

Distributing Inference Tasks over Interconnected Systems through Dynamic DNNs

*Original*

Distributing Inference Tasks over Interconnected Systems through Dynamic DNNs / Singhal, Chetna; Wu, Yashuo; Malandrino, Francesco; Levorato, Marco; Chiasserini, Carla Fabiana. - In: IEEE-ACM TRANSACTIONS ON NETWORKING. - ISSN 1063-6692. - STAMPA. - 33:4(2025), pp. 1717-1730. [10.1109/TON.2025.3543848]

*Availability:*

This version is available at: 11583/2997572 since: 2025-09-02T10:41:09Z

*Publisher:*

IEEE

*Published*

DOI:10.1109/TON.2025.3543848

*Terms of use:*

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

IEEE postprint/Author's Accepted Manuscript

©2025 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collecting works, for resale or lists, or reuse of any copyrighted component of this work in other works.

(Article begins on next page)

# Distributing Inference Tasks over Interconnected Systems through Dynamic DNNs

Chetna Singhal, *Senior Member, IEEE*, Yashuo Wu, Francesco Malandrino, *Senior Member, IEEE*, Marco Levorato, *Senior Member, IEEE*, Carla Fabiana Chiasserini, *Fellow, IEEE*

**Abstract**—An increasing number of mobile applications leverage deep neural networks (DNN) as an essential component to adapt to the operational context at hand and provide users with an enhanced experience. It is thus of paramount importance that network systems support the execution of DNN inference tasks in an efficient and sustainable way. Matching the diverse resources available at the mobile-edge-cloud network tiers with the applications requirements and the complexity of their, while minimizing energy consumption, is however challenging. A possible approach to the problem consists in exploiting the emerging concept of dynamic DNNs, characterized by multi-branched architectures with *early exits* enabling sample-based adaptation of the model depth. We leverage this concept and address the problem of deploying portions of DNNs with early exits across the mobile-edge-cloud system and allocating therein the necessary network, computing, and memory resources. We do so by developing a 3-stage graph-modeling method that allows us to represent the characteristics of the system and the applications as well as the possible options for splitting the DNN over the multi-tier network nodes. Our solution, called Feasible Inference Graph (FIN), can determine the DNN split, deployment, and resource allocation that minimizes the inference energy consumption while satisfying the nodes’ constraints and the requirements of multiple, co-existing applications. FIN closely matches the optimum and leads to over 89% energy savings with respect to state-of-the-art alternatives.

**Index Terms**—Network support to machine learning, Dynamic neural networks, Energy efficiency, Inference in the mobile-edge-cloud continuum

## I. INTRODUCTION

Deep Neural Networks (DNN) are a pervasive paradigm that empowers a wide range of applications with advanced data analysis and decision making capabilities. Examples include computer vision [1], speech recognition [2], natural language processing [3], and mobile health care [4]. However, modern DNNs are complex and incur a significant energy consumption [5], making it challenging to deploy them on mobile platforms and embedded devices that have limited computing and communication resources, memory, and energy reservoir.

Two main technical strategies attempt to address this issue: model compression [6] and edge computing [7]. In the former,

C. Singhal is with INRIA, France. Y. Wu and M. Levorato are with University of California, Irvine, USA. F. Malandrino and C.F. Chiasserini are with CNR-IEIT and CNIT, Italy. C.F. Chiasserini is with Politecnico di Torino, Italy and with Chalmers University of Technology, Sweden.

This work was supported by the Smart Networks and Services Joint Undertaking (SNS JU) under Europe Union’s Horizon Europe research and innovation programme under Grant No.101192521 (MultiX project) and Grant No.101096379 (CENTRIC project), and under the Italian National Recovery and Resilience Plan (NRRP) of NextGenerationEU, partnership on “Telecommunications of the Future” (PE0000001 - program “RESTART”). The work of Marco Levorato and Yashuo Wu was partially supported by the U.S. NSF grant CCF 2140154.

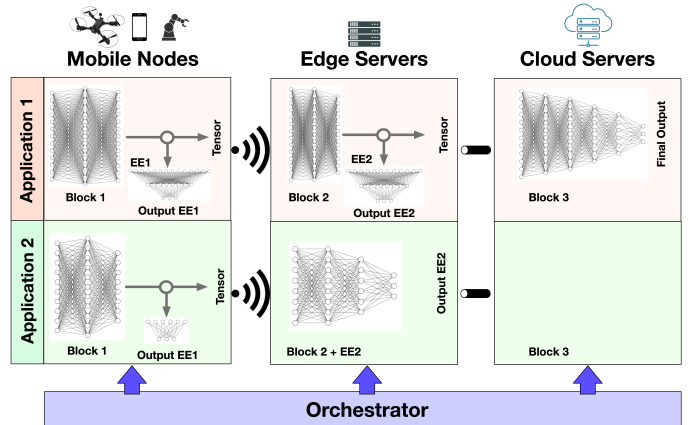


Fig. 1. Our framework allocates blocks of layers of DNNs with early-exits (EEs) over multi-tier systems so that energy consumption is minimized. An orchestrator controls the deployment of the DNN blocks that compose the intelligent applications supported by the system, as well as the information flow across them depending on the applications’ accuracy and latency requirements and the system’s bandwidth, computing and memory constraints. For Application 1, the orchestrator allocates two blocks with EE1 and EE2 to a mobile node and an edge server, while for Application 2 the application requirements are fulfilled by using the first two EEs.

the size and complexity of large DNN models are reduced using techniques such as pruning, quantization, and knowledge distillation. However, compression often results in a noticeable loss of inference performance, and in many mobile settings the execution of even compressed models necessitates a considerable amount of power. In edge computing, a compute-capable node positioned at the network edge takes over the execution of the DNN. Edge computing suffers from the need to transfer information-rich signals over a volatile wireless channel, often under tight latency constraints. Additionally, in real-world systems, edge servers likely need to support a multitude of applications and users, which may overload their available memory as well as computing and communications capabilities.

Recently, a variant of edge computing emerged where DNN models are divided into *sections* that can be distributed over mobile-edge-cloud systems. This approach is often referred to as “split computing” or “split DNN” [8]. *By controlling the splitting points defining the sections of the DNN, one can control the computing load allocated to the different devices/servers as well as the amount of data transmitted on the communication links connecting them.* In fact, each node is tasked with the execution of a portion of the DNN, and is thus in charge of a fraction of the overall operations. Furthermore, instead of the input, the *sections transmit their output tensors, whose size is a function of the splitting point.*

Relevant work in this direction includes [9]–[13], as discussed in Sec. VI. However, determining the most effective splitting, meets the performance requirements (quality and latency) whilst minimizing the cost is even more challenging when several heterogeneous applications, each requiring a DNN model, have to be deployed in the system.

In this paper, we tackle *the split DNN problem in the context of dynamic DNNs – and, specifically, multi-branched DNNs equipped with early exits [8], [14], [15] – and multi-tier mobile-edge-cloud systems* (Figure 1). The motivation behind these architectures is that, to achieve high performance in complex tasks, DNNs have to deal with the most challenging samples in a dataset. This results in models that are over-parametrized for a large set of the input samples. Early exits are model-tails attached to some layers of the original model designed to produce an analogous output as that of the full DNN, with a smaller computing effort. By deciding which branches to execute sample-by-sample, early-exit models can dynamically adapt the number of operations needed to produce an output to the input “complexity”. An overview on early-exit DNNs and related challenges can be found in [8], [16].

Our work focuses on the problem of allocating “blocks” of layers of DNNs with early exits for heterogeneous applications with diverse requirements (quality and latency) to the nodes composing the overall mobile-edge-cloud system. Notably, unlike most of the existing studies on split DNNs, we consider a setting with multiple data sources, applications, and nodes at all the tiers of the infrastructure (as exemplified in Figure 1). Also, it is worth remarking that the presence of early exits influences the flow of information throughout the blocks, and some executions may be terminated early, further complicating the allocation problem. The overall allocation problem we formulate minimizes the energy needed to complete inference under latency and accuracy application requirements, as well as bandwidth, available memory, and computing resource constraints. To resolve this challenging – NP-hard – problem, we adopt an approach based on graph optimization, where we manipulate an initial graph capturing the relationship between system nodes and DNN layers to create a specialized graph that can be used to compute feasible, low-cost paths. Each path corresponds to a possible DNN deployment and resource allocation strategy.

In summary, the contributions of this work are as follows:

- 1) We formulate an allocation problem where the blocks of DNNs layers with early exits required by mobile applications are deployed on the nodes of a mobile-edge-cloud system characterized by diverse resources availability. The objective is to minimize the energy consumed while executing inference tasks, under application-specific inference requirements, as well as system constraints in terms of computing capability, memory availability, and channel capacity.
- 2) We devise a new resolution framework – named Feasible Inference Graph (FIN) – to solve the resource and DNN block allocation problem. FIN uses multistart global optimization to allocate communication, memory, and computation resources across the intelligent applications to support. An essential component of FIN is a static

TABLE I  
NOTATION

Symbol	Definition
$s \in \mathcal{S}$	Data sources
$n \in \mathcal{N}$	Multi-tiered network nodes
$h \in \mathcal{H}$	Applications (or, equivalently, DNN models)
$\tilde{\mathcal{G}} = \{\mathcal{V}, \mathcal{E}\}$	Two-dimensional two-plane graph modeling the overall system
$\mathcal{G} = \{\mathcal{V}, E_c\}$	Single-plane extended graph
$\ell_i^h$	$i$ -th block of application $h$ 's DNN
$\phi^h(\ell_i^h)$	Fraction of input samples output after block $\ell_i^h$
$d^h(\ell_1^h, \ell_2^h)$	Size (in bits) of the data output by block $\ell_1^h$
$\sigma^h(\ell_1^h, \ell_2^h)$	Number of operations needed to execute block $\ell_1^h$
$\sigma^h$	Inference rate for application $h$
$\alpha^h$	Target inference quality of application $h$
$\delta^h$	Target inference latency of application $h$
$\phi^h(\ell_i^h)$	Fraction of input samples that are output by $\ell_i^h$
$b^h(v, v')$	Bandwidth allocated to application $h$ between nodes $v$ and $v'$
$c^h(v, v')$ , $m^h(v, v')$	Computing and memory available to application $h$ at a node $v$
$T^h(v, v')$ , $C^h(v, v')$ , $E^h(v, v')$	Data transfer time, computing time, and energy consumption weights of the edge $v \rightarrow v' \in \mathcal{E}$ , with $v = (n, \ell_i^h)$ and $v' = (n', \ell_j^h)$
$\pi^h$	Path on $\mathcal{G}$ representing a configuration of application $h$
$\gamma$	Resolution of the feasibility graph
$x_b^h, x_c^h, x_m^h$	Communication, compute, and memory resource allocation of application $h$

procedure, FIN-S, that makes near-optimal DNN block allocation decisions within each application by building a graph model amenable to optimization. This graph model indeed contains only paths representing feasible solutions and is such that the minimum-cost path corresponds to the optimal allocation.

- 3) We explore the impact on the inference energy consumption of different allocation configurations for three DNN models with early exits (B-LeNet, B-AlexNet and B-ResNet, pre-trained on multiple datasets), over nodes with different capabilities.
- 4) Our results show that models equipped with early exits can dramatically decrease the overall energy consumption when some of these exits are allocated to mobile or edge devices under system and application-level constraints, by reducing the involvement of larger-scale nodes in the completion of the inference. We also show how FIN performs close to the optimum computed by brute force, in small-scale scenarios where the optimum computation is feasible. In more complex scenarios, we demonstrate that, compared to state-of-the-art approaches [17], [18] for cost minimization, FIN yields dramatic energy savings, exceeding 71% for the computational energy and 89% for the communication one. Importantly, such gains are consistent under different branchy DNN architectures, from the small-scale B-LeNet to the much larger B-AlexNet.

## II. ENERGY-AWARE INFERENCE THROUGH DYNAMIC NNS

We now introduce the mobile-edge-cloud system, along with the structure of the DNN models and applications. We also describe the corresponding two-plane graph capturing the possible allocation of DNN blocks (Plane 1) to the system nodes (Plane-2), along with the nodes available resources (Sec.II-A). Then we translate the two-plane graph into a single-plane extended graph (Sec.II-B), used to optimally configure dynamic DNNs for energy-efficient inference tasks (Sec.II-C). We will leverage the same extended graph to develop a low complexity solution, FIN, that gives performance close to optimum, as discussed in Sec. III.

### A. System model

We consider a multi-tiered communication and computing infrastructure composed of mobile nodes, edge servers, and cloud servers. The objective of the overall system is to support a set of mobile applications whose central component is a DNN performing the analysis of the information that the mobile nodes acquire through co-located data sources. Examples include computer vision models for object detection and image classification, as well as speech recognition tasks. The edge of the infrastructure hosts a Machine Learning (ML) orchestrator that possesses knowledge of the applications, as well as some essential information about the capabilities of the mobile nodes and the data they can acquire.

Formally, the main elements defining the system are:

- 1) A set  $\mathcal{S}$  of data sources (e.g., sensors) indexed with  $s \in \{1, \dots, S\}$ ; each data source samples a physical phenomenon extracting information that serves as input to an inference task for application  $h$  at rate  $\sigma^h \geq 0$ .
- 2) A set  $\mathcal{N}$  of computationally-capable network nodes, including (i) mobile nodes connected to data sources, (ii) edge servers, and (iii) cloud servers.
- 3) A set  $\mathcal{H}$  of applications  $h \in \{1, \dots, H\}$ , each associated with a DNN model equipped with early-exits. We describe the overall architecture of the DNN associated with application  $h$  as composed of a set of layer blocks  $\mathcal{L}^h = \{\ell_1^h, \dots, \ell_N^h\}$ , where *each block corresponds to a portion of the network backbone and at most one early exit*. Although the input rate of the application is  $\sigma^h$ , early exits may “capture” some input and terminate execution; we thus define  $\phi^h(\ell_i^h)$  as the fraction of input samples that are output after block  $\ell_i^h$ . Each application  $h$  has specific requirements defined as the target inference quality (e.g., target accuracy value)  $\alpha^h$ , and maximum inference latency  $\delta^h$ . In the following, we often refer to applications and the DNN representing their essential component interchangeably.

Given application  $h$ , the ML orchestrator determines which blocks of the DNN should be executed and assigns them to network nodes in such a way that the application inference requirements are fulfilled. Specifically, according to the split computing paradigm, the network nodes can cooperatively execute the overall DNN: a node assigned a DNN block or a set of blocks will execute all the corresponding layers and send the output coefficients of their cut layer (tensor) to the nodes

hosting the subsequent block(s). Note that a node may be allocated zero, one, or multiple exits, which are all executed.

We represent the overall system by means of a directed *two-dimensional load-resource, two-plane graph* model [19]–[23], with one plane capturing the network nodes and the other the blocks of the DNN layers. We establish a relationship between these two planes to represent the possible mapping between the communication and compute resource demand of the applications in  $\mathcal{H}$  onto the resources made available by the multi-tiered network nodes in  $\mathcal{N}$ . Also, the two-dimensional load-resource model matches the communication and computing resources handled by the system, i.e., offered by the network nodes and required by the applications.

More formally, we denote such graph, illustrated in the left panel of Figure 2, with  $\tilde{\mathcal{G}} = \{\tilde{\mathcal{V}}, \tilde{\mathcal{E}}\}$  where  $\tilde{\mathcal{V}}$  are the vertices and  $\tilde{\mathcal{E}}$  are the edges. The two graph planes are as follows (the notation is summarized in Table I).

**Plane 1:** The vertices  $\tilde{\mathcal{V}}_1 \subset \tilde{\mathcal{V}}$  of this plane correspond to the system nodes in  $\mathcal{S} \cup \mathcal{N}$ . Similar to network slicing, different applications receive separate allocations of bandwidth, computing resources, and memory. Accordingly, we associate with the edges  $\tilde{\mathcal{E}}_1 \subset \tilde{\mathcal{E}}$  three-dimensional weights  $[b^h(n_1, n_2), c^h(n_1, n_2), m^h(n_1, n_2)]$ ,  $n_1, n_2 \in \mathcal{N}$ , denoting the available bandwidth, computational resources, and memory (resp.). These weights are the product between the total available resources and the fraction thereof assigned to each application  $h$ . Concerning communication, computational, and memory resources, we thus define:

- parameters  $\bar{b}(n_1, n_2), \bar{c}(n_1), \bar{m}(n_1)$  denoting (resp.) the *total* communication resources between nodes  $n_1$  and  $n_2$ , and the *total* computational and memory resources existing at  $n_1$ ;
- variables  $x_b^h(n_1, n_2), x_c^h(n_1, n_2), x_m^h(n_1, n_2) \in [0, 1]$  denoting the fraction of said resources assigned to application  $h$ ;
- auxiliary variables  $b^h(n_1, n_2) \triangleq x_b^h(n_1, n_2) \bar{b}(n_1, n_2)$  denoting the bandwidth allocated to *application*  $h$  between nodes  $n_1$  and  $n_2$ ; and  $c^h(n_1, n_2) \triangleq x_c^h(n_1, n_2) \bar{c}(n_1)$  and  $m^h(n_1, n_2) \triangleq x_m^h(n_1, n_2) \bar{m}(n_1)$  denoting the computing and the memory available to *application*  $h$  at  $n_1$ .

The existence of the edge determines whether or not two nodes can communicate, and the self loop  $n \rightarrow n$  has infinite capacity, i.e.,  $b^h(n, n) = \infty$ . Further, the edge between a network node  $n \in \mathcal{N}$  and a co-located data source  $s \in \mathcal{S}$  has weight  $[b^h(s, n), c^h(s, n), m^h(s, n)] = [\infty, 0, 0]$ , i.e., we consider the bandwidth between the two  $s$  and  $n$  as unlimited and that the data source does not have any compute capability or memory.

**Plane 2:** The vertices  $\tilde{\mathcal{V}}_2 \subset \tilde{\mathcal{V}}$  of this plane correspond to DNN layers’ blocks  $\mathcal{L} = \{\mathcal{L}^h\}_h$ . The edges  $\tilde{\mathcal{E}}_2$  Plane 2 capture the connectivity structure of the DNNs, where edges exist only between consecutive blocks of the same application. The weights  $[d^h(\ell_1^h, \ell_2^h), o^h(\ell_1^h, \ell_2^h)]$ ,  $\ell_1^h, \ell_2^h \in \mathcal{L}^h$ , represent the size (in bits) of the data output by block  $\ell_1^h$  ( $d^h(\ell_1^h, \ell_2^h)$ ) and the number of operations needed to execute block  $\ell_1^h$  ( $o^h(\ell_1^h, \ell_2^h)$ ).

**Inter-plane edges:** A set of edges  $\tilde{\mathcal{E}}_{\ell \rightarrow n}$  unidirectionally connects Plane 2 to Plane 1, with the generic edge representing that a layers’ block of a DNN is deployed at a network node.

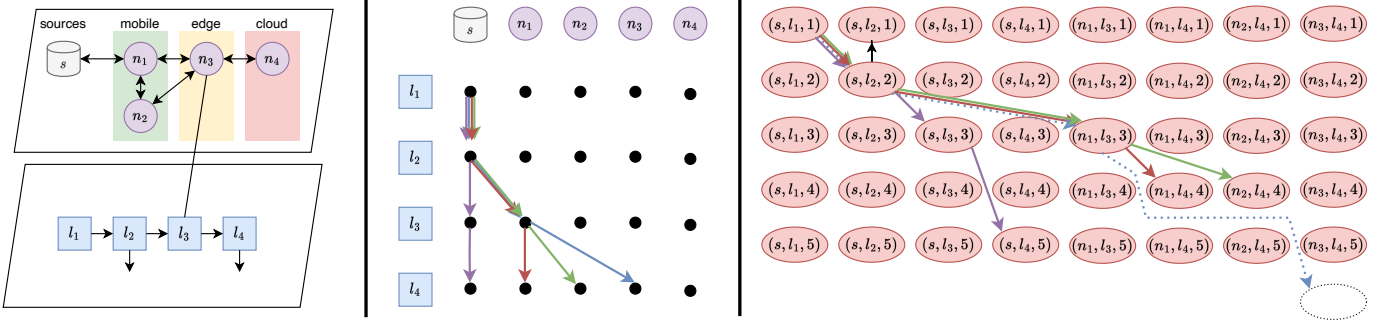


Fig. 2. The three graphs used in our solution strategy: 2-dimensional, 2-plane system model (left); single-plane extended graph (center); feasible graph (right).

### B. From the two-plane to the single-plane extended graph

We now transform the two-dimensional load-resource, two-plane graph  $\tilde{\mathcal{G}}$  into a directed *single-plane* extended graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  (second panel of Figure 2). The vertices of  $\mathcal{G}$  correspond to joint nodes and DNN blocks that are connected by inter-plane edges in  $\tilde{\mathcal{G}}$ , that is, in  $\mathcal{G}$  the two planes are collapsed following the connecting edges. Furthermore, an edge in  $\mathcal{G}$  exists only if the corresponding nodes and DNN blocks are connected in  $\tilde{\mathcal{G}}$ . Finally, vertices in the extended graph that correspond to impossible deployments do not have any outgoing or incoming edges. As an example, in Figure 2 node  $n_1$  cannot support layers  $l_1$  or  $l_2$ , hence, the corresponding vertices are not used. Consider the vertices  $v = (n, \ell_i^h)$  and  $v' = (n', \ell_j^h)$ , where  $n, n' \in \mathcal{N}$  are network nodes and  $\ell_i^h, \ell_j^h \in \mathcal{L}^h$  are DNN blocks. Notably, vertices corresponding to data sources in  $\mathcal{G}$  are not bound with any DNN block.

In  $\mathcal{G}$ , edges exist only between nodes embedding consecutive blocks of the same DNN (which may be deployed also on the same node), and are associated with the set of weights  $[T^h(v, v'), C^h(v, v'), E^h(v, v')]$ , where  $T^h(v, v')$  and  $C^h(v, v')$  are the data transfer time and the computing time, respectively, and  $E^h(v, v')$  is the per-inference energy consumption. We then define  $T^h(v, v')$  and  $C^h(v, v')$  as:

$$T^h(v, v') = \frac{d^h(\ell_i^h, \ell_j^h)}{b^h(n, n')}, \quad C^h(v, v') = \frac{o^h(\ell_i^h, \ell_j^h)}{c^h(n, n')}. \quad (1)$$

We recall that  $b^h(n, n') = \infty$  if  $n = n'$ , that is, the communication time is equal to 0 if the vertices are associated with different (and necessarily contiguous) DNN blocks allocated on the same node. The weight  $E^h(v, v')$  instead compounds the computing, transmission, and receiving energy associated with edge  $v \rightarrow v'$ , i.e.,

$$E^h(v, v') = (\xi_t + \xi_r) \frac{d^h(\ell_i^h, \ell_j^h)}{b^h(n, n')} + \xi_c \frac{o^h(\ell_i^h, \ell_j^h)}{c^h(n, n')}, \quad (2)$$

where  $\xi_t$ ,  $\xi_r$  and  $\xi_c$  are the power consumption spent, respectively, by node  $n$  to transmit, and by node  $n'$  to receive and compute, a data unit.

### C. Energy-efficient inference: Problem formulation

First, given application  $h$ , we denote with  $\pi^h$  the generic configuration indicating (i) the sources feeding data to the

application as well as (ii) *which* DNN blocks with early exits should be deployed (hence used) for inference, and (iii) *where* (i.e., on which network nodes). We then define binary selection variables  $p_\pi^h(v) \in \{0, 1\}$  based on the configuration  $\pi^h$ , where  $v \in \mathcal{V}$ . If  $p_\pi^h(v) = 1$ , then the vertex  $v$  is selected by the configuration  $\pi^h$ , i.e.,  $v \in \pi^h$ ; if  $p_\pi^h(v) = 0$ ,  $v$  is not selected, i.e.,  $v \notin \pi^h$ . The configuration  $\pi^h$  is controlled by the orchestrator. Note that the configuration needs to build a path from the sources to the DNN output for all the applications. Further, we remark that a node may be allocated multiple blocks, and even multiple exits. In the graph representation, this corresponds to a configuration that selects multiple vertices representing the same node and separate DNN blocks (like the red path in Figure 2(center)).

We adopt here a performance metric that is a function of the configuration and cannot be decoupled as a per-edge measure: the configuration inference quality. In fact, the configuration may suppress the execution of the blocks/exits after a certain index, for instance due to latency constraints. We thus define the inference quality of application  $h$  given configuration  $\pi^h$  as the quality associated with the whole sequence of DNN blocks in  $\pi^h$ , and denote it with  $a(\pi^h)$ .

The objective of the orchestrator is to minimize the overall energy consumption to support  $\sigma^h$  (tasks per second) inferences per second (3a), subject to latency, quality, available memory, network resource, and compute-resource constraints (3b)–(3i). Recalling the definition of  $b^h, c^h, m^h$ , and that  $v = (n, \ell_i^h)$  and  $v' = (n', \ell_j^h)$ , and  $\phi^h(\ell_i^h)$  is the fraction of input samples output by  $\ell_i^h$ , the resulting optimization problem is:

$$\min_{\pi^h, x_b^h, x_c^h, x_m^h} \sum_{h \in \mathcal{H}} \sum_{v, v' \in \mathcal{V}} \sigma^h \phi^h(\ell_i^h) E^h(v, v') p_\pi^h(v) p_\pi^h(v') \quad (3a)$$

$$\text{s.t.} \quad \sum_{v, v' \in \mathcal{V}} (T^h(v, v') + C^h(v, v')) p_\pi^h(v) p_\pi^h(v') \leq \delta^h, \quad \forall h \quad (3b)$$

$$a(\pi^h) \geq \alpha^h, \quad \forall h \quad (3c)$$

$$\sigma^h \phi^h(\ell_i^h) o^h(\ell_i^h, \ell_j^h) \leq c^h(v, v'), \quad \forall v, v' \in \pi^h \quad (3d)$$

$$\sigma^h \phi^h(\ell_i^h) d^h(\ell_i^h, \ell_j^h) \leq b^h(v, v'), \quad \forall v, v' \in \pi^h \quad (3e)$$

$$\sum_{\substack{n=n', \\ v, v' \in \pi^h}} \sigma^h \phi^h(\ell_i^h) d^h(\ell_i^h, \ell_j^h) \leq m^h(v, v'), \quad \forall v, v', h \quad (3f)$$

$$\sum_{h \in \mathcal{H}} x_b^h(n_1, n_2) \leq 1 \quad \forall n_1, n_2 \in \mathcal{N} \quad (3g)$$

$$\sum_{h \in \mathcal{H}} x_c^h(n_1, n_2) \leq 1 \quad \forall n_1, n_2 \in \mathcal{N} \quad (3h)$$

$$\sum_{h \in \mathcal{H}} x_m^h(n_1, n_2) \leq 1 \quad \forall n_1, n_2 \in \mathcal{N}. \quad (3i)$$

The optimization of the deployment configuration ( $\pi^h$ ) and resource allocation ( $x_b^h, x_c^h, x_m^h$ ) in the above problem minimizes the overall energy consumption for inference in applications ( $h \in \mathcal{H}$ ). It is subject to the constraints based on the requirements of the applications (inference quality- $\alpha^h$  and latency- $\delta^h$ ) and the resource availability (communication- $b^h$ , compute- $c^h$ , and memory- $m^h$ ) at the system level. Specifically, (3b) and (3c) represent the latency and inference quality constraints, respectively. The communication, compute, and memory resource allocation constraints for each application  $h$  are given by (3d), (3i), and (3f), respectively.

**Problem complexity.** The above problem exhibits high complexity, owing to its combinatorial nature and the overwhelming number of existing solutions. Specifically, we prove the following result.

*Property 1:* The problem of optimizing (3a) subject to constraints (3b)–(3i) is NP-hard.

*Proof:* To prove NP-hardness, we perform a reduction from a known NP-hard problem to the one under study. In particular, we show that any instance of the Steiner tree problem (STP) [24] can be transformed into a *simplified* instance of the problem introduced above. The STP is a generalization of the minimum spanning tree problem: given a weighted, undirected graph and a subset of nodes therein, the goal is to select the minimum-weight tree connecting all nodes in the subset. Given an instance of the STP, we build an instance of the problem in (3a)–(3i) by creating: (i) a data source for all the vertices to connect except one, and imposing that all such sources must be used for inference; (ii) one DNN layer, corresponding to the remaining vertex to connect; (iii) physical nodes for all intermediate vertices in the STP instance; (iv) only one of the physical nodes has enough capabilities to run the DNN layer. Also, the connectivity between nodes and data sources reproduces that of the STP instance, and one component of the weights in our problem instance is set to match the weights in the STP instance while all others are set to zero.

Solving our problem to optimality also yields an optimal solution to the STP instance, hence, the two problems are equivalent. Since the reduction takes polynomial (linear) time (each edge and vertex of the STP instance is processed once) and the STP problem is NP-hard [24], the thesis is proved. ■

It is also worth remarking that the instance of our problem created in the proof above is very simple (only one DNN layer, only one non-zero weight for the edges, etc.). This suggests that, on top of being NP-hard, our problem is significantly *more* complex than an already NP-hard problem like the STP. In light of the problem complexity, we propose below an algorithmic solution, leveraging a graph representation that, efficiently and very conveniently, embeds all possible decisions, the application requirements, and the system constraints.

### III. THE FIN SOLUTION

Here, we introduce our proposed heuristic, called Feasible Inference Graph (FIN). As shown in the example of Figure 2, FIN begins by creating a single-plane directed extended graph (Figure 2-center) from the two-dimensional load-resource, two-plane graph (Figure 2-left), using the approach described in Sec. II-B. Thereafter, FIN creates a feasible graph (Figure 2-right) from the single-plane extended graph and the constraints (application and system), through the FIN-S procedure. To solve the resulting problem, FIN iteratively performs a multi-start global search where, for each starting point, the FIN-S procedure is exploited to optimally allocate resources ( $x_b, x_c$ , and  $x_m$ ) in the system and achieve convergence.

We begin, in Sec. III-A, by describing the FIN-S procedure, making placement decisions under the assumption that the available resources are *statically* assigned to services, i.e., values of the  $x_b^h(n_1, n_2)$ ,  $x_c^h(n_1, n_2)$ , and  $x_m^h(n_1, n_2)$  are given and fixed. Then, in Sec. III-B we drop the static assignment assumption and describe how FIN-S is a building block of the full FIN procedure, making decisions on DNN early exiting and layers' deployment *as well as* communication, computing, and memory resource assignment. We remark that our FIN framework is applicable to a fairly broad class of alternative problems sharing the structure of (3), where a latency or quality objective is to be optimized subject to constraints such as resources, energy consumption, and quality. In those cases, the extended and feasible graphs will reflect the considered objective and constraints.

#### A. Statically-assigned resources: the FIN-S procedure

As mentioned above, the FIN-S (with the "S" stands for *static*) procedure makes decisions on the DNN early exiting and where to deploy the corresponding layers, given the values of  $x_b^h(n_1, n_2)$ ,  $x_c^h(n_1, n_2)$ , and  $x_m^h(n_1, n_2)$ . We first describe how to build a *feasible* graph – i.e., a graph summarizing all feasible solutions to the energy-aware inference problem – starting from the extended one we used to formulate the optimization problem. By construction, the feasible graph includes only those decisions (i.e., which data sources and DNN blocks are used and where such blocks should be deployed) that meet all constraints on inference latency and quality as well as on data, computational, and network resources. The second part of the section then describes how the most energy-efficient DNN configuration can be found by identifying the minimum-cost path traversing the feasible graph, as the edge weights represent the energy consumption (computation and communication) incurred by the nodes.

Specifically, the feasible graph built by FIN-S summarizes all feasible solutions through two complementary strategies:

- the *additive* constraint, namely, the inference latency, is guaranteed by the graph topology itself;
- the other constraints related to system capability and application requirements, e.g., inference quality and data requirements as well as memory, computation, and bandwidth limits, are guaranteed by pruning the edges and vertices that would violate them.

The vertices of the feasible graph are then the same as in the extended one, but each vertex is replicated a number  $\gamma$  of times. Let us denote the  $g$ -th replica of a vertex  $v=(n, l_i^h) \in \mathcal{G}$  with  $v_g, g=1, \dots, \gamma$ , and define index  $g$  as *depth* of a vertex.

Replicas of the same vertex share all the properties of the original vertex (e.g., the used model), but have different depth. This allows us to keep track of how our decisions (e.g., choosing a model) influence additive constraints like the accumulated interference latency; intuitively, vertices with a larger depth corresponds to situations in which the accumulated interference latency is closer to its target. In other words, the deeper a vertex, the closer it is to violating the additive constraint on the inference latency. Importantly, as mentioned earlier, all vertices correspond to solutions that do honor the latency limit; for vertices whose depth is exactly  $\gamma$  the constraint (3b) is met with an equality sign.

The number  $\gamma$  of replicas we generate for each vertex also corresponds to the *resolution* of our feasible graph. Intuitively, having a larger  $\gamma$  results in a finer (i.e., better) resolution; on the negative side, this also corresponds to more vertices in the graph, thus, a larger computational complexity. Accordingly, choosing  $\gamma$  allows us to select a trade-off between the quality of the solutions generated by our algorithms and the time it takes to generate them.

Once the vertices are in place, we proceed to creating the edges connecting them. Specifically, we create an edge from vertex  $v_{g_1}$  to vertex  $v'_{g_2}$ , with  $v'=(n', l_j^h)$  and  $g_2 > g_1$ , if:

- it is possible to place layer  $\ell_i^h$  at node  $n$  and layer  $\ell_j^h$  at node  $n'$  while meeting the nodes' memory and computation constraint;
- doing so incurs a combined processing time and network delay such that:

$$g_2 - g_1 = \left\lceil \gamma \cdot \frac{T^h(v, v') + C^h(v, v')}{\delta^h} \right\rceil. \quad (4)$$

In the numerator of the fraction above, the first term corresponds to the computing time, and the second to the data transfer delay. Also notice the combined effect of  $\gamma$  and the ceiling operator: the larger  $\gamma$ , the closer the input and output of the ceiling operator are, hence, the higher the quality of the solutions we will get.

Let us then define the steepness of an edge as the difference between the depth of the target and the source vertices, and the steepness of a path as the sum of the steepness values of its edges. Intuitively, steeper edges (and steeper paths) correspond to solutions with a longer inference latency; also, any path arriving to a vertex corresponding to an exit layer before depth  $\gamma$  represents, by construction, a solution conforming with the latency limit constraint.

Next, we prune all the vertices and edges that do not conform with the *local* constraints, i.e.,

- edges that would not reach the target inference quality, i.e., they belong to a configuration  $\pi^h$  s.t.  $a(\pi^h) < \alpha^h$ , or
- edges that would exceed the available bandwidth, computational capability, or available memory of the node represented by the vertex tail of the edge, i.e.,  $\sigma^h \phi^h(\ell_i^h) d^h(\ell_i^h, \ell_j^h) > b^h(v, v')$ ,

$$\sigma^h \phi^h(\ell_i^h) o^h(\ell_i^h, \ell_j^h) > c^h(v, v'), \text{ or } \sigma^h \phi^h(\ell_i^h) d^h(\ell_i^h, \ell_j^h) > m^h(v, v').$$

Surviving edges are assigned a weight corresponding to the energy consumption they incur, as defined in (2). Thanks to the way the feasible graph is built and pruned, any path going from a data source vertex to any vertex corresponding to an exit layer corresponds to a feasible solution. To find the optimal one, it is thus sufficient to compute the minimum-cost path.

For instance, the right panel of Figure 2 depicts a feasible graph with  $\gamma=5$ . The red and green paths have steepness 3 and the purple one has steepness 4; they are all feasible and one of them will be selected as optimal solution. The blue path, instead, *would* have a steepness of 5 and terminate at node  $(n_4, l_4, 6)$ ; but this node does not exist in the feasible graph. The blue path as a whole is thus infeasible (hence, it is dotted in the figure), and is dropped from the graph.

As for parameter  $\gamma$ , this indicates the *resolution* with which the inference latency is represented by the feasible graph: the smaller the  $\gamma$  values, the fewer the quantization levels (hence, the possible values of vertex depth and path steepness), and the larger the quantization error. As shown in Property 2, such an error can be arbitrarily reduced by increasing  $\gamma$ , thus getting arbitrarily close to the optimum.

*Property 2:* The competitive ratio of FIN-S (i.e., the ratio of the cost of FIN-S's solution to the cost of the optimal one) is  $1 + \frac{1}{\gamma}$ .

*Proof:* FIN-S consists of two steps: (a) building the feasible graph and (b) finding a minimum-cost path on it. Step (b) can be solved to optimality, e.g., through the Bellman-Ford algorithm. Step (a) can use the result of Theorem 4.3 in [17], proving that creating  $\gamma$  replicas of the vertices in the feasible graph results in a cost increase of at most  $\frac{1}{\gamma}$  times the optimum. ■

On the negative side, a high value of  $\gamma$  results in a higher complexity as the minimum-cost traversal will search for more edges in the feasible graph. It is thus essential to set a value of  $\gamma$  that effectively trades off quantization error with complexity. Furthermore, given  $\gamma$ , we envision a  $\lambda$ -proximity approach to reduce the complexity in searching for the vertices and edges to include in the output configuration. Specifically, as the vertices with small depth have a lower number of outgoing feasible edges, searching among vertices with depth close to 1 may not help. We thus limit the search to a  $\lambda$ -proximity ( $1 \leq \lambda \leq \gamma$ ) index on the maximum depth vertices, i.e., with index  $g \in [(\gamma - \lambda), \gamma]$ , and  $\lambda = \gamma$  corresponding to the exhaustive search among all nodes.

The steps involved in FIN-S are described in Algorithm 1 and depicted in Figure 3; the figure also shows how FIN-S is integrated with the construction of the two-plane system graph and the single-plane extended graph. The input to the FIN-S algorithm, given in Algorithm 1 Line 1, consists of the single-plane extended graph ( $\mathcal{G}$ ), application-accuracy, latency requirements ( $\alpha^h, \delta^h$ ), resolution ( $\gamma$ ), and proximity ( $\lambda$ ). The FIN-S procedure consists of two functions, namely, `Create_feasible_graph` (Lines 2–21) and `Configuration_solution` (Lines 22–34). These correspond, respectively, to the *Feasible graph* and *Shortest path* blocks depicted in Figure 3.

We generate the *feasible graph* from the single-plane extended graph  $\mathcal{G}$  in Function I (Lines 2–21) by creating replica vertices (Lines 3–8) and directed edges between vertices (Lines 9–20) such that we meet the latency constraint (Line 13), in terms of the combined processing time and network delay criterion, as given by (4). We compute the energy associated with these edges using (2) and set it as the directed edge weights in Line 16. Function II (Algorithm 1, Lines 22–34) finds the minimum cost (energy) path,  $\pi^h$ , from the *feasible graph* by performing a  $\lambda$ -proximity search (Lines 26–33) on its vertices and edges that minimizes energy during the selected path traversal.

---

**Algorithm 1:** FIN-S: Feasible graph and configuration selection

---

```

1 Input:  $\mathcal{G}, \alpha^h, \delta^h, \gamma, \lambda$ 
2 Function I: Create_feasible_graph ( $\mathcal{G}, \gamma$ ):
3   Create replica vertices
4   for each vertex  $v \in \mathcal{V}$  do
5     for  $g = 0 : 1 : \gamma$  do
6       Create  $v_g$  and include in feasible vertex set
7     end
8   end
9   Create directed edges
10  for each edge  $(v, v') \in \mathcal{E}$  do
11    for  $g_1 = 0 : 1 : \gamma$  do
12      for  $g_2 = 1 : \gamma$  do
13        if (4) holds true then
14          Create edge  $(v_{g_1}, v'_{g_2})$ 
15          Include it in feasible edges set
16          Compute edge weight  $E_h(v_{g_1}, v'_{g_2})$ 
17          using (2)
18        end
19      end
20    end
21  end
22  return feasible graph
23 Function II: Configuration_solution (feasible graph,
24  $\lambda, \gamma$ ):
25   Minimum-cost path traversal on feasible graph,
26   over  $\lambda$ -proximity vertices (i.e.,  $v_g$  with  $g \in [(\gamma - \lambda), \gamma]$ )
27   Initialize  $g_1 = 1$ 
28   Include  $v_{g_1} = (n, l_1^h)$  in  $\pi^h$ 
29   for  $i = 2 : |\mathcal{L}^h|$  do
30     for  $g_2 = (\gamma - \lambda) : \gamma$  do
31       if edge  $(v_{g_1}, v'_{g_2}) : v_{g_1} = (n, l_{i-1}^h)$  then
32         Include  $v'_{g_2}$  in  $\pi^h$  if  $\min_{v'} E^h(v, v')$ 
33         Set  $g_1 = g_2$ ;
34       end
35     end
36   end
37   return  $\pi^h$  as FIN-S output configuration

```

---

### B. Dynamically assigning resources: FIN

Our FIN solution concept integrates the FIN-S procedure described above within an optimization loop, allowing for high-quality decisions about DNN early-exiting and layers' deployment, as well as resource allocation to services.

The basic concept of FIN is presented in Figure 4, and includes two main components. The rightmost one is the FIN-S procedure presented in Sec. III-A: it makes placement and

early-exit decisions given the values of the  $x_b, x_c$ , and  $x_m$  variables. Such values come from an outer loop (leftmost box in Figure 4), running a multi-start global optimization algorithm [25], [26].

Multi-start is an iterative, global optimization technique suitable for non-convex, multi-variable problems. It iteratively improves the current solution by changing one or more variables, i.e., the  $x$ -values. The effect of such changes is evaluated in the inner loop, which (in our case) calls the FIN-S procedure. To avoid getting stuck at a local minimum, the algorithm iteratively selects a starting point for local search while updating the region of attraction each time, and it uses the negative gradient as the descent direction. Starting at multiple points and switching iteratively to updated regions for local search results in a better solution, which yields a lower cost.

FIN terminates when the multi-start algorithm converges; at that point, it returns:

- the values of  $x_b, x_c$ , and  $x_m$  the algorithm itself has converged to;
- the layer placement and early-exit decisions made by the last execution of FIN-S.

We remark that multi-start algorithms have been consistently shown to work extremely close to the optimum [26], even for complex and/or large-scale non-convex problem.

## IV. REFERENCE SCENARIO

This section first introduces the DNN models we consider for our applications, and then it describes the network system we use in our performance evaluation.

**Branched DNNs with early-exits.** To evaluate the performance of our FIN solution, we consider six different applications corresponding to three DNN models with early-exits, namely, B-LeNet, B-AlexNet, and B-ResNet, each trained with two popular datasets. We summarize such mapping in Table II and detail the three models below. We also underline that we consider accuracy as measure of the inference quality.

*B-AlexNet* is a branchy version of AlexNet [27] DNN architecture with early-exits. It has 5 convolution, 1 max-pooling, 3 fully connected, and 3 early-exit blocks. It can be used for image classification [27], video summarization [28], and human activity classification [14], [29]–[31]. The input to the B-AlexNet model is RGB format images, scaled to size  $227 \times 227 \times 3$ . The feature map and complexity of the B-AlexNet architecture is given in Table III; when trained with datasets from different sources and with different sample rate (i.e., samples/category in CIFAR10 and CIFAR100) [15], it results in an accuracy level as given in Table IV.

*B-ResNet* is a branchy version of ResNet110 [32] DNN with early-exits used for image recognition and classification tasks. ResNet110 consists of 3 stages, with each stage comprising a series of residual blocks. The first stage has 18 residual blocks with 16 filters in each block, while the second and third stages have 36 residual blocks, each with 32 and 64 filters in each block (resp.). The residual blocks consist of 2 or 3 convolutional layers, each followed by batch normalization and ReLU activation. The input is size  $32 \times 32 \times 3$ , while the

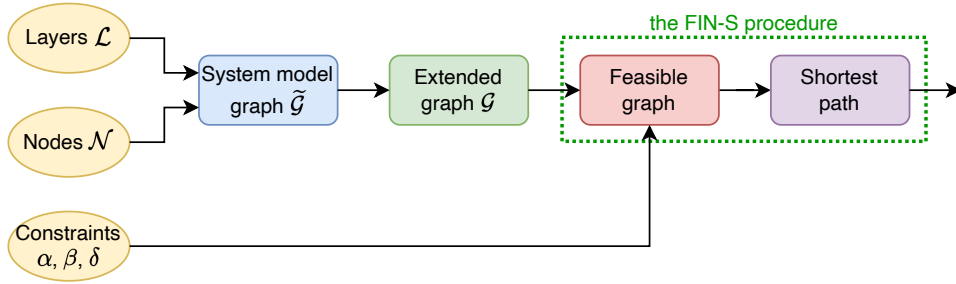


Fig. 3. Solution strategy and steps within FIN-S.

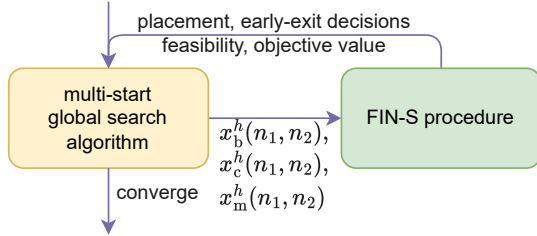


Fig. 4. High-level approach of FIN. The multi-start global optimization algorithm tries out different sets of values for the  $x_b$ ,  $x_c$ , and  $x_m$  variables, using the FIN-S procedure to obtain (i) the layer placement and early-exit decisions, and (ii) information on the feasibility and value of the objective (3a) yielded by those values.

TABLE II  
DETAILS OF THE PRE-TRAINED DNN MODELS WITH EARLY EXITS USED BY THE SIX APPLICATIONS

Application	DNN	Training dataset	#exits	Exit output $\phi^h$ [%]
$h_1$	B-AlexNet	CIFAR100	3	[65.6, 25.2, 9.2]
$h_2$		CIFAR10	3	
$h_3$	B-ResNet	CIFAR100	3	[41.5, 13.8, 44.7]
$h_4$		CIFAR10	3	
$h_5$	B-LeNet	MNIST	2	[94.3, 5.63]
$h_6$		EMNIST	2	

TABLE III  
NUMBER OF INPUT FEATURES AND COMPLEXITY OF THE DNN MODEL BLOCKS

Block	[Number of features, Complexity [MOPs]]		
	B-AlexNet	B-ResNet	B-LeNet
1	[290400, 0.043]	[16384, 0.004]	[4704, 0.118]
2	[186624, 6.711]	[16384, 0.021]	[1600, 0.040]
3	[64896, 10.145]	[16384, 0.021]	[120, 0.048]
4	[64896, 13.523]	[4096, 0.083]	-
5	[43264, 29.045]	[4096, 0.664]	-
Exit-1	[64896, 22.579]	[4096, 0.748]	[120, 0.05]
Exit-2	[43264, 9.056]	[4096, 0.665]	[10, 0.022]
Exit-3	[1000, 0.039]	[10, 0.001]	-

TABLE IV  
INFERENCE ACCURACY OF THE PRE-TRAINED DNN MODELS WITH EARLY EXITS USED BY THE APPLICATIONS

Exit block	Accuracy [%]					
	$h_1$	$h_2$	$h_3$	$h_4$	$h_5$	$h_6$
Exit-1	39.56	56.37	29.97	38.97	91.18	93.54
Exit-2	54.22	78.04	39.93	51.93	96.70	99.20
Exit-3	60.32	85.95	72.21	93.91	-	-

TABLE V  
COMMUNICATION (DOWNLINK (DL)/UPLINK (UL)) CAPACITY AND ENERGY CONSUMPTION OF THE SYSTEM NODES

Node	Power [W]		Traffic DL/UL [Gbps]	Energy DL/UL [nJ/bit]
	Idle	Max		
Mobile	3.1	3.7	0.1	30
Edge	4,096	4,550	560	37
Cloud	11,070	12,300	4,480	12.6

feature map and complexity of B-ResNet architecture layers are given in Table III. The inference accuracy of pretrained B-ResNet using CIFAR10 and CIFAR100 is given in Table IV.

*B-LeNet* is a branchy version of LeNet-5 CNN architecture [33] that can recognize handwritten digits using datasets like MNIST [34] and EMNIST [35]. B-LeNet consists of 8 layers (2 convolutional, 2 subsampling, 3 fully-connected, 1 early-exit). Its feature map and complexity are given in Table III, while the inference accuracy of pretrained B-LeNet using MNIST [34] and EMNIST [35] is given in Table IV.

**Network system.** We use three types of nodes, i.e., mobile, edge, and cloud nodes [36], respectively, associated with the following values of computational capability in trillions operations per second (TOPS) and power consumption in watt (W): [11 TOPS, 6 W]; [153.4 TOPS, 140 W], [312 TOPS, 400 W]. The capability of the communication interface of such nodes are given in Table V [37]–[39]. We recall that such parameters define the bandwidth and computing resource weights of the edges between the vertices in Plane 1 of graph  $\tilde{\mathcal{G}}$ , and the corresponding power consumption. Further, in our experiments, we initially consider one node for each network tier, and then we let the number of mobile nodes increase with the number of users.

Information about the applications to run, e.g., the DNN models to use and associated requirements, as well as network system parameters, are available at the orchestrator.

## V. PERFORMANCE EVALUATION

In this section, we first show the impact of different DNN configurations on the inference latency, quality, and energy cost. Then, we introduce the alternatives against which we compare FIN, and we present the performance obtained for the applications and network system described earlier.

### A. Impact of the DNN configurations

We start by investigating the trade-off among energy consumption, inference accuracy, and inference latency while

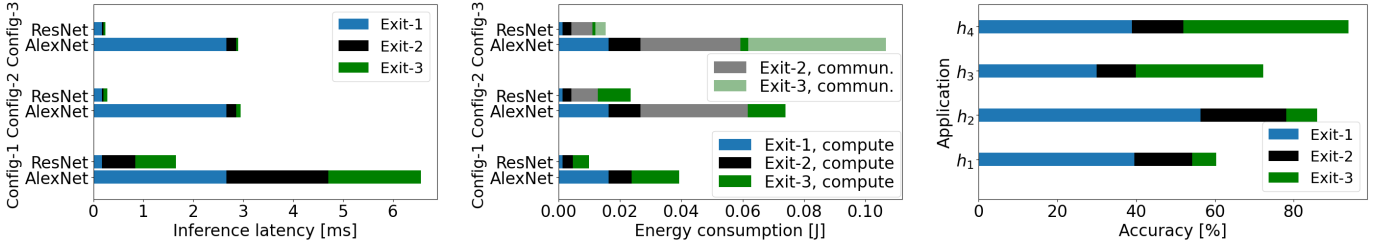


Fig. 5. Impact of the configurations listed in Table VI: inference latency (left) and energy consumption (center) for the B-AlexNet and B-ResNet; inference accuracy for B-AlexNet-based  $h_1$ ,  $h_2$  and B-ResNet-based  $h_3$ ,  $h_4$  (right).

dynamically orchestrating a DNN model on the multi-tiered network for various user applications. To this end, we consider three example configurations for B-AlexNet and B-ResNet, in which the DNN blocks are deployed on the mobile, edge, and cloud nodes as indicated in Table VI.

Figure 5 shows that, as expected, both Config-2 (involving mobile and edge) and Config-3 (involving all tiers) reduce the inference latency overall compared to Config-1 (involving the mobile only). However, it is interesting to observe that the benefit of involving the cloud in the inference task (i.e., Config-3) is negligible when compared to Config-2 where only the edge is added in support to the mobile node.

Looking instead at the energy consumption and the accuracy performance together, while Config-2 and Config-3 make the computing burden at the mobile nodes lighter, they may imply a higher overall cost due to the communication energy expenditure. Such an expenditure appears whenever exit-2 and exit-3 are enabled (i.e., a higher accuracy is required). The increase in communication energy is especially noticeable for B-AlexNet, which has indeed the largest size, resulting in a significant surge (87.6%) in the overall cost when transitioning from Config-1 to Config-2, and a further increase of 28% when shifting from Config-2 to Config-3. Instead, whenever a lower accuracy is acceptable, we can exploit the splits corresponding to exit-1 for all DNNs, and, so doing, reduce both inference latency (from 6.56 ms to 2.67 ms), and energy consumption (from 39.4 mJ with all three exits active to just 16.4 mJ when only exit-1 is activated in Config-1).

In summary, whenever the mobile nodes are used, even for a subset of the DNN blocks, the inference latency grows to such an extent that the reduction in computing time brought by cloud nodes is negligible. Furthermore, whenever we aim at the maximum accuracy, using cloud nodes may lead to very high communication energy costs. Remarkably, however, there exist configurations involving nodes from the different network tiers that can reduce the total as well as the mobile node energy expenditure. It follows that it is possible to identify DNN allocation strategies that improve the sustainability of inference tasks when the application requirements warrant it.

## B. FIN performance

1) *Benchmarks*: We compare FIN against:

- *Multi-constrained path selection (MCP)*, a solution to our problem based on the multi-constrained path selection in [17]. We select [17] because no scheme exists that specifically tackles the problem at hand. [17] finds a path

TABLE VI  
EXAMPLE TEST DEPLOYMENT CONFIGURATIONS OF B-ALEXNET AND B-RESNET

Configuration		Mobile	Edge	Cloud
B-AlexNet B-ResNet	Config-1	All blocks	-	-
	Config-2	$\ell_1^h, \text{exit-1}, \ell_2^h$	$\ell_3^h, \text{exit-2}, \ell_4^h, \ell_5^h, \text{exit-3}$	-
	Config-3	$\ell_1^h, \text{exit-1}, \ell_2^h$	$\ell_3^h, \text{exit-2}, \ell_4^h$	$\ell_5^h, \text{exit-3}$

between source and destination nodes in a graph such that it satisfies the multiple end-to-end constraints on the additive edge weights. MCP applies such an approach to our extended graph  $\mathcal{G}$  to find a feasible solution for the optimization problem in Sec. II-C. To this end, we assign to edge  $v \rightarrow v' \in \mathcal{E}$  the auxiliary weight:

$$\Omega_{\text{sum}}(v, v') = \left( \frac{T^h(v, v') + C^h(v, v')}{\delta^h} + \frac{a(v')}{\alpha^h} \right), \text{ and}$$

$$\Omega_{\text{max}}(v, v') = \max \left( \frac{T^h(v, v'), C^h(v, v')}{\delta^h}, \frac{a(v')}{\alpha^h} \right),$$

for the benchmark  $\text{MCP}_{\text{sum}}$  and  $\text{MCP}_{\text{max}}$ , respectively, where  $a(v')$  is the accuracy associated with the whole sequence of DNN blocks till  $v'$ . The auxiliary weight is computed at each node as the sum (and maximum) of observed values that are scaled by the required constraint (quality and latency) for the  $\text{MCP}_{\text{sum}}$  (and  $\text{MCP}_{\text{max}}$ ) [17]. Then, the minimum-cost path is selected using the auxiliary edge weights in the extended graph yielding the path that best meets the application requirements (quality and latency).

- *Placeto*, a reinforcement learning (RL) approach [18] to find efficient placement of distributed NNs on network nodes by iteratively tuning the policy to perform placement improvements. Specifically, it defines a reward to find a valid placement that fits the resource constraints (available memory) and a penalty if resource utilization crosses the acceptable limit. The  $\text{Placeto}_{\text{conv}}$  [18] has a predefined value of reward that is fixed, while  $\text{Placeto}_{\text{energy}}$  has the reward defined as the inverse of the inference energy consumption.
- *Optimum (Opt)*, obtained through exhaustive search.

2) *Performance of DNN deployments*: The total energy consumption of the B-AlexNet deployment configurations is obtained using  $\text{MCP}_{\text{max}}$ ,  $\text{MCP}_{\text{sum}}$ ,  $\text{Placeto}_{\text{conv}}$ ,  $\text{Placeto}_{\text{energy}}$ ,

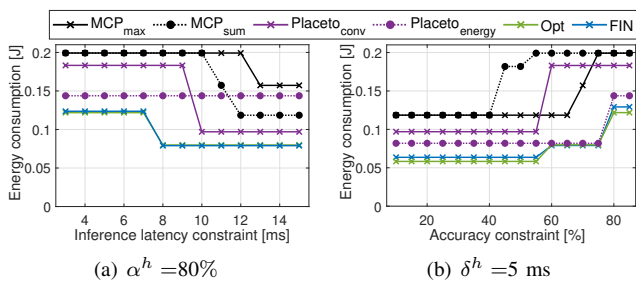


Fig. 6. Total energy consumption of the B-AlexNet configurations, with infinite (top) and limited (bottom) memory (device:edge:cloud=200kb:4Mb:8Mb) available at the network nodes, obtained through Opt, MCP<sub>max</sub>, MCP<sub>sum</sub>, Placeto<sub>conv</sub>, Placeto<sub>energy</sub>, and FIN ( $\gamma=15$ ), as the target inference latency and accuracy vary.

FIN, and Opt. It is presented in Figure 6, with the limited memory available at the network nodes, as the inference accuracy and latency constraints vary. The corresponding breakdown into communication and computation energy consumption is depicted instead in Figure 7. Figure 6 shows that the energy consumption of the DNN configurations yielded by FIN ( $\gamma=15$ ) is very close to the optimum and much less than MCP<sub>max</sub> and MCP<sub>sum</sub>, while meeting the accuracy target  $\alpha^h=80\%$  and the latency target  $\delta^h=5$  ms. We consider that the memory available at the device, edge, and cloud nodes for the applications is 200 kb, 4 Mb, and 8 Mb, respectively. Since the initial DNN block of the AlexNet model requires 290 kb, it cannot be deployed on the mobile node. Our experiments have also revealed that such performance ( $\alpha^h=80\%$ ,  $\delta^h=5$  ms) is achieved by MCP<sub>max</sub>, MCP<sub>sum</sub>, Placeto<sub>conv</sub>, Placeto<sub>energy</sub>, FIN ( $\gamma=15$ ), and Opt deploying, respectively, a set of [0,4,1], [0,4,1], [0,1,4], [0,2,3], [0,3,2], and [0,3,2] blocks on [mobile, edge, cloud] nodes, and each employs exit-3 to meet the target accuracy. In this case, the deployment of the three DNN blocks at the edge and last two on the cloud reduces the energy consumption in FIN and Opt configuration as compared to MCP<sub>max</sub>, MCP<sub>sum</sub>, Placeto<sub>conv</sub>, and Placeto<sub>energy</sub>. For a less stringent latency requirement ( $\alpha^h=80\%$ ,  $\delta^h=12$  ms), MCP<sub>max</sub> deploys [0,1,4], MCP<sub>sum</sub> deploys [0,3,2], Placeto<sub>conv</sub> deploys [0,0,5], and Placeto<sub>energy</sub> deploys [0,2,3], while FIN ( $\gamma=15$ ) and Opt both deploy [0,5,0] blocks, as the larger target latency allows keeping all blocks on the edge, thus reducing energy consumption. In summary, meeting a smaller inference latency target requires a split deployment that increases energy expenditure.

TABLE VII  
EXECUTION TIME [MS] TAKEN BY MCP<sub>MAX</sub>, MCP<sub>SUM</sub>, PLACETO<sub>CONV</sub>, PLACETO<sub>ENERGY</sub>, AND FIN ( $\gamma=15$ ) FOR FINDING THE DEPLOYMENT CONFIGURATION OF EACH DNN

Model	MCP <sub>max</sub>	MCP <sub>sum</sub>	Placeto <sub>conv</sub>	Placeto <sub>energy</sub>	FIN
B-AlexNet	0.591	0.621	1986.3	2481.5	1.025
B-ResNet	0.545	0.589	1576.4	1731.4	0.874
B-LeNet	0.243	0.301	1186.4	1279.9	0.501

This is confirmed by Figures 7(a) and 7(c): for a shorter target inference latency, the communication energy consumption grows, as a split deployment is needed. Similarly, Figures 7(b) and 7(d) underline that, with tighter accuracy constraints, the

best configurations incur higher computation and communication energy. Indeed, inspecting the resulting deployments, it emerges that they require later exit (exit-3) and split deployment. Comparing FIN to its benchmarks, we can easily discern how, even for a moderate value of  $\gamma=15$ , FIN virtually always matches the optimum and significantly outperforms MCP<sub>max</sub>, MCP<sub>sum</sub>, Placeto<sub>conv</sub>, and Placeto<sub>energy</sub> in all cases. We have similarly evaluated the B-ResNet (Figure 8) and the B-LeNet deployments (Figure 9), obtaining essentially the same results. Looking at FIN’s energy consumption, a similar effect as in Figures 8 and 9 emerges, with FIN closely matching the optimum for  $\gamma=15$  and similar to the case of B-AlexNet.

It is worth remarking that Placeto<sub>conv</sub> algorithm deploys tasks on computing nodes with the aim to minimize the inter-node communication overhead while meeting the nodes’ memory constraints. However, it does so without optimizing resource or energy, and it does not perform dynamic task selection based on quality constraints. All these result in its relatively poor energy consumption performance. The Placeto<sub>energy</sub> algorithm, on the other hand, deploys the tasks on computing nodes with the aim to minimize the energy consumption but does not have resource optimization or quality constraint-based dynamic task selection.

Finally, Table VII lists the overall execution times taken (on average) by MCP<sub>max</sub>, MCP<sub>sum</sub>, Placeto<sub>conv</sub>, Placeto<sub>energy</sub>, and FIN ( $\gamma=15$ ) for obtaining the deployment configuration of the B-AlexNet, B-ResNet, and B-LeNet models. The results have been obtained by using a Lenovo ThinkPad P1 Gen 3 with i7-10750H CPU (2.6 GHz, 32 GB RAM). The execution time of FIN is less than twice that of MCP and many orders of magnitude (200 times) lower than (high complexity, RL-based) Placeto. We also emphasize that the orchestrator will run on much more powerful hardware than ours, hence, the overhead associated with running FIN when inference rate requirements change is acceptable – and, indeed, often negligible.

3) *Multi-application scenario*: We now apply FIN to the deployment of the six applications listed in Table II. Using the pre-trained DNN models for inference, we investigate the impact of an increasing number of users on the system performance. We consider that 0.5% of the edge and cloud computing resources and 200 kb, 4 Mb, and 8 Mb, respectively, of memory are available at the device, edge, and cloud nodes for the applications inference execution. The application requirements, [inference latency [ms], accuracy [%]], are set to [5,55], [5, 55], and [0.1, 93] for  $h_{1-2}$ ,  $h_{3-4}$ , and  $h_{5-6}$  applications, respectively. The FIN framework optimizes the resource (communication, compute, and memory) allocation for each application at each node, under the given system-level resource constraints. In addition, FIN minimizes the overall energy consumption and finds the best configuration for each application that meets the user requirements (inference latency and quality).

Figure 10 shows the reduction in the energy consumption, due to computing and communication, provided by FIN ( $\gamma=15$ ) with respect to MCP<sub>max</sub>, MCP<sub>sum</sub>, Placeto<sub>conv</sub>, and Placeto<sub>energy</sub>, for the running applications  $h_{1-6}$  that require one-inference-per-second-per-user. Notice that the size of this scenario renders obtaining the optimum impractical. One can

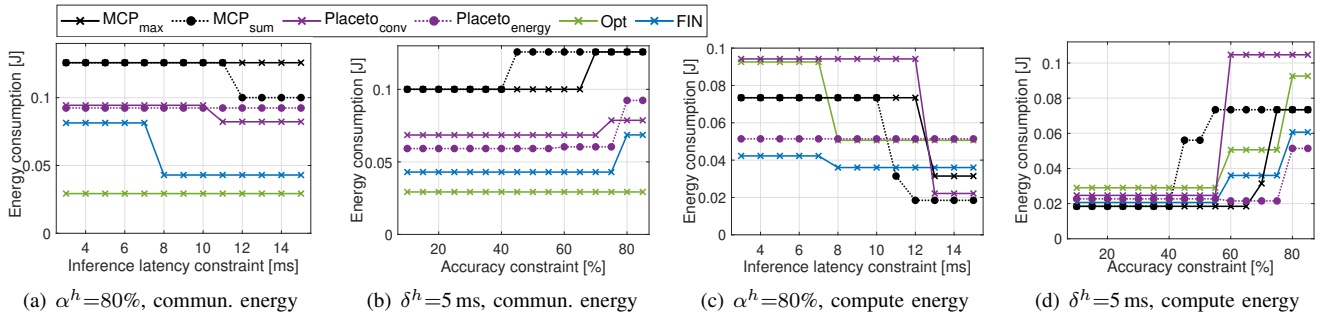


Fig. 7. Computation and communication energy consumption of  $MCP_{max}$ ,  $MCP_{sum}$ ,  $Placeto_{conv}$ ,  $Placeto_{energy}$ ,  $FIN$ , and  $Opt$  for B-AlexNet configurations, as the constraints vary.

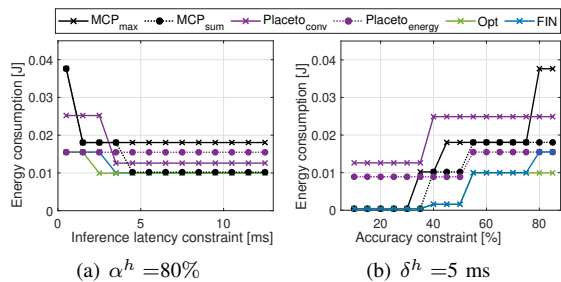


Fig. 8. Total energy consumption of the B-ResNet configurations, with infinite (top) and limited (bottom) memory (device:edge:cloud=200kb:4Mb:8Mb) available at the network nodes, obtained through  $Opt$ ,  $MCP_{max}$ ,  $MCP_{sum}$ ,  $Placeto_{conv}$ ,  $Placeto_{energy}$  and  $FIN$  ( $\gamma=15$ ), as the target inference latency and accuracy vary.

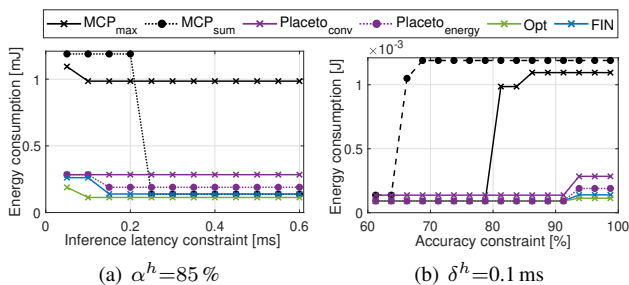


Fig. 9. Total energy consumption of the B-LeNet configurations, with infinite and limited memory (device:edge:cloud=200kb:4Mb:8Mb) available at the network nodes, obtained through  $Opt$ ,  $MCP_{max}$ ,  $MCP_{sum}$ ,  $Placeto_{conv}$ ,  $Placeto_{energy}$ , and  $FIN$  ( $\gamma=15$ ), as the target inference latency and accuracy vary.

observe how  $FIN$  deploys the DNN model configuration that entails an overall energy consumption that is 70%–85% of the benchmarks for all the considered DNNs (Figure 10).

Also, the  $MCP_{max}$ ,  $MCP_{sum}$ ,  $Placeto_{conv}$ , and  $Placeto_{energy}$  approaches lean towards the deployment of DNN layers on the mobile and cloud side more often. On the contrary, thanks to  $FIN$ 's ability to perform resource optimization in an energy-aware manner, it takes full advantage of all three tiers and, in particular, of the edge (Figure 12(left)).

Not only does  $FIN$  achieve better energy efficiency, but it also surpasses  $MCP$  and  $Placeto$  in terms of success probability across all applications (Figure 12(right)). In contrast to  $MCP$ 's and  $Placeto$ 's high failure probability (over 30% for B-Alexnet, 18% for B-ResNet, and 14% for B-Lenet),  $FIN$  sees less than 4% of users failing to meet the latency and

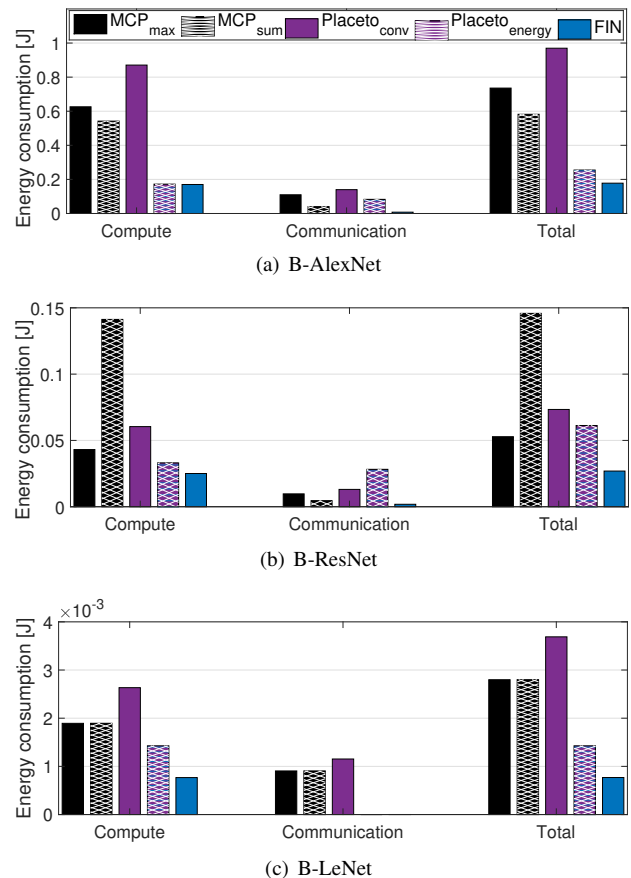


Fig. 10. Multi-application scenario: Energy consumption through  $MCP_{max}$ ,  $MCP_{sum}$ ,  $Placeto_{conv}$ ,  $Placeto_{energy}$ , and  $FIN$ . For  $FIN$ , we set  $\gamma=15$ .

accuracy constraints. Consistently, Figure 13 shows that, with high probability, and unlike  $MCP_{max}$ ,  $MCP_{sum}$ ,  $Placeto_{conv}$ , and  $Placeto_{energy}$ ,  $FIN$  can deploy the applications blocks all the way to exit-3 whenever required, while it rightfully enables the earliest exit split whenever the accuracy constraint allows it (e.g., for  $h_2$  and  $h_6$ ).

## VI. RELATED WORK

Our work lies at the intersection of two major fields, namely, model split (or, partitioning) and resource-aware ML.

**Early exit and model splitting [8]:** Early exit models have been introduced in [14], which also raises the issue of how

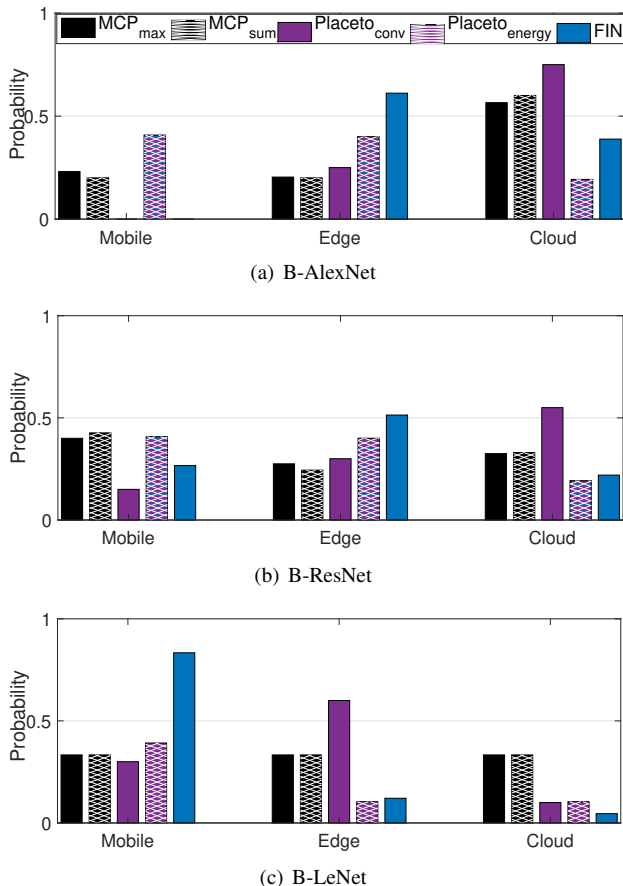


Fig. 11. Multi-application scenario: Probability of DNN block deployment on the multi-tier nodes in MCP<sub>max</sub>, MCP<sub>sum</sub>, Placeto<sub>conv</sub>, Placeto<sub>energy</sub>, and FIN. For FIN, we set  $\gamma=15$ .

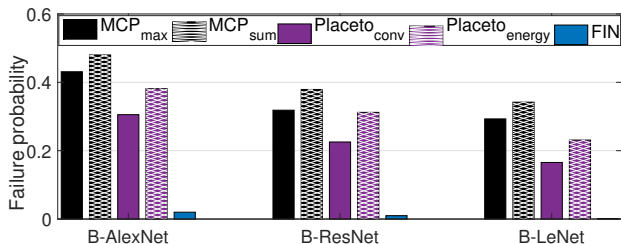


Fig. 12. Multi-application scenario: Probability that the selected configuration fails to meet the constraints. For FIN, we set  $\gamma=15$ .

to place the early-exit layers, i.e., how to make the DNN topology *branchy*. In the same context, [9] tackles distributed scenarios and seeks to adapt the placement of exit layers to the available resources in the near-edge, edge, and cloud segments of the network. [10] pursues a similar approach in IoT scenarios, minimizing the usage of edge resources. More recent work [12] performs DNN splitting in real time, with the aim of adapting to changes in channel conditions. The recent work [13] widens the focus and accounts for the location of the users that need the inference task. In a similar vein, [40] seeks to adapt the DNN behavior, i.e., whether or not early-exit branches are employed, to the specific *class* being predicted, with the goal of speeding up the detection of some more sensitive types of data.

Earlier approaches seek to split DNNs at naturally-occurring suitable locations, a.k.a. *bottlenecks* [41], [42]. If no bottleneck is available, the related problem of *bottleneck injection* arises. The goal is to change the topology of the DNN with the aim of creating suitable points to insert an early-exit layer. These techniques have been pioneered by [43], [44], and often use pairs of encoder and decode layers. Bottleneck injection can be performed in a content-aware fashion, as in [45], [46]. Also, [5] underlines that collaborative DNN partitioning and task offloading in resource-constrained edge-IoT networks can meet the DNN inference deadline requirements.

**Model compression** is the notion of extracting a smaller (hence, simpler and/or quicker to run) model from a more complex one, while retaining most of its usefulness. From the point of view of model adaptation to the nodes' resources, it can be seen as an alternative, or a complement, to model splitting. The most popular compression technique is model pruning, which consists in removing some weights from a DNN based on the so-called *lottery ticket* hypothesis [47], assuming that very few of the parameters of a model actually impact its performance. Many works seek the best way to choose the parameters to prune, following an iterative approach [48], [49] and observing the evolution of the parameter values. Knowledge distillation [50] is an alternative compression technique, where a "student" model is trained to mimic the decisions of a "teacher" model without requiring the access to the original training data. Distillation has later combined with generative models [51] and domain adaptation [52].

**Resource-aware ML** is, broadly speaking, concerned with adapting the distributed ML task to perform (whether it is training or inference) and the available resources. Works in this field often focus on selecting the best nodes to exploit, accounting for their speed [53], size of local dataset [54], and feature-richness thereof [55], as well as any communication issues they may experience [56]. A more recent trend, closer in spirit to model partitioning, consists in changing the learning task to fit the available resources, e.g., by selecting the most appropriate model [57].

A related trend is reliability in distributed ML. The main goal of this line of work is ensuring that all nodes involved in the ML task provide timely and high-quality updates, despite communication issues [58] and the presence of malicious nodes [59]. Reliability might be at odd with fairness issues, and a balance between the two goals is sought in [60].

**Novelty:** In summary, to the best of our knowledge, our work is the first to jointly investigate (i) how ML model splitting should be performed and (ii) where the different model blocks should be deployed, (iii) in the case of models with early exits as well as in the presence of inference requirements and constraints on the computational and networking resources in multi-tier systems. We mention that an early version of this work can be found in our conference paper [61] where, however, we considered that each application was operating in isolation and, hence, we did not address the optimal allocation of bandwidth, computing, and memory resources across different applications.

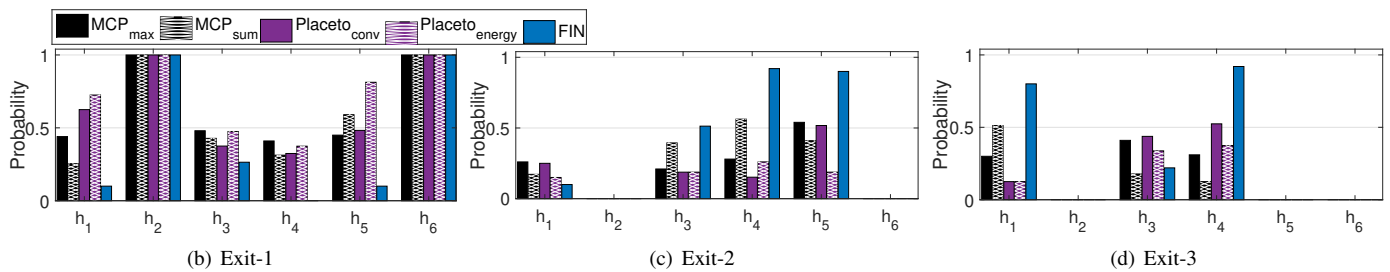


Fig. 13. Multi-application scenario: Probability of selecting a DNN exit point for inference. For FIN, we set  $\gamma=15$ .

## VII. CONCLUSIONS

We tackled the problem of deploying intelligent applications over a mobile-edge-cloud system in an energy-efficient way. As such applications require the execution of DNN inference tasks, we achieved this goal by exploiting DNNs with multi-branched architectures, and near-optimally selecting and placing the model early exists on the system nodes so as to minimize the inference energy cost. Our approach leverages a three-stage graph-modeling method that yields a compact representation of the system characteristics and constraints as well as of the applications requirements. We then used such a model to formulate an optimization problem, which resulted to be NP-hard, and to develop a low-complexity, yet effective, algorithmic solution to the problem. Extensive results obtained with different, co-existing intelligent applications demonstrate that our solution performs very closely to the optimum and, compared to the existing alternatives, it can reduce the inference energy consumption by more than 89%.

## REFERENCES

- [1] Z. Li, F. Liu, W. Yang, S. Peng, and J. Zhou, "A survey of convolutional neural networks: analysis, applications, and prospects," *IEEE Transactions on neural networks and learning systems*, vol. 33, no. 12, pp. 6999–7019, 2022.
- [2] A. B. Nassif, I. Shatin, I. Attili, M. Azzeh, and K. Shaalan, "Speech recognition using deep neural networks: A systematic review," *IEEE access*, vol. 7, no. 1, pp. 19 143–19 165, 2019.
- [3] Y. Goldberg, *Neural network methods for natural language processing*. Springer Nature, 2022.
- [4] I. Azimi, A. Anzanpour, A. M. Rahmani, T. Pahikkala, M. Levorato, P. Liljeberg, and N. Dutt, "Hich: Hierarchical fog-assisted computing architecture for healthcare iot," *ACM Transactions on Embedded Computing Systems*, vol. 16, no. 5, pp. 1–20, Sep. 2017.
- [5] X. Zhang, M. Mounesan, and S. Debroy, "EFFECT-DNN: energy-efficient edge framework for real-time DNN inference," in *Proc. IEEE WoWMoM*, Boston, MA, USA, June 2023, pp. 10–20.
- [6] 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.
- [7] Q. Luo, S. Hu, C. Li, G. Li, and W. Shi, "Resource scheduling in edge computing: A survey," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 4, pp. 2131–2165, 2021.
- [8] Y. Matsubara, M. Levorato, and F. Restuccia, "Split computing and early exiting for deep learning applications: Survey and research challenges," *ACM Computing Surveys*, vol. 55, no. 5, pp. 1–30, 2022.
- [9] 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 Comput. Archit. News*, vol. 52, no. 4, p. 615–629, Apr. 2017.
- [10] Y.-T. Yang and H.-Y. Wei, "Edge-iot computing and networking resource allocation for decomposable deep learning inference," *IEEE Internet of Things Journal*, vol. 10, no. 6, pp. 5178–5193, 2022.
- [11] W. Miao, Z. Zeng, L. Wei, S. Li, C. Jiang, and Z. Zhang, "Adaptive dnn partition in edge computing environments," in *IEEE ICPADS*, 2020.
- [12] J. Lee, H. Lee, and W. Choi, "Wireless channel adaptive dnn split inference for resource-constrained edge devices," *IEEE Communications Letters*, vol. 27, no. 6, pp. 1520–1524, 2023.
- [13] W. Fan, L. Gao, Y. Su, F. Wu, and Y. Liu, "Joint dnn partition and resource allocation for task offloading in edge-cloud-assisted iot environments," *IEEE Internet of Things Journal*, vol. 10, no. 12, pp. 10 146–10 159, 2023.
- [14] S. Teerapittayanon *et al.*, "Branchynet: Fast inference via early exiting from deep neural networks," in *IEEE ICPR*, 2016.
- [15] R. Dong, Y. Mao, and J. Zhang, "Resource-constrained edge ai with early exit prediction," *Journal of Communications and Information Networks*, vol. 7, no. 2, pp. 122–134, Jun. 2022.
- [16] S. Laskaridis, A. Kouris, and N. D. Lane, "Adaptive inference through early-exit networks: Design, challenges and directions," in *Proceedings of the 5th International Workshop on Embedded and Mobile Deep Learning*, 2021, pp. 1–6.
- [17] G. Xue, A. Sen, W. Zhang, J. Tang, and K. Thulasiraman, "Finding a path subject to many additive qos constraints," *IEEE/ACM Transactions on networking*, vol. 15, no. 1, pp. 201–211, 2007.
- [18] R. Addanki, S. B. Venkatakrishnan, S. Gupta, H. Mao, and M. Alizadeh, "Placeto: Learning generalizable device placement algorithms for distributed machine learning," in *Proc. Intern. Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2019.
- [19] L. Gouveia, M. Leitner, and M. Ruthmair, "Layered graph approaches for combinatorial optimization problems," *Computers & Operations Research*, vol. 102, pp. 22–38, Feb. 2019.
- [20] —, "Extended formulations and branch-and-cut algorithms for the black-and-white traveling salesman problem," *European Journal of Operational Research*, vol. 262, no. 3, pp. 908–928, Nov. 2017.
- [21] L. Gouveia and M. Ruthmair, "Load-dependent and precedence-based models for pickup and delivery problems," *Computers & Operations Research*, vol. 63, pp. 56–71, Nov. 2015.
- [22] S. Yang, F. Li, S. Trajanovski, X. Chen, Y. Wang, and X. Fu, "Delay-aware virtual network function placement and routing in edge clouds," *IEEE Transactions on Mobile Computing*, vol. 20, no. 2, pp. 445–459, 2021.
- [23] B. Ramamurthy, G. Rouskas, and K. Sivalingam, *Next-Generation Internet: Architectures and Protocols*. Cambridge, Feb. 2011.
- [24] F. K. Hwang *et al.*, "Steiner tree problems," *Networks*, vol. 22, no. 1, pp. 55–89, 1992.
- [25] F. Glover, "A template for scatter search and path relinking," in *Artificial Evolution, Lecture Notes in Computer Science*, J.-K. Hao, E. Lutton, E. Ronald, M. Schoenauer, and D. Snyers, Eds., vol. 1363. Berlin, Heidelberg: Springer Berlin Heidelberg, 1998, pp. 1–51.
- [26] C. Voglis and I. Lagaris, "Towards 'ideal multistart'. a stochastic approach for locating the minima of a continuous function inside a bounded domain," *Applied Mathematics and Computation*, vol. 213, no. 1, pp. 216–229, 2009.
- [27] O. Elharrouss, Y. Akbari, N. Almaadeed, and S. Al-Maadeed, "Backbones-review: Feature extraction networks for deep learning and deep reinforcement learning approaches," Jun. 2022. [Online]. Available: <https://arxiv.org/abs/2206.08016>
- [28] J. Lei, Q. Luan, X. Song, X. Liu, D. Tao, and M. Song, "Action parsing-driven video summarization based on reinforcement learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 7, pp. 2126–2137, 2019.
- [29] F. Serpush and M. Rezae, "Complex human action recognition using a hierarchical feature reduction and deep learning-based method," *SN Computer Science*, vol. 2, no. 94, pp. 1–15, Feb. 2021.
- [30] A. Ullah, J. Ahmad, K. Muhammad, M. Sajjad, and S. W. Baik, "Action recognition in video sequences using deep bi-directional LSTM with CNN features," *IEEE Access*, vol. 6, pp. 1155–1166, 2018.

- [31] S. Darafsh, S. S. Ghidary, and M. S. Zamani, "Real-time activity recognition and intention recognition using a vision-based embedded system," *CoRR*.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2015.
- [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] Y. LeCun, C. Cortes, and C. Burges, "Mnist handwritten digit database," *ATT Labs*, 2010.
- [35] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik, "EMNIST: Extending mnist to handwritten letters," in *Proc. IEEE International Joint Conference on Neural Networks (IJCNN)*, May 2017, pp. 2921–2926.
- [36] F. Malandrino, C. F. Chiasserini, and G. di Giacomo, "Efficient distributed DNNs in the mobile-edge-cloud continuum," *IEEE/ACM Transactions on Networking (early access)*, vol. 31, no. 4, pp. 1702–1716, 2023.
- [37] F. Jalali, K. Hinton, R. Ayre, T. Alpcan, and R. S. Tucker, "Fog computing may help to save energy in cloud computing," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 5, pp. 1728–1739, 2016.
- [38] Y. Li, A.-C. Orgerie, I. Rodero, B. L. Amersho, M. Parashar, and J.-M. Menaud, "End-to-end energy models for edge cloud-based IoT platforms: Application to data stream analysis in IoT," *Future Generation Computer Systems*, vol. 87, pp. 667–678, Oct. 2018.
- [39] L. Sun, H. Deng, R. K. Sheshadri, W. Zheng, and D. Koutsonikolas, "Experimental evaluation of WiFi active power/energy consumption models for smartphones," *IEEE Transactions on Mobile Computing*, vol. 16, no. 1, pp. 115–129, Mar. 2017.
- [40] M. Ayyat, T. Nadeem, and B. Krawczyk, "Classynet: Class-aware early exit neural networks for edge devices," *IEEE Internet of Things Journal*, 2023.
- [41] C. Buciluă, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *ACM SIGKDD*, 2006.
- [42] C.-H. Chiang, P. Liu, D.-W. Wang, D.-Y. Hong, and J.-J. Wu, "Optimal branch location for cost-effective inference on branchynet," in *IEEE Big Data*, 2021.
- [43] Y. Matsubara, S. Baidya, D. Callegaro, M. Levorato, and S. Singh, "Distilled split deep neural networks for edge-assisted real-time systems," in *ACM HotEdgeVideo*, 2019.
- [44] A. E. Eshratifar, A. Esmaili, and M. Pedram, "Bottlenet: A deep learning architecture for intelligent mobile cloud computing services," in *IEEE/ACM ISLPED*, 2019.
- [45] J. C. Lee, Y. Kim, S. Moon, and J. H. Ko, "A splittable dnn-based object detector for edge-cloud collaborative real-time video inference," in *IEEE AVSS*, 2021.
- [46] Y. Matsubara, D. Callegaro, S. Singh, M. Levorato, and F. Restuccia, "Bottlefit: Learning compressed representations in deep neural networks for effective and efficient split computing," in *IEEE WoWMoM*, 2022.
- [47] J. Frankle and M. Carbin, "The lottery ticket hypothesis: Finding sparse, trainable neural networks," in *ICLR*, 2018.
- [48] C. M. J. Tan and M. Motani, "Dropnet: Reducing neural network complexity via iterative pruning," in *ICML*, 2020.
- [49] T. Zhang, S. Ye, K. Zhang, J. Tang, W. Wen, M. Fardad, and Y. Wang, "A systematic dnn weight pruning framework using alternating direction method of multipliers," in *ECCV*, 2018.
- [50] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *International Journal of Computer Vision*, 2021.
- [51] Y. Huang and Y. Yu, "Distilling deep neural networks with reinforcement learning," in *IEEE ICIA*, 2018.
- [52] S. Tang, Y. Shi, Z. Ma, J. Li, J. Lyu, Q. Li, and J. Zhang, "Model adaptation through hypothesis transfer with gradual knowledge distillation," in *IEEE/RSS J IROS*, 2021.
- [53] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive federated learning in resource constrained edge computing systems," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1205–1221, 2019.
- [54] F. Malandrino and C. F. Chiasserini, "Federated learning at the network edge: When not all nodes are created equal," *IEEE Communications Magazine*, vol. 59, no. 7, p. 68–73, Jul. 2021.
- [55] H. Wu and P. Wang, "Fast-convergent federated learning with adaptive weighting," *IEEE Transactions on Cognitive Communications and Networking*, vol. 7, no. 4, pp. 1078–1088, 2021.
- [56] Y. Zhou, Q. Ye, and J. C. Lv, "Communication-Efficient Federated Learning with Compensated Overlap-FedAvg," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 01, pp. 192–205, Jan. 2022.
- [57] F. Paissan, A. Ancilotto, A. Brutti, and E. Farella, "Scalable neural architectures for end-to-end environmental sound classification," in *IEEE ICASSP*, 2022.
- [58] F. Ang, L. Chen, N. Zhao, Y. Chen, W. Wang, and F. R. Yu, "Robust federated learning with noisy communication," *IEEE Transactions on Communications*, vol. 68, no. 6, pp. 3452–3464, 2020.
- [59] S. Li, Y. Cheng, W. Wang, Y. Liu, and T. Chen, "Learning to detect malicious clients for robust federated learning," *ArXiv*, vol. abs/2002.00211, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:211011146>
- [60] T. Li, S. Hu, A. Beirami, and V. Smith, "Ditto: Fair and robust federated learning through personalization," in *ICML*, 2021.
- [61] C. Singhal, Y. Wu, F. Malandrino, M. Levorato, and C. F. Chiasserini, "Resource-aware deployment of dynamic dns over multi-tiered interconnected systems," in *IEEE INFOCOM*, 2024.