)}. \quad (4)$$

The first term of the numerator, *i.e.*, $\mathbb{1}(r_i \leq sla_i)$ quantifies SLA violation reward (one if not violated and zero otherwise). The second term, *i.e.*, p_i corresponds to the inference accuracy of the task. These two objectives have been motivated at the start of Section 4. Thus, each MAB model gets the reward function for its decisions allowing independent training of the two. The weights of the two metrics, *i.e.*, accuracy and SLA violation can be set by the user to modify the relative importance between the metrics as per application requirements. In our experiments, the weight parameters of both metrics are set to be equal based on grid-search, maximizing the average reward.

Now, for each decision context $c \in \{h, l\}$ and $d \in \{L, S\}$, we maintain a decision count $N^{c,d}$ and a reward estimate $Q^{c,d}$ which is updated using the reward functions $O^{h,d}$ or $O^{l,d}$ as follows

$$Q^{c,d} \leftarrow Q^{c,d} + \gamma(O^{c,d} - Q^{c,d}), \quad \forall d \in \{L, S\}, \quad \forall c \in \{h, l\}. \quad (5)$$

where γ is the decay parameter. Thus, each reward-estimate is updated by the corresponding reward metric.

For both these MAB models, we use a parameter-free feedback-based ϵ -greedy learning approach that is known to be versatile in adapting to diverse workload characteristics [56]. Unlike other strategies, this is known to scale asymptotically as the long-term Q estimates become exact under mild conditions [57, § 2.2]. To train the model, we take the contextual decision

$$d^i = \begin{cases} \left\{ \begin{array}{ll} \text{random decision,} & \text{with prob. } \epsilon \\ \arg \max_{d \in \{L, S\}} Q^{h,d}, & \text{otherwise} \end{array} \right. & , \quad sla_i \geq R^{a_i} \\ \left\{ \begin{array}{ll} \text{random decision,} & \text{with prob. } \epsilon \\ \arg \max_{d \in \{L, S\}} Q^{l,d}, & \text{otherwise} \end{array} \right. & , \quad sla_i < R^{a_i} \end{cases}. \quad (6)$$

Here, the probability ϵ decays using the reward feedback, starting from 1. We maintain a reward threshold ρ that is initialized as a small positive constant $k < 1$, and use aver-

age reward $O^{MAB} = \frac{1}{4} \sum_{c \in \{h, l\}} \sum_{d \in \{L, S\}} O^{c,d}$ to update ϵ and ρ using the rules

$$\epsilon \leftarrow \begin{cases} \text{decay}(\epsilon), & O^{MAB} > \rho \\ \epsilon, & \text{otherwise} \end{cases}, \quad (7)$$

$$\rho \leftarrow \begin{cases} \text{increment}(\rho), & O^{MAB} > \rho \\ \rho, & \text{otherwise} \end{cases}. \quad (8)$$

Here $\text{decay}(\epsilon) = (1 - k) \cdot \epsilon$ and $\text{increment}(\rho) = (1 + k) \cdot \rho$. Note that $O^{MAB} > \rho$ refers to the current value of ρ prior to the update. The k value controls the rate of convergence of the model. The ϵ value controls the exploration of the model at training time allowing the model to visit more states and obtain precise estimates of layer-split response times.

However, at test time we already have precise estimates of the response times; thus exploration is only required to adapt in volatile scenarios. For this, ϵ -greedy is not a suitable approach as decreasing ϵ with time would prevent exploration as time progresses. Instead, we use an Upper-Confidence-Bound (UCB) exploration strategy that is more suitable as it takes decision counts also into account [58], [59]. Thus, at test time, we take a deterministic decision using the rule

$$d^i = \begin{cases} \arg \max_{d \in \{L, S\}} Q^{h,d} + c \sqrt{\frac{\log t}{N^{h,d}}}, & sla_i \geq R^{a_i} \\ \arg \max_{d \in \{L, S\}} Q^{l,d} + c \sqrt{\frac{\log t}{N^{l,d}}}, & sla_i < R^{a_i} \end{cases}, \quad (9)$$

where t is the scheduling interval count and c is the exploration factor. An overview of the complete split-decision making workflow is shown in Figure 5. We now discuss the RL based placement module. It is worth noting that both components are independent of each other and can be improved separately in future; however, our experiments show that the MAB decision module accounts for most of the performance gains (see Section 6).

4.2 Reinforcement Learning based Placement Module

Once we have the splitting decision for each input task $i \in \mathcal{N}_t$, we now can create containers using pre-trained layer and semantic split neural networks corresponding to the application a_i . This can be done offline on a resource rich system, where the split models can be trained using existing datasets. Once we have trained models, we can generate container images corresponding to each split fragment and distribute to all worker nodes [60]. Then we need to place containers, translating to the worker node initializing a container for the corresponding image, of all active tasks C_t to workers \mathcal{H} . To do this, we use a learning model which predicts the placement matrix P_t using the state of the system S_t , decisions $\mathcal{D}_t = d^i \forall i \in \mathcal{T}_t$ and a reward signal O^P . To define O^P , we define the following metrics [50]:

- 1) *Average Energy Consumption* (AEC) is defined for any interval I_t as the mean energy consumption of all edge workers in the system.
- 2) *Average Response Time* (ART) is defined for any interval I_t as mean response time (in scheduling intervals) of all leaving tasks E_t .

The choice of these two objectives for the placement sub-problem has been motivated at the start of Section 4. Using these metrics, for any interval I_t , O^P is defined as

$$O^P = O^{MAB} - \alpha \cdot AEC_t - \beta \cdot ART_t. \quad (10)$$

Algorithm 1 SplitPlace Decision and Placement Module**Require:**

-
- Pre-trained MAB models MAB^h, MAB^l
 Discounting factor $\gamma \in (0, 1)$
- 1: **procedure** SPLITPLACE(scheduling interval I_t)
 - 2: Get new tasks \mathcal{N}_t
 - 3: Get leaving tasks E_t
 - 4: Calculate $O^{c,d} \forall c \in \{h, l\}, d \in \{L, S\}$ using (3) and (4)
 - 5: Update gain estimates $Q^{c,d} \forall c, d$ using (5) ▷ Q update
 - 6: Update decision counts $N^{c,d} \forall c, d$ ▷ Count update
 - 7: **for** $i \in \mathcal{N}_t$ **do**
 - 8: $i = \{b_i, sla_i, a_i\}$
 - 9: $d^i \leftarrow \begin{cases} \arg \max_{d \in \{L, S\}} Q^{h,d} + c\sqrt{\frac{\log t}{N^{h,d}}}, & sla_i \geq R^{a_i} \\ \arg \max_{d \in \{L, S\}} Q^{l,d} + c\sqrt{\frac{\log t}{N^{l,d}}}, & sla_i < R^{a_i} \end{cases}$ ▷ UCB based decision
 - 10: $\mathcal{D}_t \leftarrow \{d^i\}_{i \in \mathcal{N}_t}$
 - 11: Get S_t from resource-monitor ▷ Current State
 - 12: $P_t \leftarrow \text{DASO}([S_t, P_{t-1}, \mathcal{D}_t])$ ▷ Placement Decision
 - 13: $O^{MAB} \leftarrow \frac{1}{4} \sum_{c \in \{h, l\}} \sum_{d \in \{L, D\}} O^{c,d}$
 - 14: Fine-tune DASO using O^P calculated using (10)
 - 15: **return** P_t
-

Here, α and β (such that $\alpha + \beta = 1$) are hyper-parameters that can be set by users as per the application requirements. Higher α aims to optimize energy consumption at the cost of higher response times, whereas low α aims to reduce average response time. Thus, a RL model f , parameterized by θ takes a decision P_i , where the model uses the reward estimate as the output of the function $f([S_t, P_t, \mathcal{D}_t]; \theta)$, where the parameters θ are updated based on the reward signal O^P . We call this learning approach “decision-aware” as part of the input is the split-decision taken by the MAB model.

Clearly, the proposed formulation is agnostic to the underlying implementation of the learning approach. Thus, any policy like Q-learning or Actor-Critic Learning could be used in the SplitPlace model [57]. However, recently developed techniques like GOBI [50] use gradient-based optimization of the reward to quickly converge to a local-maximum of the objective function. GOBI uses a neural-network based surrogate model to estimate the reward from a given input state, which is then used to update the state by calculating the gradients of the reward estimates with respect to the input. Moreover, advances like momentum, annealing and restarts allow such models to quickly reach a global optima [50].

DASO placement module. In the proposed framework, we use decision-aware surrogate based optimization method (termed as DASO) to place containers in a distributed mobile edge environment. This is motivated from prior neural network based surrogate optimization methods [50]. Here, we consider a Fully-Connected-Network (FCN) model $f(x; \theta)$ that takes an x as a tuple of input state S_t , split-decision \mathcal{D}_t and placement decision P_t , and outputs an estimate of the QoS objective score O_t . This is because FCNs are agnostic to the structure of the input and hence a suitable choice for modeling dependencies between QoS metrics and model inputs like resource utilization and placement decision [49], [50]. Exploration of other styles

of neural models, such as graph neural networks that can take the network topology graph as an input are part of future work. Now, using existing execution trace dataset, $\Lambda = \{[S_t, P_t, \mathcal{D}_t], O_t\}_b$, the FCN model is trained to optimize its network parameters θ such that the Mean-Square-Error (MSE) loss

$$\mathcal{L}(f(x; \theta), y) = \frac{1}{b} \sum_{t=0}^b (y - f(x; \theta))^2, \text{ where } (x, y) \in \Lambda. \quad (11)$$

is minimized as in [50]. To do this, we use AdamW optimizer [61] and update θ up till convergence. This allows the surrogate model f to predict an QoS objective score for a given system state S_t , split-decisions \mathcal{D}_t and task placement P_t . Once the surrogate model is trained, starting from the placement decision from the previous interval $P_t = P_{t-1}$, we leverage it to optimize the placement decision using the following rule

$$P_t \leftarrow P_t - \eta \cdot \nabla_{P_t} f([S_t, P_t, \mathcal{D}_t]; \theta), \quad (12)$$

for a given state and decision pair S_t, \mathcal{D}_t . Here, η is the learning rate of the model. The above equation is iterated till convergence, *i.e.*, the L_2 norm between the placement matrices of two consecutive iterations is lower than a threshold value. Thus, at the start of each interval I_t , using the output of the MAB decision module, the DASO model gives us a placement decision P_t .

4.3 SplitPlace Algorithm

An overview of the SplitPlace approach is given in Algorithm 1. Using pre-trained MAB models, *i.e.*, Q-estimates $Q^{c,d}$ and decision counts $N^{c,d}$, the model decides the optimal splitting decision using the UCB metric (line 9). To adapt the model in non-stationary scenarios, we dynamically update the Q-estimates and decision counts (lines 5 and 6). Using the current state and the split-decisions of all active tasks, we use the DASO approach to take a placement

Table 3: Edge Worker characteristics of Azure Edge Environment.

Name	Qty	Core count	MIPS	RAM	RAM Bandwidth	Ping time	Network Bandwidth	Disk Bandwidth	Cost Model
Worker Nodes									
B2ms	20	2	4029	4295 MB	372 MB/s	2 ms	1000 MB/s	13.4 MB/s	0.0944 \$/hr
E2asv4	10	2	4019	4172 MB	412 MB/s	2 ms	1000 MB/s	10.3 MB/s	0.148 \$/hr
B4ms	10	4	8102	7962 MB	360 MB/s	3 ms	2500 MB/s	10.6 MB/s	0.189 \$/hr
E4asv4	10	4	7962	7962 MB	476 MB/s	3 ms	2500 MB/s	11.64 MB/s	0.296 \$/hr
Broker Node									
L8sv2	1	8	16182	17012 MB	945 MB/s	1 ms	4000 MB/s	17.6 MB/s	0.724 \$/hr

decision for the active containers (line 12). Again, we fine-tune the DASO’s surrogate model using the reward metric to adapt to changes in the environment, for instance the changes in the latency of mobile edge nodes and their consequent effect on the reward metrics (line 14). However, the placement decision must conform to the allocation constraints as described in Section 3.2. To relax the constraint of having only feasible placement decisions, in SplitPlace we allocate or migrate only those containers for which it is possible. Those containers that could not be allocated in a scheduling interval are placed to nodes corresponding to the highest output of the neural network f . If no worker placement is feasible the task is added to a wait queue, which are considered again for allocation in the next interval.

5 IMPLEMENTATION

To implement and evaluate the SplitPlace policy, we need a framework that we can use to deploy containerized neural network split fragments on an edge computing environment. One such framework is COSCO [50]. It enables the development and deployment of integrated edge-cloud environments with structured communication and platform independent execution of applications. It connects various IoT sensors, which can be healthcare sensors with gateway devices, to send data and tasks to edge computing nodes, including edge or cloud workers. The resource management and task initiation is undertaken on edge nodes in the broker layer. The framework uses HTTP RESTful APIs for communication and seamlessly integrates a Flask based web-environment to deploy and manage containers in a distributed setup [62].

We use only the edge-layer deployment in the framework and use the Docker container engine to containerize and execute the split-neural networks in various edge workers [60]. We use the Checkpoint/Restore In Userspace (CRIU) [63] tool for container migration. Further, the DASO approach is implemented using the Autograd package in the PyTorch module [64].

To implement SplitPlace in the COSCO framework, we extend the Framework class to allow constraints for sequential execution of layer-splits. The function `getPlacementPossible()` was modified to also check for containers of layer-split partitioning scheme to be scheduled sequentially. Moreover, we implemented data transferring pipeline for broadcasting inputs in semantic-split decision and forwarding the outputs in layer-split decision. Finally, the inference outputs were synchronized and brought

to the broker to calculate the performance accuracy and measure the workflow response time. For synchronization of outputs and execution of network splits, we use the HTTP Notification API.

6 PERFORMANCE EVALUATION

To test the efficacy of the SplitPlace approach and compare it against the baseline methods, we perform experiments on a heterogeneous edge computing testbed. To do this we emulate a setting with mobile edge devices mounted on self-driving cars, that execute various image-recognition tasks.

6.1 Experiment Setup

As in prior work [50], [49], [65], we use $\alpha = \beta = 0.5$ in (10) for our experiments (we consider other value pairs in Appendix A.2). Also, we use the exploration factor $c = 0.5$ for the UCB exploration and the exponential moving average parameter $\phi = 0.9$, chosen using grid-search using the cumulative reward as the metric to maximize. We create a testbed of 50 resource-constrained VMs located in the same geographical location of London, United Kingdom using Microsoft Azure. The worker resources are shown in Table 3. All machines use Intel i3 2.4 GHz processor cores with processing capacity of no more than a Raspberry Pi 4B device. To keep storage costs consistent, we keep Azure P15 Managed disk with 125 MB/s disk throughput and 256 GB size³. The worker nodes have 4-8 GB of RAM, whereas the broker has 16 GB RAM. To factor in the mobility of the edge nodes, we use the NetLimiter tool to tweak the communication latency with the broker node using the mobility model described in [66]. Specifically, we use the latency and bandwidth parameters of workers from the traces generated using the Simulation of Urban Mobility (SUMO) tool [67] that emulates mobile vehicles in a city like environment. SUMO gives us the parameters like ping time and network bandwidth to simulate in our testbed using NetLimiter. The moving averages and periodic fine-tuning allow our approach to be robust towards any kind of dynamism in the edge environment, including the one arising from mobility of worker nodes.

Our Azure environment is such that all devices are in the same LAN with 10 MBps network interface cards to avoid

³. Azure Managed Disks <https://docs.microsoft.com/en-us/azure/virtual-machines/disks-types#premium-ssd>.

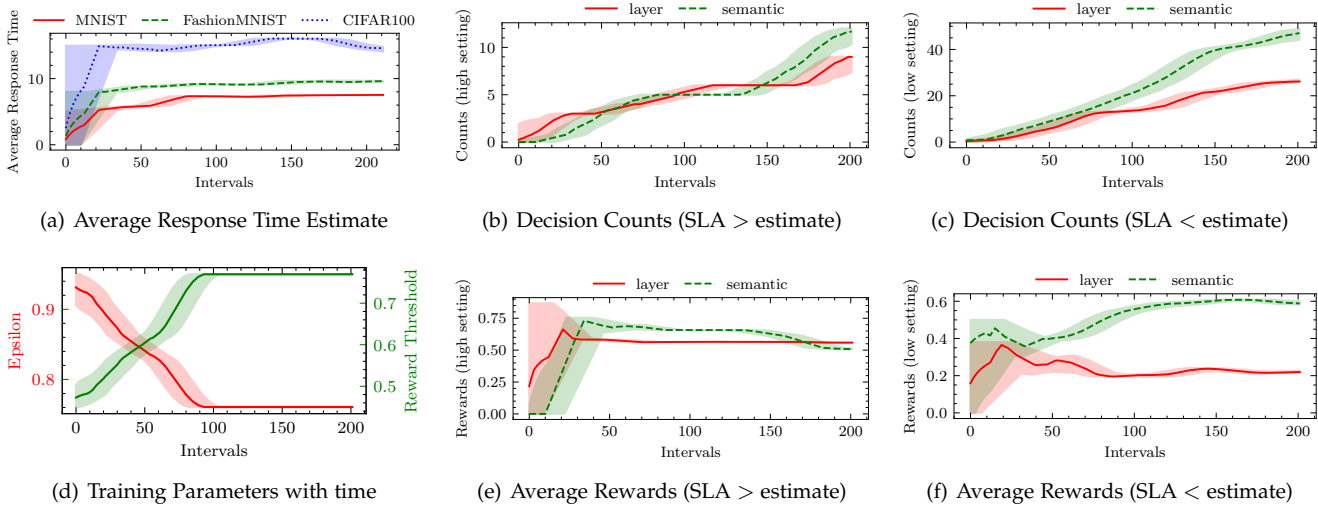


Figure 6: MAB training curves

network bottlenecks while transferring inputs, outputs and intermediate results across neural network splits. Even if the nodes are in the same LAN with high bandwidth connections, the SUMO model would emulate the affects of mobility as is common in prior work on mobile edge computing [65], [68]. Further, we use the `cPickle`⁴ Python module to save the intermediate results using `bzip2` compression and `rsync`⁵ file-transfer utility to minimize the communication latency. For containers corresponding to a layer-split workload that are deployed in different nodes, the intermediate results are forwarded using the `scp` utility to the next container in the neural network pipeline. Similarly, for semantic splitting, the `cPickle` outputs are collected using `rsync` and concatenated using the `torch.cat` function.

We use the Microsoft Azure pricing calculator to obtain the cost of execution per hour (in US Dollars)⁶. The power consumption models are taken from the Standard Performance Evaluation Corporation (SPEC) benchmarks repository⁷. The Million-Instruction-per-Second (MIPS) of all VMs are computed using the `perf-stat`⁸ tool on the SPEC benchmarks. We run all experiments for 100 scheduling intervals, *i.e.*, $\Gamma = 100$, with each interval being 300 seconds long, giving a total experiment time of 8 hours 20 minutes. We average over five runs and use diverse workload types to ensure statistical significance in our experiments. We consider variations of the experimental setup in Appendix A.3.

6.2 Workloads

Motivated from prior work [32], we use three families of popular DNNs as the benchmarking models: ResNet50-

V2 [69], MobileNetV2 [70] and InceptionV3 [71]. Each family has many variants of the model, each having a different number of layers in the neural model. For instance, the ResNet model has 34 and 50 layers. We use three image-classification data sets: MNIST, FashionMNIST and CIFAR100 [33], [34], [35]. MNIST is a hand-written digit recognition dataset with 28×28 gray-scale images to 10-dimensional output. FashionMNIST has 28×28 RGB images with 10 dimensional output. CIFAR100 has 32×32 RGB images with 100-dimensional output. Thus the application set \mathcal{A} becomes $\{\text{MNIST, FashionMNIST, CIFAR100}\}$. These models have been taken directly from the AIOtBench workloads [72]. This is a popular suite of AI benchmark applications for IoT and Edge computing solutions. The three specific datasets used in our experiments are motivated from the vertical use case of self-driving cars, which requires DNN-based applications to continuously recognize images with low latency requirements. Herein, an image recognition software is deployed that reads speed signs (digit recognition, MNIST), recognizes humans (through apparel and pose [73], FashionMNIST), identifies other objects like cars and barriers (object detection, CIFAR100). We use the implementation of neural network splitting from prior work [32], [16].

We use the inference deadline from the work [32] as our SLA. To create the input tasks, we use batch sizes sampled uniformly from 16,000 – 64,000. At the beginning of each scheduling interval, we create $Poisson(\lambda)$ tasks with $\lambda = 6$ tasks for our setup, sampled uniformly from one of the three applications [50]. We consider other λ values and single workload type (from MNIST, FashionMNIST and CIFAR100) in Appendices A.1 and A.4. The split fragments for MNIST, FashionMNIST and CIFAR100 lead to container images of sizes 8-14 MB, 34-56 MB and 47-76 MB, respectively. To calculate the inference accuracy to feed in the MAB models and perform UCB exploration, we also share the ground-truth labels of all datasets with all worker nodes at the time of sharing the neural models as Docker container images. We also compare edge and cloud setups

4. `cPickle` module <https://docs.python.org/2/library/pickle.html#module-cPickle>.

5. `rsync` tool <https://linux.die.net/man/1/rsync>.

6. Microsoft Azure pricing calculator for South UK <https://azure.microsoft.com/en-gb/pricing/calculator/>.

7. SPEC benchmark repository https://www.spec.org/cloud_iaas2018/results/.

8. `perf-stat` tool <https://man7.org/linux/man-pages/man1/perf-stat.1.html>.

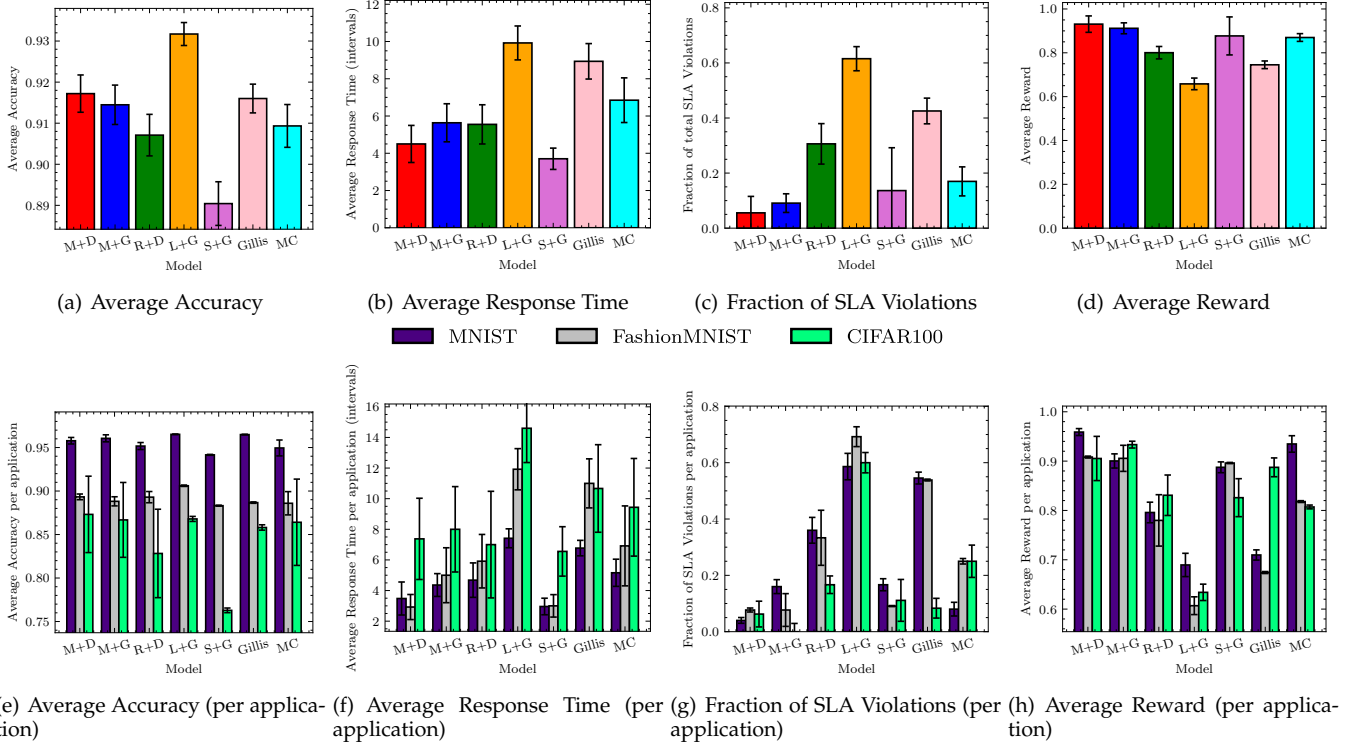


Figure 7: Comparison of SplitPlace against baselines and ablated models on physical setup with 50 edge workers

in Appendix A.5 to establish the need for edge devices for latency critical workloads.

6.3 MAB Training

To train our MAB models, we execute the workloads on the test setup for 200 intervals and use feedback-based ϵ -greedy exploration to update the layer-split decision response time estimates, Q -estimates and decision counts. Figure 6 shows the training curves for the two models.

Figure 6(a) shows how the response time estimates for the layer-split decision are learned starting from zero using moving averages. Figure 6(d) shows how the reward-threshold ρ and decay parameter ϵ change with time. We use the decay and increment multipliers as 0.9 and 1.1 ($k = 0.1$ in (7)) for ϵ and ρ respectively, as done in [56]. Figures 6(b) and 6(c) show the decision counts for high and low SLA settings for both decisions. Figures 6(e) and 6(f) show the Q -estimates for high and low SLA settings. The dichotomy between the two settings is reflected here. When the sla_i of the input task i is less than the estimate R^{a_i} (low setting) there is a clear distinction between the rewards of the two decisions as layer-split is likely to lead to SLA violation and hence lower rewards. However, when sla_i is greater than the estimate R^{a_i} (high setting), both decisions give relatively high rewards with layer-split decision slightly surpassing the semantic-split due to higher average accuracy as discussed in Section 2.

The feedback-based ϵ -greedy training allows us to obtain close estimates of the average response times of the layer-split executions for each application type and average rewards for both decisions in high and low SLA settings. Thus, in our experiments, we initialize the expected reward

(Q) and layer-split response time (R) estimates by the values we get from this training approach. At test time, we dynamically update these estimates using (2) and (5).

6.4 Performance Metrics

We use the following evaluation metrics in our experiments as motivated from prior works [2], [50], [49]. We also use AEC and ART as discussed in Section 4.

- 1) *Average Accuracy* is defined for an execution trace as the average accuracy of all tasks run in an experiment, *i.e.*,

$$Accuracy = \frac{\sum_t \sum_{i \in E_t} p_i}{\sum_t |E_t|}. \tag{13}$$

- 2) *Fraction of SLA violation* is defined for an execution trace as the fraction of all tasks run in an experiment for which the response time is higher than the SLA deadline, *i.e.*,

$$SLA\ Violations = \frac{\sum_t \sum_{i \in E_t} \mathbb{1}(sla_i \geq r_i)}{\sum_t |E_t|}. \tag{14}$$

- 3) *Average Reward* is defined for an execution trace as follows

$$Reward = \frac{\sum_t \sum_{i \in E_t} \mathbb{1}(sla_i \geq r_i) + p_i}{2 \cdot \sum_t |E_t|}. \tag{15}$$

- 4) *Execution Cost* is defined for an execution trace as the total cost incurred during the experiment, *i.e.*,

$$Cost = \sum_{h \in \mathcal{H}} \int_x C_h(x) dx. \tag{16}$$

where $C_h(x)$ is the cost function for worker h with time.

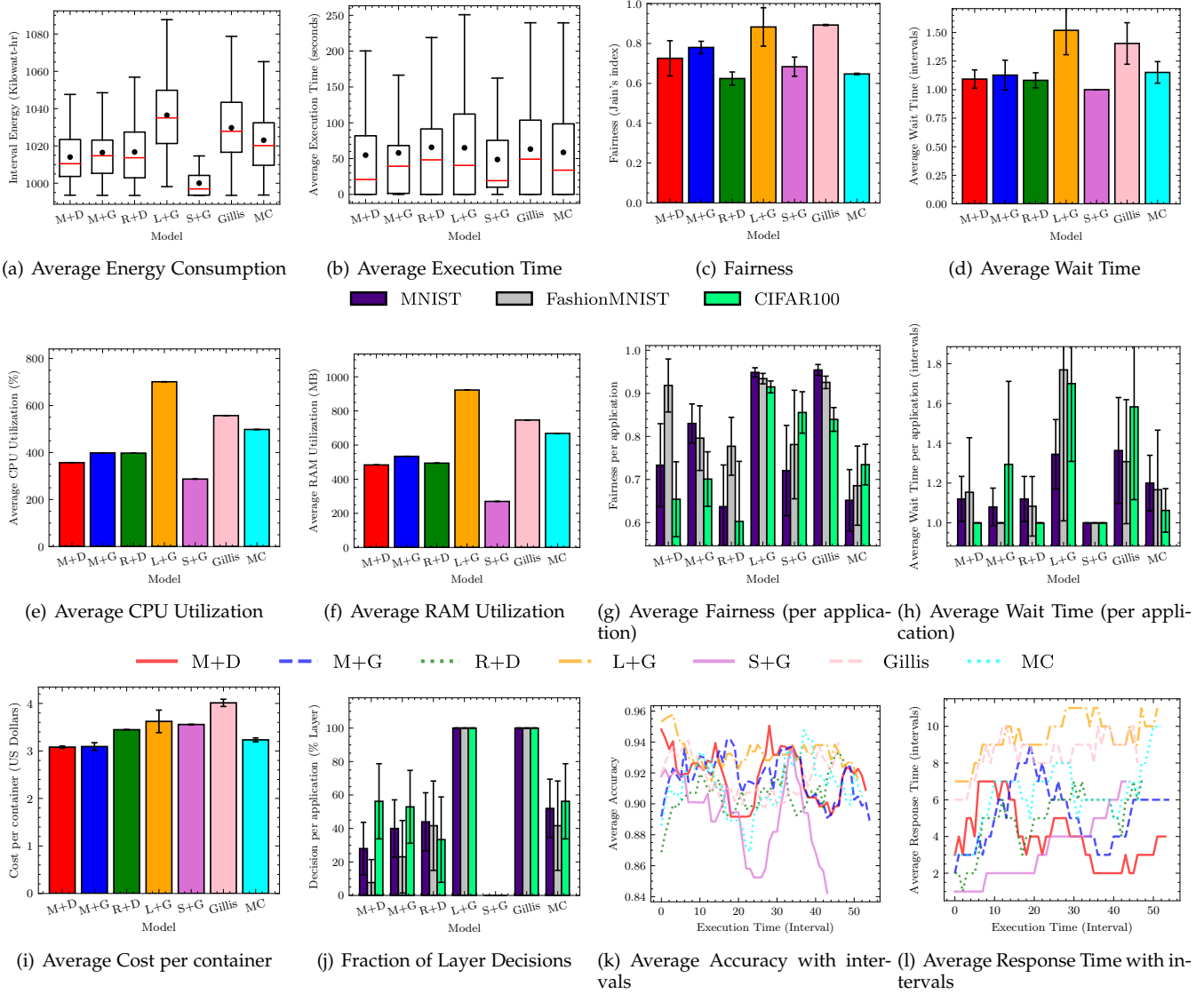


Figure 8: Additional results comparing SplitPlace with baselines and ablated models.

- 5) *Average Wait Time* is the average time a task had to wait in the wait queue till it could be allocated to a worker for execution.
- 6) *Average Execution Time* is the response time minus the wait time, averaged for all tasks run in an experiment.
- 7) *Fairness* is defined as the Jain’s fairness index for execution on tasks over the edge workers [50].

6.5 Baselines and Ablated Models

We compare the performance of the SplitPlace approach against the state-of-the-art baselines *Gillis* and *BottleNet++* Model Compression (denoted as *MC* in our graphs) [32], [37], [38]. *Gillis* refers to the reinforcement learning method proposed in [32] that leverages both layer-splitting and compression models to achieve optimal response time and inference accuracy. Note that contrary to the original *Gillis*’ work, our implementation does not leverage serverless functions. *MC* is a model-compression approach motivated from *BottleNet++* that we implement using the `PyTorch`

`Prune` library.⁹ Further details in Section 2. We do not include results for other methods discussed in Section 2 as *MC* and *Gillis* give better results empirically for all comparison metrics. We also compare SplitPlace with ablated models, where we replace one or both of the MAB or DASO components with simpler versions as described below.

- *Semantic+GOBI (S+G)*: Semantic-split decision only with vanilla GOBI placement module.
- *Layer+GOBI (L+G)*: Layer-split decision only with vanilla GOBI placement module.
- *Random+DASO (R+D)*: Random split decision with DASO placement module.
- *MAB+GOBI (M+G)*: MAB based split decider with vanilla GOBI placement module.

The final SplitPlace approach is represented as *MAB+DASO* or *M+D* in shorthand notation in the graphs. These ablated baselines help us determine the relative im-

⁹ `PyTorch Prune`. https://pytorch.org/docs/stable/generated/torch.nn.utils.prune.ln_structured.html. Accessed 10 October 2021.

Table 4: Comparison of SplitPlace with baseline and ablated models. The best achieved value for each metric is shown in bold. Units: Energy (MW-hr), Scheduling Time (seconds), Fairness (Jain’s Index), Wait Time (Intervals), Response Time (Intervals).

Model	Energy	Scheduling Time	Fairness	Wait Time	Response Time	SLA Violations	Accuracy	Average Reward
Baselines								
Model Compression	1.1368	8.84±0.02	0.65±0.01	1.15±0.09	6.85±1.20	0.26±0.02	89.93	83.98
Gillis	1.1442	8.22±0.01	0.89±0.03	1.40±0.18	8.39±0.95	0.22±0.03	91.90	84.17
Ablation								
Semantic+GOBI	1.1112	8.68±0.03	0.68±0.04	1.08±0.00	3.70±0.57	0.14±0.04	89.04	83.91
Layer+GOBI	1.1517	8.72±0.01	0.88±0.03	1.52±0.21	9.92±0.91	0.62±0.07	93.17	64.87
Random+DASO	1.1297	8.86±0.01	0.62±0.05	1.00±0.07	5.55±1.05	0.29±0.09	90.71	81.62
MAB+GOBI	1.1290	9.12±0.02	0.78±0.08	1.13±0.13	5.64±1.02	0.10±0.03	91.45	90.18
SplitPlace Model								
MAB+DASO	1.0867	9.32±0.02	0.73±0.01	1.09±0.08	4.50±1.00	0.08±0.02	92.72	94.18

improvements in performance by the two components of MAB and DASO separately.

6.6 Results and Ablation Analysis

We now provide comparative results showing the performance of the proposed SplitPlace approach against the baseline models and argue the importance of the MAB and decision-aware placement using ablation analysis. We train the *GOBI* and *DASO* models using the execution trace dataset used to train the MAB models. The learning rate (η) was set to 10^{-3} from [50].

Figure 7 shows the average reward and related performance metrics, *i.e.*, accuracy, response time and SLA violation rate. As expected, the L+G policy gives the highest accuracy of 93.17% as all decisions are layer-wise only with a higher inference performance than semantic-split execution. The S+G policy gives the least accuracy of 89.04%. However, due to layer-splits only the L+G policy also has the highest average response time, subsequently giving the highest SLA violation rate. On the other hand, S+G policy has the least average response time. However, due to the intelligent decision making in SplitPlace, it is able to get the highest total reward of 0.9418. Similar trends are also seen when comparing across models for each application. The accuracy is the highest for the MNIST dataset and lowest for CIFAR100. Average response time is highest for the CIFAR100 and lowest for MNIST in general. Among the baselines, the Gillis approach has the lowest SLA violation rate of 22% and SplitPlace improves upon this by giving 14% lower SLA violations (only 8%). Gillis has higher accuracy between the baselines of 91.9%, with SplitPlace giving an average improvement of 0.82%. Overall, the total reward of SplitPlace is higher than the baselines by at least 10.1%, giving the reward of 94.18%.

Figure 8 shows the performance of all models for other evaluation metrics like energy, execution time and fairness. Compared to the baselines, SplitPlace can reduce energy consumption by up to 4.41% – 5.03% giving an average energy consumption of 1.0867 MW-hr. However, the SplitPlace approach has higher scheduling time and lower fairness index (Table 4). The Gillis baseline has the highest fairness

index of 0.89, however this index for SplitPlace is 0.73. SplitPlace has a higher overhead of 11.8% compared to the Gillis baseline in terms of scheduling time. Figure 8(i) compares the average execution cost (in USD) for all models. As SplitPlace is able to run the maximum number of containers in the 100 intervals, it has the least cost of 3.07 USD/container. The main advantage of SplitPlace is the intelligent splitting decisions facilitate overcoming the memory bottlenecks in edge environments, giving up to 32% lower RAM utilization compared to Gillis and Model Compression.

In terms of the initial communication time of the Docker container images, the SplitPlace method takes 30 seconds at the start of an execution. Gillis and MC have such communication times of 20 and 18 seconds, respectively. This demonstrates that SplitPlace has a low one-time overhead (up to 12 seconds) compared to the baselines when compared to the gains in response time (up to 46%) that linearly scales as the number of workloads increase.

A summary of comparisons with values of main performance metrics for all models is given in Table 4. The best values achieved for each metric are highlighted in bold.

7 CONCLUSIONS

In this work, we present SplitPlace, a novel framework for efficiently managing demanding neural network based applications. SplitPlace exploits the trade-off between layer and semantic split models where the former gives higher accuracy, but the latter gives much lower response times. This allows SplitPlace to not only manage tasks to maintain high inference accuracy on average, but also reduce SLA violation rate. The proposed model uses a Multi-Armed-Bandits based policy to decide which split strategy to use according to the SLA deadline of the incoming task. Moreover, it uses a decision-aware learning model to take appropriate placement decisions for those neural fragments on mobile edge workers. Further, both MAB and learning models are dynamically tuned to adapt to volatile scenarios. All these contributions allow SplitPlace to out-perform the baseline models in terms of average response time, SLA violation rate, inference accuracy and total reward by up to 46.3%,

69.2%, 3.1% and 12.1% respectively in a heterogeneous edge environment with real-world workloads.

We propose the following future directions for this work. An extension of the current work may be developed that dynamically updates the splitting configuration to adapt to more heterogeneous and non-stationary edge environments [74]. Moreover, the current model assumes that all neural models are divisible into independent layers. This may be hard for deep learning models like attention based neural networks or transformer models [75]. Finally, the model only considers splits and their placement as containers, more fine-grained methods involving Neural Architecture Search and cost efficient deployment methods may be explored like serverless frameworks [76]. Other considerations such as privacy concerns and non-stationary number of active edge nodes with extreme levels of heterogeneity such that the placement decision has a significant impact on response time is also part of future work.

SOFTWARE AVAILABILITY

The code is available at <https://github.com/imperial-qore/SplitPlace>. The Docker images used in the experiments are available at <https://hub.docker.com/u/shreshthtuli>.

ACKNOWLEDGMENTS

Shreshth Tuli is grateful to the Imperial College London for funding his Ph.D. through the President's Ph.D. Scholarship scheme. We thank Feng Yan for helpful discussions.

REFERENCES

- [1] S. Tuli, "SplitPlace: Intelligent Placement of Split Neural Nets in Mobile Edge Environments," *SIGMETRICS Perform. Eval. Rev.*, 2021.
- [2] S. S. Gill, S. Tuli, M. Xu, I. Singh, K. V. Singh, D. Lindsay, S. Tuli, D. Smirnova, M. Singh, U. Jain *et al.*, "Transformative effects of IoT, Blockchain and Artificial Intelligence on cloud computing: Evolution, vision, trends and open challenges," *Internet of Things*, vol. 8, pp. 100–118, 2019.
- [3] H. Zhu, M. Akrouf, B. Zheng, A. Pelegris, A. Jayarajan, A. Phanishayee, B. Schroeder, and G. Pekhimenko, "Benchmarking and analyzing deep neural network training," in *2018 IEEE International Symposium on Workload Characterization (IISWC)*. IEEE, 2018, pp. 88–100.
- [4] E. Li, L. Zeng, Z. Zhou, and X. Chen, "Edge AI: On-demand accelerating deep neural network inference via edge computing," *IEEE Transactions on Wireless Communications*, vol. 19, no. 1, pp. 447–457, 2019.
- [5] A. Khanna, A. Sah, and T. Choudhury, "Intelligent mobile edge computing: A deep learning based approach," in *International Conference on Advances in Computing and Data Sciences*. Springer, 2020, pp. 107–116.
- [6] F. A. Kraemer, A. E. Braten, N. Tamkittikhun, and D. Palma, "Fog computing in healthcare—a review and discussion," *IEEE Access*, vol. 5, pp. 9206–9222, 2017.
- [7] L. Zhang and L. Zhang, "Deep learning-based classification and reconstruction of residential scenes from large-scale point clouds," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 4, pp. 1887–1897, 2017.
- [8] M. Roopaei, P. Rad, and M. Jamshidi, "Deep learning control for complex and large scale cloud systems," *Intelligent Automation & Soft Computing*, vol. 23, no. 3, pp. 389–391, 2017.
- [9] J. R. Gunasekaran, C. S. Mishra, P. Thinakaran, M. T. Kandemir, and C. R. Das, "Implications of public cloud resource heterogeneity for inference serving," in *Proceedings of the 2020 Sixth International Workshop on Serverless Computing*, 2020, pp. 7–12.
- [10] S. Laskaridis, S. I. Venieris, M. Almeida, I. Leontiadis, and N. D. Lane, "Spinn: synergistic progressive inference of neural networks over device and cloud," in *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking*, 2020, pp. 1–15.
- [11] Q. Liang, P. Shenoy, and D. Irwin, "AI on the edge: Characterizing AI-based IoT applications using specialized edge architectures," in *2020 IEEE International Symposium on Workload Characterization (IISWC)*. IEEE, 2020, pp. 145–156.
- [12] N. Abbas, Y. Zhang, A. Taherkordi, and T. Skeie, "Mobile edge computing: A survey," *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 450–465, 2017.
- [13] Y. Mao, J. Zhang, and K. B. Letaief, "Dynamic computation offloading for mobile-edge computing with energy harvesting devices," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 12, pp. 3590–3605, 2016.
- [14] J. Shao and J. Zhang, "Communication-computation trade-off in resource-constrained edge inference," *IEEE Communications Magazine*, vol. 58, no. 12, pp. 20–26, 2020.
- [15] J. Liu, M. L. Curry, C. Maltzahn, and P. Kufeldt, "Scale-out edge storage systems with embedded storage nodes to get better availability and cost-efficiency at the same time," in *3rd {USENIX} Workshop on Hot Topics in Edge Computing (HotEdge 20)*, 2020.
- [16] J. Kim, Y. Park, G. Kim, and S. J. Hwang, "Splitnet: Learning to semantically split deep networks for parameter reduction and model parallelization," in *Proceedings of the 34th International Conference on Machine Learning—Volume 70*. JMLR. org, 2017, pp. 1866–1874.
- [17] Y. Shi, K. Yang, T. Jiang, J. Zhang, and K. B. Letaief, "Communication-efficient edge ai: Algorithms and systems," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 4, pp. 2167–2191, 2020.
- [18] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, "Federated learning in mobile edge networks: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2031–2063, 2020.
- [19] Y. Siriwardhana, P. Porambage, M. Liyanage, and M. Ylianttila, "A Survey on Mobile Augmented Reality With 5G Mobile Edge Computing: Architectures, Applications, and Technical Aspects," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 2, pp. 1160–1192, 2021.
- [20] S. Tuli, N. Basumatary, S. S. Gill, M. Kahani, R. C. Arya, G. S. Wander, and R. Buyya, "Healthfog: An ensemble deep learning based smart healthcare system for automatic diagnosis of heart diseases in integrated iot and fog computing environments," *Future Generation Computer Systems*, vol. 104, pp. 187–200, 2020.
- [21] W. Xiao, R. Bhardwaj, R. Ramjee, M. Sivathanu, N. Kwatra, Z. Han, P. Patel, X. Peng, H. Zhao, Q. Zhang *et al.*, "Gandiva: Introspective cluster scheduling for deep learning," in *13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*, 2018, pp. 595–610.
- [22] A. Gujarati, R. Karimi, S. Alzayat, W. Hao, A. Kaufmann, Y. Vigfusson, and J. Mace, "Serving dnns like clockwork: Performance predictability from the bottom up," in *14th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 20)*, 2020, pp. 443–462.
- [23] J. Chen and X. Ran, "Deep learning with edge computing: A review." *Proceedings of the IEEE*, vol. 107, no. 8, pp. 1655–1674, 2019.
- [24] A. Capotondi, M. Rusci, M. Fariselli, and L. Benini, "Cmix-nn: Mixed low-precision cnn library for memory-constrained edge devices," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 67, no. 5, pp. 871–875, 2020.
- [25] J. Huang, C. Samplawski, D. Ganesan, B. Marlin, and H. Kwon, "CLIO: Enabling automatic compilation of deep learning pipelines across IoT and Cloud," in *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking*, 2020, pp. 1–12.
- [26] Q. Le, L. Miralles-Pechuán, S. Kulkarni, J. Su, and O. Boydell, "An overview of deep learning in industry," *Data Analytics and AI*, pp. 65–98, 2020.
- [27] Y. Matsubara, S. Baidya, D. Callegaro, M. Levorato, and S. Singh, "Distilled split deep neural networks for edge-assisted real-time systems," in *Proceedings of the 2019 Workshop on Hot Topics in Video Analytics and Intelligent Edges*, 2019, pp. 21–26.
- [28] Y. A. Ushakov, P. N. Polezhaev, A. E. Shukhman, M. V. Ushakova, and M. Nadezhda, "Split neural networks for mobile devices," in

- 2018 26th Telecommunications Forum (TELFOR). IEEE, 2018, pp. 420–425.
- [29] V. S. Gordon and J. Crouson, "Self-splitting modular neural network-domain partitioning at boundaries of trained regions," in *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*. IEEE, 2008, pp. 1085–1091.
- [30] E. Ahmed and M. H. Rehmani, "Mobile edge computing: Opportunities, solutions, and challenges," pp. 59–63, 2017.
- [31] D. Bouneffouf, I. Rish, and C. Aggarwal, "Survey on applications of multi-armed and contextual bandits," in *2020 IEEE Congress on Evolutionary Computation (CEC)*. IEEE, 2020, pp. 1–8.
- [32] M. Yu, Z. Jiang, H. C. Ng, W. Wang, R. Chen, and B. Li, "Gillis: Serving large neural networks in serverless functions with automatic model partitioning," in *International Conference on Distributed Computing Systems*, 2021.
- [33] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [34] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [35] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [36] A. Goli, O. Hajjehassani, H. Khazaei, O. Ardakanian, M. Rashidi, and T. Dauphinee, "Migrating from monolithic to serverless: A fintech case study," in *Companion of the ACM/SPEC International Conference on Performance Engineering*, 2020, pp. 20–25.
- [37] A. E. Eshratifar, A. Esmaili, and M. Pedram, "Bottlenet: A deep learning architecture for intelligent mobile cloud computing services," in *2019 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*. IEEE, 2019, pp. 1–6.
- [38] J. Shao and J. Zhang, "Bottlenet++: An end-to-end approach for feature compression in device-edge co-inference systems," in *2020 IEEE International Conference on Communications Workshops (ICC Workshops)*. IEEE, 2020, pp. 1–6.
- [39] S. Gao, F. Huang, J. Pei, and H. Huang, "Discrete model compression with resource constraint for deep neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1899–1908.
- [40] S. Teerapittayanon, B. McDanel, and H.-T. Kung, "Distributed deep neural networks over the cloud, the edge and end devices," in *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2017, pp. 328–339.
- [41] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang, "Neurosurgeon: Collaborative intelligence between the cloud and mobile edge," *ACM SIGARCH Computer Architecture News*, vol. 45, no. 1, pp. 615–629, 2017.
- [42] S. Zhang, S. Zhang, Z. Qian, J. Wu, Y. Jin, and S. Lu, "DeepSlicing: collaborative and adaptive cnn inference with low latency," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 9, pp. 2175–2187, 2021.
- [43] L. Deng, G. Li, S. Han, L. Shi, and Y. Xie, "Model compression and hardware acceleration for neural networks: A comprehensive survey," *Proceedings of the IEEE*, vol. 108, no. 4, pp. 485–532, 2020.
- [44] P. S. Chandakkar, Y. Li, P. L. K. Ding, and B. Li, "Strategies for re-training a pruned neural network in an edge computing paradigm," in *2017 IEEE International Conference on Edge Computing (EDGE)*. IEEE, 2017, pp. 244–247.
- [45] L. Wang, L. Xiang, J. Xu, J. Chen, X. Zhao, D. Yao, X. Wang, and B. Li, "Context-aware deep model compression for edge cloud computing," in *2020 IEEE 40th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2020, pp. 787–797.
- [46] F. Yu, L. Cui, P. Wang, C. Han, R. Huang, and X. Huang, "Easiedge: A novel global deep neural networks pruning method for efficient edge computing," *IEEE Internet of Things Journal*, vol. 8, no. 3, pp. 1259–1271, 2020.
- [47] A. Kaplunovich and Y. Yesha, "Automatic tuning of hyperparameters for neural networks in serverless cloud," in *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 2020, pp. 2751–2756.
- [48] S. Tuli, R. Mahmud, S. Tuli, and R. Buyya, "Fogbus: A blockchain-based lightweight framework for edge and fog computing," *Journal of Systems and Software*, vol. 154, pp. 22–36, 2019.
- [49] D. Basu, X. Wang, Y. Hong, H. Chen, and S. Bressan, "Learn-as-you-go with megh: Efficient live migration of virtual machines," *IEEE Transactions on Parallel and Distributed Systems*, vol. 30, no. 8, pp. 1786–1801, 2019.
- [50] S. Tuli, S. R. Poojara, S. N. Srirama, G. Casale, and N. R. Jennings, "COSCO: Container Orchestration Using Co-Simulation and Gradient Based Optimization for Fog Computing Environments," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 1, pp. 101–116, 2022.
- [51] S.-H. Park, O. Simeone, and S. Shamai, "Joint cloud and edge processing for latency minimization in fog radio access networks," in *2016 IEEE 17th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 2016, pp. 1–5.
- [52] R. A. da Silva and N. L. da Fonseca, "Resource allocation mechanism for a fog-cloud infrastructure," in *2018 IEEE International Conference on Communications (ICC)*. IEEE, 2018, pp. 1–6.
- [53] R. Kohavi and R. Longbotham, "Online controlled experiments and a/b testing," *Encyclopedia of machine learning and data mining*, vol. 7, no. 8, pp. 922–929, 2017.
- [54] M. Wiering and M. Van Otterlo, "Reinforcement learning," *Adaptation, learning, and optimization*, vol. 12, no. 3, 2012.
- [55] X. Chen, H. Zhang, C. Wu, S. Mao, Y. Ji, and M. Bennis, "Optimized computation offloading performance in virtual edge computing systems via deep reinforcement learning," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 4005–4018, 2018.
- [56] A. Maroti, "RBED: Reward Based Epsilon Decay," *arXiv preprint arXiv:1910.13701*, 2019.
- [57] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [58] Y. Zhang, P. Cai, C. Pan, and S. Zhang, "Multi-agent deep reinforcement learning-based cooperative spectrum sensing with upper confidence bound exploration," *IEEE Access*, vol. 7, pp. 118 898–118 906, 2019.
- [59] S. Gupta, G. Joshi, and O. Yağın, "Correlated multi-armed bandits with a latent random source," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 3572–3576.
- [60] A. Ahmed and G. Pierre, "Docker container deployment in fog computing infrastructures," in *2018 IEEE International Conference on Edge Computing (EDGE)*. IEEE, 2018, pp. 1–8.
- [61] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *International Conference on Learning Representations*, 2018.
- [62] M. Grinberg, *Flask web development: developing web applications with python*. " O'Reilly Media, Inc.", 2018.
- [63] R. S. Venkatesh, T. Smejkal, D. S. Milojevic, and A. Gavrilovska, "Fast in-memory criu for docker containers," in *The International Symposium on Memory Systems*, 2019, pp. 53–65.
- [64] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.
- [65] S. Tuli, S. Ilager, K. Ramamohanarao, and R. Buyya, "Dynamic scheduling for stochastic edge-cloud computing environments using a3c learning and residual recurrent neural networks," *IEEE Transactions on Mobile Computing*, 2020.
- [66] K. Gilly, S. Alcaraz, N. Akinin, S. Filiposka, and A. Mishev, "Modelling edge computing in urban mobility simulation scenarios," in *2020 IFIP Networking Conference (Networking)*. IEEE, 2020, pp. 539–543.
- [67] D. Krajzewicz, J. Erdmann, M. Behrisch, and L. Bieker, "Recent development and applications of sumo-simulation of urban mobility," *International journal on advances in systems and measurements*, vol. 5, no. 3&4, 2012.
- [68] C. Wang, Y. Xu, J. Zhang, and B. Ran, "Integrated traffic control for freeway recurrent bottleneck based on deep reinforcement learning," *IEEE Transactions on Intelligent Transportation Systems*, 2022.
- [69] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [70] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.
- [71] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 2017.
- [72] C. Luo, F. Zhang, C. Huang, X. Xiong, J. Chen, L. Wang, W. Gao, H. Ye, T. Wu, R. Zhou *et al.*, "AIoT bench: towards comprehensive

- benchmarking mobile and embedded device intelligence,” in *International Symposium on Benchmarking, Measuring and Optimization*. Springer, 2018, pp. 31–35.
- [73] K. Hoang, “Image classification with fashion-mnist and cifar10,” *California State University*, 2018.
- [74] K. A. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. M. Kiddon, J. Konečný, S. Mazzocchi, B. McMahan, T. V. Overveldt, D. Petrou, D. Ramage, and J. Roselander, “Towards federated learning at scale: System design,” in *SysML 2019*, 2019.
- [75] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in Neural Information Processing Systems*, vol. 30, pp. 5998–6008, 2017.
- [76] G. Casale, M. Artáč, W.-J. van den Heuvel, A. van Hoorn, P. Jakovits, F. Leymann, M. Long, V. Papanikolaou, D. Presentza, A. Russo *et al.*, “Radon: rational decomposition and orchestration for serverless computing,” *SICS Software-Intensive Cyber-Physical Systems*, vol. 35, no. 1, pp. 77–87, 2020.

APPENDIX A ADDITIONAL EXPERIMENTS

A.1 Sensitivity with λ

We perform sensitivity analysis with the number of workloads being generated at the start of each scheduling interval to test the scalability of the SplitPlace model. As described in Section 6.2, we create *Poisson*(λ) tasks at the beginning of each scheduling interval with $\lambda = 6$ for our setup. The tasks are sampled uniformly from the three applications: MNIST, FashionMNIST and CIFAR100 as done in our experiments in Section 6. We vary the λ parameter from 2 to 50 to test how the performance varies with this parameter. Beyond $\lambda = 50$, the baseline models diverge and the system breaks down due to computational overload for these models. Figure 9 demonstrates the average accuracy, response time, fraction of SLA violations, reward and energy consumption for all models with varying the λ parameter. Again, all experiments are performed for 100 scheduling intervals with the same setup as described in Section 6.1. Clearly, as λ increases, the number of workloads and hence the resource requirements increase.

The average accuracy observed with policies that take the splitting decision based on external reward signals (M+D, M+G, MC and Gillis) drops as the λ parameter increases. This is not seen for policies R+D, L+G or S+G, as their splitting decision (which affects the classification accuracy) are unaffected. The other policies see a gradual decline of accuracy by 2.61% – 4.83% as the ratio of semantic-split decisions increases to reduce resource consumption (see Figure 11). However, we observe that the response times increase with λ . This is due to the increasing wait times and contention affects in congested edge environments. Still, the SplitPlace model (M+D) maintains a low average response time of 4.76 intervals, 33.61% – 45.31% lower than the baseline models (MC and Gillis). All models see a sharp increase in violation rate as λ increases to 30 tasks in terms of SLA violations. The L+G, MC and Gillis models show average violation rates of nearly 1 for $\lambda = 50$. However, the average slope is least for the SplitPlace model, increasing from 8% to 35.1% (increase of 27%), whereas for others this change is much higher 34.87% – 76.46%. Naturally, as SLA violation rates increase and accuracies decrease, the average reward would decrease (Figure 9(d)). The average

rewards of the M+D model (88.44) are still higher than the baselines and ablated models (64.69 – 85.66). Figure 9(e) shows the change of energy consumption as λ increases. This metric sharply increases as λ changes from 2 tasks to 6 tasks. Upon further increase, the energy does not increase as significantly as edge environment reaches a saturation point and no more tasks can be scheduled. All models have similar trends in terms of the energy consumption metric.

A.2 Sensitivity with α, β

We also perform sensitivity analysis on the α, β parameters introduced in Section 4.2 for the SplitPlace and its ablated versions. α and β are the weights corresponding to the average energy consumption (AEC) and average response time (ART) metrics in Equation 10. Note that it is a convex combination of the two metrics, so $\alpha + \beta = 1$. These experiments help us test the robustness of the SplitPlace model in the presence of diverse QoS objectives.

Figure 10 shows the change in performance metrics when α varies from 0 to 1. Some observations overlap with Figure 9. In terms of energy consumption, all models show a decline in this metric as α increases. This is because of the energy conserving scheduling decisions by the GOBI model, allocating a higher fraction of tasks to the low-power edge nodes (B2ms and E2asv4). This leads to higher resource contention in these devices, leading to higher SLA violations, for which the MAB model tends to take a higher ratio of semantic-split decisions (due to their lower response times, see Figure 12). This leads to the decline in the accuracy of the M+D and M+G models as α increases, whereas the average accuracy for reward-free models (R+D, L+G, S+G) remains nearly unchanged. The average response times and SLA violation rates increase steadily for all models. However, the increase is not as significant for the MAB based models (M+D and M+G) due to the UCB based external signal. The average reward is still the highest for the SplitPlace model (91.93) compares to the ablated models (65.91 – 90.39).

A.3 Constraining Environments

We now perform the same experiments with various constraints in the original edge computing environment. These experiments demonstrate the ability of the MAB based models to adapt to varied deployment scenarios efficiently.

Changes in Setup: For each host in Table 3, we limit the core count (compute constrained), network bandwidth (network constrained) and RAM size (memory constrained) to half the original value¹⁰. These setups show us the performance of the SplitPlace model with respect to the baselines techniques in more diverse experimental scenarios.

Results: Figure 13 shows the performance metrics for the three constrained environments for all models. The primary y -axis gives the value of the metric in the new setting and the secondary y -axis gives us the ratio with respect to metric value in the original setup. Figure 14 shows the distribution of the response time to the average time taken for scheduling, waiting, data transfer, task execution and

10. To limit the CPU core usage in Linux kernel we modify the `grubconfig`. To limit the RAM usage, we use the `ulimit` tool. To limit the network bandwidth, we appropriately modify the parameters of the `Netlimiter` tool.

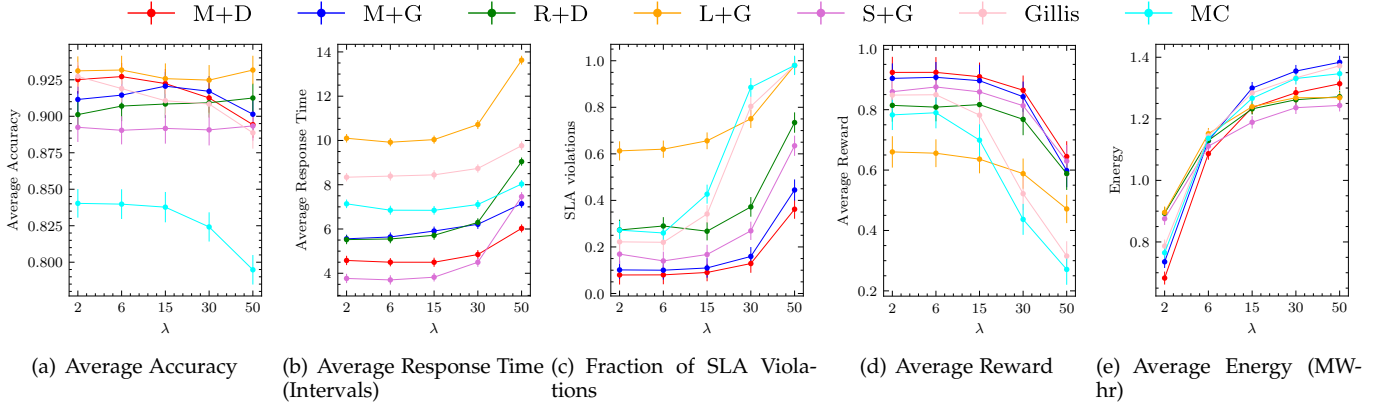


Figure 9: Sensitivity of λ on all models.

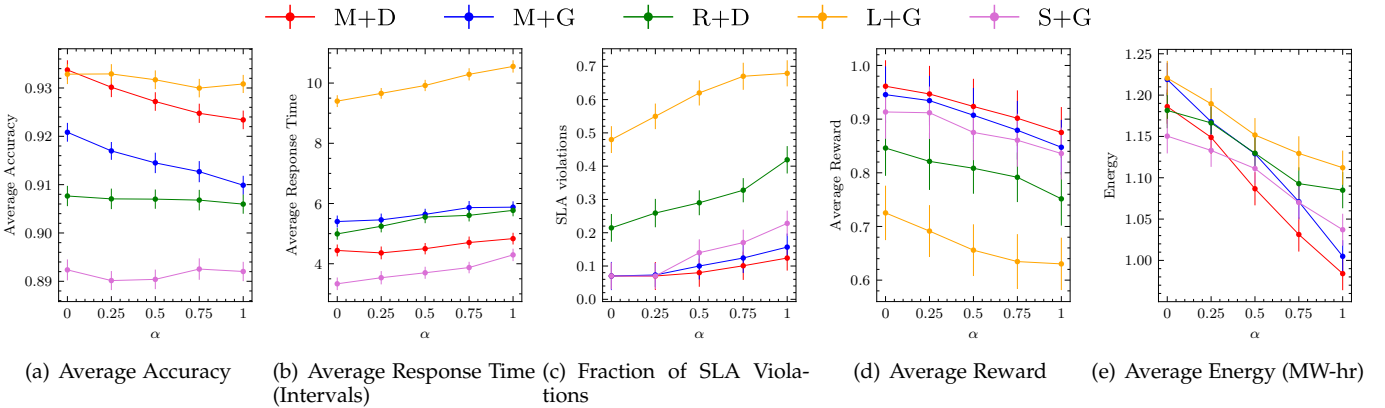


Figure 10: Sensitivity of α, β on ablated SplitPlace models.

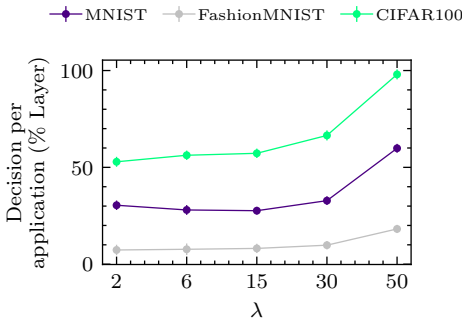


Figure 11: Fraction of layer-split decisions of the SplitPlace model with increasing λ .

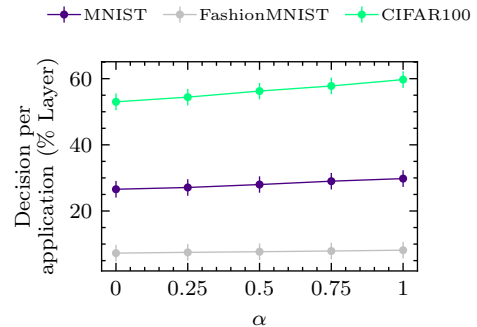


Figure 12: Fraction of layer-split decisions of the SplitPlace model with increasing α .

container migration. Figure 15 shows the distribution of the SLA violations for each application. Similar to what we observed earlier, for reward free models (R+D, L+G and S+G), the average accuracy is nearly unchanged, whereas for M+D, M+G, MC and Gillis, the accuracy drop is up to 40%. This drop is more pronounced in network and compute constrained environments as compared with the memory constrained one. This is because of the increasing average response times compared to the original setup due to lower computational resource availability, higher migration and data response times or low memory resources. This

leads to higher resource contention in these devices, again for which the MAB model tends to take a higher ratio of semantic-split decisions to avoid too many SLA violations. This leads to the decline in accuracy of the M+D, M+G and baseline models in the constrained environments, whereas the average accuracy for reward free models (R+D, L+G, S+G) remains nearly unchanged. However, the response times and fraction of SLA violations are much higher for the reward free models. The relative increase in the violation rates are much lower for the M+D and M+G policies, thanks to the intelligent adaptability of the MAB model to the

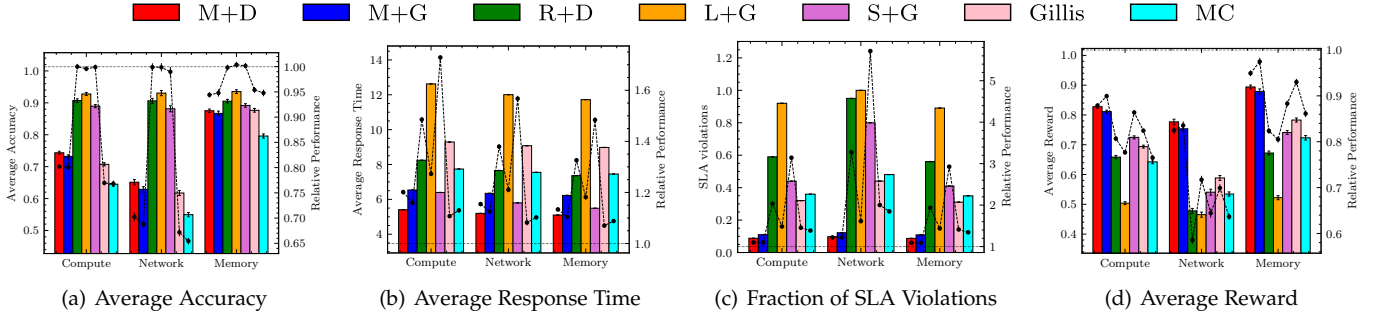


Figure 13: Comparison of SplitPlace against baselines and ablated models in Compute, Network and Memory constrained environments.

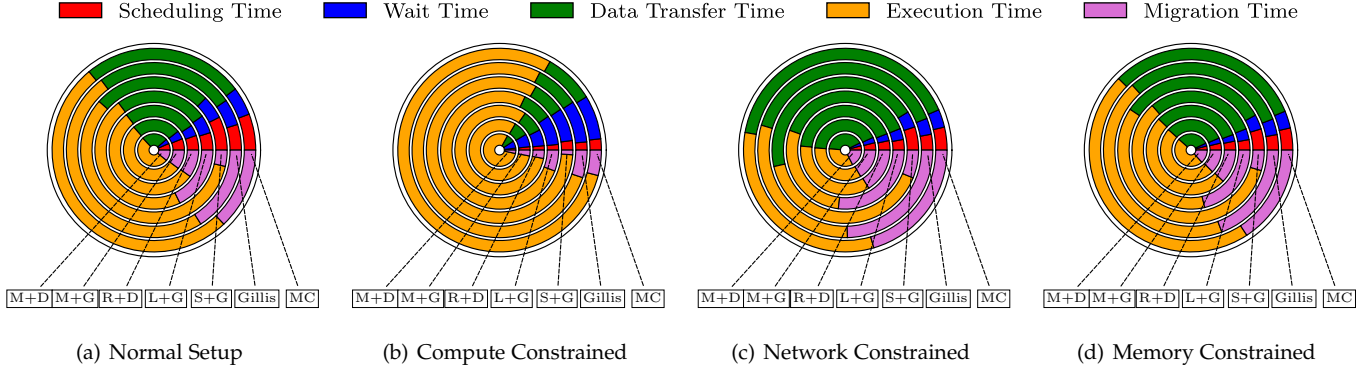


Figure 14: Distribution of response time in Normal and Compute, Network and Memory constrained environments.

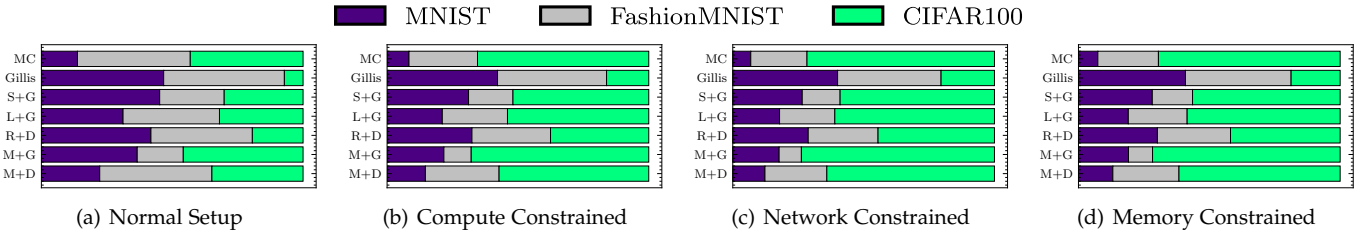


Figure 15: Distribution of SLA violations in Normal and Compute, Network and Memory constrained environments.

changing response times in the new setups. The increase in the response times in the case of computational constraints is predominantly due to the increase in the execution times due to low compute resources available (Figure 14(b)). For network constrained setup, the increase in response time is majorly due to the increase in data transfer and container migration times (due to the lower bandwidth availability, see Figure 14(c)). In memory constrained setup, low memory availability causes the edge devices to use the disk-swap space increasing both execution and data transfer times (Figure 14(d)). These constraints have the highest impact on the resource hungry CIFAR100 workloads which leads to an increase in its SLA violation rate compared to other application types (Figure 15).

The average rewards for all models are lower than the original setup (Figure 13(d)). The relative rewards for MAB based models are the highest (0.84 – 0.95) compared to the baseline models (0.63 – 0.86). The reward drop is the least for the SplitPlace model 11.61%. The reward drops are higher for the reward-free models (14.05% – 22.91), due to

their significantly higher SLA violation rate, even without any accuracy change. For all the three cases, the SplitPlace model has the highest average reward (0.77 – 0.90), whereas the Gillis and MC models have rewards in the range 0.58–0.78 and 0.53–0.71 respectively. The increase in the SLA violation rates are also the lowest for the SplitPlace model (9.24% – 22.01%) compared to the baselines with 34.73% – 103.11%. Similarly, the decrease in the average accuracy for the SplitPlace model is 9.31%, which for the baseline models is 16.98%–32.13%. The average increase in the response time in the constrained setup for the SplitPlace model is 14.56%, whereas for the baselines it is 12.31%.

A.4 Constraining Workloads

We now perform the same experiments, but now changing the workloads to MNIST only, FashionMNIST only and CIFAR100 only applications in lieu of sampling workloads uniformly at random from one of these. This helps us illustrate the robustness of the SplitPlace model against diverse application domains.

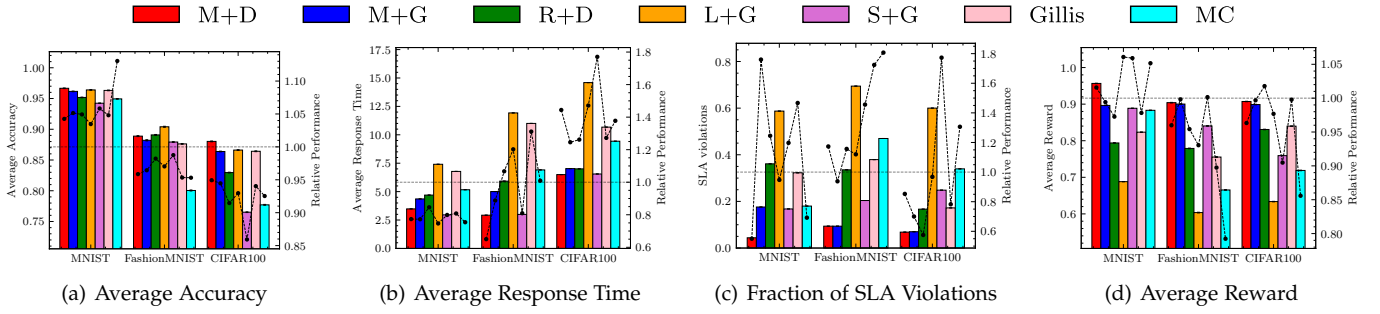


Figure 16: Comparison of SplitPlace against baselines and ablated models in settings with only MNIST, FashionMNIST and CIFAR100 workloads.

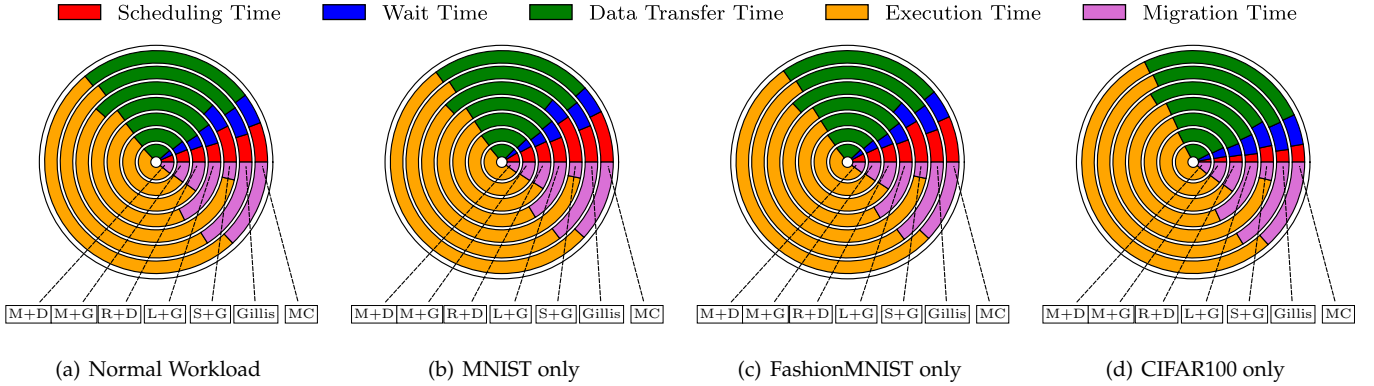


Figure 17: Distribution of response time in Normal and MNIST only, FashionMNIST only, CIFAR100 only workloads.

Results: Figure 16 shows the performance parameters for all models in the three workload settings. The graphs show contrasting workload specific trends. For instance, the accuracy of the MNIST application is higher and response time lower than the other two application types. Thus, for all models, the average accuracy is higher compared to the original workload and response time lower. However, due to lower resource requirements of the MNIST workload, the MC model seldom takes the decision to use compressed models, giving a high accuracy uplift compared to the original setup. Figure 17 shows the distribution of the response times of tasks for each workload setting. There are no drastic changes in the time distribution compared to the original workloads, the only difference being the different execution and data transfer times for various workloads. For smaller workloads like MNIST, the execution time and data transfer times are lower than larger workloads like CIFAR100. Considering the average response time and SLA violation rates, the SplitPlace model gives a lower average violations in all three cases. This is due to the limited workload complexity and simpler decision making for the MAB model.

Even with these constrained workloads, the SplitPlace model has the highest average accuracy with the lowest average response times and SLA violation rates. Overall, we see that the average reward of the SplitPlace model (89.83 – 95.36) is highest among all methods. The baseline models, Gillis and MC have rewards in the range 75.54 – 83.98 and 66.57 – 88.31 respectively.

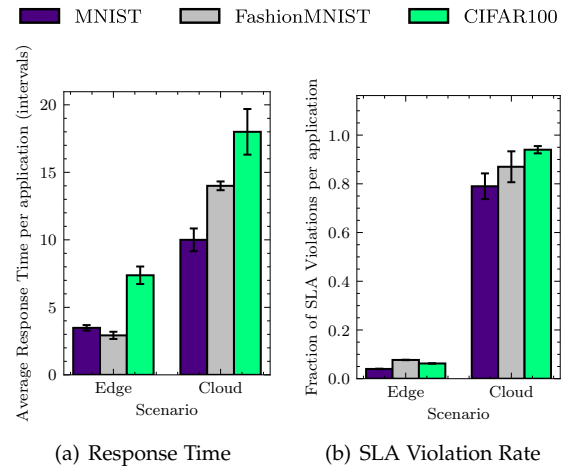


Figure 18: Comparing the response time and SLA violation rates for SplitPlace for edge and cloud setups.

A.5 Comparing Edge with Cloud

We also perform experiments to compare the performance of cloud versus edge platforms. As motivated in Section 1, for latency-critical applications, it is crucial to resort to edge devices that are in close proximity with the users. Relying on cloud nodes, which may be at multi-hop distance from the users lead to higher latency and subsequently the average response times. Thus, for latency critical applications, utilizing memory abundant cloud nodes, with typically 32-64 GB of memory available, is not feasible. This makes the

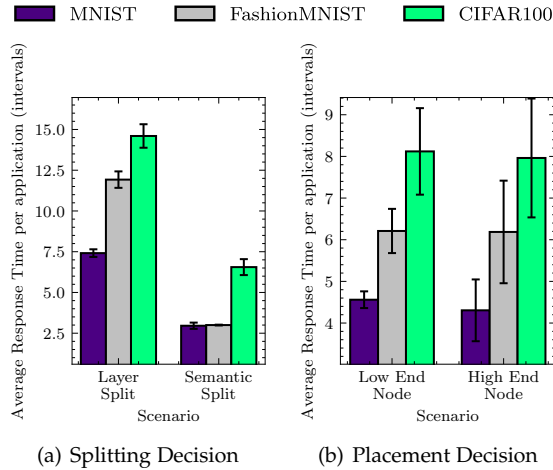


Figure 19: Comparing the response time across splitting and placement decisions.

problem challenging as now we need to ensure that we can run large-scale DNNs on edge, without any cloud backend whatsoever. Thus, in our formulation and experiments, we consider an environment with only memory-constrained edge nodes with 4-8 GB RAM and all in proximity with the users. To empirically verify this, we consider a “Cloud” setup wherein the broker nodes remains in the Azure UK-South datacenter, whereas the worker nodes are now initialized in the East-US datacenter. For a fair comparison, we consider a setup where no splitting is required and end-to-end models are used in the cloud nodes. The response time and SLA violation rates are shown in Figure 18. The figure demonstrates that in a Cloud-AI setup, the same worker nodes are unable to conform to the strict deadlines when compared to the distributed Edge-AI setup. In terms of the initial one-time communication for transferring the containers, SplitPlace method on edge takes 30 seconds, whereas the cloud-only approach takes 72 seconds. The single-hop communication topology among the broker-worker nodes foregoes the need for sending the container images via the Internet.

A.6 Comparing the Impact of Splitting and Placement Decisions on Response Time

The decomposition of decision making in SplitPlace is based on the hypothesis that the MAB models do not need split placement decisions to decide the splitting decisions. The major contributing factor for response time of a task is the split decision type as shown in Figure 2. As the figure demonstrates, the deviation in the response time is much higher when we compare layer vs semantic. The deviation is low when comparing across placement decisions. This is demonstrated by Figure 19. Even though we consider a heterogeneous edge computing environment, with edge devices having different resource capacities, the difference in resources is not sufficient to give significant deviation in terms of response time as per the placement decisions.