

Intra-DP: A High Performance Distributed Inference System on Robotic IoT

Anonymous Author(s)

Submission Id: 445*

Abstract

The rapid advancements in machine learning (ML) techniques have led to significant achievements in various robotic tasks, where deploying these ML approaches on real-world robots requires fast and energy-efficient inference of their deep neural network (DNN) models. Distributed inference, which involves inference across multiple powerful GPU devices, has emerged as a promising optimization to improve inference performance in modern data centers. However, when deployed over Internet of Things of these real-world robots (robotic IoT), all existing parallel methods (data parallelism, tensor parallelism, pipeline parallelism) fail to simultaneously meet the robots' latency and energy requirements and raise significant challenges due to the failure of data parallelism, the unacceptable communication overhead of tensor parallelism, and the significant transmission bottlenecks in pipeline parallelism due to the limited bandwidth of robotic IoT.

We present Intra-DP, the first high-performance distributed inference system optimized for model inference on robotic IoT. Intra-DP introduces a fine-grained approach to transmission and computation by confining them to each local operator of DNN layers (i.e., operators that can be computed independently without the complete input, such as the convolution kernel in the convolution layer). By adaptively scheduling the computation and transmission of each local operator based on various hardware conditions and network bandwidth, Intra-DP enables different local operators of different layers to be computed and transmitted concurrently, and overlap the computation and transmission phases within the same inference task to achieve fast and energy-efficient distributed inference on robotic IoT. The evaluation shows that, Intra-DP reduces inference time by 20% and energy consumption by 10% compared with the state-of-the-art baselines.

Keywords: Distributed inference, Robotic IoT, Distributed system and network

1 Introduction

The rapid progress in machine learning (ML) techniques has led to remarkable achievements in various fundamental robotic tasks, such as object detection [18, 28, 31], robotic control [22, 46, 56], and environmental perception [4, 23, 51]. However, deploying these ML applications on real-world robots requires fast and energy-efficient inference of their deep neural network (DNN) models, given the need for swift

environmental responses and the limited battery capacity of robots. Placing the entire model on robots not only requires additional computing accelerators on robots (e.g., GPU [35], FPGA [36], SoC [15]), but also introduce additional energy consumption (e.g., 162% more for [31] in our experiments) due to the computationally intensive nature of DNN models, while placing the entire model in the cloud brings an extended response delay.

Distributed inference, which involves inference across multiple GPU devices, has emerged as a promising approach to meet the latency requirements of robotic applications and extend the battery lifetime of robots. This paradigm has been widely adopted in data centers [17, 52, 57], where numerous GPUs are utilized to speed large model inference, such as in the case of ChatGPT [50]. Adopting distributed inference across robots and other powerful GPU devices through the Internet of Things for these robots (robotic IoT) not only accelerates the inference process by leveraging the high computing capabilities of powerful GPUs but also alleviates the local computational burden, thereby reducing energy consumption, making it an ideal solution for robotic applications.

However, all existing parallel methods for distributed inference in the data center are ill-suited for robotic IoT. In data centers, there are mainly three kinds of parallel methods: Data parallelism (DP) replicates the model across devices, and lets each replica handle one mini-batch (i.e., a subset that slices out of an input data set); Tensor parallelism (TP) splits a single DNN layer over devices; Pipeline parallelism (PP) places different layers of a DNN model over devices (layer partitioning) and pipelines the inference to reduce devices' idling time (pipeline execution).

For DP, the small batch sizes inherent to robotic IoT applications (typically 1) hinder the mini-batch computation, rendering DP inapplicable for robotic IoT. In the data center, DP is feasible due to the large batch sizes employed (e.g., 16 images), allowing for the division of inputs into mini-batches that still contain several complete inputs (e.g., 2 images). However, in robotic IoT, real-time performance is crucial, necessitating immediate inference upon receiving inputs, which typically have smaller batch sizes (e.g., 1 image). Further splitting these inputs would result in mini-batches containing incomplete inputs (e.g., 1/4 of an image), which cannot be computed parallel to speed up inference.

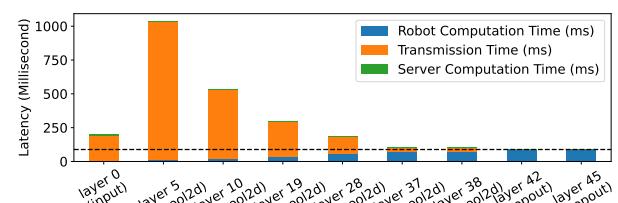
TP requires frequent synchronization among devices, leading to unacceptable communication overhead in robotic IoT.

111 By partitioning parameter tensors of a layer across GPUs,
 112 TP allows concurrent computation on different parts of this
 113 tensor but requires an all-reduce communication [57] to
 114 combine computation results from different devices, which
 115 entails significant communication overhead. Consequently,
 116 TP is used mainly for large layers that are too large to fit in
 117 one device in data centers and require dedicated high-speed
 118 interconnects (e.g., 400 Gbps for NVLink [21]) even within
 119 data centers. On the contrary, robots must prioritize seam-
 120 less mobility and primarily depend on wireless connections,
 121 which inherently possess limited bandwidth, as described in
 122 Sec. 2.1, making all-reduce synchronization an unacceptable
 123 overhead (e.g., the inference time with TP was up to 143.9X
 124 slower than local computation in Sec.2.3.).

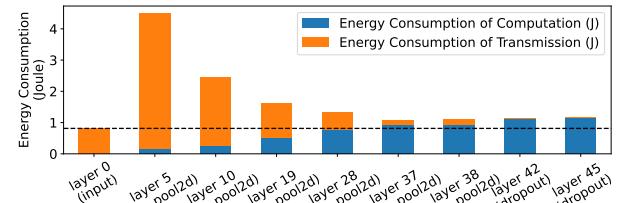
125 Consequently, existing distributed inference approaches [5,
 126 24] in robotic IoT primarily adopt the PP paradigm and fo-
 127 cus on layer partitioning of PP, aiming to achieve fast and
 128 energy-efficient inference. This is because the PP paradigm
 129 in data centers consists of layer partitioning and pipeline
 130 execution, where the pipeline execution of PP enhances infer-
 131 ence throughput rather than reducing the completion time of
 132 a single inference [6], which is the most critical requirement
 133 in robotic IoT. Based on the fact that the amounts of output
 134 data in some intermediate layers of a DNN model are sig-
 135 nificantly smaller than that of its raw input data [16], DNN
 136 layer partitioning strategies constitute various trade-offs be-
 137 tween computation and transmission, taking into account
 138 application-specific inference speed requirements and en-
 139 ergy consumption demands, as shown in Fig. 1.

140 However, existing methods based on PP face significant
 141 transmission bottlenecks in robotic IoT due to the inherent
 142 scheduling mechanism. PP paradigm on robotic IoT consists
 143 of three sequential phases: computing earlier DNN layers on
 144 robots, transmitting intermediate results, and completing in-
 145 ference on the GPU server. Despite optimal layer partitioning
 146 strategies from [5, 24], the transmission overhead becomes
 147 a substantial bottleneck due to limited bandwidth of robotic
 148 IoT, accounting for up to 69% of inference time in our ex-
 149 periments. While the pipeline execution of PP can mitigate
 150 this overhead by overlapping computation and transmis-
 151 sion phases across multiple inference tasks, it cannot reduce
 152 the completion time of a single inference task [6], which is
 153 crucial for robotic applications. Furthermore, the prolonged
 154 transmission phase not only slows down inference speed but
 155 also consumes significantly more energy.

156 The key reason for the problem of the above methods is
 157 that existing methods conduct layer-granulated scheduling,
 158 whose transmission time is typically longer than the computa-
 159 tion time due to the limited bandwidth of real-world robotic
 160 IoT networks. As transmission time constitutes a substantial
 161 portion of the total inference time (approximately half) in
 162 existing methods, a novel parallel method with finer gran-
 163 ularity that overlaps computation and transmission within
 164 the same inference task has the potential to address this



(a) Inference Latency



(b) Energy consumption on robot

Figure 1. Existing distributed inference approaches on VGG19 [42] in our experiments, which adopt PP paradigm with various layer partitioning scheduling strategies. The X-axis of the graph represents different layer partitioning strategies, where ‘layer i’ indicates that all layers up to and including the i-th layer are computed on the robot, while the subsequent layers are processed on the GPU server.

shortcoming, achieving fast and energy-efficient inferences. Note that the robot can not enter low-power sleep mode during the transmission phase due to the need to promptly continue working upon receiving inference results, but can only enter standby mode, when chips like CPU, GPU, and memory consume non-negligible power even when not computing (e.g., 95% power consumption in our experiments). Such a parallel method would reduce the robot’s standby time without significantly increasing energy consumption during the computation phase, thereby decreasing overall energy consumption.

In this paper, we present Intra-DP (Intra-Data Parallel), a high performance distributed inference system optimized for real-world robotic IoT networks. We discovered that operators for each DNN layer (e.g., convolution, ReLU, softmax) can be categorized into two types: local operators and global operators, depending on whether they can be computed independently without the complete input. For instance, softmax [29] requires the complete input vector to calculate the corresponding probability distribution, referring to it as a global operator, while ReLU [7] and convolution [32] can be computed with part of the input tensor (the elements in the input vector for ReLU and the blocks in the input tensor for convolution), referring to them as local operators. Local

operators are widely used in robotic applications, especially convolution layers in computer vision [31] and point cloud tasks [46]. The local operator granularity provides a finer granularity for Intra-DP, allowing different local operators of different layers to be computed and transmitted concurrently, enabling the overlapping of computation and transmission phases within the same inference task to achieve fast and energy-efficient inference.

The design of Intra-DP is confronted with two major challenges. The first one is how to guarantee the correctness of inference results based on local operator. We propose Local Operator Parallelism (LOP), which reduces the granularity of calculation from each layer to each local operator. LOP determines the correct input required for different local operators based on their calculation characteristics and processes at first. When a part of the local operators in a layer completes the calculation and the tensor composed of these local operators satisfies the input requirements of the local operators in the subsequent layer, the local operators in the subsequent layer can be calculated in advance, without waiting for all local operators of the current layer to be computed in LOP. For global operator layers, Intra-DP enforces a synchronization before these layers to combine the complete input for them, as TP's all-reduce communications do. In this way, Intra-DP only change the execution sequence of local operators among local operator layers and ensures the calculation correctness of local operator layers through LOP and global operator layers through synchronization. We formally prove that LOP guarantees the correctness of inference results in Sec. 4.1.

The second challenge is under LOP, how to properly schedule the computation and transmission of each local operator to achieve fast and energy-efficient inference under various hardware conditions and network bandwidths. Intra-DP places part of the local operator execution on GPU servers and transmits the corresponding part of the input tensor based on LOP, while computing the rest of the local operators with a novel Local Operator Scheduling Strategy (LOSS). LOSS formulates the problem of determining which part of the local operators should be executed on robots and which part should be executed on GPU servers as a nonlinear optimization problem (see Sec. 4.2), and schedules the computation and transmission of each local operator based on the solution obtained via the differential evolution algorithm [39].

We implemented Intra-DP in PyTorch [38] and evaluated Intra-DP on our real-world robots under two typical real-world robotic applications [31, 46] and several models common to mobile devices on a larger scale [42–44, 48, 53]. We compared Intra-DP with two SOTA pipeline parallelism methods as baselines: DSCCS [24], aimed at accelerating inference, and SPSO-GA [5], focused on optimizing energy consumption, under different real-world robotic IoT networks

environments (namely indoors and outdoors). Evaluation shows that:

- Intra-DP is fast. Intra-DP reduced inference time by ~80% ~98% compared to baselines under indoors and outdoors environments.
- Intra-DP is energy-efficient. Intra-DP reduced 91% ~98% energy consumption per inference compared to baselines, due to faster inference speed and no-increased power consumption against time.
- Intra-DP is robust in various robotic IoT environments. When the robotic IoT environment changed (from indoors to outdoors), Intra-DP's superior performance remained consistent.
- Intra-DP is easy to use. It took only two lines of code to apply Intra-DP to existing ML applications.

Our main contribution are LOP, a fine-grained parallel method based on local operators, and LOSS, a new scheduling strategy based on LOP optimized for distributed inference over real-world robotic IoT networks. By leveraging these contributions, Intra-DP dramatically reduces the transmission overhead in existing distributed inference on robotic IoT by overlapping the computation and transmission phases within the same inference task, achieving fast and energy-efficient distributed inference on robotic IoT. We envision that the fast and energy-efficient inference of Intra-DP will foster the deployment of diverse robotic tasks on real-world robots in the field. Intra-DP's code is released on <https://github.com/xxx/xxx>.

In the rest of this paper, we introduce the background of this paper in Sec. 2, give an overview of Intra-DP in Sec. 3, present the detailed design of Intra-DP in Sec. 4, evaluate Intra-DP in Sec. 6, and finally conclude in Sec. 7.

2 background

2.1 Characteristics of Robotic IoT

In real-world robotic IoT scenarios, devices often navigate and move around for tasks like search and exploration. While wireless networks provide high mobility, they also have limited bandwidth. For instance, Wi-Fi 6, the most advanced Wi-Fi technology, offers a maximum theoretical bandwidth of 1.2 Gbps for a single stream [27]. However, not only the limited hardware resources on the robot can not fully play the potential of Wi-Fi 6 [55], but also the actual available bandwidth of wireless networks is often reduced in practice due to factors such as movement of the devices [30, 37], occlusion from physical barriers [9, 41], and preemption of the wireless channel by other devices [2, 40].

To demonstrate the instability of wireless transmission in real-world situations, we conducted a robot surveillance experiment using four-wheel robots navigating around several given points at 5–40 cm/s speed in our lab (indoors) and campus garden (outdoors), with hardware and wireless network settings as described in Sec. ???. We believe our setup

represents robotic IoT devices' state-of-the-art computation and communication capabilities. We saturated the wireless network connection with iperf [1] and recorded the average bandwidth capacity between these robots every 0.1s for 5 minutes.

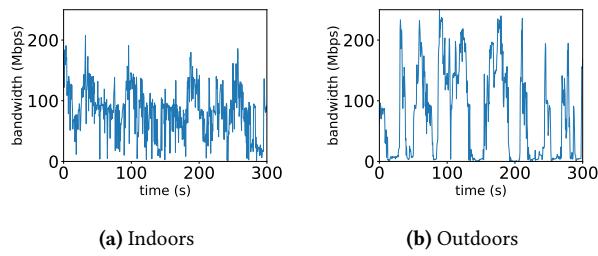


Figure 2. The instability of wireless transmission in robotic IoT networks.

The results in Fig. 2 show average bandwidth capacities of 93 Mbps and 73 Mbps for indoor and outdoor scenarios, respectively. The outdoor environment exhibited higher instability, with bandwidth frequently dropping to extremely low values around 0 Mbps, due to the lack of walls to reflect wireless signals and the presence of obstacles like trees between communicating robots, resulting in fewer received signals compared to indoor environments.

In summary, robotic IoT systems' wireless transmission is constrained by limited bandwidth, both due to the theoretical upper limit of wireless transmission technologies and the practical instability of wireless networks. Moreover, the unstable bandwidth in robotic IoT wireless networks can cause the transmission time to vary dramatically, sometimes changing by hundreds of times. In our experiments, even a relatively small model with only 0.84M parameters still suffers from its significant transmission overhead.

2.2 Characteristics of Data Center Networks

Data center networks, which are used for large model inference (e.g., ChatGPT [50]), are wired and typically exhibit higher bandwidth capacity and lower fluctuation compared to robotic IoT networks. GPU devices in data centers are interconnected using high-speed networking technologies such as InfiniBand [45] or PCIe [21], offering bandwidths ranging from 40 Gbps to 500 Gbps. The primary cause of bandwidth fluctuation in these networks is congestion on intermediate switches, which can be mitigated through traffic scheduling techniques implemented on the switches [34]. The stable and high-bandwidth nature of data center networks makes them well-suited for demanding tasks like large model inference, in contrast to the more variable and resource-constrained environments found in robotic IoT networks.

2.3 Existing distributed inference methods in the data center

Data parallelism. Data parallelism [52] is a widely used technique in distributed inference that partitions input data across multiple devices, such as GPUs, to perform parallel inference. Each device maintains a complete model replica and independently processes a subset of the input data (mini-batch), aggregating results to generate the final output. Data parallelism enhances throughput by distributing workload across devices, leveraging their combined computational power.

However, data parallelism's scalability is constrained by the total batch size [33], which is particularly problematic in robotic IoT applications where smaller batch sizes are inherent due to the need for swift environmental responses. In robotic applications, immediate inference upon receiving inputs is crucial for obtaining real-time target points quickly. For example, in our experiments, the robot constantly obtains the latest images from the camera for inference, with a batch size of only 1. These small batches cannot be further split into mini-batches, a fundamental requirement for effective data parallelism.

Tensor parallelism. Tensor parallelism [57] is a distributed inference technique that divides a model's layer parameters across multiple devices, each storing and computing a portion of the weights. This approach requires an all-reduce communication step after each layer to combine results from different devices, introducing significant overhead, especially for large DNN layers. To mitigate this, TP is typically deployed across GPUs within the same server in data centers, using fast intra-server GPU-to-GPU links like NVLink [21], which is beneficial when the model is too large for a single device.

In contrast to data center networks, the limited bandwidth in robotic IoT renders the communication cost of TP prohibitively high, as evidenced by our evaluation of DINA [32], a state-of-the-art TP method. Table 1 reveals that transmission time constitutes 49% to 94% of the total inference time due to all-reduce communication for each layer, resulting in TP's inference time being 45.2X to 143.9X longer than local computation. Although Table 2 indicates lower power consumption with TP (13.4% to 67.3% less than local computation), the extended transmission times significantly increase energy consumption per inference, ranging from 28.5X to 62.7X. Since TP significantly extends inference time, making it impractical for real-world robotic applications, we did not further consider TP in this paper.

Pipeline parallelism. Pipeline parallelism [17] is a distributed inference technique that partitions DNN model layers across multiple devices (layer partitioning), forming an inference pipeline for concurrent processing of multiple tasks. While PP can increase throughput and resource utilization via pipeline execution, it primarily focuses on

Model(number of parameters)	Local computation time(s)	Environment	Transmission time (s) with TP	Inference time (s) with TP	Percentage(%) with TP
MobileNet_V3_Small(2M)	0.031(± 0.004)	indoors	0.698(± 0.135)	1.400(± 0.232)	49.85
		outdoors	0.901(± 0.778)	1.775(± 1.370)	51.23
ResNet101(44M)	0.065(± 0.005)	indoors	7.156(± 3.348)	8.106(± 3.403)	87.95
		outdoors	8.470(± 6.337)	9.356(± 6.328)	90.46
VGG19_BN(143M)	0.063(± 0.002)	indoors	5.152(± 4.873)	5.444(± 4.831)	70.18
		outdoors	5.407(± 6.673)	5.759(± 6.635)	93.70

Table 1. Average transmission time (Second), inference time (Second), percentage that transmission time accounts for of inference time and their standard deviation ($\pm n$) with TP on different models in different environments. “Local computation” refers to placing the whole layers on the robot.

Model(number of parameters)	Environment	Power consumption(W)		Energy consumption(J) per inference	
		Local	TP	Local	TP
MobileNet_V3_Small(2M)	indoors	6.05(± 0.21)	5.24(± 0.19)	0.3(± 0.09)	7.33(± 1.21)
	outdoors	6.05(± 0.21)	5.11(± 0.28)	0.3(± 0.09)	9.08(± 7.0)
ResNet101(44M)	indoors	11.27(± 0.51)	4.97(± 0.16)	0.93(± 0.19)	40.28(± 16.91)
	outdoors	11.27(± 0.51)	4.9(± 0.23)	0.93(± 0.19)	45.8(± 30.98)
VGG19_BN(143M)	indoors	14.86(± 0.43)	4.88(± 0.29)	1.19(± 0.18)	26.55(± 23.56)
	outdoors	14.86(± 0.43)	4.87(± 0.27)	1.19(± 0.18)	28.06(± 32.33)

Table 2. Power consumption against time (Watt) and energy consumption per inference (Joule) with standard deviation ($\pm n$) with TP on different models in different environments. “Local” represents “Local computation”

enhancing overall throughput rather than reducing single-inference latency [6], which is crucial in robotic IoT. As a result, existing distributed inference approaches [5, 24] in robotic IoT primarily adopt PP paradigm and focus on layer partitioning to achieve fast and energy-efficient inference, with two main categories based on their optimization goals: accelerating inference for diverse DNN structures [16, 19, 24, 32, 54] and optimizing robot energy consumption during inference [5, 25, 49]. Both two kinds of methods suffer from the transmission bottleneck inherent to PP’s scheduling mechanism, which can be eliminated by Intra-DP.

2.4 Other methods to speed up DNN Models Inference on Robotic IoT

Compressed communication. Compressed communication is crucial for efficient distributed inference in wireless networks, as it significantly reduces communication overhead through techniques such as quantization and model distillation. Quantization [8, 12, 13] is a technique that reduces the numerical precision of model weights and activations, thereby minimizing the memory footprint and computational requirements of deep learning models. This process typically involves converting high-precision (e.g., 32-bit) floating-point values to lower-precision (e.g., 8-bit) floating-point representations, with minimal loss of model

accuracy. Model distillation [14, 26, 47], on the other hand, is an approach that involves training a smaller, more efficient “student” model to mimic the behavior of a larger, more accurate “teacher” model by minimizing the difference between the student model’s output and the teacher model’s output. The distilled student model retains much of the teacher model’s accuracy while requiring significantly fewer resources. These model compression methods complement distributed inference by achieving faster inference speed through model modifications, potentially sacrificing some accuracy with smaller models, while distributed inference realizes fast inference without loss of accuracy by intelligently scheduling computation tasks across multiple devices.

Inference Job scheduling. Significant research efforts have been devoted to exploring inference parallelism and unleashing the potential of layer partition to accelerate DNN inference, such as inference job scheduling, aiming to accelerate multiple DNN inference tasks by optimizing their execution on various devices under different network bandwidths while considering application-specific inference speed requirements and energy consumption demands. For instance, [3, 10] support online scheduling of offloading inference tasks based on the current network and resource status of mobile systems while meeting user-defined energy constraints. [11] focused on optimizing DNN inference workloads in cloud

computing using a deep reinforcement learning based scheduler for QoS-aware scheduling of heterogeneous servers, aiming to maximize inference accuracy and minimize response delay. While these methods focus on overall optimization in multi-task scenarios involving multi-robots, they do not address the optimization of single inference tasks and are thus orthogonal to distributed inference for a single inference, where improved distributed inference can provide faster and more energy-efficient inference for these scenarios.

3 Overview

3.1 Workflow of Intra-DP

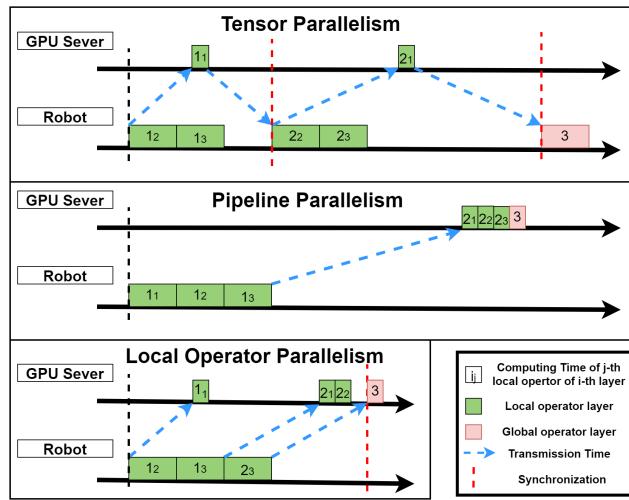


Figure 3. Workflow of Intra-DP. Each local operator layer have to complete the calculation of three local operators, and the same local operator in the three cases has the same computation time on robots and GPU servers, as well as the corresponding transmission time. The output tensor volume of layer 2 is larger than that of layer 1, resulting in longer transmission times for local operators in layer 2, and PP selects a layer partition strategy at layer 1 [24].

Fig. 3 presents the workflow of Intra-DP and compares it with TP and PP under robotic IoT networks with limited bandwidth, illustrating why existing methods suffer from transmission overhead and how Intra-DP solves this issue via its LOP.

While TP can place some local operator execution on the GPU server, it requires an all-reduce communication [57] to combine computation results from different devices, which entails significant communication overhead (as shown by the red dotted lines for synchronization Fig. 3). Although the layer partition algorithm [24] can be used to minimize overall inference time, the transmission time still becomes a significant bottleneck, as illustrated by the extremely long transmission in Fig. 3.

To alleviate the transmission overhead in distributed inference, Intra-DP overlaps the computation and transmission of different local operators from different local operator layers, as shown in Fig. 3. Compared with TP, Intra-DP cancels the synchronization of the all-reduce communication for local operator layers and ensures that each local operator can get the required input in time and obtain the correct calculation results through LOP, rather than relying on all-reduce communication for local operator layers. Intra-DP maintains synchronization for global operator layers, which do not have local operators and require the complete input (the entire output tensor of the previous layer) to perform the calculation, ensuring the correctness of the global operator layers' inference. Compared with PP, Intra-DP starts transmitting some local operators and the corresponding input in advance, without waiting for all local operators in the current local operator layer to complete the calculation. In this way, Intra-DP achieves much faster inference compared with existing distributed inference methods.

Moreover, the idle time on the robot (when the robot is not computing, as shown in Fig. 3) consumes significant energy. This is because the robot cannot enter a low-power sleep mode while waiting for the final inference result from the GPU server, as it has to promptly continue working when it receives the inference results. During the standby phase (idle time), chips like CPU, GPU, and memory consume non-negligible power even when not computing, due to the static power consumption rooted in transistors' leakage current [20]. Meanwhile, we found that wireless network cards consume only 0.21 Watt for transmission during the idle time, while the robot consumes 13.35 Watt during computing. In this way, Intra-DP dramatically reduces the idle time on the robot, alleviating the energy wasted by standby mode, and increases a negligible amount of network card transmission power consumption during computing, thereby reducing the overall energy consumption for each inference.

To achieve the workflow shown in Fig. 3, the design of Intra-DP must tackle two problems: guaranteeing the correctness of inference results based on local operators and scheduling the computation and transmission of each local operator. In Sec.4.1, we will explain how Intra-DP ensures that each local operator can still obtain the correct calculation result via LOP, and in Sec.4.2, we will discuss how Intra-DP achieves fast and energy-efficient inference through its LOSS.

3.2 Architecture of Intra-DP

Fig. 5 shows the architecture of Intra-DP.

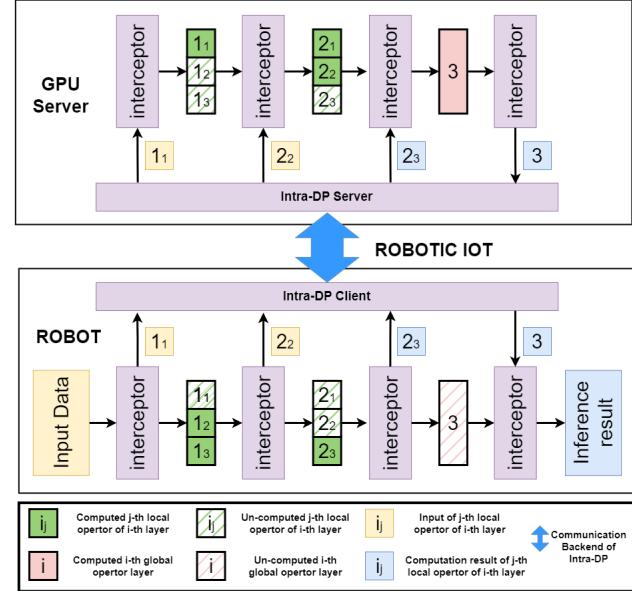


Figure 4. Architecture of Intra-DP. Intra-DP adopts the same scheduling scheme as in Fig. 3.

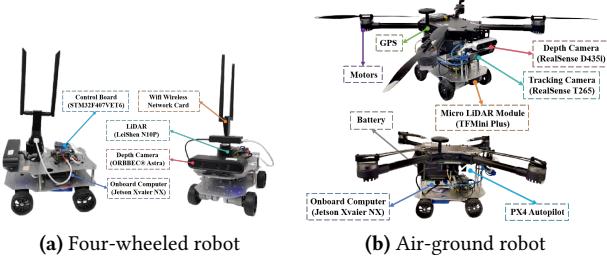


Figure 5. The detailed composition of the robot platforms

4 Detailed Design

4.1 Local Operator Parallelism

4.2 Local Operator Scheduling Strategy

4.3 Algorithms of Intra-DP

5 Implementation

6 Evaluation

Testbed. The evaluation was conducted on a custom four-wheeled robot (Fig 5a), and a custom air-ground robot(Fig 5b). They are equipped with a Jetson Xavier NX [35] 8G onboard computer that is capable of AI model inference with local computation resources. The system runs Ubuntu 20.04 with ROS Noetic and a dual-band USB network card (MediaTek MT76x2U) for wireless connectivity. The Jetson Xavier NX interfaces with a Leishen N10P LiDAR, ORBBEC Astra depth camera, and an STM32F407VET6 controller via USB serial

Algorithm 1: Local_Operator_Scheduling_Strategy

```

716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770

```

Input: Cuda kernel function called by the corresponding operator *func* and the required parameters *args*

Output: The execution result *ret*

Data: inference operator sequence *IOS* = \emptyset

- if** *IOS.empty()* **then**
 - // recorder
 - SendRPCtoServer(func, args)*
 - IOS = DataDependencySearch(func, args)*
 - ret = GetRPCExecutionResult()*
 - RecordReturn(ret)*
- end**
- else**
 - // replayer on robot
 - if** *func == IOS.start()["func"]* **then**
 - ret = StartRRTO(args, IOS)*
 - // start a new inference
 - end**
 - else if** *func == IOS.end()["func"]* **then**
 - ret = WaitingForRRTO()*
 - // Waiting for the final computation result
 - end**
 - else**
 - ret = IOS.find(func)["ret"]*
 - end**
- end**
- return** *ret*

	inference	transmission	standby
Power (W)	13.35	4.25	4.04

Table 3. Power consumption (Watt) of our robot in different states.

ports. Both LiDAR and depth cameras facilitate environmental perception, enabling autonomous navigation, obstacle avoidance, and SLAM mapping. The GPU server accepting offloaded computation tasks from the robot is a PC equipped with an Intel(R) i5 12400f CPU @ 4.40GHz and an NVIDIA GeForce GTX 2080 Ti 11GB GPU, connected to our robot via Wi-Fi 6 over 80MHz channel at 5GHz frequency in our experiments.

Tab. 3 presents the overall on-board energy consumption (excluding motor energy consumption for robot movement) of the robot in various states: inference (model inference with full GPU utilization, including CPU and GPU energy consumption), transmission (communication with the GPU server, including wireless network card energy consumption), and standby (robot has no tasks to execute). Notice that different models, due to varying numbers of parameters,

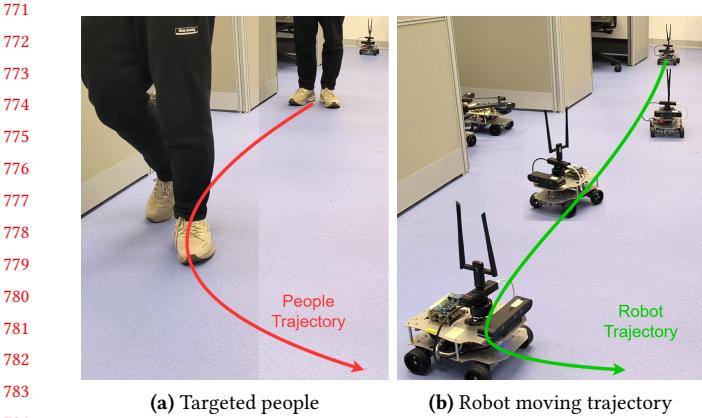


Figure 6. A real-time people-tracking robotic application on our robot based on a well-known human pose estimation ML model, Kapao [31].

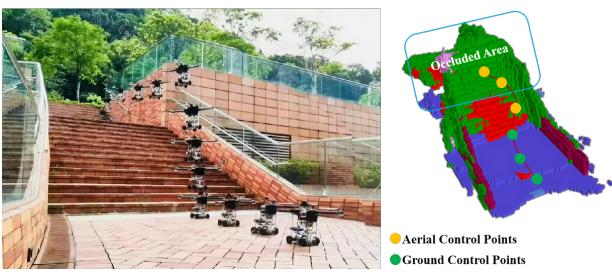


Figure 7. By predicting occlusions in advance, AGRNav [46] gains an accurate perception of the environment and avoids collisions, resulting in efficient and energy-saving paths.

exhibit distinct GPU utilization rates and power consumption during inference.

Experiment Environments. We evaluated two real-world environments: indoors (robots move in our laboratory with desks and separators interfering with wireless signals) and outdoors (robots move in our campus garden with trees and bushes interfering with wireless signals, resulting in lower bandwidth). The corresponding bandwidths between the robot and the GPU server in indoors and outdoors scenarios are shown in Fig. 2.

Workload. We evaluated two typical real-world robotic applications on our testbed: Kapao, a real-time people-tracking application on our four-wheeled robot (Fig 6), and AGRNav, an autonomous navigation application on our air-ground robot (Fig 7). These applications feature different model input and output size patterns: Kapao takes RGB images as input and outputs key points of small data volume. In contrast, AGRNav takes point clouds as input and outputs predicted point clouds and semantics of similar data volume as input, implying that AGRNav needs to transmit more data during offloading. And we have verified several models common

to mobile devices on a larger scale to further corroborate our observations and findings: MobileNet [43], ResNet [44], VGGNet [42], ConvNeXt [48], RegNet [53].

Baselines. We selected two SOTA pipeline parallelism methods as baselines: DSCCS [24], aimed at accelerating inference, and SPSO-GA [5], focused on optimizing energy consumption. We set SPSO-GA’s deadline constraints to 1 Hz, the minimum frequency required for robot movement control. Given our primary focus on inference time and energy consumption per inference, we disabled pipeline execution to concentrate solely on assessing the performance of various layer partitioning methods.

6.1 Inference Time

Kapao. From the results in the upper part of Tab. 4, both SPSO-GA and DSCCS reduced Kapao’s inference time by 39.69% and 56.92% indoors and 28.67% and 47.46% outdoors, with DSCCS achieving 28.57% (indoors) and 26.34% (outdoors) lower inference time than SPSO-GA. While both systems significantly reduced inference time via offloading, transmission time accounts for 49.69% to 69.46% of the whole inference time, indicating that even with SOTA layer partitioning, the transmission bottleneck inherent to PP’s scheduling mechanism cannot be mitigated. The difference between DSCCS and SPSO-GA can be attributed to their optimization goals: DSCCS minimizes inference latency, while SPSO-GA minimizes power consumption under deadline constraints.

AGRNav. The performance gain of the two offloading systems varied for AGRNav, as shown in the lower part of Tab. 4. DSCCS still reduced inference time by 18.34% and 12.43% in indoors and outdoors. However, SPSO-GA achieved similar inference time (3.65% and 3.06% reduction) as local computation both indoors and outdoors. We will explain and analyze this phenomenon in Sec.6.2.

Notice that the large standard deviation in transmission time in outdoors in both offloading systems indicates that bandwidth fluctuated more frequently and more fiercely outdoors compared with indoors, which complies with Fig. 2. Additionally, the lower average bandwidth for outdoors scenarios (see Sec.2.1) results in increased transmission and inference times relative to indoor scenarios.

6.2 Breakdown

Both SPSO-GA and DSCCS automatically adapt to available bandwidth, transitioning to edge computation (placing all DNN layers on the GPU server) when bandwidth is sufficient, and to local computation (placing all DNN layers on robots) when bandwidth is low. To better understand how their layer partitioning scheduling varies with different network conditions and models, we recorded and analyzed the Categories and percentages of various layer partitioning schedules under different baselines and environments, as detailed in Fig. 8.

Model(number of parameters)	Local computation time (s)	Environment	Transmission time (s)			Inference time (s)		
			DSCCS	SPSO-GA	Intra-DP	DSCCS	SPSO-GA	Intra-DP
kapao(77M)	0.78(± 0.23)	indoors	0.228(± 0.176)	0.235(± 0.164)	0.259(± 0.181)	0.343(± 0.192)	0.311(± 0.168)	0.264(± 0.189)
		outdoors	0.087(± 0.178)	0.35(± 1.045)	0.385(± 1.15)	0.696(± 0.125)	0.434(± 1.046)	0.369(± 0.140)
agrnav(0M)	1.12(± 0.11)	indoors	0.266(± 0.166)	0.239(± 0.149)	0.263(± 0.164)	0.433(± 0.134)	0.427(± 0.12)	0.365(± 0.130)
		outdoors	0.263(± 0.843)	0.218(± 0.796)	0.24(± 0.875)	0.512(± 0.779)	0.536(± 0.72)	0.456(± 0.74)

Table 4. Average transmission time (Second), inference time (Second), percentage that transmission time accounts for of inference time and their standard deviation ($\pm n$) of Kapao and AGRNav with different pipeline parallelism offloading systems and different environments. “Local computation” refers to placing the whole layers on the robot.

Model(number of parameters)	Environment	Power consumption(W)			Energy consumption(J) per inference			
		Local	DSCCS	SPSO-GA	Intra-DP	Local	DSCCS	SPSO-GA
kapao(77M)	indoors	15.04(± 0.64)	7.03(± 3.57)	5.92(± 2.18)	5.03(± 1.87)	15.03(± 0.63)	2.41(± 1.35)	1.84(± 1.0)
	outdoors	15.04(± 0.64)	14.15(± 1.71)	5.89(± 2.3)	5.02(± 1.97)	1.56(± 0.85)	9.85(± 1.77)	2.56(± 6.17)
agrnav(0.84M)	indoors	10.26(± 1.58)	6.6(± 1.95)	6.74(± 2.03)	5.74(± 1.78)	10.82(± 1.44)	2.86(± 0.88)	2.88(± 0.81)
	outdoors	10.26(± 1.58)	7.36(± 2.34)	7.91(± 2.45)	6.71(± 2.1)	10.82(± 1.44)	3.77(± 5.74)	4.25(± 5.7)

Table 5. The power consumption against time (Watt) and energy consumption per inference (Joule) with standard deviation ($\pm n$) of Kapao and AGRNav at different baselines and environments. “Local” represents “Local computation”

Local computation and edge computation are special cases of layer partitioning, with the bandwidth conditions required for each model to reach these cases varying based on the model structure and partitioning method used. Analyzing Fig. 8a and Fig. 8b, both SPSO-GA and DSCCS tend to allocate more layers on the robot for AGRNav. When comparing indoor and outdoor scenarios in Fig. 8, it is evident that higher bandwidth leads to more layers being scheduled on GPU server. Additionally, when comparing SPSO-GA and DSCCS in Fig. 8, DSCCS, which focuses on optimizing energy consumption, tends to place fewer layers on the robot to reduce computation energy consumption.

In summary, the conditions under which layer partitioning schemes make these special cases are influenced by multiple factors: model structure, and the trade-offs between inference delay and energy consumption. And the higher the bandwidth, the more layers are scheduled to be placed on GPU server.

6.3 Energy Consumption

Kapao. From the results in the upper part of Tab. 5, DSCCS consumed 3.38% and 2.02% more power per second than SPSO-GA indoors and outdoors due to more layers placed on robots shown in Fig. 8a. However, SPSO-GA consumed 58.54% and 49.74% more energy overall to process a frame than DSCCS because it only aims at minimizing the power consumption against time at the cost of possibly prolonged inference time.

AGRNav. From the results in the lower part of Tab. 5, DSCCS consumed 61.21% and 22.92% more energy per second than SPSO-GA indoors and outdoors (Tab. 5), while DSCCS consumed 34.15% and 5.43% more energy to process a frame than SPSO-GA. SPSO-GA’s advantages in power consumption aganist time shrinks in energy consumption per inference due to prolonged inference time.

6.4 Validation on a larger range of models

We evaluated PP across a broad range of models with varying parameter counts (from 0.84M to 644M, as detailed in Tab. 6 and Tab. 7), which are commonly used in mobile devices. Our findings confirm that transmission time constitutes a significant portion of the total inference time in robotic IoT when using PP. The inherent transmission overhead of PP’s scheduling mechanism significantly wastes both inference time and energy.

7 Conclusion

In this paper, we explored the problems that hinder the application of existing parallel methods for distributed inference on robotic IoT, including the failure of data parallelism due to small batch sizes, the unacceptable communication overhead of tensor parallelism caused by all-reduce communication, and the significant transmission bottlenecks inherent to pipeline parallelism’s scheduling mechanism. By raising awareness of these issues, we aim to stimulate research efforts towards developing novel parallel methods that address these problems. We envision that fast and energy-efficient

Model(number of parameters)	Local computation time (s)	Environment	Transmission time (s)		Inference time (s)		Percentage(%)	
			SPSO-GA	DSCCS	SPSO-GA	DSCCS	SPSO-GA	DSCCS
MobileNet_V3_Small (2M)	0.033(± 0.019)	indoors	0.035(± 0.019)	0.016(± 0.005)	0.044(± 0.020)	0.031(± 0.008)	79.79	1048 53.24
		outdoors	0.035(± 0.044)	0.017(± 0.005)	0.047(± 0.037)	0.033(± 0.018)	50.04	1049 51.49
RegNet_X_3_2GF (15M)	0.060(± 0.022)	indoors	0.049(± 0.026)	0.033(± 0.011)	0.065(± 0.028)	0.049(± 0.016)	76.25	1051 64.17
		outdoors	0.049(± 0.055)	0.032(± 0.032)	0.069(± 0.050)	0.051(± 0.030)	53.23	1052 44.50
ResNet101 (44M)	0.060(± 0.023)	indoors	0.054(± 0.451)	0.033(± 0.010)	0.072(± 0.453)	0.050(± 0.016)	75.64	1053 57.37
		outdoors	0.052(± 0.064)	0.033(± 0.036)	0.077(± 0.059)	0.054(± 0.034)	51.54	1054 42.48
ConvNeXt_Base (88M)	0.047(± 0.004)	indoors	0.033(± 0.018)	0.020(± 0.006)	0.044(± 0.019)	0.032(± 0.009)	75.39	1055 49.37
		outdoors	0.032(± 0.038)	0.020(± 0.022)	0.045(± 0.033)	0.034(± 0.019)	52.82	1056 35.63
ConvNeXt_Large (197M)	0.051(± 0.005)	indoors	0.033(± 0.017)	0.023(± 0.008)	0.046(± 0.019)	0.035(± 0.013)	72.96	1057 62.68
		outdoors	0.032(± 0.038)	0.023(± 0.024)	0.054(± 0.040)	0.041(± 0.028)	48.94	1058 43.96
RegNet_Y_128GF (644M)	0.139(± 0.016)	indoors	0.076(± 0.289)	0.041(± 0.024)	0.305(± 0.382)	0.100(± 0.035)	23.58	1059 40.76
		outdoors	0.171(± 0.602)	0.016(± 0.055)	0.432(± 0.615)	0.117(± 0.242)	32.39	1060 9.41

Table 6. Average transmission time (Second), inference time (Second), percentage that transmission time accounts for of inference time and their standard deviation ($\pm n$) of common AI models in different environments with different offloading systems. “Local computation” refers to placing the whole layers on the robot.

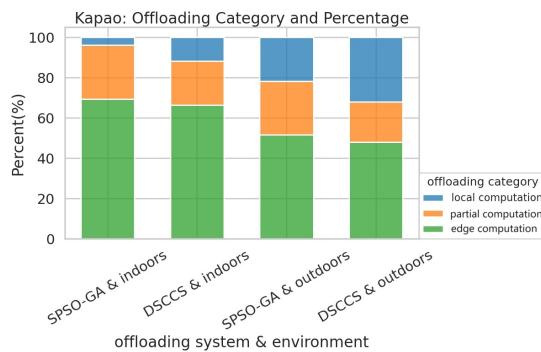
Model(number of parameters)	Environment	Power consumption(W)			Energy consumption(J) per inference			1066 1067
		Local	SPSO-GA	DSCCS	Local	SPSO-GA	DSCCS	
MobileNet_V3_Small (2M)	indoors	6.131(± 0.061)	5.448(± 0.168)	5.658(± 0.085)	0.202(± 0.002)	0.241(± 0.107)	0.174(± 0.046)	1068 1069
	outdoors	6.131(± 0.061)	5.567(± 0.273)	5.557(± 0.186)	0.202(± 0.002)	0.260(± 0.204)	0.185(± 0.099)	1070 1071
RegNet_X_3_2GF (15M)	indoors	8.208(± 0.140)	5.490(± 0.178)	5.714(± 0.342)	0.492(± 0.008)	0.356(± 0.156)	0.278(± 0.091)	1072 1073
	outdoors	8.208(± 0.140)	5.878(± 0.659)	6.041(± 0.624)	0.492(± 0.008)	0.406(± 0.295)	0.311(± 0.184)	1074 1075
ResNet101 (44M)	indoors	11.851(± 0.404)	5.457(± 0.240)	5.953(± 0.789)	0.711(± 0.024)	0.390(± 2.471)	0.298(± 0.094)	1076 1077
	outdoors	11.851(± 0.404)	6.179(± 1.083)	6.431(± 1.060)	0.711(± 0.024)	0.478(± 0.364)	0.349(± 0.216)	1078 1079
ConvNeXt_Base (88M)	indoors	15.335(± 0.273)	5.507(± 0.358)	7.713(± 2.613)	0.721(± 0.013)	0.241(± 0.103)	0.250(± 0.069)	1073 1074
	outdoors	15.335(± 0.273)	7.638(± 3.297)	9.148(± 3.338)	0.721(± 0.013)	0.346(± 0.254)	0.307(± 0.171)	1075 1076
ConvNeXt_Large (197M)	indoors	15.403(± 0.082)	5.518(± 0.638)	6.604(± 2.860)	0.786(± 0.004)	0.251(± 0.104)	0.230(± 0.088)	1077 1078
	outdoors	15.403(± 0.082)	8.400(± 4.345)	8.895(± 4.505)	0.786(± 0.004)	0.452(± 0.339)	0.366(± 0.248)	1079 1080
RegNet_Y_128GF (644M)	indoors	15.430(± 0.020)	5.384(± 1.071)	6.151(± 2.155)	2.145(± 0.003)	1.642(± 0.327)	0.615(± 0.216)	1081 1082
	outdoors	15.430(± 0.020)	6.361(± 2.349)	9.127(± 4.724)	2.145(± 0.003)	2.748(± 1.015)	1.068(± 0.553)	1083 1084

Table 7. The power consumption against time (Watt) and energy consumption per inference (Joule) with standard deviation ($\pm n$) of common AI models at different baselines and environments. “Local” represents “Local computation”

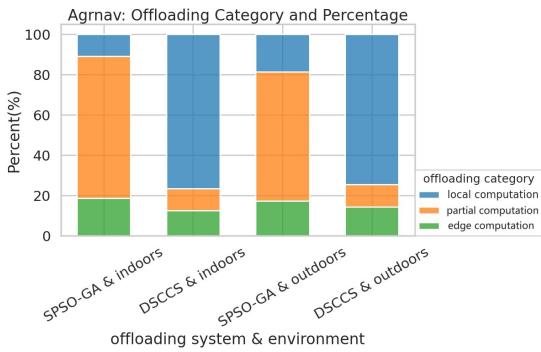
inference will foster the deployment of diverse robotic tasks on real-world robots in the field.

References

- [1] [n. d.]. iPerf - Download iPerf3 and original iPerf pre-compiled binaries. <https://iperf.fr/iperf-download.php>
- [2] Toni Adame, Marc Carrascosa-Zamacois, and Boris Bellalta. 2021. Time-sensitive networking in IEEE 802.11 be: On the way to low-latency WiFi 7. *Sensors* 21, 15 (2021), 4954.
- [3] Majid Altamimi, Atef Abdabou, Kshirasagar Naik, and Amiya Nayak. 2015. Energy cost models of smartphones for task offloading to the cloud. *IEEE Transactions on Emerging Topics in Computing* 3, 3 (2015), 384–398.
- [4] Anh-Quan Cao and Raoul de Charette. 2022. Monoscene: Monocular 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3991–4001.
- [5] Xing Chen, Jianshan Zhang, Bing Lin, Zheyi Chen, Katinka Wolter, and Geyong Min. 2021. Energy-efficient offloading for DNN-based smart IoT systems in cloud-edge environments. *IEEE Transactions on Parallel and Distributed Systems* 33, 3 (2021), 683–697.
- [6] Daniel Crankshaw, Gur-Eyal Sela, Xiangxi Mo, Corey Zumar, Ion Stoica, Joseph Gonzalez, and Alexey Tumanov. 2020. InferLine: latency-aware provisioning and scaling for prediction serving pipelines. In *Proceedings of the 11th ACM Symposium on Cloud Computing*. 477–491.
- [7] Ingrid Daubechies, Ronald DeVore, Simon Foucart, Boris Hanin, and Guergana Petrova. 2022. Nonlinear approximation and (deep) ReLU networks. *Constructive Approximation* 55, 1 (2022), 127–172.
- [8] Alexandre Défossez, Yossi Adi, and Gabriel Synnaeve. 2021. Differentiable model compression via pseudo quantization noise. *arXiv preprint arXiv:2104.09987* (2021).
- [9] Ming Ding, Peng Wang, David López-Pérez, Guoqiang Mao, and Zihuai Lin. 2015. Performance impact of LoS and NLoS transmissions in dense cellular networks. *IEEE Transactions on Wireless Communications* 15, 1045



(a) Categories and percentage of various scheduling for Kapao



(b) Categories and percentage of various scheduling for AGRNav

Figure 8. The layer partitioning scheduling under different baselines and environments. “Local computation” refers to placing the whole layers on the robot when the bandwidth is too low, “edge computation” means placing the whole layers on GPU server when the bandwidth is sufficient, and “partial computation” means placing part of the layers on the robot and part on GPU server.

3 (2015), 2365–2380.

- [10] Khalid Elgazzar, Patrick Martin, and Hossam S Hassanein. 2014. Cloud-assisted computation offloading to support mobile services. *IEEE Transactions on Cloud Computing* 4, 3 (2014), 279–292.
- [11] Zhou Fang, Tong Yu, Ole J Mengshoel, and Rajesh K Gupta. 2017. Qos-aware scheduling of heterogeneous servers for inference in deep neural networks. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. 2067–2070.
- [12] Stefan Gheorghe and Mihai Ivanovici. 2021. Model-based weight quantization for convolutional neural network compression. In *2021 16th International Conference on Engineering of Modern Electric Systems (EMES)*. IEEE, 1–4.
- [13] Cheng Gong, Yao Chen, Ye Lu, Tao Li, Cong Hao, and Deming Chen. 2020. VecQ: Minimal loss DNN model compression with vectorized weight quantization. *IEEE Trans. Comput.* 70, 5 (2020), 696–710.
- [14] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision* 129, 6 (2021), 1789–1819.
- [15] Vinayak Honkote, Dileep Kurian, Sriram Muthukumar, Dibyendu Ghosh, Satish Yada, Kartik Jain, Bradley Jackson, Ilya Klotchkov, Mallikarjuna Rao Nimmagadda, Shreela Dattawadkar, et al. 2019. 2.4

a distributed autonomous and collaborative multi-robot system featuring a low-power robot soc in 22nm cmos for integrated battery-powered minibots. In *2019 IEEE International Solid-State Circuits Conference-(ISSCC)*. IEEE, 48–50.

- [16] Chuang Hu, Wei Bao, Dan Wang, and Fengming Liu. 2019. Dynamic adaptive DNN surgery for inference acceleration on the edge. In *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*. IEEE, 1423–1431.
- [17] Yang Hu, Connor Imes, Xuanang Zhao, Souvik Kundu, Peter A Beerel, Stephen P Crago, and John Paul Walters. 2022. Pipeedge: Pipeline parallelism for large-scale model inference on heterogeneous edge devices. In *2022 25th Euromicro Conference on Digital System Design (DSD)*. IEEE, 298–307.
- [18] KJ Joseph, Salman Khan, Fahad Shahbaz Khan, and Vineeth N Balasubramanian. 2021. Towards Open World Object Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 5830–5840.
- [19] Yiping Kang, Johann Hauswald, Cao Gao, Austin Rovinski, Trevor Mudge, Jason Mars, and Lingjia Tang. 2017. Neurosurgeon: Collaborative intelligence between the cloud and mobile edge. *ACM SIGARCH Computer Architecture News* 45, 1 (2017), 615–629.
- [20] Nam Sung Kim, Todd Austin, David Baauw, Trevor Mudge, Krisztian Flautner, Jie S Hu, Mary Jane Irwin, Mahmut Kandemir, and Vijaykrishnan Narayanan. 2003. Leakage current: Moore’s law meets static power. *computer* 36, 12 (2003), 68–75.
- [21] Ang Li, Shuaiwen Leon Song, Jieyang Chen, Jiajia Li, Xu Liu, Nathan R Tallent, and Kevin J Barker. 2019. Evaluating modern gpu interconnect: Pcie, nvlink, nv-sli, nvswitch and gpudirect. *IEEE Transactions on Parallel and Distributed Systems* 31, 1 (2019), 94–110.
- [22] Qingbiao Li, Fernando Gama, Alejandro Ribeiro, and Amanda Prorok. 2020. Graph neural networks for decentralized multi-robot path planning. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 11785–11792.
- [23] Yiming Li, Zhiding Yu, Christopher Choy, Chaowei Xiao, Jose M Alvarez, Sanja Fidler, Chen Feng, and Anima Anandkumar. 2023. Voxformer: Sparse voxel transformer for camera-based 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9087–9098.
- [24] Huanghuang Liang, Qianlong Sang, Chuang Hu, Dazhao Cheng, Xiaobo Zhou, Dan Wang, Wei Bao, and Yu Wang. 2023. DNN surgery: Accelerating DNN inference on the edge through layer partitioning. *IEEE transactions on Cloud Computing* (2023).
- [25] Bing Lin, Yiniao Huang, Jianshan Zhang, Junqin Hu, Xing Chen, and Jun Li. 2019. Cost-driven off-loading for DNN-based applications over cloud, edge, and end devices. *IEEE Transactions on Industrial Informatics* 16, 8 (2019), 5456–5466.
- [26] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. 2020. Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems* 33 (2020), 2351–2363.
- [27] Ruofeng Liu and Nakjung Choi. 2023. A First Look at Wi-Fi 6 in Action: Throughput, Latency, Energy Efficiency, and Security. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* 7, 1 (2023), 1–25.
- [28] Shuai Liu, Xin Li, Huchuan Lu, and You He. 2022. Multi-Object Tracking Meets Moving UAV. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 8876–8885.
- [29] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. 2016. Large-margin softmax loss for convolutional neural networks. *arXiv preprint arXiv:1612.02295* (2016).
- [30] Antoni Masiukiewicz. 2019. Throughput comparison between the new HEW 802.11 ax standard and 802.11 n/ac standards in selected distance windows. *International Journal of Electronics and Telecommunications* 65, 1 (2019), 79–84.

- 1211 [31] William McNally, Kanav Vats, Alexander Wong, and John McPhee. 2022. Rethinking keypoint representations: Modeling keypoints and 1266
1212 poses as objects for multi-person human pose estimation. In *European 1267
1213 Conference on Computer Vision*. Springer, 37–54. 1268
1214
- 1215 [32] Thaha Mohammed, Carlee Joe-Wong, Rohit Babbar, and Mario 1269
1216 Di Francesco. 2020. Distributed inference acceleration with adaptive 1270
1217 DNN partitioning and offloading. In *IEEE INFOCOM 2020-IEEE 1271
1218 Conference on Computer Communications*. IEEE, 854–863. 1272
1219
- 1220 [33] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, 1273
1221 Mostofa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al. 2021. Efficient 1274
1222 large-scale language model training on gpu clusters using megatron-lm. In *Proceedings of the International Conference for High Performance 1275
1223 Computing, Networking, Storage and Analysis*. 1–15. 1276
1224
- 1225 [34] Mohammad Noormohammadpour and Cauligi S Raghavendra. 2017. 1277
1226 Datacenter traffic control: Understanding techniques and tradeoffs. 1278
1227 *IEEE Communications Surveys & Tutorials* 20, 2 (2017), 1492–1525. 1279
1228
- 1229 [35] NVIDIA. 2024. The World’s Smallest AI Supercomputer. 1280
1230 [https://www.nvidia.com/en-us/autonomous-machines/embedded- 1281
1231 systems/jetson-xavier-series/](https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-xavier-series/). 1282
- 1232 [36] Takeshi Ohkawa, Kazushi Yamashina, Hitomi Kimura, Kanemitsu 1283
1233 Ootsu, and Takashi Yokota. 2018. FPGA components for integrating 1284
1234 FPGAs into robot systems. *IEICE TRANSACTIONS on Information and 1285
1235 Systems* 101, 2 (2018), 363–375. 1286
1236
- 1237 [37] Yuanteng Pei, Matt W Mutka, and Ning Xi. 2013. Connectivity and 1287
1238 bandwidth-aware real-time exploration in mobile robot networks. 1288
1239 *Wireless Communications and Mobile Computing* 13, 9 (2013), 847– 1289
1240 863. 1290
- 1241 [38] pytorch. 2024. pytorch. <https://pytorch.org/>. 1291
- 1242 [39] A Kai Qin, Vicky Ling Huang, and Ponnuthurai N Suganthan. 2008. 1292
1243 Differential evolution algorithm with strategy adaptation for global 1293
1244 numerical optimization. *IEEE transactions on Evolutionary Computation* 13, 2 (2008), 398–417. 1294
- 1245 [40] Yi Ren, Chih-Wei Tung, Jyh-Cheng Chen, and Frank Y Li. 2018. 1295
1246 Proportional and preemption-enabled traffic offloading for IP flow mobility: 1296
1247 Algorithms and performance evaluation. *IEEE Transactions on Vehicular 1297
1248 Technology* 67, 12 (2018), 12095–12108. 1298
- 1249 [41] Nurul I Sarkar and Osman Mussa. 2013. The effect of people movement 1299
1250 on Wi-Fi link throughput in indoor propagation environments. In *IEEE 1300
1251 2013 Tencon-Spring*. IEEE, 562–566. 1301
- 1252 [42] Karen Simonyan and Andrew Zisserman. 2015. Very Deep 1302
1253 Convolutional Networks for Large-Scale Image Recognition. 1303
1254 arXiv:1409.1556 [cs.CV]. 1304
- 1255 [43] Debjyoti Sinha and Mohamed El-Sharkawy. 2019. Thin mobilenet: An 1305
1256 enhanced mobilenet architecture. In *2019 IEEE 10th annual ubiquitous 1306
1257 computing, electronics & mobile communication conference (UEMCON)*. 1307
1258 IEEE, 0280–0285. 1308
- 1259 [44] Sasha Targ, Diogo Almeida, and Kevin Lyman. 2016. Resnet in resnet: 1309
1260 Generalizing residual architectures. *arXiv preprint arXiv:1603.08029* 1310
1261 (2016). 1311
- 1262 [45] Hao Wang, Sreeram Potluri, Miao Luo, Ashish Kumar Singh, Sayantan 1312
1263 Sur, and Dhabaleswar K Panda. 2011. MVAPICH2-GPU: optimized 1313
1264 GPU to GPU communication for InfiniBand clusters. *Computer Science- 1314
1265 Research and Development* 26, 3 (2011), 257–266. 1315
- 1266 [46] Junming Wang, Zekai Sun, Xiuxian Guan, Tianxiang Shen, Zongyuan 1316
1267 Zhang, Tianyang Duan, Dong Huang, Shixiong Zhao, and Heming Cui. 1317
1268 2024. AGRNav: Efficient and Energy-Saving Autonomous Navigation 1318
1269 for Air-Ground Robots in Occlusion-Prone Environments. In *IEEE 1319
1270 International Conference on Robotics and Automation (ICRA)*. 1320
- 1271 [47] Lin Wang and Kuk-Jin Yoon. 2021. Knowledge distillation and student- 1320
1272 teacher learning for visual intelligence: A review and new outlooks. 1321
1273 *IEEE transactions on pattern analysis and machine intelligence* 44, 6 1322
1274 (2021), 3048–3068.