


An Offloading Algorithm for Maximizing Inference Accuracy on Edge Device in an Edge Intelligence System

Andrea Fresa 

Edge Networks Group, IMDEA Networks Institute
University Carlos III of Madrid
Madrid, Spain
andrea.fresa@imdea.org

Jaya Prakash Varma Champati 

Edge Networks Group, IMDEA Networks Institute
Madrid, Spain
jaya.champati@imdea.org

ABSTRACT

With the emergence of edge computing, the problem of offloading jobs between an Edge Device (ED) and an Edge Server (ES) received significant attention in the past. Motivated by the fact that an increasing number of applications are using Machine Learning (ML) inference from the data samples collected at the EDs, we study the problem of offloading *inference jobs* by considering the following novel aspects: 1) in contrast to a typical computational job an inference job has *accuracy* measure, 2) both inference accuracy and processing time of an inference job increases with the size of the ML model, and 3) recently proposed Deep Neural Networks (DNNs) for resource-constrained EDs provide the choice of scaling down the model size by trading off the inference accuracy. Therefore, we consider a new system with multiple small-size ML models at the ED and a powerful large-size ML model at the ES and study the problem of offloading inference jobs with the objective of maximizing the total inference accuracy at the ED subject to a time constraint T on the makespan. Noting that the problem is NP-hard, we propose an approximation algorithm: Accuracy Maximization using LP-Relaxation and Rounding (AMR²), and prove that it results in a makespan at most $2T$, and achieves a total accuracy that is lower by a small constant (less than 1) from the optimal total accuracy. As proof of concept, we implemented AMR² on a Raspberry Pi, equipped with MobileNets, that is connected via LAN to a server equipped with ResNet, and studied the total accuracy and makespan performance of AMR² for image classification.



CCS CONCEPTS

• **Theory of computation** → **Scheduling algorithms; Numeric approximation algorithms; Scheduling algorithms.**

KEYWORDS

Approximation ratio algorithm, Machine Learning Inference, Edge Computing

ACM Reference Format:

Andrea Fresa  and Jaya Prakash Varma Champati . 2022. An Offloading Algorithm for Maximizing Inference Accuracy on Edge Device in an Edge Intelligence System. In *Proceedings of the International Conference on Modeling Analysis and Simulation of Wireless and Mobile Systems (MSWiM '22)*, October 24–28, 2022, Montreal, QC, Canada. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3551659.3559044>

1 INTRODUCTION

Edge computing is seen as a key component of future networks that augments the computation, memory, and battery limitations of Edge Devices (EDs) (e.g., IoT devices, mobile phones, etc.), by allowing the devices to offload computational jobs to nearby Edge Servers (ESs) [28]. Since the *offloading decision*, i.e., which jobs to offload, is the key to minimizing the execution delay of the jobs and/or the energy consumption at the ED, it received significant attention in the past [18]. Recently, an increasing number of applications are using Machine Learning (ML) inference from the data samples collected at the EDs, and there is a major thrust for deploying pre-trained Deep Neural Networks (DNNs) on the EDs as this, thus paving the way for Edge Intelligence. Performing inference on an ED has, among other advantages, a reduced latency. Thanks to the development of DNN models with reduced computation and storage requirements (at the cost of reduced inference accuracy) and the advancements in the hardware of EDs [11], ML frameworks such as Tensorflow Lite [1] and PyTorch Mobile [2] are now able to support the deployment of DNNs on EDs. In this context, we study the offloading decision between an ED and an ES for *inference jobs*, where an inference job refers to the execution of a pre-trained ML model on a data sample.

In comparison to the fixed processing time requirement of a generic computational job (typically represented by a directed acyclic task graph), the processing time requirement of an inference job depends on the ML model size: a larger model size results in longer processing time and may provide higher inference accuracy. For example, on Pixel 3 smartphone, ResNet [13] has size 178 MB, requires 526 ms, and provides 76.8% accuracy (top-1 accuracy) for the ImageNet dataset [10], while a small-size DNN model of MobileNet [26] has size 1.9 MB, requires 1.2 ms, but provides 41.4% accuracy [1]. Furthermore, recently developed DNNs for EDs allow for scaling the model size by simply setting a few hyperparameters (cf. [4, 26, 31]), enabling the EDs to choose between multiple model sizes. However, as we explain in Section 2, the offloading decision for inference jobs considering the above novel aspects has received little attention in the literature.

Taking into account the novel aspects for inference jobs, or simply *jobs* in the sequel, we consider the system in Figure 1, where the ED has m ML models to choose from, and the ES is equipped with a state-of-the-art ML model for a given application. Consider that n jobs (corresponding to n data samples) are available at the ED. It may offload them all to the ES to maximize the inference accuracy. However, offloading each job incurs a communication time to upload the data sample in addition to the processing time at the ES. This may result in a large *makespan*, i.e., the total time to finish all

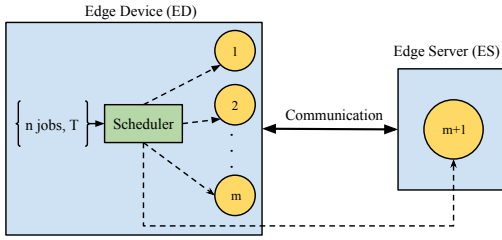


Figure 1: Scheduling inference in an edge intelligence system with an ED and an ES.

the jobs. On the other hand, executing all the jobs on the smallest ML model at the ED may result not only in a smaller makespan, but also the lowest inference accuracy. Thus, a scheduler at the ED needs to strike a trade-off between the accuracy and the makespan. Toward this end, we formulate the following problem: *given n data samples at time zero, find a schedule that offloads a partition of the jobs to the ES and assigns the remaining jobs to m models on the ED, such that the total accuracy is maximized and the makespan is within a time constraint T .* Solution to this problem will be beneficial to applications such as Google Photos where a set of photos selected by a user need to be classified into multiple categories in real time. Also, the problem has relevance to applications which do periodic scheduling, i.e., the ED periodically collects all the data samples arrived in a time period T and aims to finish their processing within the next time period T .

Since the true label provided by a model for a given data sample can only be inferred after the job is executed, for analytical tractability, we consider the top-1 accuracy of the model is its accuracy for any data sample. Given the processing and communication times of the jobs, we formulate the problem as an Integer Linear Program (ILP). We note that the ILP is agnostic to the actual ML models used on the ED and the ES. Different ML models on the ED may correspond to different instantiations of the same DNN with different hyperparameter values (cf. [26]), or they may correspond to different ML algorithms such as DNN, State Vector Machines, Random Forests, etc.

Solving the formulated ILP is challenging due to the following reasons. Partitioning the set of jobs between the ED and the ES is related to scheduling jobs on parallel machines [22], and assigning the jobs to the models on the ED is related to the knapsack problem [16], both are known to be NP-hard. A special case of our problem, where the ED has a single model ($m = 1$), is the Generalized Assignment Problem (GAP) with two machines [24]. GAP is known to be APX-hard¹, and the best-known approximation algorithm provides a solution that has makespan at most $2T$ [29]. However, the algorithms for GAP and their performance guarantees are not directly applicable to our problem due to the additional aspect that, on the ED there are multiple models to choose from. We propose a novel algorithm that solves a Linear Programming relaxation (LP-relaxation) of the ILP, uses a counting argument to bound the number of fractional solutions, and uses a simple rounding rule to round the fractional solution.

Our main contributions are summarized below:

- We formulate the total accuracy maximization problem subject to a constraint T on the makespan as an ILP. Noting that the ILP is NP-hard, we propose an approximation algorithm Accuracy Maximization using LP-Relaxation and Rounding (AMR²) to solve it. The runtime of AMR² is $O(n^3(m+1)^3)$, where n is the number of jobs and m is the number of local ML models.
- We prove that the total accuracy achieved by AMR² is at most a small constant (less than 1) lower than the optimum total accuracy, and its makespan is at most $2T$.
- As proof of concept, we perform experiments using a Raspberry Pi, equipped with MobileNets, and a server, equipped with ResNet50, that are connected over a Local Area Network (LAN). Our application is image classification for the images from ImageNet [10]. We estimate processing and communication times for different sizes of images, and implemented AMR² and a greedy algorithm on Raspberry Pi. Our results indicate that the total test accuracy achieved by AMR² is close to that of the LP-relaxed solution.

The rest of the paper is organized as follows. In Section 2, we present the related work. The system model is presented in Section 3. In Sections 4 and 5 we present AMR² and its performance bounds, respectively. In Section 6, we present the experimental results and finally conclude in Section 7.

2 RELATED WORKS

In this section, we first present the related works for computation offloading problem and then discuss closely related classical job scheduling problems.

2.1 Offloading and ML Inference Jobs

Since the initial proposal of edge computing in [27], significant attention was given to the computational offloading problem, wherein the ED needs to decide which jobs to offload, and how to offload them to an ES [18]. The objectives that were considered for optimizing the offloading decision are, 1) minimize the total execution delay of the jobs, see for example [6, 7, 17, 19, 30], and 2) minimize the energy of the ED spent in computing and/or transmitting the jobs, subject to a constraint on the execution delay, see for example [9, 15, 32]. However, the above works consider generic computation jobs, and the aspect of accuracy, which is relevant for the case of inference jobs, has not been considered.

Recently, a few works considered the problem of maximizing accuracy for inference jobs on the ED [33],[21],[20]. In [33], the authors studied the problem of maximizing the accuracy within a deadline for each frame of a video analytics application. They do not consider offloading to the edge and their solution is tailored to the DNNs that use early exits [31]. A similar problem was studied in [21], where offloading between a mobile device and a cloud is considered. The authors account for the time-varying communication times by using model selection at the cloud and by allowing the duplication of processing the job at the mobile device. A heuristic solution was proposed in [20] for offloading inference jobs for maximizing inference accuracy subject to a maximum energy constraint. In contrast to the above works, we consider multiple models on the ED and provide performance guarantees for AMR².

¹An APX-hard problem has no polynomial-time approximation scheme unless $P = NP$.

2.2 Job Scheduling

As noted in Section 1, our problem is related to the knapsack problem [16]. To see this, note that if it is not feasible to schedule on the ES and all jobs have to be assigned to the ED, then maximizing the total accuracy is equivalent to maximizing profit, and the constraint T is equivalent to the capacity of knapsack. In this case, our problem turns out to be a generalization of the CCKP [12]. Another special case of our problem, where the ED has only a single model, can be formulated as a GAP [5, 24], with two machines. In GAP, n jobs (or items) have to be assigned to r machines (or knapsacks). Each job-machine pair is characterized by two parameters: processing time and cost. The objective is to minimize the total cost subject to a time constraint T on the makespan. It is known that GAP is APX-hard [8].

In their seminal work [29], the authors proposed an algorithm for GAP that achieves minimum total cost and has makespan at most $2T$. Their method involves solving a sequence of LP feasibility problems, in order to tackle the processing times that are greater than T , and compute the minimum total cost using bisection search. Their algorithm can also be used for solving a related extension of GAP, where the cost of scheduling a job on a machine increases linearly with decrease in the processing time of the job. In comparison to this setting, the accuracies (equivalent to negative costs) are not linearly related to the processing times of the jobs and thus the proposed method in [29] is not directly applicable to the problem at hand. Our proposed algorithm AMR² is different from their method in that it does not require to solve LP feasibility problems and the use of bisection search. Further, we prove the performance bounds using a different analysis technique which is based on a counting argument for the LP-relaxation and rounding of the values of the basic variable used to obtain the optimal solution.

3 SYSTEM MODEL

Consider an ED and an ES connected over a network and the ED enlists the help of the ES for computation offloading. At time zero, n inference jobs, each representing the processing requirement of a data sample on a pre-trained ML model, are available to a scheduler at the ED. Let j denote the job index and $J = \{1, 2, \dots, n\}$ denote the set of job indices.

3.1 ML Models and Accuracy

The ED is equipped with m pre-trained ML models, or simply models. Note that these may correspond to m instantiations of the same DNN with different hyperparameter values resulting in different model sizes; see for example [26, 31], or the models may correspond to different ML algorithms such as logistic regression, support vector machines, DNN, etc. Since the ES is a computationally powerful machine, we consider that it is equipped with a state-of-the-art model. We note that our problem formulation and the solution are applicable to any family of ML models deployed on the ED and the ES.

Let $a_i \in [0, 1]$ denote the top-1 accuracy of model i . Since we do not know if a job is classified correctly by a model without first processing it on that model, for analytical tractability, we consider that the accuracy of a model i for any job is a_i . Note that assigning a job to a model with higher top-1 accuracy increases its probability of

correct classification. WLOG, we assume that $a_1 \leq a_2 \leq \dots \leq a_m$, and also assume that the model $m+1$ is a state-of-the-art model with a higher top-1 accuracy than the models on the ED, i.e., $a_m \leq a_{m+1}$. In the sequel, the term ‘accuracy’ refers to the top-1 accuracy, unless otherwise specified.

3.2 Processing and Communication Times

The processing time of job j on model $i \in M \setminus \{m+1\}$ is denoted by p_{ij} , and on model $m+1$ it is denoted by $p'_{(m+1)j}$. In several applications, the data samples may need pre-processing before they are input to the ML model. For example, in computer vision tasks, images require pre-processing and the time required for pre-processing varies with the size of the image [25]. In our experiments with the images from the ImageNet dataset, the pre-processing stage only involves reshaping the images to input to the DNN models. Let τ_{ij} denote the pre-processing time of job j on model i . We consider the pre-processing times are part of the processing times defined above.

Let c_j denote the communication time for offloading job j . It is determined by the data size of the job, i.e., the size of the data sample in bits, and the data rate of the connection between the ED and the ES. Given $p'_{(m+1)j}$ and c_j , the *total time* to process job j on the ES, denoted by $p_{(m+1)j}$, is given by $p_{(m+1)j} = c_j + p'_{(m+1)j}$. We deliberately use similar notation for the processing times p_{ij} on the ED and the total times $p_{(m+1)j}$ on the ES because it simplifies the expressions in the sequel. We consider that the communication times c_j are fixed and are known a priori. This is possible in the scenarios where the ED and the ES are connected in a LAN or in a private network with fixed bandwidth. In our experiments, the ED and the ES are connected via our institute’s LAN, and the communication times have negligible variance. We also consider that the processing times of the jobs are known a priori and that they can be estimated from the historical job executions.

3.3 Optimization Problem

Given the set of jobs J at time zero, the *makespan* is defined as the time when the processing of the last job in J is complete. The objective of the scheduler at the ED is to assign the set of jobs J to the set of models M such that the total accuracy, denote by A , is maximized and the makespan is within the time constraint T . Note that a schedule involves the partitioning of the set J between the ED and the ES, and the constraint on the makespan implies that the completion time of all the jobs should be within T on both ED and ES. The above objective is relevant in applications where the ED periodically collects the data samples in a period T and aims to finish their processing within the next period. By choosing a small period T , a real-time application can aim for fast ML inference at a reduced total accuracy.

Let x_{ij} denote a binary variable such that $x_{ij} = 1$, if the scheduler assigns job j to model i , and $x_{ij} = 0$, otherwise. Note that, if $x_{(m+1)j} = 1$, then job j is offloaded to the ES. Therefore, a schedule is determined by the matrix $\mathbf{x} = [x_{ij}]$. We impose the following constraints on \mathbf{x} :

$$\sum_{i=1}^m \sum_{j=1}^n p_{ij} x_{ij} \leq T \quad (1)$$

$$\sum_{j=1}^n p_{(m+1)j} x_{(m+1)j} \leq T \quad (2)$$

$$\sum_{i=1}^{m+1} x_{ij} = 1, \forall j \in J \quad (3)$$

$$x_{i,j} \in \{0, 1\}, \forall i \in M, \forall j \in J, \quad (4)$$

where constraints (1) and (2) ensure that the total processing times on the ED and the ES, respectively, are within T , and thus the makespan is within T . Constraints in (3) imply that each job is assigned to only one model and no job should be left unassigned, and (4) are integer constraints. We are interested in the following accuracy maximization problem \mathcal{P} :

$$\begin{aligned} & \underset{\mathbf{x}}{\text{maximize}} && A = \sum_{i=1}^{m+1} \sum_{j=1}^n a_{ij} x_{ij} \\ & \text{subject to} && (1), (2), (3), \text{ and } (4). \end{aligned}$$

Note that \mathcal{P} is an ILP. We will show later that a special case of \mathcal{P} reduces to CCKP, which is NP-hard, and thus \mathcal{P} is NP-hard. Let A^* denote the optimal total accuracy for \mathcal{P} .

4 ACCURACY MAXIMIZATION USING LP-RELAXATION AND ROUNDING (AMR²)

In this section, we first present the the LP-relaxation of \mathcal{P} and a result that guides the design of AMR². Later, we present AMR² in Algorithm 1.

4.1 LP-Relaxation

Given \mathcal{P} , we proceed with solving the LP-relaxation of \mathcal{P} , where the integer constraints in (4) are replaced using the following non-negative constraints:

$$x_{ij} \geq 0, \forall i \in M \text{ and } \forall j \in J. \quad (5)$$

Note that, the constraints $x_{ij} \leq 1$ are not required as this is ensured by the constraints in (3). Let the matrix $\bar{x} = [\bar{x}_{ij}]$ and A_{LP}^* denote the schedule and the total accuracy, respectively, output by the LP-relaxation. The LP-relaxed solution provides an upper bound on the total accuracy achieved by an optimal schedule, and thus we have $A_{LP}^* \geq A^*$.

Note that the solution to the LP-relaxation may contain x_{ij} values that are fractional, and the rounding procedure is critical to proving the performance bounds. To design a rounding procedure, we first refer to a key result in [23], where the author studied the problem of assigning N jobs to K parallel machines with the objective of minimizing the makespan. For the LP-relaxation of this problem, the author presented the following counting argument: there exists an optimal basic solution in which there can be at most $K - 1$ fractional jobs, i.e., the jobs that are divided between machines, and all the other jobs are fully assigned. Further, the simplex algorithm outputs such a basic optimal solution. In our problem, there are two parallel machines, the ED and the ES, but in contrast to [23], the ED has multiple models and the jobs assigned to the ED are processed in sequence. Taking this new aspect into account, we extend the counting argument for the problem at hand and show that solving the LP-relaxation of \mathcal{P} results in at most two fractional jobs. This structural result is stated in the following lemma.

LEMMA 1. *For the LP-relaxation of \mathcal{P} , there exists an optimal basic solution with at most two fractional jobs.*

PROOF. Since LP-relaxation of \mathcal{P} has $n+2$ constraints, apart from the non-negative constraints in (5), one can show using LP theory that there exists an optimal basic solution with $n+2$ basic variables that may take positive values and all the non-basic variables take value zero. Under such an optimal basic solution, for the n constraints in (3) to be satisfied, at least one positive basic variable should belong to each of those n constraints. The remaining 2 basic variables may belong to at most two equations. This implies that at least $n-2$ equations should have exactly one positive basic variable whose value should be 1 in order to satisfy the constraint. Therefore, there can be at most two equations with multiple basic variables whose values are in $(0, 1)$, and the two jobs that correspond to these equations are the fractional jobs. \square

Given the basic optimal solution to the LP-relaxation, the result in Lemma 1 reduces the rounding procedure to assigning at most two fractional jobs. WLOG, we re-index the jobs and refer to the fractional jobs by job 1 and job 2. We define the set $I = J \setminus \{1, 2\}$ and refer to the assignment of I under \bar{x} as the *integer solution of the LP-relaxation*. We define:

$$P_1 = \sum_{i=1}^m \sum_{j \in I} p_{i,j} \bar{x}_{i,j}, \quad (6)$$

$$P_2 = \sum_{j \in I} p_{m+1,j} \bar{x}_{m+1,j}. \quad (7)$$

With a slight abuse in notation, we use i_j and k_j to denote the indices of the machines on which the fractional job $j \in \{1, 2\}$ is scheduled. We have

$$\bar{x}_{i_1} + \bar{x}_{k_1} = 1, \quad (8)$$

$$\bar{x}_{i_2} + \bar{x}_{k_2} = 1. \quad (9)$$

4.2 AMR² Description

The main steps of AMR² are summarized in Algorithm 1. In the first step AMR² solves the LP-relaxation. In the second step, if there is one fractional job, it is assigned to model with largest accuracy such that the makespan does not exceed $2T$. If there are two fractional jobs, we use the simple rounding rule, assign the job to the model on which it has higher fraction. Though the algorithm is not sophisticated, we will later see that proving the performance bounds is involved. We use \mathbf{x}^\dagger and A^\dagger to denote the schedule and the total accuracy, respectively, output by AMR².

Computational complexity. The computational complexity of solving an LP with l variables is $O(l^3)$ (cf.[3]). In the LP-relaxation, the number of variables are $n(m+1)$ and thus its runtime is $O(n^3(m+1)^3)$. The rounding technique has negligible complexity when compared to the complexity of solving the LP-relaxation. In conclusion the runtime is $O(n^3(m+1)^3)$.

5 ANALYSIS OF AMR²

In this section, we analyse AMR² and present a $2T$ bound for its makespan and show that its total accuracy is at most $a_{m+1} - a_1$ lower than the optimal accuracy.

Algorithm 1: AMR²

```
1: Input:  $p_{ij}$ , for all  $i \in M$  and  $j \in J$ .
2: Solve the LP-relaxation of  $\mathcal{P}$ .
3: if One fractional job then
4:   if  $P_2 + p_{m+1,1} \leq 2T$  then
5:     Assign job 1 to model  $m + 1$ 
6:   else
7:     Assign job 1 to model
        $\arg \max_{i \in M \setminus \{m+1\}} \{a_i : p_{i1} + P_1 \leq 2T\}$ .
8:   end if
9: end if
10: if Two fractional jobs then
11:   for all  $j \in \{1, 2\}$  do
12:     if  $\bar{x}_{ij} > \bar{x}_{kj}$  then
13:        $x_{ij}^\dagger = 1$ 
14:     else
15:        $x_{kj}^\dagger = 1$ 
16:     end if
17:   end for
18: end if
19: Output: Assignment matrix  $\mathbf{x}^\dagger$  and total accuracy  $A^\dagger$ 
```

THEOREM 5.1. *If \mathcal{P} is feasible, then the makespan of the system under AMR² is at most $2T$.*

PROOF. For the case of one fractional job, the result follows by the construct of AMR². For the case of two fractional jobs, based on the fractional job assignment output by the LP-relaxation solution we consider three cases and for each case, we consider sub-cases based on the schedule of AMR². For the proof, we assume (10) and (11) are true. The proof steps are similar for other cases.

$$\bar{x}_{k_1} < \bar{x}_{i_1}, \quad (10)$$

$$\bar{x}_{i_2} < \bar{x}_{k_2}. \quad (11)$$

Case 1: Both jobs are assigned as fractional on the ED, i.e., i_1 and k_2 are in $\{1, 2, \dots, m\}$. Clearly, in this case the makespan on the ES is at most T , same as that given in the LP-relaxed solution. Suppose that after rounding, the completion time on the ED under AMR² violates $2T$, i.e.,

$$P_1 + p_{i_1} + p_{k_2} > 2T. \quad (12)$$

From the LP-relaxed solution, we obtain

$$T - P_1 = p_{i_1} \bar{x}_{i_1} + p_{k_1} \bar{x}_{k_1} + p_{i_2} \bar{x}_{i_2} + p_{k_2} \bar{x}_{k_2}. \quad (13)$$

Substituting (13) in (12), we obtain

$$p_{i_1} + p_{k_2} > T + p_{i_1} \bar{x}_{i_1} + p_{k_1} \bar{x}_{k_1} + p_{i_2} \bar{x}_{i_2} + p_{k_2} \bar{x}_{k_2}. \quad (14)$$

Using (8) and (9) in (14), we obtain

$$p_{i_1} \bar{x}_{k_1} + p_{k_2} \bar{x}_{i_2} - p_{k_1} \bar{x}_{k_1} - p_{i_2} \bar{x}_{i_2} > T. \quad (15)$$

The inequality in (15) implies that

$$p_{i_1} \bar{x}_{k_1} + p_{k_2} \bar{x}_{i_2} + p_{k_1} \bar{x}_{k_1} + p_{i_2} \bar{x}_{i_2} > T. \quad (16)$$

Given (10) and (11), (16) should hold if we substitute \bar{x}_{i_1} in place of \bar{x}_{k_1} and \bar{x}_{k_2} in place of \bar{x}_{i_2} , i.e.,

$$p_{i_1} \bar{x}_{i_1} + p_{k_1} \bar{x}_{k_1} + p_{i_2} \bar{x}_{i_2} + p_{k_2} \bar{x}_{k_2} > T. \quad (17)$$

However, the left hand side (LHS) of (17) is equal to $T - P_1$ (cf. (13)), which is smaller than T . Therefore, (15) is false and by contradiction $P_1 + p_{i_1} + p_{k_2} \leq 2T$ is true.

Case 2: One job is assigned as fractional between the ED and the ES and the other job is assigned as fractional between two models on the ED. WLOG, we consider job 1 is assigned to models on the ED and job 2 is assigned between the ED and the ES. We consider the following sub cases.

Case 2a: Job 2 is scheduled on the ES and from (11) we must have $k_2 = m + 1$. From the LP-relaxed solution we have

$$P_1 + p_{i_1} \bar{x}_{i_1} + p_{k_1} \bar{x}_{k_1} + p_{i_2} \bar{x}_{i_2} \leq T, \quad (18)$$

$$P_2 + p_{k_2} \bar{x}_{k_2} \leq T. \quad (19)$$

For the ES, consider that

$$P_2 + p_{k_2} > 2T. \quad (20)$$

Because of (11), $p_{k_2} \bar{x}_{k_2} > T$ is true. If $p_{k_2} \bar{x}_{k_2} > T$ then also $P_2 + p_{k_2} \bar{x}_{k_2} > T$, which contradicts (19), and thus $P_2 + p_{k_2} \leq 2T$.

For the ED consider that

$$P_1 + p_{i_1} > 2T. \quad (21)$$

Substituting (18) in (21), we obtain:

$$p_{i_1} > T + p_{i_1} \bar{x}_{i_1} + p_{k_1} \bar{x}_{k_1} + p_{i_2} \bar{x}_{i_2}.$$

Using (9) we arrive to

$$p_{i_1} \bar{x}_{k_1} - p_{k_1} \bar{x}_{k_1} - p_{i_2} \bar{x}_{i_2} > T. \quad (22)$$

We consider (18) and (22), we have

$$\begin{aligned} P_1 + p_{i_1} \bar{x}_{i_1} + p_{k_1} \bar{x}_{k_1} + p_{i_2} \bar{x}_{i_2} &< p_{i_1} \bar{x}_{k_1} - p_{k_1} \bar{x}_{k_1} - p_{i_2} \bar{x}_{i_2} \\ \implies P_1 + p_{i_1} (\bar{x}_{i_1} - \bar{x}_{k_1}) + 2p_{k_1} \bar{x}_{k_1} + 2p_{i_2} \bar{x}_{i_2} &< 0, \end{aligned}$$

which is false as all quantities on LHS are positive. This means that (22) is false, thus (21) too.

Case 2b: Job 2 is scheduled on the ED. When using the basic solution of the LP problem, the completion time of the ED is

$$P_1 + p_{i_1} \bar{x}_{i_1} + p_{k_2} \bar{x}_{k_2} + p_{k_1} \bar{x}_{k_1} \leq T. \quad (23)$$

We claim that $P_1 + p_{i_1} + p_{k_2} \leq 2T$.

Consider

$$P_1 + p_{i_1} + p_{k_2} > 2T. \quad (24)$$

Substituting (23) in (24), and using (8) and (9), we obtain

$$p_{i_1} \bar{x}_{k_1} + p_{k_2} \bar{x}_{i_2} - p_{k_1} \bar{x}_{k_1} > T. \quad (25)$$

Substituting (23) in (25), we obtain

$$-P_1 - p_{i_1} (\bar{x}_{k_1} - \bar{x}_{i_1}) + p_{k_2} (\bar{x}_{i_2} - \bar{x}_{k_2}) \geq 0.$$

The LHS is the sum of negative terms because of (10) and (11), thus the inequality is false.

Case 3: Both jobs are assigned as fractional between ED and ES by the LP-relaxation. We have three different sub cases based on the fractional assignment of the LP-relaxation problem.

Case 3a: Both jobs are scheduled on the ES. The completion time of the ES, when considering the basic solution \bar{x} is:

$$P_2 + p_{i_1} \bar{x}_{i_1} + p_{k_2} \bar{x}_{k_2} \leq T \quad (26)$$

We claim that $P_2 + p_{i_1} + p_{k_2} \leq 2T$. Suppose that

$$P_2 + p_{i_1} + p_{k_2} > 2T. \quad (27)$$

We substitute (26) in (27) and obtain

$$p_{i_1} \bar{x}_{k_1} + p_{k_2} \bar{x}_{i_2} > T. \quad (28)$$

Substituting (26) in (28) we have

$$-P_2 - p_{i_1} (\bar{x}_{i_1} - \bar{x}_{k_1}) - p_{k_2} (\bar{x}_{k_2} - \bar{x}_{i_2}) \geq 0$$

which is false because of hypotheses (10) and (11).

Case 3b: One job is scheduled on the ED and the other on the ES. The completion time of the ED under the LP-relaxed solution satisfies

$$P_1 + p_{i_1} \bar{x}_{i_1} + p_{i_2} \bar{x}_{i_2} \leq T, \quad (29)$$

while the completion on the ES satisfies:

$$P_2 + p_{k_2} \bar{x}_{k_2} + p_{k_1} \bar{x}_{k_1} \leq T. \quad (30)$$

We claim $P_1 + p_{i_1} \leq 2T$, and $P_2 + p_{k_2} \leq 2T$. Suppose that

$$P_1 + p_{i_1} > 2T. \quad (31)$$

We substitute (29) in (31) and obtain

$$p_{i_1} \bar{x}_{k_1} - p_{i_2} \bar{x}_{i_2} > T. \quad (32)$$

Using (32) and (29) we have:

$$\begin{aligned} P_1 + p_{i_1} \bar{x}_{i_1} + p_{i_2} \bar{x}_{i_2} &< p_{i_1} \bar{x}_{k_1} - p_{i_2} \bar{x}_{i_2} \\ \implies P_1 + p_{i_1} (\bar{x}_{i_1} - \bar{x}_{k_1}) + 2p_{i_2} \bar{x}_{i_2} &< 0, \end{aligned}$$

which is false, thus (31) is not true.

On the ES, suppose that

$$P_2 + p_{k_2} > 2T. \quad (33)$$

Substituting (30) in (33) we obtain

$$p_{k_2} \bar{x}_{i_2} - p_{k_1} \bar{x}_{k_1} > T. \quad (34)$$

Using (34) and (30) we have

$$\begin{aligned} P_2 + p_{k_2} \bar{x}_{k_2} + p_{k_1} \bar{x}_{k_1} &< p_{k_2} \bar{x}_{i_2} - p_{k_1} \bar{x}_{k_1} \\ \implies P_2 + p_{k_2} (\bar{x}_{k_2} - \bar{x}_{i_2}) + 2p_{k_1} \bar{x}_{k_1} &< 0 \end{aligned}$$

which is false, thus (33) is not true.

Case 3c: Both jobs are scheduled on the ED. We claim

$$P_1 + p_{i_1} + p_{k_2} \leq 2T. \quad (35)$$

The completion time equation on the ED using the basic solution \bar{x} is

$$P_1 + p_{i_1} \bar{x}_{i_1} + p_{k_2} \bar{x}_{k_2} \leq T. \quad (36)$$

We negate (35):

$$P_1 + p_{i_1} + p_{i_2} > 2T. \quad (37)$$

We substitute (36) in (37) and obtain

$$p_{i_1} \bar{x}_{k_1} + p_{k_2} \bar{x}_{i_2} > T. \quad (38)$$

Substituting (36) in (38), we obtain

$$-P_1 - p_{i_1} (\bar{x}_{i_1} - \bar{x}_{k_1}) - p_{k_2} (\bar{x}_{k_2} - \bar{x}_{i_2}) \geq 0,$$

which is false because of hypotheses (10) and (11). \square

THEOREM 5.2. *The difference between the optimal total accuracy A^* and A^\dagger , the total accuracy achieved by AMR², is upper bounded by $a_{m+1} - a_1$.*

PROOF. Since $A^* \leq A_{LP}^*$, we prove the result with respect to A_{LP}^* . WLOG, we consider that i_1 and i_2 as the indices of the models with lower accuracy, respectively, for jobs 1 and 2, and k_1 and k_2 are the index of the models which provide higher accuracy. To prove the performance bound we distinguish the following three cases.

Case 1: $\bar{x}_{k_1} \geq \frac{1}{2}$, $\bar{x}_{k_2} \geq \frac{1}{2}$. In this case, AMR² will schedule job 1 on model k_1 and job 2 on model k_2 . The contribution of the following jobs to the A^\dagger is $a_{k_1} + a_{k_2}$. The contribution of the same jobs to the optimal solution A_{LP}^* is: $a_{i_1} \bar{x}_{i_1} + a_{k_1} \bar{x}_{k_1} + a_{i_2} \bar{x}_{i_2} + a_{k_2} \bar{x}_{k_2}$. However, $a_{k_1} > a_{i_1}$ and $a_{k_2} > a_{i_2}$: thus it is trivial that $A^\dagger > A^*$.

Case 2: $\bar{x}_{k_1} \geq \frac{1}{2}$ and $\bar{x}_{k_2} < \frac{1}{2}$. Here, AMR² schedules job 1 on k_1 and job 2 on i_2 .

$$\begin{aligned} A_{LP}^* - A^\dagger &= a_{i_1} \bar{x}_{i_1} + a_{k_1} \bar{x}_{k_1} + a_{i_2} \bar{x}_{i_2} + a_{k_2} \bar{x}_{k_2} - a_{i_2} - a_{k_1} \\ &= a_{i_1} \bar{x}_{i_1} + a_{i_2} (\bar{x}_{i_2} - 1) + a_{k_1} (\bar{x}_{k_1} - 1) + a_{k_2} \bar{x}_{k_2} \\ &= a_{i_1} \bar{x}_{i_1} - a_{i_2} \bar{x}_{k_2} - a_{k_1} \bar{x}_{i_1} + a_{k_2} \bar{x}_{k_2} \\ &= \bar{x}_{i_1} (a_{i_1} - a_{k_1}) + \bar{x}_{k_2} (a_{k_2} - a_{i_2}). \end{aligned}$$

Substituting $\bar{x}_{k_1} \geq \frac{1}{2}$ and $\bar{x}_{k_2} < \frac{1}{2}$ in the above equation we have $a_{i_1} - a_{k_1} < 0$. In conclusion

$$A_{LP}^* - A^\dagger \leq \frac{1}{2} (a_{k_2} - a_{i_2}). \quad (39)$$

Proof for the case $\bar{x}_{k_1} < \frac{1}{2}$, $\bar{x}_{k_2} \geq \frac{1}{2}$ is similar to **Case 2**.

Case 3: $\bar{x}_{k_1} < \frac{1}{2}$ and $\bar{x}_{k_2} < \frac{1}{2}$. AMR² schedules job 1 on i_1 , and job 2 on i_2 .

$$\begin{aligned} A_{LP}^* - A^\dagger &= a_{i_1} \bar{x}_{i_1} + a_{k_1} \bar{x}_{k_1} + a_{i_2} \bar{x}_{i_2} + a_{k_2} \bar{x}_{k_2} - a_{i_1} - a_{i_2} \\ &= \bar{x}_{k_1} (a_{k_1} - a_{i_1}) + \bar{x}_{k_2} (a_{k_2} - a_{i_2}) \leq a_{m+1} - a_1. \end{aligned}$$

In the last equation above, we used $\bar{x}_{k_1} < \frac{1}{2}$ and $\bar{x}_{k_2} < \frac{1}{2}$, and the fact that $a_{k_1} - a_{i_1}$ and $a_{k_2} - a_{i_2}$ are upper bounded by $a_{m+1} - a_1$. \square

COROLLARY 1. *If the processing times of all jobs on the ES are at most T , then $A^* \leq A^\dagger$.*

PROOF. Considering AMR² when all jobs have processing time less than T on all the models, we consider three cases. In the first case there is no job assigned as fractional in the LP-relaxed solution which implies that we obtain $A^* = A_{LP}^* = A^\dagger$. In the second case there is a single fractional job and we have $A_{LP}^* \leq A_{LP,I}^* + a_{m+1}$. In this case, AMR² schedules the fractional job on the ES and we obtain $A^\dagger = A_{LP,I}^* + a_{m+1}$. Thus, $A_{LP}^* \leq A^\dagger$. In the third case, the number of fractional jobs is two. AMR² schedules or either both jobs on the ES or one on the ES and the other on model m of the ED, achieving a total accuracy that is at least $A^\dagger = A_{LP,I}^* + a_{m+1} + a_m$. The solution of AMR² will be always greater or at most equal to the solution of A_{LP}^* , as it will have T seconds to schedule the two jobs, meanwhile the LP-relaxation will have a time that is less or equal to T to schedule both of them. In all the three cases $A^\dagger \geq A_{LP}^* \geq A^*$. \square

Remark 1: The schedule \mathbf{x}^\dagger given by AMR² may result in a makespan greater than T . In our experimental results, we show that the percentage of violation on an average is at most 40% for the considered application. As noted before, a special case of our

problem is GAP for which the best known approximation algorithm, proposed in [29], has the makespan bound $2T$ and produces a schedule that may exceed T . The algorithm in [29] achieves the optimal cost for GAP. However, it requires a bisection search to find this optimal cost and each step in the bisection search requires solving an LP-relaxed feasibility problem. In contrast, in AMR² we solve an LP-relaxed problem only once and thus it has a lower computational complexity, which is important because, as we will see later in Section 6, computing the schedule itself cannot take significant time when T is small.

6 EXPERIMENTAL RESULTS

In this section, we first present the experimental setup. We then present the implementation details for estimating the processing and communication times. As explained in Section 2, the aspect of multiple models on the ED has not been considered in computation offloading literature and there are no existing algorithms that are applicable for the problem at hand for a performance comparison. Therefore, we present the performance comparison between AMR² and a baseline: Greedy Round Robin Algorithm (Greedy-RRA). Given the list of jobs, Greedy-RRA offloads them from the start of the list to the ES until the constraint T is met. The remaining jobs are assigned in a round robin fashion to the models on the ED until the constraint T is met. Any further remaining jobs are assigned to model 1. Note that Greedy-RRA solution may violate the time constraint T and its runtime is $O(n)$.

6.1 Experimental Setup

Our experimental setup comprises a Raspberry Pi device (the ED) and a local server (the ES) that are connected and located in the same LAN. Raspberry Pi has 4 cores, 1.5 GHz CPU frequency, and 4 GB RAM, with the operating system Raspbian 10, while the server has 512 cores, 1.4 GHz CPU frequency, and 504 GB RAM, with the operating system Debian 11. All the functions on Raspberry Pi and on the server are implemented using Python 3. We used HTTP protocol to offload images from Raspberry Pi to the ES and implemented HTTP Client and Server using Requests and Flask, respectively.

The data samples are images from the ImageNet dataset for which we use DNN models for inference. On Raspberry Pi we import, from the TensorFlow Lite library, two pre-trained MobileNets corresponding to two values 0.25 and 0.75 for the hyperparameter α , which is a width multiplier for the DNN [14]. Both the models are quantized and require input images of dimensions 128×128 . On the ES, we import a pre-trained ResNet50 model [13] from the Tensorflow library. The ResNet50 model requires input images of dimensions 224×224 . Images of different dimensions need to be reshaped to the respective dimensions on the ES and the ED. The top-1 accuracies for the three models are presented in Table 1.

We implemented both AMR² and Greedy-RRA on Raspberry Pi in Python 3. AMR² takes up to 50 ms for computing a schedule for 40 jobs. The runtime of AMR² is dominated by the runtime of the solver from the Python library for solving the LP-relaxation. In future, we plan to reduce this runtime by implementing AMR² in C.

Model	Top-1 Accuracy
MobileNet $\alpha = 0.25$ (model 1)	0.395
MobileNet $\alpha = 0.75$ (model 2)	0.559
ResNet50 (model 3)	0.771

Table 1: Test accuracies of the considered DNN models [1].

6.2 Estimation of Processing and Communication Times

In our experiments, we consider images of dimensions 333×500 , 375×500 , and 480×640 , for which we estimate the processing and communication times using the following procedure. On Raspberry Pi, we run 30 samples of same image dimensions and use the median processing times as our estimate. Note that median is an unbiased estimate, and unlike the mean, it is not affected by cold start. We note that the estimates for the processing times include the reshape times.

In order to estimate the total time on the ES, we use the HTTP client/server connection to send 30 images of same image dimensions from Raspberry Pi to the server. For each image we measure the time till the reception of an inference for the image from the ES, and finally use the median. At the server we also measure the reshape time and the processing time, and the estimate for the communication time is obtained by subtracting the reshape time and the processing time from the total time. Since Raspberry Pi and the dedicated local server are in the same LAN, the observed communication times are almost constant with negligible variance. This is also true for the observed processing times, and we will later verify this when implementing the schedules using these estimates.

The estimates for the processing times are presented in Table 2. Observe that the processing times increase with the model size. On Raspberry Pi, the variance in the processing times on a model is small. In contrast, the total times on the ES vary with the dimensions of the image and are an order of magnitude higher than the processing times on Raspberry Pi. In Figure 2, we present the communication, reshape, and processing times on the ES. It is worth noting that, as the dimensions of the image increases, both communication time and the reshape times increase. Thus, it is more advantageous to offload images with smaller dimensions.

Model	Location	333×500	375×500	480×640
MobileNet $\alpha = 0.25$	ED	0.01	0.011	0.011
MobileNet $\alpha = 0.75$	ED	0.04	0.04	0.043
ResNet50	ES	0.28	0.32	0.38

Table 2: Estimated processing times (in seconds).

6.3 Performance of AMR²

In Figure 3, we examine the number of jobs assigned to different models under AMR². Observe that as T increases the number of jobs assigned to larger models increases. Also, note that MobileNet $\alpha = 0.25$ is only being used when T is small. In all the subsequent

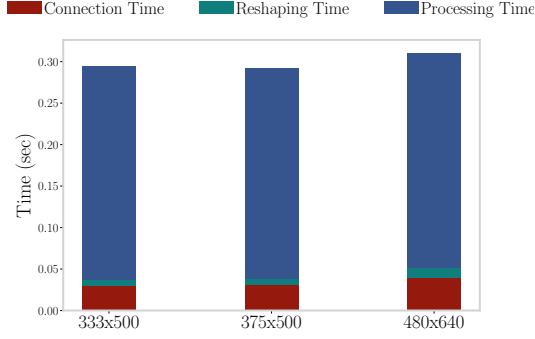


Figure 2: Estimated total time for inference on the ES.

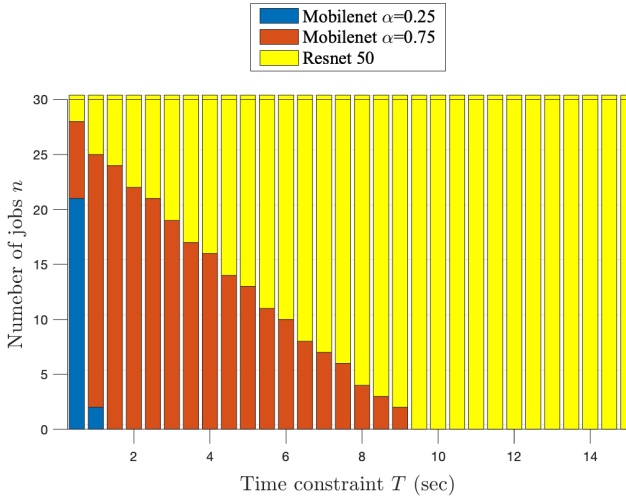


Figure 3: Job assignment under AMR^2 for varying T .

figures, for each point, we run 30 experiments and compute the average. Recall that the total accuracy A^\dagger is based on the top-1 accuracy of the models.

In Figures 4 and 5, we compare total accuracy achieved under different schedules, by varying T and n , respectively. For $n = 60$, no LP-relaxed solution exists for $T = 2$ sec. From both figures, we observe that A^\dagger overlaps with, and in some cases exceeds, the total accuracy of the LP-relaxed solution A_{LP}^* . This is because all the processing times (cf. Table 2) are less than 2 sec, the minimum value used for T , and therefore, from Corollary 1, A^\dagger exceeds A^* . Furthermore, in some cases, where T is large enough, AMR^2 may assign both the fractional jobs to the server and A^\dagger exceeds A_{LP}^* . In the above cases, however, the makespan under AMR^2 exceeds T .

From Figure 4, we observe that AMR^2 always has higher total accuracy than Greedy-RRA with a percentage gain between 20–60% averaging at 40%, but the percentage gains are lower at smaller T . The latter fact is also confirmed in Figure 5 when $T = 2$ sec. This is expected, because for $T = 2$ sec, not many jobs can be offloaded to the server as the processing times are around 0.3 seconds. For $T = 4$ sec we see significant gains of around 40–50%.

In Figure 6, we present the makespan achieved by AMR^2 and Greedy-RRA for varying n . The actual makespan, i.e., the time elapsed at Raspberry Pi from the start of scheduling the jobs till the

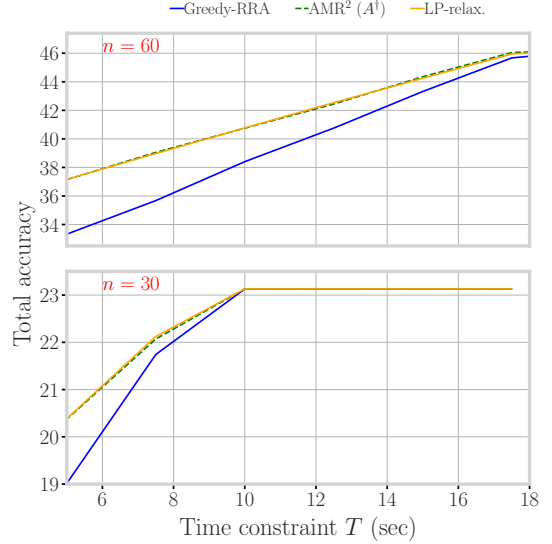


Figure 4: Total accuracy varying T for $n = (30, 60)$.

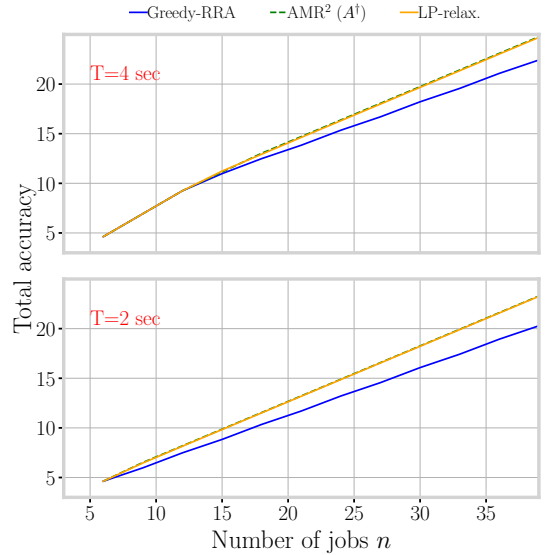


Figure 5: Total accuracy varying n with $T = (2, 4)$ sec.

finishing time of the last job is indicated by AMR^2 in the legend. The estimated makespan that is numerically computed using the schedule \mathbf{x}^\dagger and the estimated processing and communication times is indicated by AMR^2 (estd. proc. time). Observe that both these makespans have negligible difference asserting that the variances in our estimates for both communication and processing times are small. For $T = 4$, AMR^2 violates T for $n \geq 17$, but then it saturates at a makespan with a maximum percentage of violation of 15%. This is expected, because from Lemma 1 there cannot be more than two fractional jobs irrespective of n value and thus, the constraint violation due to the reassignment of the fractional jobs do not increase beyond $n = 30$. This saturation effect can also be observed for $T = 2$. In this case, the percentage of violation under AMR^2 is higher because the processing times on the server are comparable to $T = 2$ sec and reassigning a fractional job to the server results in higher percentage of violation.

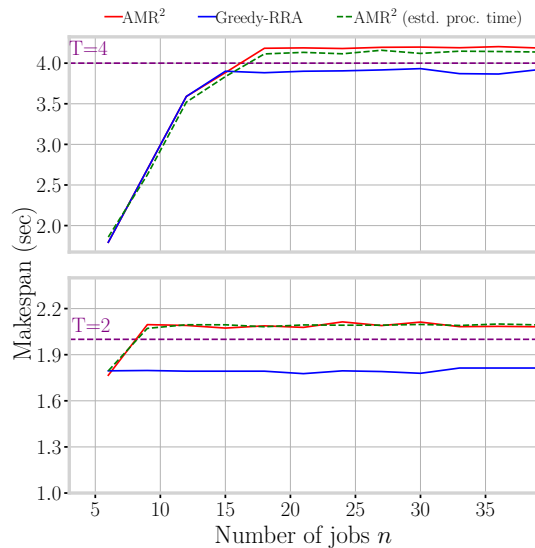


Figure 6: Makespan under AMR² and Greedy-RRA for varying n , and $T = 2$ sec and $T = 4$ sec.

7 CONCLUSION

We have studied the offloading decision for inference jobs between an ED and an ES, where the ED has m models and the ES has a state-of-the-art model. Given n data samples at the ED, we proposed an approximation algorithm AMR² for maximizing the total accuracy for the inference jobs subject to a time constraint T on the makespan. We proved that the makespan under AMR² is at most $2T$, and its total accuracy is lower than the optimal total accuracy by at most 1, and for typical problem instances its total accuracy is at least the optimal total accuracy. We have implemented AMR² on Raspberry Pi and demonstrated its efficacy in improving the inference accuracy: under the considered scenarios AMR² provides, on average, 40% higher total accuracy than that of Greedy-RRA.

In our problem model, the communication times are deterministic, which is applicable when ED and the ES are connected via Ethernet. If the ED and ES are connected over a wireless channel, where the communication times are random, the problem is more difficult and is open. We plan to study the performance analysis of AMR² that uses estimated mean communication times over the wireless channel. Also, we aim to study application scenarios where inference jobs arrive dynamically over time.

8 ACKNOWLEDGEMENT

This research is funded by the European Union through MSCA-PF project “DIME: Distributed Inference for Energy-efficient Monitoring at the Network Edge” under Grant Agreement No. 101062011.

REFERENCES

- [1] [n.d.]. Image Classification using TensorFlow Lite. https://www.tensorflow.org/lite/guide/hosted_models.
- [2] [n.d.]. PyTorch Mobile. <https://pytorch.org/mobile/home/>.
- [3] Jan van den Brand. 2019. A Deterministic Linear Program Solver in Current Matrix Multiplication Time. <https://arxiv.org/abs/1910.11957>
- [4] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. 2019. Once-for-All: Train One Network and Specialize it for Efficient Deployment. <https://arxiv.org/abs/1908.09791>

- [5] Dirk G. Cattrysse and Luk N. Van Wassenhove. 1992. A survey of algorithms for the generalized assignment problem. *European Journal of Operational Research* 60, 3 (1992), 260–272.
- [6] Jaya Prakash Champati and Ben Liang. 2017. Semi-Online Algorithms for Computational Task Offloading with Communication Delay. *IEEE Transactions on Parallel and Distributed Systems* 28, 4 (2017), 1189–1201.
- [7] Jaya Prakash Champati and Ben Liang. 2020. Single Restart with Time Stamps for Parallel Task Processing with Known and Unknown Processors. *IEEE Transactions on Parallel and Distributed Systems* 31, 1 (2020), 187–200.
- [8] Chandra Chekuri and Sanjeev Khanna. 2000. A PTAS for the Multiple Knapsack Problem. In *Proc. ACM SODA*. 213–222.
- [9] Meng-Hsi Chen, Ben Liang, and Min Dong. 2015. A semidefinite relaxation approach to mobile cloud offloading with computing access point. In *Proc. IEEE SPAWC*. 186–190. <https://doi.org/10.1109/SPAWC.2015.7227025>
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *Proc. IEEE CVPR*. 248–255. <https://doi.org/10.1109/CVPR.2009.5206848>
- [11] Lei Deng, Guoqi Li, Song Han, Luping Shi, and Yuan Xie. 2020. Model Compression and Hardware Acceleration for Neural Networks: A Comprehensive Survey. *Proc. IEEE* 108, 4 (2020), 485–532. <https://doi.org/10.1109/JPROC.2020.2976475>
- [12] Krzysztof Dudzinski. 1989. On a cardinality constrained linear programming knapsack problem. *Operations Research Letters* 8, 4 (1989), 215–218.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep Residual Learning for Image Recognition. <https://arxiv.org/abs/1512.03385>
- [14] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. 2017. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. <https://arxiv.org/abs/1704.04861>
- [15] Mohamed Kamoun, Wael Labidi, and Mireille Sarkiss. 2015. Joint resource allocation and offloading strategies in cloud enabled cellular networks. In *Proc. IEEE ICC*. 5529–5534.
- [16] H. Kellerer, U. Pferschy, and D. Pisinger. 2004. *Knapsack Problems*. Springer, Berlin, Germany.
- [17] Juan Liu, Yuyi Mao, Jun Zhang, and Khaled B. Letaief. 2016. Delay-optimal computation task scheduling for mobile-edge computing systems. In *Proc. IEEE ISIT*.
- [18] Pavel Mach and Zdenek Becvar. 2017. Mobile Edge Computing: A Survey on Architecture and Computation Offloading. *IEEE Communications Surveys Tutorials* 19, 3 (2017), 1628–1656.
- [19] Yuyi Mao, Jun Zhang, and Khaled B. Letaief. 2016. Dynamic Computation Offloading for Mobile-Edge Computing With Energy Harvesting Devices. *IEEE Journal on Selected Areas in Communications* 34, 12 (2016), 3590–3605.
- [20] Ivana Nikoloska and Nikola Zlatanov. 2021. Data Selection Scheme for Energy Efficient Supervised Learning at IoT Nodes. *IEEE Communications Letters* 25, 3 (2021), 859–863.
- [21] Samuel S. Ogden and Tian Guo. 2020. MDInference: Balancing Inference Accuracy and Latency for Mobile Applications. <https://arxiv.org/abs/2002.06603>
- [22] Michael L. Pinedo. 2008. *Scheduling: Theory, Algorithms, and Systems* (3rd ed.). Springer Publishing Company, Incorporated.
- [23] C.N. Potts. 1985. Analysis of a linear programming heuristic for scheduling unrelated parallel machines. *Discrete Applied Mathematics* 10, 2 (1985), 155–164.
- [24] G. Ross and R. Soland. 1975. A branch and bound algorithm for the generalized assignment problem. *Mathematical Programming* 8 (1975), 91–103.
- [25] Philipp Ross and Andre Luckow. 2019. EdgeInsight: Characterizing and Modeling the Performance of Machine Learning Inference on the Edge and Cloud. In *2019 IEEE International Conference on Big Data (Big Data)*. 1897–1906.
- [26] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. 2018. MobileNetV2: Inverted Residuals and Linear Bottlenecks. (2018). <https://arxiv.org/abs/1801.04381>
- [27] Mahadev Satyanarayanan, Paramvir Bahl, Ramon Caceres, and Nigel Davies. 2009. The Case for VM-Based Cloudlets in Mobile Computing. *IEEE Pervasive Computing* 8, 4 (2009), 14–23.
- [28] Weisong Shi, Jie Cao, Quan Zhang, Youhuizi Li, and Lanyu Xu. 2016. Edge Computing: Vision and Challenges. *IEEE Internet of Things Journal* 3, 5 (2016), 637–646.
- [29] David B. Shmoys and Éva Tardos. 1993. An Approximation Algorithm for the Generalized Assignment Problem. *Math. Program.* 62, 1–3 (feb 1993), 461–474.
- [30] Sowndarya Sundar, Jaya Prakash Varma Champati, and Ben Liang. 2020. Multi-user Task Offloading to Heterogeneous Processors with Communication Delay and Budget Constraints. *IEEE Transactions on Cloud Computing* (2020), 1–1.
- [31] Surat Teerapittayanon, Bradley McDanel, and H. T. Kung. 2017. BranchyNet: Fast Inference via Early Exiting from Deep Neural Networks. <https://arxiv.org/abs/1709.01686>
- [32] Yanting Wang, Min Sheng, Xijun Wang, Liang Wang, and Jiandong Li. 2016. Mobile-Edge Computing: Partial Computation Offloading Using Dynamic Voltage Scaling. *IEEE Transactions on Communications* 64, 10 (2016), 4268–4282.
- [33] Zizhao Wang, Wei Bao, Dong Yuan, Liming Ge, Nguyen H. Tran, and Albert Y. Zomaya. 2019. SEE: Scheduling Early Exit for Mobile DNN Inference during Service Outage. In *Proc. ACM MSWIM*. 279–288.