

# 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.

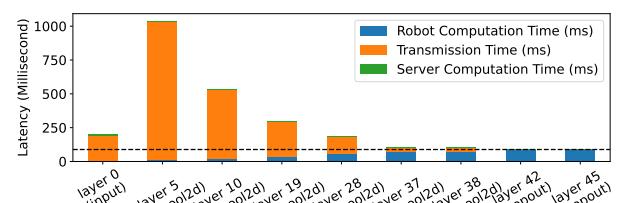
TP requires frequent synchronization among devices, leading to unacceptable communication overhead in robotic IoT.

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.).

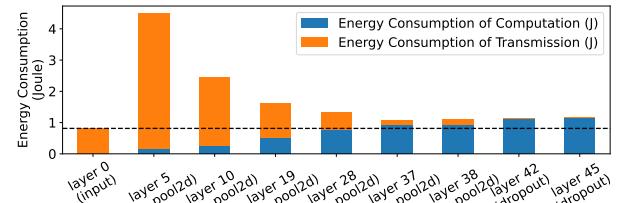
125 Consequently, existing distributed inference approaches [5,  
 126 24] in robotic IoT primarily adopt the PP paradigm and fo-  
 127 cus on layer partitioning of PP, aiming to achieve fast and  
 128 energy-efficient inference. This is because the PP paradigm  
 129 in data centers consists of layer 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 Parallel (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. 3.4.

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. 4.1), 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.

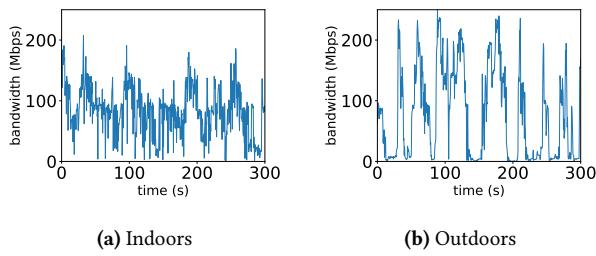
## 2 background

### 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( $\pm 0.004$ )	indoors outdoors	0.698( $\pm 0.135$ ) 0.901( $\pm 0.778$ )	1.400( $\pm 0.232$ ) 1.775( $\pm 1.370$ )	49.85 51.23
ResNet101(44M)	0.065( $\pm 0.005$ )	indoors outdoors	7.156( $\pm 3.348$ ) 8.470( $\pm 6.337$ )	8.106( $\pm 3.403$ ) 9.356( $\pm 6.328$ )	87.95 90.46
VGG19_BN(143M)	0.063( $\pm 0.002$ )	indoors outdoors	5.152( $\pm 4.873$ ) 5.407( $\pm 6.673$ )	5.444( $\pm 4.831$ ) 5.759( $\pm 6.635$ )	70.18 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( $\pm 0.21$ )	5.24( $\pm 0.19$ )	0.3( $\pm 0.09$ )	7.33( $\pm 1.21$ )
	outdoors	6.05( $\pm 0.21$ )	5.11( $\pm 0.28$ )	0.3( $\pm 0.09$ )	9.08( $\pm 7.0$ )
ResNet101(44M)	indoors	11.27( $\pm 0.51$ )	4.97( $\pm 0.16$ )	0.93( $\pm 0.19$ )	40.28( $\pm 16.91$ )
	outdoors	11.27( $\pm 0.51$ )	4.9( $\pm 0.23$ )	0.93( $\pm 0.19$ )	45.8( $\pm 30.98$ )
VGG19_BN(143M)	indoors	14.86( $\pm 0.43$ )	4.88( $\pm 0.29$ )	1.19( $\pm 0.18$ )	26.55( $\pm 23.56$ )
	outdoors	14.86( $\pm 0.43$ )	4.87( $\pm 0.27$ )	1.19( $\pm 0.18$ )	28.06( $\pm 32.33$ )

**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,

551 where improved distributed inference can provide faster and  
 552 more energy-efficient inference for these scenarios.

### 553 3 Overview

#### 556 3.1 Existing distributed inference on robotic IoT

557 However, the increasing complexity of DNN structures, now  
 558 evolved into directed acyclic graphs (DAGs), poses new chal-  
 559 lenges, potentially leading to NP-hardness in performance  
 560 optimization [24]. This issue is addressed by graph theory  
 561 techniques [24, 54] and varying hardware and network con-  
 562 ditions further complicate the problem.

563 In summary, these two categories primarily adopt the  
 564 PP paradigm but suffer from the transmission bottleneck  
 565 inherent to PP’s scheduling mechanism (see Sec. 3.2). Con-  
 566 sequently, achieving fast and energy-efficient inference on  
 567 robotic IoT remains an open issue.

#### 569 3.2 Dilemma on Inference Time and Energy 570 Consumption

572 Regardless of the complexity of DNN models, layer partition-  
 573 ing methods consist of three phases: computing earlier DNN  
 574 layers on robots, transmitting intermediate results, and com-  
 575 pleting inference on the GPU device. Since the GPU device’s  
 576 computation time is negligible compared to the other two  
 577 phases (see Fig. 1) due to the high computing capabilities of  
 578 GPU devices, this paper focuses on the computation phase  
 579 of robots and the data transmission phase via robotic IoT.

580 The data transmission phase can only begin after obtain-  
 581 ing the calculation result of the intermediate layer when the  
 582 computation phase on robots is completed, preventing over-  
 583 lap for a single inference task. PP can only overlap computa-  
 584 tion and data transmission phases from different inference  
 585 tasks, not from the same task [6]. However, the transmission  
 586 cost inherent to the PP’s scheduling mechanism becomes a  
 587 bottleneck in robotic IoT due to limited bandwidth. In our  
 588 experiments, even with optimal layer partitioning [5, 24],  
 589 such communication cost takes up to 63% of inference time.

590 To make matters worse, such transmission overhead not  
 591 only leads to prolonged inference time but also to high en-  
 592 ergy consumption during the data transmission phase, re-  
 593 ferred to as transmission energy consumption. Our find-  
 594 ings reveal that such transmission energy consumption ac-  
 595 counts for nearly one-third of the total energy consumed  
 596 during inference (see Sec.??). This is because the device can-  
 597 not be put into low-power sleep mode while waiting for  
 598 the final inference result from the GPU device, as it has to  
 599 promptly continue working when it receives the inference re-  
 600 sults. Moreover, chips like CPU, GPU, and memory consume  
 601 non-negligible power even when not computing, due to the  
 602 static power consumption rooted in transistors’ leakage cur-  
 603 rent [20]. Consequently, both the energy consumed during  
 604 the execution of DNN layers on robots, referred to as robot  
 605

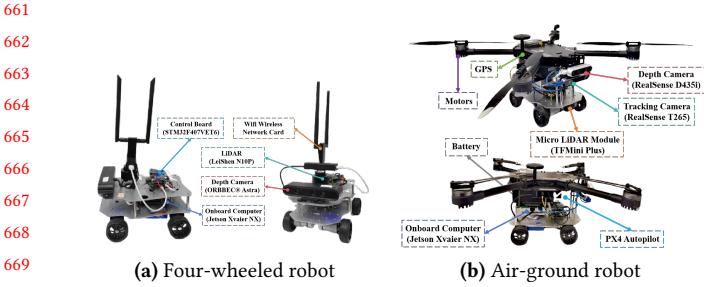
606 computation energy consumption, and the transmission en-  
 607 ergy consumption resulting from prolonged transmission  
 608 times substantially impact the overall power consumption  
 609 of the inference process in robotic IoT.

610 Only models with limited transmission overhead can mit-  
 611 igate the impact of these shortcomings on inference per-  
 612 formance. However, the unstable bandwidth in robotic IoT  
 613 wireless networks can cause the transmission time for layer  
 614 partitioning to vary dramatically, sometimes changing by  
 615 hundreds of times (see Fig. 2). In our experiments, even a rel-  
 616 atively small model with only 0.84M parameters still suffers  
 617 from its significant transmission overhead. The significant  
 618 impact of transmission overhead on both inference time  
 619 and energy consumption highlights the need for innovative  
 620 approaches that can effectively mitigate the transmission  
 621 bottleneck in robotic IoT.

#### 623 3.3 Special Cases

625 Since layer partitioning methods schedule at the granularity  
 626 of model layers, “local computation” and “edge computation”  
 627 are special cases of layer partitioning. “Local computation”  
 628 refers to placing the whole layers on the robot when the  
 629 bandwidth is too low, while “edge computation” means plac-  
 630 ing the whole layers on GPU devices when the bandwidth is  
 631 sufficient. Local computation avoids the impact of network  
 632 transmission on inference time but consumes the maximum  
 633 computation energy consumption. On the other hand, edge  
 634 computation minimizes computation energy consumption  
 635 but requires a high enough bandwidth to ensure the lowest  
 636 possible transmission energy consumption and overall in-  
 637 ference time. These two special cases are indispensable for  
 638 existing methods to cope with different network conditions,  
 639 when they are too low or sufficient, and to address the need  
 640 for various trade-offs between inference delay and energy  
 641 consumption.

642 In our experiments, we found that the bandwidth condi-  
 643 tions under which the layer partitioning scheme of different  
 644 models becomes these special cases vary, and the higher  
 645 the bandwidth, the more layers are scheduled to be placed  
 646 on GPU devices. We explain the reasons causing different  
 647 bandwidth conditions for different models in Sec. ?? with  
 648 some detailed real-world cases. The existence of these special  
 649 cases highlights the importance of considering the relation-  
 650 ship between bandwidth, model structure, and the resulting  
 651 trade-offs between inference delay and energy consump-  
 652 tion.

**Figure 3.** The detailed composition of the robot platforms

	inference	transmission	standby
Power (W)	13.35	4.25	4.04

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

### 3.4 Local Operator Parallel

## 4 Detailed Design

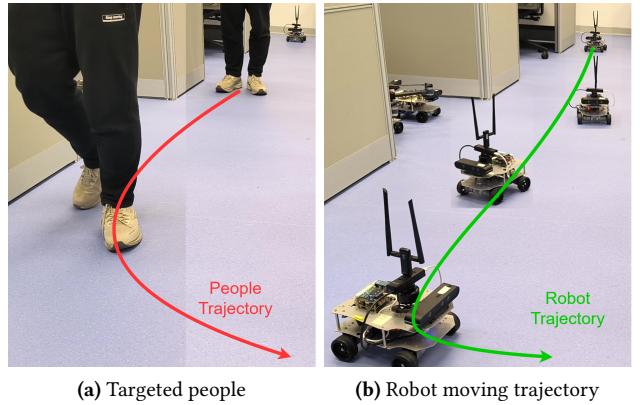
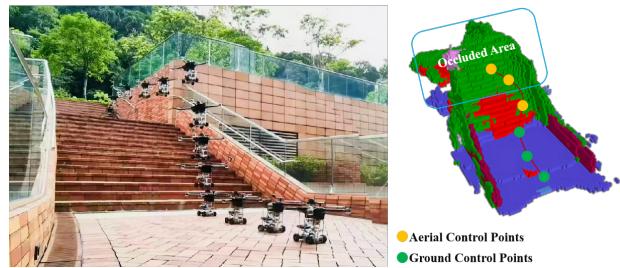
### 4.1 Local Operator Scheduling Strategy

## 5 Implementation

## 6 Evaluation

**Testbed.** The evaluation was conducted on a custom four-wheeled robot (Fig 3a), and a custom air-ground robot(Fig 3b). 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 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, exhibit distinct GPU utilization rates and power consumption during inference.

**Figure 4.** A real-time people-tracking robotic application on our robot based on a well-known human pose estimation ML model, Kapao [31].**Figure 5.** By predicting occlusions in advance, AGRNav [46] gains an accurate perception of the environment and avoids collisions, resulting in efficient and energy-saving paths.

**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 4), and AGRNav, an autonomous navigation application on our air-ground robot (Fig 5). 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

661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
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

771 our observations and findings: MobileNet [43], ResNet [44],  
 772 VGGNet [42], ConvNeXt [48], RegNet [53].

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

## 783 6.1 Inference Time

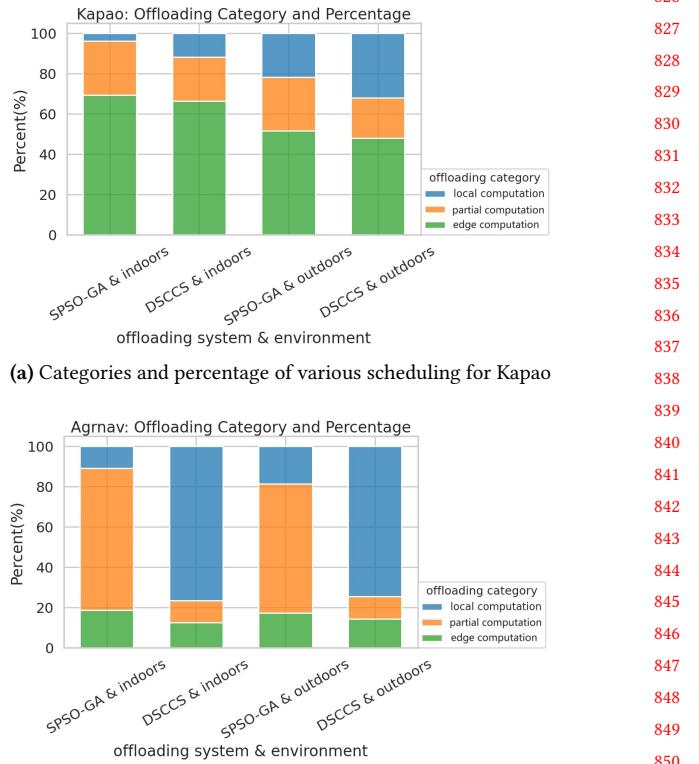
785 **Kapao.** From the results in the upper part of Tab. 4, both  
 786 SPSO-GA and DSCCS reduced Kapao's inference time by  
 787 39.69% and 56.92% indoors and 28.67% and 47.46% outdoors,  
 788 with DSCCS achieving 28.57% (indoors) and 26.34% (out-  
 789 doors) lower inference time than SPSO-GA. While both  
 790 systems significantly reduced inference time via offloading,  
 791 transmission time accounts for 49.69% to 69.46% of the whole  
 792 inference time, indicating that even with SOTA layer parti-  
 793 tioning, the transmission bottleneck inherent to PP's schedul-  
 794 ing mechanism cannot be mitigated. The difference between  
 795 DSCCS and SPSO-GA can be attributed to their optimization  
 796 goals: DSCCS minimizes inference latency, while SPSO-GA  
 797 minimizes power consumption under deadline constraints.

798 **AGRNav.** The performance gain of the two offloading  
 799 systems varied for AGRNav, as shown in the lower part of  
 800 Tab. 4. DSCCS still reduced inference time by 18.34% and  
 801 12.43% in indoors and outdoors. However, SPSO-GA achieved  
 802 similar inference time (3.65% and 3.06% reduction) as local  
 803 computation both indoors and outdoors. We will explain and  
 804 analyze this phenomenon in Sec.6.2.

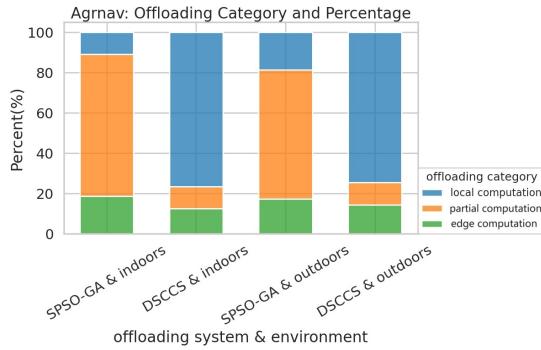
805 Notice that the large standard deviation in transmission  
 806 time in outdoors in both offloading systems indicates that  
 807 bandwidth fluctuated more frequently and more fiercely out-  
 808 doors compared with indoors, which complies with Fig. 2.  
 809 Additionally, the lower average bandwidth for outdoors sce-  
 810 narios (see Sec.2.1) results in increased transmission and  
 811 inference times relative to indoor scenarios.

## 813 6.2 Breakdown

815 Both SPSO-GA and DSCCS automatically adapt to available  
 816 bandwidth, transitioning to edge computation (placing all  
 817 DNN layers on the GPU server) when bandwidth is suffi-  
 818 cient, and to local computation (placing all DNN layers on  
 819 robots) when bandwidth is low. To better understand how  
 820 their layer partitioning scheduling varies with different net-  
 821 work conditions and models, we recorded and analyzed the  
 822 Categories and percentages of various layer partitioning  
 823 schedules under different baselines and environments, as  
 824 detailed in Fig. 6.



(a) Categories and percentage of various scheduling for Kapao



(b) Categories and percentage of various scheduling for AGRNav

853 **Figure 6.** The layer partitioning scheduling under difference  
 854 baselines and environments. “Local computation” refers to  
 855 placing the whole layers on the robot when the bandwidth is  
 856 too low, “edge computation” means placing the whole layers  
 857 on GPU server when the bandwidth is sufficient, and “partial  
 858 computation” means placing part of the layers on the robot  
 859 and part on GPU server.

862 Local computation and edge computation are special cases  
 863 of layer partitioning, with the bandwidth conditions required  
 864 for each model to reach these cases varying based on the  
 865 model structure and partitioning method used. Analyzing  
 866 Fig. 6a and Fig. 6b, both SPSO-GA and DSCCS tend to allocate  
 867 more layers on the robot for AGRNav. When comparing  
 868 indoor and outdoor scenarios in Fig. 6, it is evident that  
 869 higher bandwidth leads to more layers being scheduled on  
 870 GPU server. Additionally, when comparing SPSO-GA and  
 871 DSCCS in Fig. 6, DSCCS, which focuses on optimizing energy  
 872 consumption, tends to place fewer layers on the robot to  
 873 reduce computation energy consumption.

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

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

**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( $\pm 0.64$ )	7.03( $\pm 3.57$ )	5.92( $\pm 2.18$ )	5.03( $\pm 1.87$ )	15.03( $\pm 0.63$ )	2.41( $\pm 1.35$ )	1.84( $\pm 1.0$ )
	outdoors	15.04( $\pm 0.64$ )	14.15( $\pm 1.71$ )	5.89( $\pm 2.3$ )	5.02( $\pm 1.97$ )	1.56( $\pm 0.85$ )	9.85( $\pm 1.77$ )	2.56( $\pm 6.17$ )
agrnav(0.84M)	indoors	10.26( $\pm 1.58$ )	6.6( $\pm 1.95$ )	6.74( $\pm 2.03$ )	5.74( $\pm 1.78$ )	10.82( $\pm 1.44$ )	2.86( $\pm 0.88$ )	2.88( $\pm 0.81$ )
	outdoors	10.26( $\pm 1.58$ )	7.36( $\pm 2.34$ )	7.91( $\pm 2.45$ )	6.71( $\pm 2.1$ )	10.82( $\pm 1.44$ )	3.77( $\pm 5.74$ )	4.25( $\pm 5.7$ )

**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”

### 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. 6a. 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 against 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 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 Abdrabou, 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).

	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
991	MobileNet_V3_Small (2M)	0.033( $\pm 0.019$ )	indoors	0.035( $\pm 0.019$ )	0.016( $\pm 0.005$ )	0.044( $\pm 0.020$ )	0.031( $\pm 0.008$ )	79.79	53.24
992			outdoors	0.035( $\pm 0.044$ )	0.017( $\pm 0.005$ )	0.047( $\pm 0.037$ )	0.033( $\pm 0.018$ )	50.04	51.49
993	RegNet_X_3_2GF (15M)	0.060( $\pm 0.022$ )	indoors	0.049( $\pm 0.026$ )	0.033( $\pm 0.011$ )	0.065( $\pm 0.028$ )	0.049( $\pm 0.016$ )	76.25	64.17
994			outdoors	0.049( $\pm 0.055$ )	0.032( $\pm 0.032$ )	0.069( $\pm 0.050$ )	0.051( $\pm 0.030$ )	53.23	44.50
995	ResNet101 (44M)	0.060( $\pm 0.023$ )	indoors	0.054( $\pm 0.451$ )	0.033( $\pm 0.010$ )	0.072( $\pm 0.453$ )	0.050( $\pm 0.016$ )	75.64	57.37
996			outdoors	0.052( $\pm 0.064$ )	0.033( $\pm 0.036$ )	0.077( $\pm 0.059$ )	0.054( $\pm 0.034$ )	51.54	42.48
997	ConvNeXt_Base (88M)	0.047( $\pm 0.004$ )	indoors	0.033( $\pm 0.018$ )	0.020( $\pm 0.006$ )	0.044( $\pm 0.019$ )	0.032( $\pm 0.009$ )	75.39	49.37
998			outdoors	0.032( $\pm 0.038$ )	0.020( $\pm 0.022$ )	0.045( $\pm 0.033$ )	0.034( $\pm 0.019$ )	52.82	35.63
999	ConvNeXt_Large (197M)	0.051( $\pm 0.005$ )	indoors	0.033( $\pm 0.017$ )	0.023( $\pm 0.008$ )	0.046( $\pm 0.019$ )	0.035( $\pm 0.013$ )	72.96	62.68
1000			outdoors	0.032( $\pm 0.038$ )	0.023( $\pm 0.024$ )	0.054( $\pm 0.040$ )	0.041( $\pm 0.028$ )	48.94	43.96
1001	RegNet_Y_128GF (644M)	0.139( $\pm 0.016$ )	indoors	0.076( $\pm 0.289$ )	0.041( $\pm 0.024$ )	0.305( $\pm 0.382$ )	0.100( $\pm 0.035$ )	23.58	40.76
1002			outdoors	0.171( $\pm 0.602$ )	0.016( $\pm 0.055$ )	0.432( $\pm 0.615$ )	0.117( $\pm 0.242$ )	32.39	9.41
1003									1059
1004									

**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			1064
			Local	SPSO-GA	DSCCS	Local	SPSO-GA	DSCCS	
1009	MobileNet_V3_Small (2M)	indoors	6.131( $\pm 0.061$ )	5.448( $\pm 0.168$ )	5.658( $\pm 0.085$ )	0.202( $\pm 0.002$ )	0.241( $\pm 0.107$ )	0.174( $\pm 0.046$ )	1066
1010		outdoors	6.131( $\pm 0.061$ )	5.567( $\pm 0.273$ )	5.557( $\pm 0.186$ )	0.202( $\pm 0.002$ )	0.260( $\pm 0.204$ )	0.185( $\pm 0.099$ )	1065
1011	RegNet_X_3_2GF (15M)	indoors	8.208( $\pm 0.140$ )	5.490( $\pm 0.178$ )	5.714( $\pm 0.342$ )	0.492( $\pm 0.008$ )	0.356( $\pm 0.156$ )	0.278( $\pm 0.088$ )	1068
1012		outdoors	8.208( $\pm 0.140$ )	5.878( $\pm 0.659$ )	6.041( $\pm 0.624$ )	0.492( $\pm 0.008$ )	0.406( $\pm 0.295$ )	0.311( $\pm 0.184$ )	1069
1013	ResNet101 (44M)	indoors	11.851( $\pm 0.404$ )	5.457( $\pm 0.240$ )	5.953( $\pm 0.789$ )	0.711( $\pm 0.024$ )	0.390( $\pm 2.471$ )	0.298( $\pm 0.094$ )	1070
1014		outdoors	11.851( $\pm 0.404$ )	6.179( $\pm 1.083$ )	6.431( $\pm 1.060$ )	0.711( $\pm 0.024$ )	0.478( $\pm 0.364$ )	0.349( $\pm 0.216$ )	1071
1015	ConvNeXt_Base (88M)	indoors	15.335( $\pm 0.273$ )	5.507( $\pm 0.358$ )	7.713( $\pm 2.613$ )	0.721( $\pm 0.013$ )	0.241( $\pm 0.103$ )	0.250( $\pm 0.069$ )	1072
1016		outdoors	15.335( $\pm 0.273$ )	7.638( $\pm 3.297$ )	9.148( $\pm 3.338$ )	0.721( $\pm 0.013$ )	0.346( $\pm 0.254$ )	0.307( $\pm 0.171$ )	1071
1017	ConvNeXt_Large (197M)	indoors	15.403( $\pm 0.082$ )	5.518( $\pm 0.638$ )	6.604( $\pm 2.860$ )	0.786( $\pm 0.004$ )	0.251( $\pm 0.104$ )	0.230( $\pm 0.088$ )	1073
1018		outdoors	15.403( $\pm 0.082$ )	8.400( $\pm 4.345$ )	8.895( $\pm 4.505$ )	0.786( $\pm 0.004$ )	0.452( $\pm 0.339$ )	0.366( $\pm 0.248$ )	1074
1019	RegNet_Y_128GF (644M)	indoors	15.430( $\pm 0.020$ )	5.384( $\pm 1.071$ )	6.151( $\pm 2.155$ )	2.145( $\pm 0.003$ )	1.642( $\pm 0.327$ )	0.615( $\pm 0.216$ )	1075
1020		outdoors	15.430( $\pm 0.020$ )	6.361( $\pm 2.349$ )	9.127( $\pm 4.724$ )	2.145( $\pm 0.003$ )	2.748( $\pm 1.015$ )	1.068( $\pm 0.553$ )	1076
1021									1078

**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”

- [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, 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] K J 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.

- 1101 [19] Yiping Kang, Johann Hauswald, Cao Gao, Austin Rovinski, Trevor  
1102 Mudge, Jason Mars, and Lingjia Tang. 2017. Neurosurgeon: Collabora-  
1103 tive intelligence between the cloud and mobile edge. *ACM SIGARCH Computer Architecture News* 45, 1 (2017), 615–629.
- 1104 [20] Nam Sung Kim, Todd Austin, David Baauw, Trevor Mudge, Krisztián  
1105 Flautner, Jie S Hu, Mary Jane Irwin, Mahmut Kandemir, and Vijaykrishnan Narayanan. 2003. Leakage current: Moore’s law meets static  
1106 power. *computer* 36, 12 (2003), 68–75.
- 1107 [21] Ang Li, Shuaiwen Leon Song, Jieyang Chen, Jiajia Li, Xu Liu, Nathan R  
1108 Tallent, and Kevin J Barker. 2019. Evaluating modern gpu interconnect:  
1109 Pcie, nvlink, nv-sli, nvswitch and gpudirect. *IEEE Transactions on Parallel and Distributed Systems* 31, 1 (2019), 94–110.
- 1110 [22] Qingbiao Li, Fernando Gama, Alejandro Ribeiro, and Amanda Prorok.  
1111 2020. Graph neural networks for decentralized multi-robot path plan-  
1112 ning. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 11785–11792.
- 1113 [23] Yiming Li, Zhiding Yu, Christopher Choy, Chaowei Xiao, Jose M Al-  
1114 varez, Sanja Fidler, Chen Feng, and Anima Anandkumar. 2023. Vox-  
1115 former: Sparse voxel transformer for camera-based 3d semantic scene  
1116 completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9087–9098.
- 1117 [24] Huanghuang Liang, Qianlong Sang, Chuang Hu, Dazhao Cheng, Xi-  
1118 aobo Zhou, Dan Wang, Wei Bao, and Yu Wang. 2023. DNN surgery:  
1119 Accelerating DNN inference on the edge through layer partitioning.  
1120 *IEEE transactions on Cloud Computing* (2023).
- 1121 [25] Bing Lin, Yinhao Huang, Jianshan Zhang, Junqin Hu, Xing Chen, and  
1122 Jun Li. 2019. Cost-driven off-loading for DNN-based applications  
1123 over cloud, edge, and end devices. *IEEE Transactions on Industrial  
1124 Informatics* 16, 8 (2019), 5456–5466.
- 1125 [26] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. 2020.  
1126 Ensemble distillation for robust model fusion in federated learning.  
1127 *Advances in Neural Information Processing Systems* 33 (2020), 2351–  
1128 2363.
- 1129 [27] Ruofeng Liu and Nakjung Choi. 2023. A First Look at Wi-Fi 6 in Action:  
1130 Throughput, Latency, Energy Efficiency, and Security. *Proceedings  
1131 of the ACM on Measurement and Analysis of Computing Systems* 7, 1  
1132 (2023), 1–25.
- 1133 [28] Shuai Liu, Xin Li, Huchuan Lu, and You He. 2022. Multi-Object Track-  
1134 ing Meets Moving UAV. In *Proceedings of the IEEE/CVF Conference on  
1135 Computer Vision and Pattern Recognition (CVPR)*. 8876–8885.
- 1136 [29] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. 2016. Large-  
1137 margin softmax loss for convolutional neural networks. *arXiv preprint arXiv:1612.02295* (2016).
- 1138 [30] Antoni Masiukiewicz. 2019. Throughput comparison between the new  
1139 HEW 802.11 ax standard and 802.11 n/ac standards in selected distance  
1140 windows. *International Journal of Electronics and Telecommunications*  
1141 65, 1 (2019), 79–84.
- 1142 [31] William McNally, Kanav Vats, Alexander Wong, and John McPhee.  
1143 2022. Rethinking keypoint representations: Modeling keypoints and  
1144 poses as objects for multi-person human pose estimation. In *European  
1145 Conference on Computer Vision*. Springer, 37–54.
- 1146 [32] Thaha Mohammed, Carlee Joe-Wong, Rohit Babbar, and Mario  
1147 Di Francesco. 2020. Distributed inference acceleration with adap-  
1148 tive DNN partitioning and offloading. In *IEEE INFOCOM 2020-IEEE  
1149 Conference on Computer Communications*. IEEE, 854–863.
- 1150 [33] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGres-  
1151 ley, Mostafa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi  
1152 Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al. 2021. Efficient  
1153 large-scale language model training on gpu clusters using megatron-  
1154 lm. In *Proceedings of the International Conference for High Performance  
Computing, Networking, Storage and Analysis*. 1–15.
- 1155 [34] Mohammad Noormohammadpour and Cauligi S Raghavendra. 2017.  
1156 Datacenter traffic control: Understanding techniques and tradeoffs.  
1157 *IEEE Communications Surveys & Tutorials* 20, 2 (2017), 1492–1525.
- [35] NVIDIA. 2024. The World’s Smallest AI Supercomputer.  
<https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-xavier-series/>.
- [36] Takeshi Ohkawa, Kazushi Yamashina, Hitomi Kimura, Kanemitsu  
Ootsu, and Takashi Yokota. 2018. FPGA components for integrating  
FPGAs into robot systems. *IEICE TRANSACTIONS on Information and  
Systems* 101, 2 (2018), 363–375.
- [37] Yuanteng Pei, Matt W Mutka, and Ning Xi. 2013. Connectivity and  
bandwidth-aware real-time exploration in mobile robot networks.  
*Wireless Communications and Mobile Computing* 13, 9 (2013), 847–  
863.
- [38] pytorch. 2024. pytorch. <https://pytorch.org/>.
- [39] A Kai Qin, Vicky Ling Huang, and Ponnuthurai N Suganthan. 2008.  
Differential evolution algorithm with strategy adaptation for global  
numerical optimization. *IEEE transactions on Evolutionary Computation* 13, 2 (2008), 398–417.
- [40] Yi Ren, Chih-Wei Tung, Jyh-Cheng Chen, and Frank Y Li. 2018. Proportional  
and preemption-enabled traffic offloading for IP flow mobility:  
Algorithms and performance evaluation. *IEEE Transactions on Vehicular  
Technology* 67, 12 (2018), 12095–12108.
- [41] Nurul I Sarkar and Osman Mussa. 2013. The effect of people movement  
on Wi-Fi link throughput in indoor propagation environments. In *IEEE  
2013 Tencon-Spring*. IEEE, 562–566.
- [42] Karen Simonyan and Andrew Zisserman. 2015. Very Deep  
Convolutional Networks for Large-Scale Image Recognition.  
*arXiv:1409.1556 [cs.CV]*
- [43] Debjyoti Sinha and Mohamed El-Sharkawy. 2019. Thin mobilenet: An  
enhanced mobilenet architecture. In *2019 IEEE 10th annual ubiquitous  
computing, electronics & mobile communication conference (UEMCON)*.  
IEEE, 0280–0285.
- [44] Sasha Targ, Diogo Almeida, and Kevin Lyman. 2016. Resnet in resnet:  
Generalizing residual architectures. *arXiv preprint arXiv:1603.08029*  
(2016).
- [45] Hao Wang, Sreeram Potluri, Miao Luo, Ashish Kumar Singh, Sayantan  
Sur, and Dhabaleswar K Panda. 2011. MVAPICH2-GPU: optimized  
GPU to GPU communication for InfiniBand clusters. *Computer Science-  
Research and Development* 26, 3 (2011), 257–266.
- [46] Junming Wang, Zekai Sun, Xiuxian Guan, Tianxiang Shen, Zongyuan  
Zhang, Tianyang Duan, Dong Huang, Shixiong Zhao, and Heming Cui.  
2024. AGRNav: Efficient and Energy-Saving Autonomous Navigation  
for Air-Ground Robots in Occlusion-Prone Environments. In *IEEE  
International Conference on Robotics and Automation (ICRA)*.
- [47] Lin Wang and Kuk-Jin Yoon. 2021. Knowledge distillation and student-  
teacher learning for visual intelligence: A review and new outlooks.  
*IEEE transactions on pattern analysis and machine intelligence* 44, 6  
(2021), 3048–3068.
- [48] Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen,  
Zhuang Liu, In So Kweon, and Saining Xie. 2023. Convnext v2: Co-  
designing and scaling convnets with masked autoencoders. In *Pro-  
ceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
Recognition*. 16133–16142.
- [49] Huaming Wu, William J Knottenbelt, and Katinka Wolter. 2019. An  
efficient application partitioning algorithm in mobile environments.  
*IEEE Transactions on Parallel and Distributed Systems* 30, 7 (2019),  
1464–1480.
- [50] Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long  
Han, and Yang Tang. 2023. A brief overview of ChatGPT: The history,  
status quo and potential future development. *IEEE/CAA Journal of  
Automatica Sinica* 10, 5 (2023), 1122–1136.
- [51] Zhaoyang Xia, Youquan Liu, Xin Li, Xinge Zhu, Yuexin Ma, Yikang Li,  
Yuenan Hou, and Yu Qiao. 2023. SCPNet: Semantic Scene Completion  
on Point Cloud. In *Proceedings of the IEEE/CVF Conference on Computer  
Vision and Pattern Recognition*. 17642–17651.

1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210

- 1211 [52] Yecheng Xiang and Hyoseung Kim. 2019. Pipelined data-parallel  
1212 CPU/GPU scheduling for multi-DNN real-time inference. In *2019 IEEE  
1213 Real-Time Systems Symposium (RTSS)*. IEEE, 392–405.
- 1214 [53] Jing Xu, Yu Pan, Xinglin Pan, Steven Hoi, Zhang Yi, and Zenglin Xu.  
1215 2022. RegNet: self-regulated network for image classification. *IEEE  
1216 Transactions on Neural Networks and Learning Systems* (2022).
- 1217 [54] Min Xue, Huaming Wu, Guang Peng, and Katinka Wolter. 2021.  
1218 DDPQN: An efficient DNN offloading strategy in local-edge-cloud  
1219 collaborative environments. *IEEE Transactions on Services Computing*  
1220 15, 2 (2021), 640–655.
- 1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265 [55] Xinlei Yang, Hao Lin, Zhenhua Li, Feng Qian, Xingyao Li, Zhiming He,  
1266 Xudong Wu, Xianlong Wang, Yunhao Liu, Zhi Liao, et al. 2022. Mobile  
1267 access bandwidth in practice: Measurement, analysis, and implications.  
1268 In *Proceedings of the ACM SIGCOMM 2022 Conference*. 114–128.
- [56] Yang Yang, Li Juntao, and Peng Lingling. 2020. Multi-robot path  
1269 planning based on a deep reinforcement learning DQN algorithm.  
1270 *CAAI Transactions on Intelligence Technology* 5, 3 (2020), 177–183.
- [57] Yonghao Zhuang, Hexu Zhao, Lianmin Zheng, Zhuohan Li, Eric Xing,  
1271 Qirong Ho, Joseph Gonzalez, Ion Stoica, and Hao Zhang. 2023. On  
1272 optimizing the communication of model parallelism. *Proceedings of  
1273 Machine Learning and Systems* 5 (2023).
- 1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320