

Dynamic DNN Decomposition for Lossless Synergistic Inference

Beibei Zhang, Tian Xiang, Hongxuan Zhang, Te Li, Shiqiang Zhu, Jianjun Gu
Intelligent Robotics Research Center, Zhejiang Lab, Hangzhou, China
{beibei, txiang, hongxuan, lite, zhusq, jgu}@zhejianglab.com

Abstract—Deep neural networks (DNNs) sustain high performance in today’s data processing applications. DNN inference is resource-intensive thus is difficult to fit into a mobile device. An alternative is to offload the DNN inference to a cloud server. However, such an approach requires heavy raw data transmission between the mobile device and the cloud server, which is not suitable for mission-critical and privacy-sensitive applications such as autopilot. To solve this problem, recent advances unleash DNN services using the edge computing paradigm. The existing approaches split a DNN into two parts and deploy the two partitions to computation nodes at two edge computing tiers. Nonetheless, these methods overlook collaborative device-edge-cloud computation resources. Besides, previous algorithms demand the whole DNN re-partitioning to adapt to computation resource changes and network dynamics. Moreover, for resource-demanding convolutional layers, prior works do not give a parallel processing strategy without loss of accuracy at the edge side. To tackle these issues, we propose D^3 , a dynamic DNN decomposition system for synergistic inference without precision loss. The proposed system introduces a heuristic algorithm named horizontal partition algorithm to split a DNN into three parts. The algorithm can partially adjust the partitions at run time according to processing time and network conditions. At the edge side, a vertical separation module separates feature maps into tiles that can be independently run on different edge nodes in parallel. Extensive quantitative evaluation of five popular DNNs illustrates that D^3 outperforms the state-of-the-art counterparts up to $3.4\times$ in end-to-end DNN inference time and reduces backbone network communication overhead up to $3.68\times$.

Index Terms—Distributed computing, Edge computing, DNN inference acceleration.

I. INTRODUCTION

The proliferation of mobile devices such as smartphones and smart robotics brings a tremendous amount of data generated from users. On mobile devices, pervasive data processing applications including machine translation [1], object detection [2], and many others process the data and share the data over a communication network. The results of these applications should be exceedingly accurate [3]. As such, deep neural networks (DNNs) become one of the de-facto solutions in the data processing applications due to its high accuracy. However, DNN inference requires abundant computation resources and consumes considerable energy [4]. Therefore, it is not suitable to deploy DNNs on mobile devices that have restricted computation power and limited energy.

One of the popular approaches to tackle this issue is to offload the intensive DNN inference to a cloud server [5]. However, this approach involves transferring raw data collected by the data processing applications to the cloud server

through a backbone network, which incurs long transmission delays and privacy concerns. With the observation that the intermediate result size of a DNN is significantly smaller than the raw data size, the idea of only transferring the intermediate result to the cloud server comes into existence. Recent research employs the *edge computing paradigm* that comprises three computing tiers (i.e., device¹, edge, cloud) to solve the problem. Specifically, previous studies propose to split a DNN into two parts according to the processing time of DNN layers and the data transmission delay between two layers [6]–[9]. Generally, a mobile device collects raw data and passes the input to the first DNN partition located at an edge node. Next, the edge node processes the input and transmits the intermediate results to a cloud server. Finally, the cloud server handles the rest of the DNN inference that requires more processing capability. The collaborative computation leverages the resources provided by the edge and the cloud, reducing DNN inference latency and communication overhead over the network core.

Nonetheless, the rapid development of hardware makes it possible to perform partial DNN inference on mobile devices [10]. For instance, the latest smartphone has an octa-core CPU running up to 3.1 GHz and a GPU with 1.37 TFLOPS [11]. In this context, current solutions fail to leverage the synergistic device-edge-cloud computation power [8]. Besides, prior works require re-partitioning the whole DNN to accommodate dynamics of computation resources and network bandwidth [12]. Moreover, in a joint DNN inference pipeline, the node with the most processing time becomes the bottleneck of the overall inference. For an edge node with limited resources compared with a cloud node, if it is the bottleneck of the collaborative inference, the state-of-the-art does not provide a parallel processing strategy for convolutional layers assigned to the edge node without loss of accuracy [13].

To address these limitations, in this paper, we present a **dynamic DNN decomposition system** named D^3 for lossless synergistic inference. D^3 accelerates DNN inference by leveraging the synergistic device-edge-cloud computation power without precision loss. In D^3 , we employ a regression model that takes computation resources and DNN layer configurations as input and estimates the processing time of DNN layers. According to the per-layer execution time and the trans-

¹To avoid ambiguity, we use the term “device” to represent the “device tier of edge computing paradigm” and the term “node” to denote the “computing device” throughout this paper.

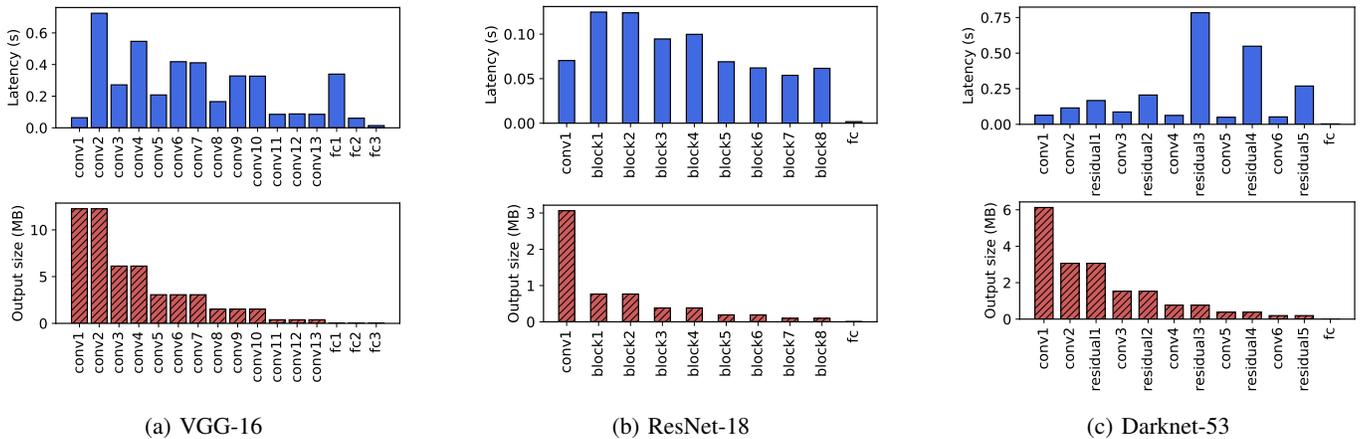


Fig. 1: With the input size of $3 \times 224 \times 224$, we measure the layer-wise inference latency and the per-layer output size of VGG-16, ResNet-18, and Darknet-53 on a Raspberry Pi 4 model B running at 1.5GHz with 4 GB system memory. Each block or residual contains several convolutional layers.

mission delay between layers, a heuristic algorithm, which we refer to as **horizontal partition algorithm** (HPA), partitions a DNN into three parts, each of which runs on one computing tier. In the case of changes in the DNN layer processing time or the transmission delay, HPA can partially adjust the DNN segmentation, adapting to the changes via local updates. At a finer granularity, a **vertical separation module** (VSM) further splits a stack of feature maps of convolutional layers into blocks spatially. We assign a block of correlated feature maps to an edge node. In this way, we can run feature maps of convolutional layers on the edge nodes in parallel independently, utilizing the edge resources.

We implement D^3 and evaluate its performance based on ImageNet dataset [14] using various real-world DNNs including AlexNet [15], VGG-16 [16], ResNet-18 [17], Darknet-53 [18], and Inception-v4 [19]. Our experimental results show that compared with the state-of-the-art counterparts, D^3 accelerates the DNN inference time up to $3.4\times$ and reduces the communication overhead up to $3.68\times$.

The rest of this paper is organized as follows. First, we review related work in section II. Then, we describe the proposed system in section III. Section IV introduces the implementation and the test-bed of D^3 . Extensive comparisons and evaluations between D^3 and its counterparts follow in section V. Finally, we conclude in section VI.

II. RELATED WORK

Before delving into the D^3 system, we first provide an overview of the existing DNN inference acceleration mechanisms. The status-quo approaches that are related to our method can be classified into two categories. The first group aims to split a DNN model into partitions according to per-layer processing time and inter-layer transmission delay. The partitions are distributed to multiple computation nodes. We denote such approaches as horizontal partition. The second group, which we refer to as vertical separation, focuses on

dividing the feature map of a convolutional layer spatially to multiple tiles hosted by multiple computation nodes.

To accelerate DNN inference without sacrificing accuracy, a few previous works focus on the DNN horizontal partition. Neurosurgeon [7] offloads computation from resource-constrained mobile devices to cloud servers. It splits a DNN of chain topology at a layer granularity to minimize processing latency and energy consumption. IONN [20] models a chain topology DNN as an auxiliary DAG and finds the optimal incremental offloading with the shortest path algorithm on the DAG. DADS [8] extends the layer-wise partition to multi-branch DNNs represented by DAGs and exploits the min-cut algorithm to find the optimal cutting points. It employs edge computing and deploys DNN layers to an edge node and a cloud server. However, DADS cannot generalize the min-cut approach to separate a DNN into more than two parts. DINA [21] presents an adaptive partition algorithm to divide DNN layers into pieces that can be smaller than a layer. It offloads the pieces to fog nodes based on a swap-matching algorithm. DDNN [6] introduces the idea of early-exit that sacrifices accuracy in exchange for reducing inference delay and communication overhead. However, this requires a specific training process to ensure the prediction confidence. Thus, it is not suitable for accelerating the lossless inference of the trained DNN discussed in this paper. Edgent [22] combines the chain topology DNN segmentation and the early-exit method and performs collaborative device-edge DNN inference. SPINN [23] applies the early-exit strategy to progressive inference and designs a scheduler to flexibly handle service-level agreements.

The convolutional layer, which is commonly deployed in current data processing applications, is one of the most resource-intensive components of a DNN [24]. Fig. 1 illustrates the per-layer inference latency and the inter-layer output size of three widely used DNNs. We notice that some convolutional layers require substantial computation resources.

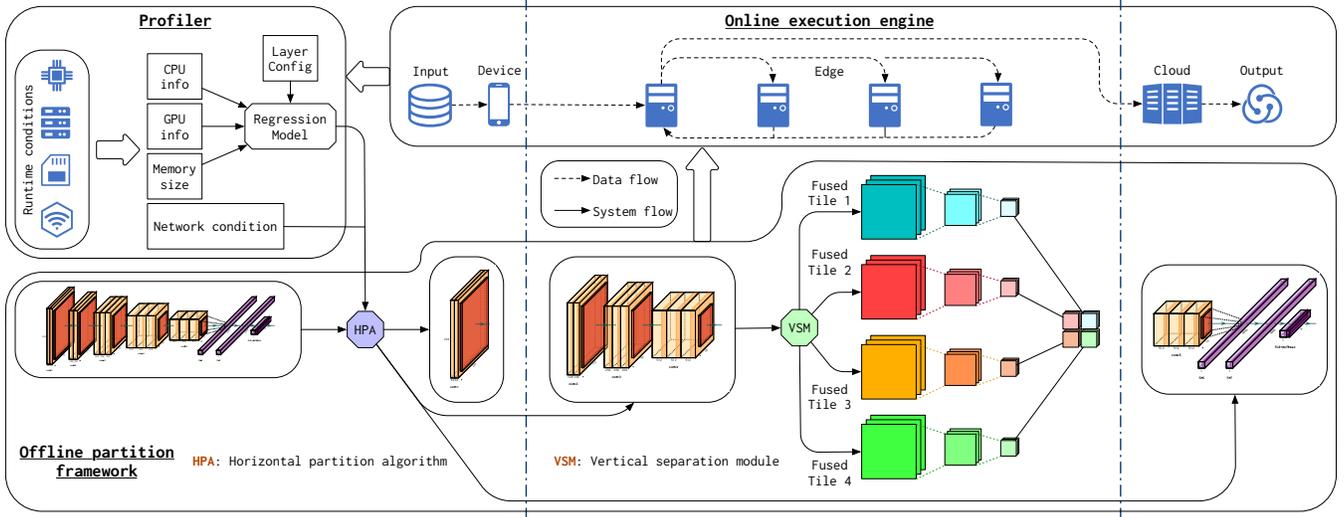


Fig. 2: D^3 system architecture.

To explore inference parallelism for convolutional layers, prior works aim to separate feature maps spatially. MoDNN [25] proposes layer-wise parallelism that divides one feature map of a convolutional layer into pieces. Each computation node executes a part of the feature map and generates an output. A host node gathers the output and re-partitions the feature map for parallel processing of the next convolutional layer. This procedure results in significant communication overhead. DeepThings [13] removes the communication overhead by introducing a fused tile partition (FTP) that slices a stack of correlated feature maps spatially and distributes the stack to a computation node. AOFL [9] extends the idea of fused tiles and offers an algorithm to find the optimal tile partition according to resources of each computation node.

III. PROPOSED SYSTEM

This section formally describes D^3 , a system that dynamically decomposes a DNN to segments for collaborative inference over device, edge, and cloud with no precision loss. We first discuss the edge computing framework, based on which our system is designed. We then give an overview of D^3 , followed by the modeling of our system. A regression model is provided to estimate the per-layer execution latency of a DNN. Next, we propose our horizontal partition algorithm that splits a DNN into three parts. Finally, we describe the vertical separation module that enables parallel convolutional layer inference.

A. Edge Computing Paradigm

The edge computing architecture comprises three tiers that are device, edge, and cloud. Each layer consists of multiple computation nodes. From a computation perspective, edge nodes provide high computation capabilities compared with device nodes. Even so, as the edge nodes are often heterogeneous, the computation power of the edge nodes is still bounded compared with the cloud servers. Consequently, we

say that the computation resources are gradually increasing over device, edge, and cloud [26]. From a communication perspective, since we deploy edge nodes close to the data source, network bandwidth maintains high between device and edge. Nonetheless, as device nodes and edge nodes connect to cloud servers through a backbone network (e.g., the Internet backbone), the bandwidth to the cloud nodes remains limited [27]. Compared with the transmission delay between the computing tiers, the in-memory transmission delay of a node or the transmission delay between two computation nodes within the same computing tier is negligible. To simplify the problem, without loss of generality, we assume that the transmission delay within each computing tier is infinitesimal.

B. System Overview

We show the D^3 system architecture in Fig. 2. D^3 comprises a profiler, an offline partition framework, and an online execution engine. The profiler collects the operating conditions of computation nodes at device, edge, and cloud as well as the network status between tiers. A regression model takes the computation statistics as input and estimates the inference time of each DNN layer processed at different computing tiers. The offline partition framework comprises two components. A horizontal partition algorithm splits a DNN model into three parts according to the per-layer inference time and the transmission delay between layers. D^3 distributes the three DNN segments to three computation nodes located at device, edge, and cloud respectively. Considered that the computation resource of a single edge node is limited, if convolutional layers locate at an edge node, to further accelerate the inference, a vertical separation module divides a sequence of correlated feature maps to multiple feature map stacks, each of which is processed on an edge node independently. Fig. 2 illustrates the condition that D^3 breaks a sequence of feature maps into four stacks. The online execution engine orchestrates the distributed and parallelism processing, handling communication among

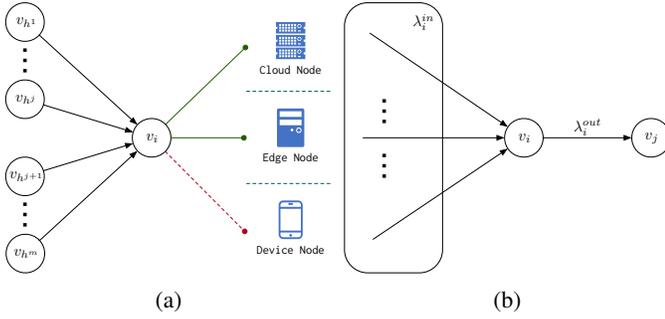


Fig. 5: The potential tier of a vertex depends on its direct predecessors. The optimal tier of a vertex depends on its input and output size.

E. Horizontal Partition Algorithm

Our main objective is to minimize the DNN inference time by leveraging the collaborative computation provided by three computing tiers. To achieve this goal, we propose to split the DNN model into three parts that are executed over device, edge, and cloud. Given a DNN layer, the total latency includes the DNN layer processing time and its input data transmission delay. Theoretically, for $v_i \in \mathcal{V}$ and a set of its direct predecessors \mathcal{V}_i^p , the total latency is $t_i^{l_i} + \sum_{v_h \in \mathcal{V}_i^p} t_{hi}^{[l_h, l_i]}$. The optimal tier of a vertex v_i is decided by comparing the latencies of attaching v_i to different computing tiers. We regard the tier that is possible to yield the smallest latency as a *potential tier*. A set of potential tiers of v_i is denoted by Γ_i , where $\Gamma_i \subset \{d, e, c\}$.

To this end, mathematically, our major goal is to split the DAG in (1) into three sub-graphs by assigning each vertex to one of the three computing tiers and minimizing the total latency

$$\Theta(v_0, v_1, \dots, v_n) = \sum_{v_i \in \mathcal{V}} t_i^{l_i} + \sum_{(v_i, v_j) \in \mathcal{L}} t_{ij}^{[l_i, l_j]}.$$

Each sub-graph contains a subset of \mathcal{V} that have the same optimal tier. However, partitioning a DAG according to multiple vertex weights and link weights falls into an NP-hard problem [30], [31]. Thus, heuristics are essential to partition the DAG.

With the DAG in (1) representing a DNN, the **horizontal partition algorithm** (HPA) first computes the longest distance from v_0 to v_i , denoted by $\delta(v_i)$, $\forall v_i \in \mathcal{V}$. We get the longest distance with the dynamic programming method mentioned in [32]. The time complexity of the method is $\mathcal{O}(|\mathcal{V}| + |\mathcal{L}|)$. Subsequently, we define a partition of \mathcal{V} by

$$\mathcal{Z}_q := \{v_i : \delta(v_i) = q, v_i \in \mathcal{V}\}, \quad q = 0, 1, \dots, n-1,$$

and HPA arranges v_i to the graph layer \mathcal{Z}_q . To illustrate, in Fig. 3b, HPA assigns the vertices to 7 graph layers that are $\mathcal{Z}_0 = \{v_0\}$, $\mathcal{Z}_1 = \{v_1\}$, $\mathcal{Z}_2 = \{v_2, v_3, v_4, v_5\}$, $\mathcal{Z}_3 = \{v_6, v_7, v_8, v_9\}$, $\mathcal{Z}_4 = \{v_{10}\}$, $\mathcal{Z}_5 = \{v_{11}, v_{12}\}$, $\mathcal{Z}_6 = \{v_{13}\}$. In graph layer \mathcal{Z}_q where $q \in \mathbb{Z}^+$, to decide the optimal tier of each vertex in \mathcal{Z}_q , HPA must assign all vertices in \mathcal{Z}_{q-1} to

TABLE I: The total latencies of processing v_i and v_j .

Location of v_i	Location of v_j	Total Latency
device	device	$t_i^d + t_j^d$
device	edge	$t_i^d + t_j^e + \lambda_i^{out}/\sigma_{de}$
edge	edge	$t_i^e + t_j^e + \lambda_i^{in}/\sigma_{de}$
edge	cloud	$t_i^e + t_j^c + \lambda_i^{in}/\sigma_{de} + \lambda_i^{out}/\sigma_{ec}$
cloud	cloud	$t_i^c + t_j^c + \lambda_i^{in}/\sigma_{dc}$
device	cloud	$t_i^d + t_j^c + \lambda_i^{out}/\sigma_{dc}$

their optimal tiers. By mathematical induction, we know that HPA starts from graph layer \mathcal{Z}_0 and calculates the optimal tiers layer by layer.

For a vertex v_i , its potential tiers and the tiers of its direct predecessors have the relation described in the Proposition 1. As an example, if the direct predecessors of vertex v_i are assigned to an edge node, then the potential tiers for v_i are edge and cloud.

Proposition 1. For a vertex $v_i \in \mathcal{V}$ and a set of its direct predecessors $\mathcal{V}_i^p = \{v_{h^1}, v_{h^2}, \dots, v_{h^m}\}$, the potential tier l_i of v_i has the subsequent relation with the tiers of its direct predecessors, i.e., $\max\{l_{h^1}, l_{h^2}, \dots, l_{h^m}\} \succeq l_i$.

Proof. We prove the proposition 1 by contradiction. Consider a vertex v_i with multiple direct predecessors $\{v_{h^1}, \dots, v_{h^j}, v_{h^{j+1}}, \dots, v_{h^m}\}$ shown in Fig. 5a. We assume that the DNN layers represented by $\{v_{h^1}, \dots, v_{h^j}\}$ are assigned to a cloud node and the DNN layers denoted by $\{v_{h^{j+1}}, \dots, v_{h^m}\}$ are assigned to an edge node. Suppose that we assign the DNN layer described by v_i to a device node, which means $l_i \succ \max\{l_{h^1}, l_{h^2}, \dots, l_{h^m}\}$. The total latency of v_i is $t_i^d + \sum_{x=1}^j t_{hx}^{[c, d]} + \sum_{y=j+1}^m t_{hy}^{[e, d]}$, which is larger than the delay $t_i^e + \sum_{x=1}^j t_{hx}^{[c, e]}$ when the DNN layer signified by v_i is assigned to an edge node. Assigning v_i to the device tier cannot yield the smallest latency. By contradiction, we prove $\max\{l_{h^1}, l_{h^2}, \dots, l_{h^m}\} \succeq l_i$. \square

We now derive the optimal tier selection strategy for a vertex $v_i \in \mathcal{V}$. Intuitively, for a single vertex $v_i \in \mathcal{V}$, the optimal tier l_i^{opt} is calculated as:

$$l_i^{opt} = \arg \min_{l_i \in \Gamma_i} (t_i^{l_i} + \sum_{v_h \in \mathcal{V}_i^p} t_{hi}^{[l_h, l_i]}). \quad (2)$$

Specially, $l_0^{opt} = d$ for the virtual input vertex v_0 . However, Equation (2) only selects the local optimal location for v_i . To further optimize the selection, we take the input and output size of a DNN layer into consideration. Assuming that a vertex v_i has multiple direct successors, if we place the DNN layers represented by these successors on the same computation node and give them the same input, we call the successor representing the DNN layer with the longest processing time as the *largest direct successor* of v_i . Fig. 5b depicts two vertices $v_i \in \mathcal{V}$ and $v_j \in \mathcal{V}$, where v_j is the largest direct successor of v_i and $\Gamma_i = \{d, e, c\}$. The total input size of v_i is λ_i^{in} , and its output size is λ_i^{out} . The network bandwidth

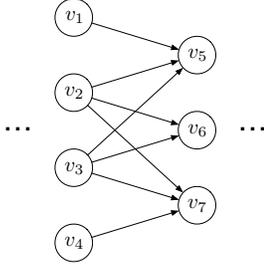


Fig. 6: Illustration of SIS relation between vertices.

between a device node and an edge node, an edge node and a cloud node, a device node and a cloud node are $\sigma_{de}, \sigma_{ec}, \sigma_{dc}$ respectively. We list the total latencies of allocating the DNN layers indicated by v_i and v_j to computation nodes at different tiers in TABLE I. The inputs of v_i are from the device tier. Heuristically, HPA selects the optimal tier of v_i via the following mechanism. On one hand, if $\lambda_i^{in} > \lambda_i^{out}$, v_i 's optimal tier l_i^{opt} is computed via Equation (2). On the other hand, if $\lambda_i^{in} \leq \lambda_i^{out}$, HPA chooses the largest direct successor of v_i and computes the total latencies of processing v_i and its largest direct successor as in TABLE I. According to the smallest value among these total latencies, HPA selects the optimal tier of v_i .

Given a vertex $v_i \in \mathcal{Z}_q$, we refer to a vertex v_j as the *subset input sibling* (SIS) vertex of v_i if $\mathcal{V}_j^p \subset \mathcal{V}_i^p$. We clarify the SIS relation with an example in Fig. 6 where 7 vertices are connected via directed links. In the example, v_6 is the SIS vertex of v_5 since $\mathcal{V}_6^p \subset \mathcal{V}_5^p$, whereas v_7 is not the SIS vertex of v_5 in that $\mathcal{V}_7^p \not\subset \mathcal{V}_5^p$. After deciding the optimal tier for all vertices in layer \mathcal{Z}_q , for each $v_i \in \mathcal{Z}_q$, if the optimal tier of v_i 's SIS vertex v_j is before l_i^{opt} (i.e., $l_j^{opt} \succ l_i^{opt}$), HPA updates the optimal tier of the SIS vertex to l_i^{opt} . We regard this approach as *SIS update* that follows the Proposition 2.

Proposition 2. *SIS update optimizes the processing delay and the transmission latency of a SIS vertex.*

Proof. Given a vertex v_i and one of its SIS vertices v_j where $l_j^{opt} \succ l_i^{opt}$. Since the inputs of v_i are already transmitted to the tier l_i^{opt} , therefore relocating the SIS vertex that is at the previous tier l_j^{opt} to l_i^{opt} reduces the processing time and brings no transmission delay overhead. \square

We now show the HPA in Algorithm 1. The algorithm calls `get_longest_path()` to calculate the longest path from v_0 to all vertices in \mathcal{G} . With the result of the longest path, HPA constructs the graph layers by invoking the function `get_graph_layer()`. In each graph layer \mathcal{Z}_q , HPA computes the optimal tiers for all vertices. For every vertex v_i in the layer, the algorithm uses the function `get_pred_loc()` to get the optimal tiers of v_i 's direct predecessors. Next, it leverages the idea presented in Proposition 1 and employs `get_loc_choice()` to get all potential tiers of v_i . In particular, if the potential tier of v_i is c , the optimal tier l_i^{opt} is c . Otherwise, HPA computes the optimal tier for v_i

Algorithm 1: Horizontal Partition Algorithm: HPA()

Data: DAG: $\mathcal{G} = (\mathcal{V}, \mathcal{L})$;
Vertex weights: \mathcal{T}_* ;
Link weights: \mathcal{T}_\dagger .
Result: Optimal tiers l_i^{opt} where $i = 1, 2, \dots, |\mathcal{V}|$.

- 1 $\mathcal{Q} \leftarrow \text{get_longest_path}(\mathcal{G})$;
- 2 $\Delta \leftarrow \text{get_graph_layer}(\mathcal{Q}, \mathcal{V})$;
- 3 **foreach** \mathcal{Z}_q **in** Δ **do**
- 4 **foreach** v_i **in** \mathcal{Z}_q **do**
- 5 $\Phi_i \leftarrow \text{get_pred_loc}(v_i)$;
- 6 $\Gamma_i \leftarrow \text{get_loc_choice}(\Phi_i)$;
- 7 **if** $\Gamma_i = \{c\}$ **then**
- 8 $l_i^{opt} \leftarrow c$;
- 9 **else**
- 10 $l_i^{opt} \leftarrow \text{get_opt_loc}(v_i, \Gamma_i, \mathcal{T}_*, \mathcal{T}_\dagger)$;
- 11 **end**
- 12 **end**
- 13 $\text{sis_update}(\mathcal{Z}_q)$;
- 14 **end**

via function `get_opt_loc()` which leverages the heuristics in the optimal tier selection strategy. After the calculation of all vertices in the graph layer \mathcal{Z}_q , `sis_update()` performs SIS update for all vertices in \mathcal{Z}_q . HPA stops when it finishes processing all the graph layers.

Resource changes and network dynamics lead to variations of DNN layer processing time and input data transmission latencies, which further affect the optimal locations to process DNN layers. Assuming that the optimal tier of a vertex changes, HPA can accommodate the modification by locally adjusting the optimal tiers of its SIS vertices, its direct successors, and the SIS vertices of its direct successors. For instance, in Fig. 3b, supposing that l_6^{opt} of v_6 changes to a different value, HPA recalculates l_{10}^{opt} of v_{10} since v_{10} is the direct successor of v_6 . To avoid constantly calculating the optimal tier to respond to the resource and network fluctuation, we can set upper and lower thresholds to limit the scope of the optimal tier alteration. HPA only recalculates the optimal tiers when DNN layer processing time or the network bandwidth is outside of the threshold range.

F. Vertical Separation Module

The preliminary experimental results in TABLE II show the inference latencies of DNNs after HPA given an input image of $3 \times 224 \times 224$. The device node is an NVIDIA Jetson Nano 2GB Developer Kit [33], the edge node is a Linux machine with Intel Core i7-8700 CPU and 8 GB system memory, and the cloud node is a Linux server with NVIDIA GeForce RTX 2080 Ti GPU and 256 GB system memory. We observe that the processing time of the edge node is longer than that of the cloud node, causing the cloud node to be idle and waiting for the results from the edge node in the inference pipeline. The edge node becomes the bottleneck of the synergistic inference. To accelerate the inference at the edge tier, we

TABLE II: The synergistic inference time at three nodes.

DNNs	Device Node (millisecond)	Edge Node (millisecond)	Cloud Node (millisecond)
AlexNet	2.2	3.6	1.4
VGG-16	5.7	46.7	0.5
ResNet-18	6.1	7.5	0.5
Darknet-53	27.9	48.1	0.1
Inception-v4	21.4	46.4	16.7

design a parallel processing strategy to avoid the bottleneck condition. The parallelization option is not suitable for the resource stringent device node since processing raw input in parallel incurs privacy concerns and amount of communication overhead. Next, we introduce our parallel processing method.

The parameters of a convolutional layer are a set of learnable filters [34]. Each filter is a weighted tensor spatially defined by its hyper-parameters that are width, height, and depth. An input feature map is the input activation for a given filter. The number of input feature maps is the same as the filter depth in a convolutional layer. Besides, a convolutional layer contains two hyper-parameters that are filter stride and padding for input feature maps. If padding exists, the convolutional layer adds entries to the borders of an input feature map, resulting in an input feature map with paddings. The convolutional layer systematically performs a dot product between the entries of the padded input feature maps and the filter. The results are assembled to an output feature map. We call this process a convolution operation. In order to optimize the inference latency at the resource constrained edge node, we leverage the idea of separating a sequence of correlated input feature maps to multiple feature map stacks. This idea is firstly proposed in DeepThings [9]. However, DeepThings does not consider input feature maps with paddings, leading to the precision loss that affects the inference accuracy. To settle this issue, we derive a parallel convolutional layer inference module without loss of accuracy referred to as vertical separation module (VSM).

Given a sequence of k convolutional layers, we describe each convolutional layer as c_i where $i = 1, 2, \dots, k$. The input feature maps of layer c_i have the dimension of $\mathcal{W}_i \times \mathcal{H}_i \times \mathcal{D}_i$ (width \times height \times depth). For the filter of c_i , we designate the filter size as $\mathcal{F}_i^w \times \mathcal{F}_i^h \times \mathcal{D}_i$ (width \times height \times depth) with a horizontal stride of \mathcal{S}_i^w and a vertical stride of \mathcal{S}_i^h . The padding is \mathcal{P}_i^w horizontally and \mathcal{P}_i^h vertically. Consequently, the input feature map size of layer c_i has the coming relation with the input feature map size of layer c_{i-1} :

$$\begin{aligned} \mathcal{W}_i &= \frac{\mathcal{W}_{i-1} - \mathcal{F}_{i-1}^w + 2 \times \mathcal{P}_{i-1}^w}{\mathcal{S}_{i-1}^w} + 1, \\ \mathcal{H}_i &= \frac{\mathcal{H}_{i-1} - \mathcal{F}_{i-1}^h + 2 \times \mathcal{P}_{i-1}^h}{\mathcal{S}_{i-1}^h} + 1. \end{aligned} \quad (3)$$

We divide the input feature maps of layer c_i into $\mathcal{A} \times \mathcal{B}$ non-overlapping continuous *tiles* whose depth is \mathcal{D}_i . Generally, we index an entry of the input feature map with two-dimensional coordinates. Therefore, we use the coordinates of the top left corner and the bottom right cor-

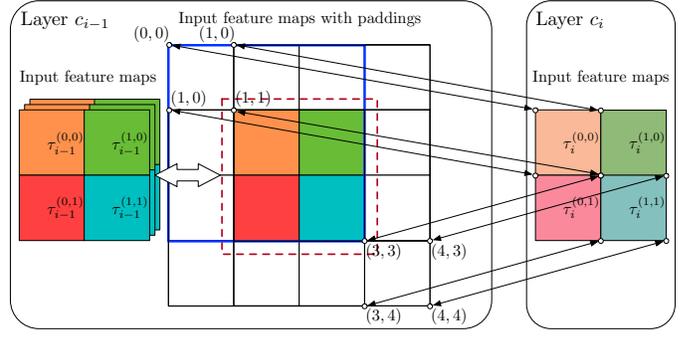


Fig. 7: Layer c_{i-1} has an input feature map of $2 \times 2 \times 3$ with a padding of $\mathcal{P}_{i-1}^w = \mathcal{P}_{i-1}^h = 1$ and a stride of $\mathcal{S}_{i-1}^w = \mathcal{S}_{i-1}^h = 1$.

ner of a tile to locate the tile in the input feature maps. Mathematically, the tile $\tau_i^{(a,b)} = (\alpha_i^{(a,b)}, \beta_i^{(a,b)})$ at layer c_i where $a = 0, 1, \dots, \mathcal{A} - 1$; $b = 0, 1, \dots, \mathcal{B} - 1$. Specifically, the tile at the top left corner of the layer c_i 's input feature map is $\tau_i^{(0,0)}$. We represent the top left coordinate of $\tau_i^{(a,b)}$ as $\alpha_i^{(a,b)} = (x_i^{[\alpha,(a,b)]}, y_i^{[\alpha,(a,b)]})$ and the bottom right coordinate as $\beta_i^{(a,b)} = (x_i^{[\beta,(a,b)]}, y_i^{[\beta,(a,b)]})$. Particularly, $\alpha_i^{(0,0)} = (0,0)$. In the light of Equation (3), given the coordinates of $\tau_i^{(a,b)}$ ($i > 1$), we compute the coordinates of the correlated tile with paddings at layer c_{i-1} denoted by $\hat{\tau}_{i-1}^{(a,b)} = (\hat{\alpha}_{i-1}^{(a,b)}, \hat{\beta}_{i-1}^{(a,b)})$ where $\hat{\alpha}_{i-1}^{(a,b)} = (\hat{x}_{i-1}^{[\alpha,(a,b)]}, \hat{y}_{i-1}^{[\alpha,(a,b)]})$ and $\hat{\beta}_{i-1}^{(a,b)} = (\hat{x}_{i-1}^{[\beta,(a,b)]}, \hat{y}_{i-1}^{[\beta,(a,b)]})$ in following manner:

$$\begin{aligned} \hat{x}_{i-1}^{[\alpha,(a,b)]} &= \mathcal{S}_{i-1}^w \times x_i^{[\alpha,(a,b)]}, \\ \hat{y}_{i-1}^{[\alpha,(a,b)]} &= \mathcal{S}_{i-1}^h \times y_i^{[\alpha,(a,b)]}, \\ \hat{x}_{i-1}^{[\beta,(a,b)]} &= \mathcal{S}_{i-1}^w \times (x_i^{[\beta,(a,b)]} - 1) + \mathcal{F}_{i-1}^w, \\ \hat{y}_{i-1}^{[\beta,(a,b)]} &= \mathcal{S}_{i-1}^h \times (y_i^{[\beta,(a,b)]} - 1) + \mathcal{F}_{i-1}^h. \end{aligned} \quad (4)$$

Considering the paddings added to the input feature maps of layer c_{i-1} , given the coordinates of a padded tile $\hat{\tau}_{i-1}^{(a,b)}$ ($i > 1$), to compute the coordinates of tile $\tau_{i-1}^{(a,b)}$, we need to remove the paddings from $\hat{\tau}_{i-1}^{(a,b)}$. This alters the coordinates of the padded tile, transforming $\hat{\tau}_{i-1}^{(a,b)}$ to $\tau_{i-1}^{(a,b)}$. The coordination of $\tau_{i-1}^{(a,b)}$ is described as follows:

$$\begin{aligned} x_{i-1}^{[\alpha,(a,b)]} &= \max(0, \hat{x}_{i-1}^{[\alpha,(a,b)]} - \mathcal{P}_{i-1}^w), \\ y_{i-1}^{[\alpha,(a,b)]} &= \max(0, \hat{y}_{i-1}^{[\alpha,(a,b)]} - \mathcal{P}_{i-1}^h), \\ x_{i-1}^{[\beta,(a,b)]} &= \begin{cases} \mathcal{W}_{i-1}, & \text{if } \hat{x}_{i-1}^{[\beta,(a,b)]} = \mathcal{W}_{i-1} + 2 \times \mathcal{P}_{i-1}^w, \\ \max(0, \hat{x}_{i-1}^{[\beta,(a,b)]} - \mathcal{P}_{i-1}^w), & \text{otherwise,} \end{cases} \\ y_{i-1}^{[\beta,(a,b)]} &= \begin{cases} \mathcal{H}_{i-1}, & \text{if } \hat{y}_{i-1}^{[\beta,(a,b)]} = \mathcal{H}_{i-1} + 2 \times \mathcal{P}_{i-1}^h, \\ \max(0, \hat{y}_{i-1}^{[\beta,(a,b)]} - \mathcal{P}_{i-1}^h), & \text{otherwise.} \end{cases} \end{aligned} \quad (5)$$

Fig. 7 depicts the process of generating layer c_{i-1} 's input feature maps from the input feature maps of layer c_i . When the stride is 1×1 at layer c_{i-1} , the $3 \times 3 \times 3$ filter moves one entry at a time, producing the input feature map of layer c_{i-1} .

Algorithm 2: Vertical Separation Module: VSM()

Data: k correlated convolutional layers: c_i where $i = 1, 2, \dots, k$;

Decision of separation: $\mathcal{A} \times \mathcal{B}$ tiles;

The tiles at layer c_{k+1} : $T_{k+1} = \{\tau_{k+1}^{(a,b)}\}$
($a = 0, 1, \dots, \mathcal{A} - 1$; $b = 0, 1, \dots, \mathcal{B} - 1$).

Result: Coordinates of the tiles at c_1 .

```
1 foreach  $\tau_{k+1}^{(a,b)}$  in  $T_{k+1}$  do
2   foreach  $i \leftarrow k$  to 1 do
3      $\tau_i^{(a,b)} \leftarrow \text{RTC}(c_i, \tau_{i+1}^{(a,b)})$ ;
4   end
5 end
```

Reversely, given the input feature maps of layer c_i , we separate the feature maps into 2×2 tiles. With the coordinates of the four tiles at layer c_i , we find the corresponding padded tiles at layer c_{i-1} with Equation (4). Then, we obtain the tiles at layer c_{i-1} by offsetting the paddings from the padded tiles via Equation (5). We regard the procedure of computing tile $\tau_{i-1}^{(a,b)}$ from $\tau_i^{(a,b)}$ as **reverse tile calculation (RTC)**.

To accelerate the inference at the edge tier without loss of accuracy, we develop VSM whose overall procedure is shown in Algorithm 2. Assuming that a sequence of k convolutional layers are assigned to an edge node, to help to describe our algorithm, we add a virtual convolutional layer c_{k+1} , whose input feature maps are the output feature maps of layer c_k . VSM separates the input feature maps of layer c_{k+1} into $\mathcal{A} \times \mathcal{B}$ non-overlapping continuous tiles. For each tile at layer c_{k+1} , VSM finds the coordinates of corresponding tiles from layer c_k to layer c_1 via a series of RTC operations. We refer to a stack of correlated tiles from layer c_1 to c_k as *fused tiles*. Next, VSM locates the tiles at layer c_1 and splits the input feature maps of c_1 according to the coordinates of the tiles. VSM then feeds the tiles at layer c_1 to $\mathcal{A} \times \mathcal{B}$ edge nodes, each of which holds the inference parameters and hyper-parameters of the k convolutional layers. Each edge node possesses one tile and processes a fused tile stack independently via convolution operations, generating the tile of the output feature maps at layer c_k . Finally, an edge node gathers and combines the $\mathcal{A} \times \mathcal{B}$ tiles of the output feature maps at layer c_k . The results are the output feature maps of layer c_k , which are transferred to a cloud node. We neglect batch normalization layers and activation layers in the middle of two convolutional layers, since their operations do not change the volume of input feature maps. Moreover, pooling layers that are used to sub-sample feature maps between two convolutional layers are separated and fused by VSM in the same way as the convolutional layers. Fig. 8 shows four fused tile stacks of three consecutive convolutional layers. We assign the four fused tile stacks to four edge nodes.

IV. IMPLEMENTATION

We introduce our implementation details in this section. We consider a scenario where the device node is mobile. The

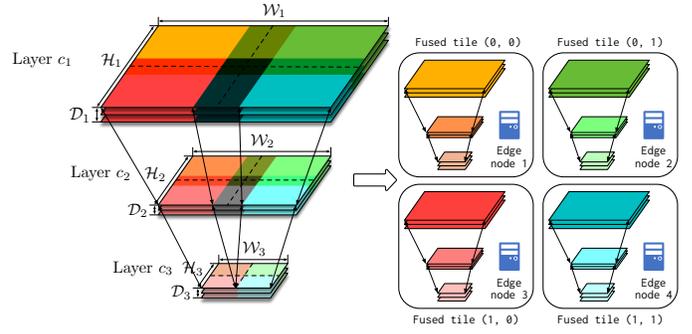


Fig. 8: A sequence of input feature maps of three consecutive convolutional layers are divided into 2×2 fused tile stacks. Each fused tile stack is assigned to an edge node.

edge node and the device node are in the same local area network (LAN), both of which connect to the cloud node via the Internet [26]. Hence, we use a Raspberry Pi 4 model B with 4 GB system memory as the device node. The edge node is a Linux machine with Intel Core i7-8700 CPU and 8 GB system memory. We employ a remote server with NVIDIA GeForce RTX 2080 Ti GPU and 256 GB system memory.

Responsibility. We implement the profiler and the offline partition framework on a dedicated computation node, which monitors running conditions, performs dynamic partitioning according to the conditions, and distributes partitions to the on-line execution nodes. For the online execution procedure, the device node is responsible for collecting the input, processing through the DNN layers allocated to it, and transferring the output to an edge node or a cloud node. The edge node handles the output from the device node and passes the intermediate inference results to the cloud node for further computation. Depending on the partition decision, the inference pipeline may involve only part of the computing tiers.

Software stack. We implement a client-server interface using gRPC [35] in each node to accommodate inter-process communication. The dedicated computation node loads a trained DNN in the ONNX [36] format. We process the DNN models with PyTorch [37] and build the DAG representation of the DNN with NetworkX [38]. After applying HPA, we rebuild the computation graph from the partitioned DAG, transform the computation graph into a partial DNN by reconstructing the `forward()` function in PyTorch, and store the partial DNN in the ONNX format.

Datasets and models. We evaluate D³ on the ImageNet (ILSVRC2012) [14] dataset that contains 150000 validation and test images of 1000 categories. We compress the size of the input images to $3 \times 224 \times 224$. Then, we feed the images to the device node at 30 FPS for 100 seconds and test the per-image average end-to-end latency. We evaluate our methods across 5 widely used DNNs: AlexNet, VGG-16, ResNet-18, Darknet-53, and Inception-v4, all of which are trained before deployment.

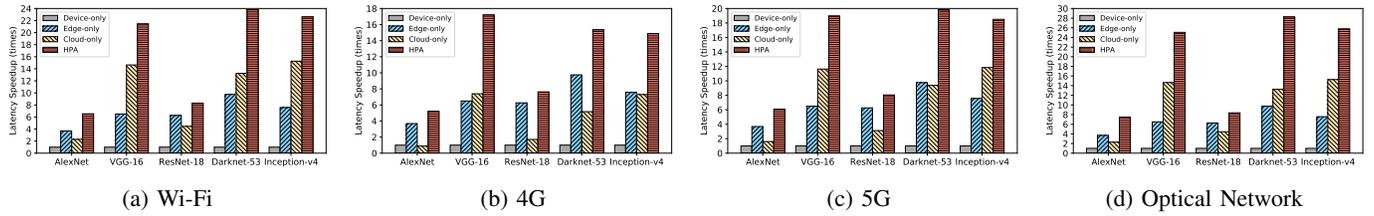


Fig. 9: End-to-end latency speedup comparison among HPA, device-only, edge-only, and cloud-only under different network conditions.

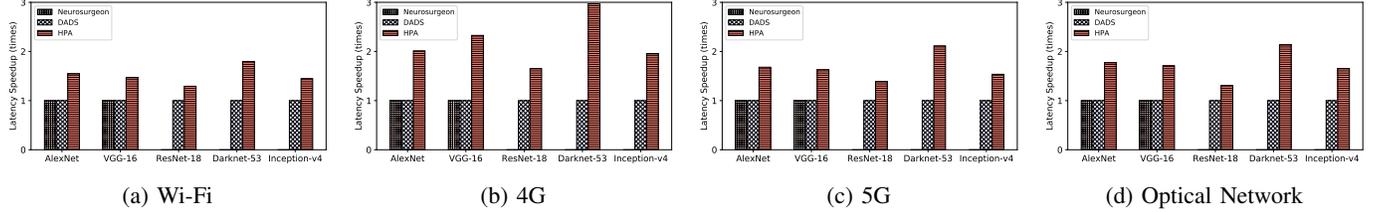


Fig. 10: End-to-end latency speedup comparison among HPA, Neurosurgeon [7], and DADS [8] under different network conditions.

V. EVALUATION

This section presents the performance comparison of D^3 with its counterparts. In addition to executing a DNN model solely on a device node, an edge node, or a cloud node, we choose two state-of-the-art precision lossless DNN offloading systems Neurosurgeon [7] and DADS [8] for comparison. We mainly examine the effectiveness of HPA and VSM from two performance metrics: end-to-end inference speedup and per-image communication overhead. To investigate the performance of D^3 in various network settings, we test the performance metrics under different network conditions shown in TABLE III, which shows the average uplink rate. Particularly, we consider the following scenarios. The communication link between the device node and the edge node is a Wi-Fi running at 5 GHz (Gigabit Ethernet 802.3, Wi-Fi 5 IEEE 802.11ac). Both the device node and the edge node connect to the cloud node with the same type of communication link, which is Wi-Fi, 4G, or 5G. Besides, when the edge node employs optical network to bridge to the cloud node, the device node connects to the cloud node via the 5 GHz Wi-Fi. In our experiment, we alter the type of communication link between the LAN and the cloud node (i.e., Wi-Fi, 4G, 5G, Optical Network), and examine the performance of our algorithms.

A. End-to-end Inference Speedup

This subsection compares the end-to-end inference speedup of D^3 with device-only, edge-only, cloud-only, Neurosurgeon, and DADS under different network conditions. In the edge-only and the cloud-only method, input data is collected by the device node and transmitted to the edge node and cloud node respectively for processing. The experimental results validate the effectiveness of HPA and VSM.

Fig. 9 demonstrates the end-to-end latency speedup of HPA over device-only approach, edge-only approach, and cloud

TABLE III: The average uplink rate (Mbps) between two nodes.

DNNs	Wi-Fi	4G	5G	Optical Network
device to edge	84.95	N.A.	N.A.	N.A.
edge to cloud	31.53	13.79	22.75	50.23
device to cloud	18.75	6.12	11.64	N.A.

only approach. We set the device-only as the baseline for our comparison. The result shows that HPA accelerates the end-to-end latency up to $28.2\times$, $3.85\times$, and $5.90\times$ compared with device-only, edge-only, and cloud-only approaches respectively. We discover that the device-only is constrained by the computation resources. When the model demands more computation power, offloading part of the model to the edge node obtains a higher inference speedup. Meanwhile, the cloud-only method is limited by the low bandwidth between the device node and the cloud node. By increasing the bandwidth between the device node and the cloud node, the cloud-only method achieves a lower inference latency in that the input data transmission time between the device node and the cloud node is smaller. Fig. 10 illustrates the end-to-end latency speedup among HPA, Neurosurgeon, and DADS. Since Neurosurgeon can only partition the DNN of chain topology, it is not applicable for ResNet-18, Darknet-53, and Inception-v4, which are of DAG topology. The experimental results show that HPA outperforms Neurosurgeon up to $2.33\times$ in the chain topology DNNs. Compared with DADS, HPA accelerates the inference latency up to $2.97\times$ in DNNs of DAG topology. Notably, we observe that HPA attains more inference speedup when the model gets larger, which requires more computation resources. To show the impact of network bandwidth variations, we apply HPA to Inception-v4 and measure its

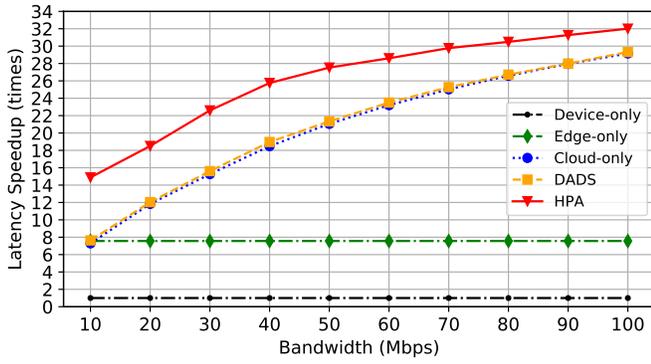


Fig. 11: We measure the latency speedup of Inception-v4 under various bandwidth between the LAN and the cloud node.

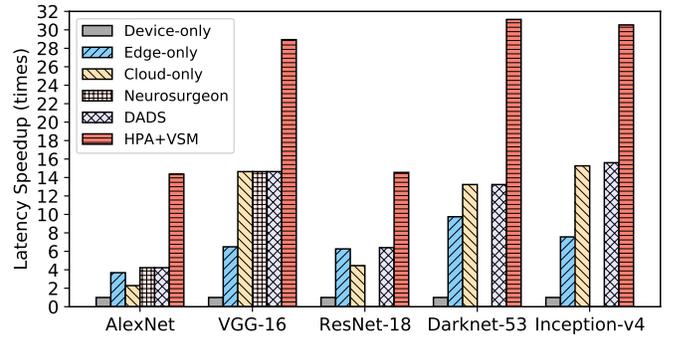


Fig. 12: The latency speedup when applying both HPA and VSM. The device node and the edge nodes connect to the cloud node via Wi-Fi.

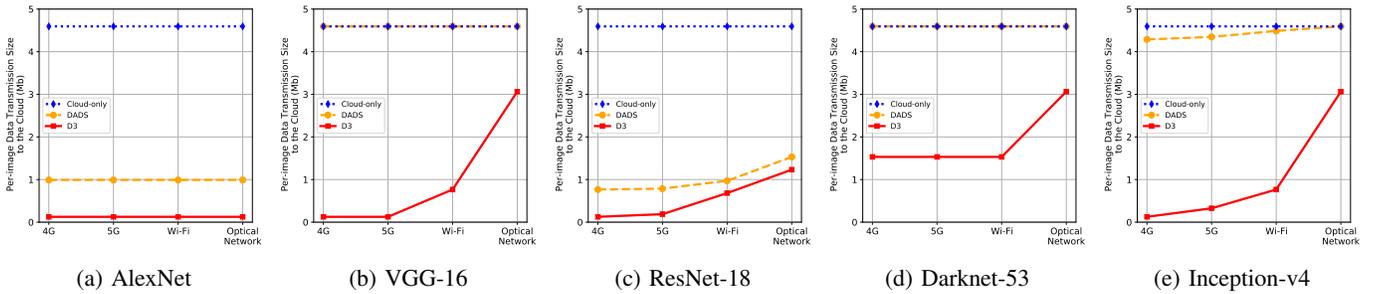


Fig. 13: Per-image communication overhead comparison among D^3 , Cloud-only, and DADS [8] for different models under different network conditions.

end-to-end latency speedup under different network bandwidth between the LAN and the cloud node. From Fig. 11, we monitor that when the network bandwidth between the LAN and the cloud node increases, HPA tends to offload more DNN layers to the cloud to optimize the end-to-end inference latency.

We apply VSM to the convolutional layers that are assigned to the edge node under HPA. We employ four Linux machines with Intel Core i7-8700 CPU and 8 GB system memory as the edge nodes. Both the device node and the edge nodes connect to the cloud node via Wi-Fi. From Fig. 12, we can see the latency speedup when both HPA and VSM are applied. The D^3 system surpasses the device-only, edge-only, cloud-only, Neurosurgeon, and DADS up to $31.13\times$, $4.46\times$, $6.28\times$, $3.4\times$, and $3.4\times$ respectively. When VSM is deployed, the processing time of convolutional layers at the edge tier does not shrink to $1/4$ of the original processing time compared with the HPA-only approach, since there are spatial overlaps among the fused tile stacks, which in turn leads to computational redundancy.

Under the condition that both the device node and the edge nodes connect to the cloud node via Wi-Fi, HPA improves the end-to-end latency to $1.29\times - 1.8\times$, and HPA+VSM accelerates the latency to $1.96\times - 3.4\times$, where the state-of-the-art is the baseline. Overall, the end-to-end latency measurements verify the effectiveness of D^3 .

B. Per-image Communication Overhead

Typically, a cloud node locates at a remote place, which is accessible from a LAN through the Internet. A lower data transmission between the LAN and the cloud server reduces the data transmission over the network core hence mitigating the Internet congestion [39]. Transferring the intermediate results of a DNN to the cloud curtails the data transmission between the LAN and the cloud. We test the per-image communication overhead to the cloud node of D^3 and its counterparts and show the outcomes in Fig. 13. D^3 shrinks the per-image data transmission size on the Internet backbone to 27.21% - 66.67% of the cloud-only approach. D^3 reduces the per-image data transmission size on the Internet backbone to 27.21% - 80.42% when DADS is the baseline. When the network bandwidth between the LAN and the cloud server becomes larger, D^3 tends to offload more DNN layers and transmit more intermediate data to the cloud server.

VI. CONCLUSION

With the increasing computation capability of mobile devices and the growing need of accelerating DNN inference, there is an expansion demand of performing synergistic DNN inference across device, edge, and cloud without precision loss. In this work, we introduce D^3 , a system that consists of HPA and VSM. HPA partitions a DNN model into three parts according to the per-layer processing time and the inter-layer

transmission delay of the DNN. Besides, VSM further divides the feature maps of DNN layers assigned to the edge node into multiple fused tile stacks for parallel processing. Under our experiment settings, D³ system provides up to 3.4× end-to-end inference speedup on chain topology DNNs and accelerates the end-to-end inference up to 2.35× on DAG topology DNNs compared to the state-of-the-art. In addition, the per-image data transmission size reduces to 27.21% of the state-of-the-art, which relieves network congestion and communication cost.

REFERENCES

- [1] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [2] Z. He, J. Li, D. Liu, H. He, and D. Barber, "Tracking by animation: Unsupervised learning of multi-object attentive trackers," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1318–1327.
- [3] D. He, Y. Xia, T. Qin, L. Wang, N. Yu, T.-Y. Liu, and W.-Y. Ma, "Dual learning for machine translation," in *Advances in neural information processing systems*, 2016, pp. 820–828.
- [4] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, "Efficient processing of deep neural networks," *Synthesis Lectures on Computer Architecture*, vol. 15, no. 2, pp. 1–341, 2020.
- [5] X. Wang, Y. Han, V. C. Leung, D. Niyato, X. Yan, and X. Chen, "Convergence of edge computing and deep learning: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, pp. 869–904, 2020.
- [6] S. Teerapittayanon, B. McDanel, and H.-T. Kung, "Distributed deep neural networks over the cloud, the edge and end devices," in *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2017, pp. 328–339.
- [7] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang, "Neurosurgeon: Collaborative intelligence between the cloud and mobile edge," in *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '17, 2017, pp. 615–629.
- [8] C. Hu, W. Bao, D. Wang, and F. Liu, "Dynamic adaptive dnn surgery for inference acceleration on the edge," in *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*. IEEE, 2019, pp. 1423–1431.
- [9] L. Zhou, M. H. Samavatian, A. Bacha, S. Majumdar, and R. Teodorescu, "Adaptive parallel execution of deep neural networks on heterogeneous edge devices," in *Proceedings of the 4th ACM/IEEE Symposium on Edge Computing*, 2019, pp. 195–208.
- [10] S. Wang, A. Pathania, and T. Mitra, "Neural network inference on mobile socs," *IEEE Design & Test*, 2020.
- [11] *Qualcomm Snapdragon 865+ 5G mobile platform*, 2020 (accessed September 30, 2020). [Online]. Available: <https://www.qualcomm.com/media/documents/files/qualcomm-snapdragon-865-5g-mobile-platfor-m-product-brief.pdf>
- [12] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge intelligence: Paving the last mile of artificial intelligence with edge computing," *Proceedings of the IEEE*, vol. 107, no. 8, pp. 1738–1762, 2019.
- [13] Z. Zhao, K. M. Barijough, and A. Gerstlauer, "Deepthings: Distributed adaptive deep learning inference on resource-constrained iot edge clusters," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 11, pp. 2348–2359, 2018.
- [14] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [16] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [18] A. Farhadi and J. Redmon, "Yolov3: An incremental improvement," *Computer Vision and Pattern Recognition*, cite as, 2018.
- [19] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, ser. AAAI'17, 2017, pp. 4278–4284.
- [20] H.-J. Jeong, H.-J. Lee, C. H. Shin, and S.-M. Moon, "Ionn: Incremental offloading of neural network computations from mobile devices to edge servers," in *Proceedings of the ACM Symposium on Cloud Computing*, 2018, pp. 401–411.
- [21] T. Mohammed, C. Joe-Wong, R. Babbar, and M. Di Francesco, "Distributed inference acceleration with adaptive dnn partitioning and offloading," in *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*. IEEE, 2020, pp. 854–863.
- [22] E. Li, L. Zeng, Z. Zhou, and X. Chen, "Edge ai: On-demand accelerating deep neural network inference via edge computing," *IEEE Transactions on Wireless Communications*, vol. 19, no. 1, pp. 447–457, 2019.
- [23] S. Laskaridis, S. I. Venieris, M. Almeida, I. Leontiadis, and N. D. Lane, "Spinn: synergistic progressive inference of neural networks over device and cloud," in *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking*, 2020, pp. 1–15.
- [24] B. Cao, L. Zhang, Y. Li, D. Feng, and W. Cao, "Intelligent offloading in multi-access edge computing: A state-of-the-art review and framework," *IEEE Communications Magazine*, vol. 57, no. 3, pp. 56–62, 2019.
- [25] J. Mao, X. Chen, K. W. Nixon, C. Krieger, and Y. Chen, "Modnn: Local distributed mobile computing system for deep neural network," in *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2017. IEEE, 2017, pp. 1396–1401.
- [26] N. Abbas, Y. Zhang, A. Taherkordi, and T. Skeie, "Mobile edge computing: A survey," *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 450–465, 2017.
- [27] C.-J. Wu, D. Brooks, K. Chen, D. Chen, S. Choudhury, M. Dukhan, K. Hazelwood, E. Isaac, Y. Jia, B. Jia *et al.*, "Machine learning at facebook: Understanding inference at the edge," in *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 2019, pp. 331–344.
- [28] *Intel Core i7-8700 Processor*, 2020 (accessed October 10, 2020). [Online]. Available: <https://www.intel.com/content/www/us/en/products/processors/core/core-vpro/i7-8700.html>
- [29] *NVIDIA GeForce RTX 2080 Ti GPU*, 2020 (accessed October 12, 2020). [Online]. Available: <https://www.nvidia.com/en-us/geforce/graph-ics-cards/rtx-2080-ti/>
- [30] J. Hartmanis, "Computers and intractability: a guide to the theory of np-completeness (michael r. garey and david s. johnson)," *Siam Review*, vol. 24, no. 1, p. 90, 1982.
- [31] J. Nossack and E. Pesch, "A branch-and-bound algorithm for the acyclic partitioning problem," *Computers & operations research*, vol. 41, pp. 174–184, 2014.
- [32] R. Sedgewick and K. Wayne, "Algorithms," 2011.
- [33] *NVIDIA Jetson Nano 2GB Developer Kit*, 2020 (accessed October 22, 2020). [Online]. Available: <https://developer.nvidia.com/embedded/jets-on-nano-2gb-developer-kit>
- [34] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*. MIT press Cambridge, 2016, vol. 1.
- [35] Google, LLC, "grpc." [Online]. Available: <https://grpc.io/>
- [36] Facebook, Inc. & Microsoft Corporation, "Onnx." [Online]. Available: <https://onnx.ai/>
- [37] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, "Pytorch." [Online]. Available: <https://pytorch.org/>
- [38] Aric Hagberg, Pieter Swart, Dan Schult, "Networkx." [Online]. Available: <https://networkx.org/>
- [39] Y. Harchol, A. Mushtaq, V. Fang, J. McCauley, A. Panda, and S. Shenker, "Making edge-computing resilient," in *Proceedings of the 11th ACM Symposium on Cloud Computing*, 2020, pp. 253–266.