

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

Submission Id: xxx

Abstract

The rapid advancements in machine learning techniques have led to significant achievements in various real-world robotic tasks. These tasks heavily rely on fast and energy-efficient inference of deep neural network (DNN) models when deployed on robots. To enhance inference performance, distributed inference has emerged as a promising approach, parallelizing inference across multiple powerful GPU devices in modern data centers using techniques such as data parallelism, tensor parallelism, and pipeline parallelism. However, when deployed on real-world robots, existing parallel methods fail to provide low inference latency and meet the energy requirements due to the limited bandwidth of robotic IoT.

We present Intra-DP, a high-performance distributed inference system optimized for robotic IoT. Intra-DP employs a fine-grained approach to parallelize inference at the granularity of local operators within DNN layers (i.e., operators that can be computed independently with the partial input, such as the convolution kernel in the convolution layer). By doing so, Intra-DP enables different operators of different layers to be computed and transmitted concurrently, and overlap the computation and transmission phases within the same inference task. The evaluation demonstrate that Intra-DP reduces inference time by 14.9% ~41.1% and energy consumption per inference by up to 35.3% compared to the state-of-the-art baselines.

1 Introduction

The rapid progress in machine learning (ML) techniques has led to remarkable achievements in various fundamental robotic tasks, such as object detection [21, 31, 34], robotic control [25, 51, 61], and environmental perception [4, 26, 56]. 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 [38], FPGA [39], SoC [16]), but also introduce additional energy consumption (e.g., 162% more for [34] 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 [18, 57, 65], where numerous GPUs are utilized to speed up large model inference, such as in the case of ChatGPT [55]. 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 in parallel to speed up inference.

TP requires frequent synchronization among devices, leading to unacceptable communication overhead in robotic IoT. By partitioning parameter tensors of a layer across GPUs, TP allows concurrent computation on different parts of this tensor but requires an all-reduce communication [65] to combine computation results from different devices, which entails significant communication overhead. Consequently, TP is used mainly for large layers that are too large to fit in one device in data centers and require dedicated high-speed interconnects (e.g., 400 Gbps for NVLink [24]) even within data centers. On the contrary, robots must prioritize seamless mobility and primarily depend on wireless connections, which inherently possess limited bandwidth, as described in Sec. 2.1, making all-reduce synchronization an unacceptable overhead (e.g., the inference time with TP was up to 143.9X slower than local computation in Sec. 2.3).

Consequently, existing distributed inference approaches [5, 27] in robotic IoT are constrained to the PP paradigm. Since the PP paradigm in data centers consists of layer partitioning and pipeline execution, where the pipeline execution of PP enhances inference throughput rather than reducing the completion time of a single inference [6], the most critical requirement in robotic IoT, existing methods on robotic IoT concentrate on optimizing the layer partitioning aspect of PP to achieve fast and energy-efficient inference. Based on the fact that the amounts of output data in some intermediate layers of a DNN model are significantly smaller than that of its raw input data [17], DNN layer partitioning strategies constitute various trade-offs between computation and transmission, taking into account application-specific inference speed requirements and energy consumption demands, as shown in Fig. 1.

However, existing methods based on the PP paradigm face significant transmission bottlenecks in robotic IoT due to the inherent scheduling mechanism. The PP paradigm on robotic IoT involves three sequential phases: computing early DNN layers on robots, transmitting intermediate results, and completing inference on a GPU server, where the limited bandwidth of real-world networks often results in transmission time exceeding computation time. Despite optimal layer partitioning strategies [5, 27], the transmission overhead becomes a substantial bottleneck, accounting for up to 70.45% of inference time in our experiments, due to the limited bandwidth of robotic IoT. This overhead not only slows down inference speed and consumes significantly more energy but also cannot be effectively mitigated by overlapping computation and transmission phases across multiple inference tasks via pipeline execution, which still fails to reduce the completion time of a single inference task [6], a crucial aspect for robotic applications.

The key reason for the problem of the above methods is that existing methods conduct layer-granulated scheduling, which

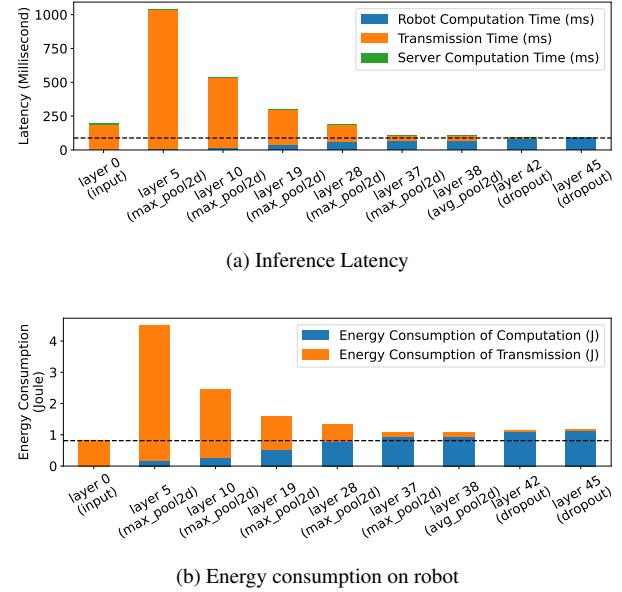


Figure 1: Existing distributed inference approaches on VGG19 [46] in our experiments, which adopt PP paradigm with various layer partitioning scheduling strategies. The X-axis of the graph represents different layer partitioning strategies, where ‘layer i ’ indicates that all layers up to and including the i^{th} layer are computed on the robot, while the subsequent layers are processed on the GPU server.

divides a single inference task into multiple sequential phases, thereby precluding parallel execution within the scope of an individual inference task. As transmission time constitutes a substantial portion of the total inference time (approximately half) in existing methods, a novel parallel method that can efficiently overlap computation and transmission within the same inference task has the potential to address this shortcoming, achieving fast 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 also 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 with partial input. For instance, softmax [32] requires the complete input vector to calculate the correspond-

ing probability distribution, referring to it as a global operator, while ReLU [7] and convolution [35] can be computed with partial 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. Since a single local operator like convolution kernel may require multiple calculations per layer, we treat each calculation of the local operator as an independent local operator in this article for easy discussion. Local operators are widely used in robotic applications, especially convolution layers in computer vision [34] and point cloud tasks [51]. 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 overlap 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.

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 on robot 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 [43].

We implemented Intra-DP in PyTorch [41] and evaluated Intra-DP on our real-world robots under two typical real-world robotic applications [34, 51] and several models common to mobile devices on a larger scale [46, 47, 49, 53, 58]. We

compared Intra-DP with two SOTA pipeline parallelism methods as baselines: DSCCS [27], 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 14.9% ~41.1% compared to baselines under indoors and outdoors environments.
- Intra-DP is energy-efficient. Intra-DP reduced up to 35.3% energy consumption per inference compared to baselines, due to faster inference speed and limited-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 three 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/eurosys25paper445/intraDP>.

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 [30]. However, not only the limited hardware resources on the robot can not fully play the potential of Wi-Fi 6 [60], but also the actual available bandwidth of wireless networks is often reduced in practice

due to factors such as movement of the devices [33, 40], occlusion from by physical barriers [9, 45], and preemption of the wireless channel by other devices [2, 44].

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 5–40cm/s speed in our lab (indoors) and campus garden (outdoors), with hardware and wireless network settings as described in Sec. 6. 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.

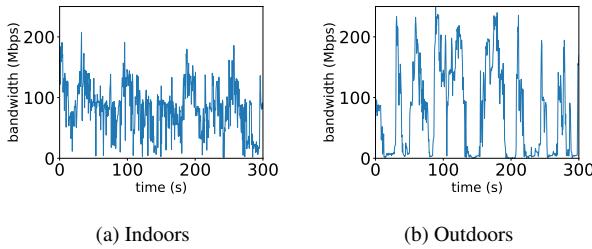


Figure 2: The instability of wireless transmission between our robot and a base station in robotic IoT networks.

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.

2.2 Characteristics of Data Center Networks

Data center networks, which are used for large model inference (e.g., ChatGPT [55]), 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 [50] or PCIe [24], 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 [37]. 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. DP [57] is a widely used technique in distributed inference that splits input data across multiple GPU devices to perform parallel inference. Each device has a complete copy of the model and processes a portion of the input data independently, combining the results to produce the final output. This approach improves throughput by distributing the workload across devices. However, data parallelism's scalability is limited by the total batch size [36], which is especially challenging in robotic IoT applications. In these applications, smaller batch sizes are common because of the need for quick responses to the environment and immediate inference when inputs are received. For example, in our experiments, the robot continuously receives the latest images from the camera for inference, with a batch size of only 1, which cannot be further divided into mini-batches, a crucial requirement for effective data parallelism.

Tensor parallelism. TP [65] is a distributed inference technique that splits a model’s layer parameters across multiple devices, each storing and computing a portion of the parameter tensors. This approach requires an all-reduce communication step after each layer to combine results from different devices, which introduces significant overhead, especially for large DNN layers. To mitigate this, TP is typically used across GPUs within the same server in data centers, using fast intra-server GPU-to-GPU links like NVLink [24], which is helpful when the model is too large for a single device. However, in robotic IoT, the limited bandwidth makes the communication cost of TP prohibitively high. Our evaluation on the same testbed as Sec. 6 of DINA [35], a state-of-the-art TP method, shows that transmission time takes up 49% to 94% of the total inference time due to all-reduce communication for each layer, making TP’s inference time 45.2X to 143.9X longer than local computation in Table 1. Although TP has lower power consumption (13.4% to 67.3% less than local computation), the extended transmission times significantly increase energy consumption per inference by 28.5X to 62.7X in Table 2. As TP greatly extends inference time, making it impractical for real-world robotic applications, we did not consider it further in this paper.

Pipeline parallelism. PP [18] 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 enhancing overall throughput rather than reducing single-inference latency [6], which is crucial in robotic IoT. As a result, existing distributed inference approaches [5, 27] in robotic IoT primarily adopt PP

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 the total inference time and their standard deviation ($\pm n$) with TP on different models in different environments. “Local computation” refers to inference the entire model locally 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”

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 [17, 22, 27, 35, 59] and optimizing robot energy consumption during inference [5, 28, 54]. However, both 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 essential for efficient distributed inference in wireless networks, as it significantly reduces communication overhead through techniques such as quantization and model distillation. Quantization [8, 13, 14] reduces the numerical precision of model weights and activations, minimizing the memory footprint and computational requirements of deep learning models by converting high-precision floating-point values (e.g., 32-bit) to lower-precision representations (e.g., 8-bit) with minimal loss of model accuracy. Model distillation [15, 29, 52], on the other hand, involves training a smaller, more efficient “student” model to mimic the behavior of a larger, more accurate “teacher” model by minimizing the difference between their outputs, allowing the distilled student model to retain 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, which aims 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, while [11] focus 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. However, 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

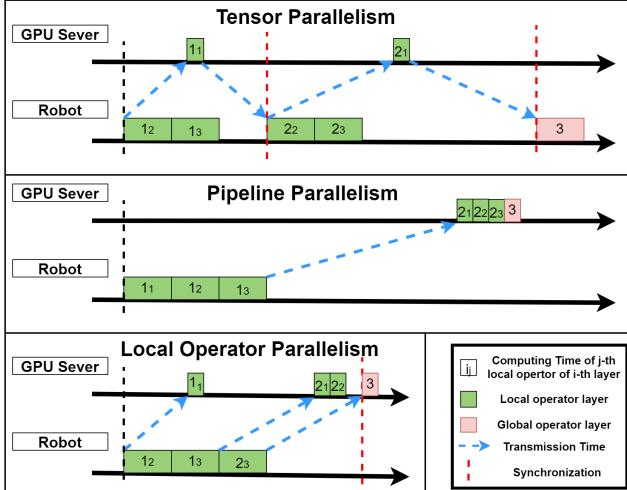


Figure 3: Workflow of Intra-DP. Each local operator layer have to complete the calculation of three local operators, and the same local operator in the three cases has the same computation time on robots and GPU servers, as well as the corresponding transmission time. The output tensor volume of layer 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 [27].

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 [65] 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 [27] 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 [23]. 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. 4 shows the architecture of Intra-DP, which adds an interceptor for each DNN layer to flexibly split the input tensor and combine the output tensor for each operator. Compared with the original model inference process on the robot, Intra-DP only increases the time cost of interceptors, which is the time cost of splitting the input tensor and combining the output tensor. The time cost of splitting the input tensor is negligible because the data transfer can be completed through the backend processes of the Intra-DP client and server while the local operators assigned to be executed on the robot and GPU server continue to perform subsequent layer calculations. The time cost of combining the output tensor is mainly bound by the time when the device on the other side completes the corresponding computation and transmission, causing prolonged waiting time. Intra-DP formulates such waiting time into the nonlinear optimization problem in its LOSS, minimizing the

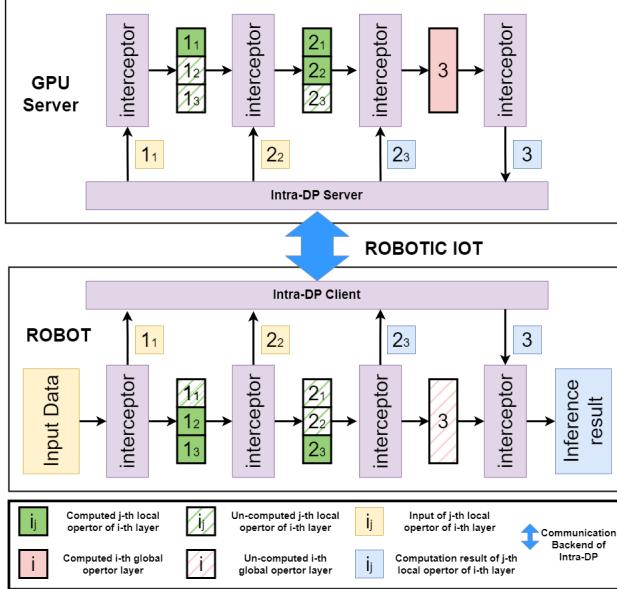


Figure 4: Architecture of Intra-DP. The core components of Intra-DP are highlighted in purple. Intra-DP adopts the same scheduling scheme as in Fig. 3.

waiting time and implementing scheduling schemes on local operators with a higher degree of parallelism. In this way, Intra-DP only increases negligible extra time on system cost and achieves faster inference via LOP and LOSS.

To address frequent fluctuations in real-world wireless networks of robotic IoT, Intra-DP generates optimal local operator schedule plans for the DNN model under different bandwidth conditions in advance. During inference, Intra-DP predicts the network bandwidth using mature tools [63] in the field of wireless transmission and adopts the corresponding schedule plan based on the predicted bandwidth. To ensure that Intra-DP can flexibly switch among various schedule plans, it keeps a copy of the model on the robot at the GPU server (Fig. 4), avoiding unnecessary transmission when migrating the parameters of local operators between robots and the GPU server. It is important to note that the model inference time, typically tens or hundreds of milliseconds, is finer (smaller) than the granularity (or frequency) of bandwidth fluctuation in real-world robotic IoT networks, as shown in Fig. 2. Therefore, we assume that the network bandwidth of robotic IoT during each inference is stable, while the network bandwidth for different inferences may differ.

4 Detailed Design

4.1 Local Operator Parallelism

LOP guarantees the correctness of inference results by determining the correct input required for different local operators based on their calculation characteristics and processes. We

summarize three classes of local operators common in models used on mobile devices:

- Element-wise local operator. This class of operators compute each element of the input tensor separately, requiring only the corresponding element from the input tensor to perform the calculation. They are widely used in activation functions such as ReLU [7], Sigmoid [62], SiLU [20]. However, it is important to note that some activation functions, like softmax [32], require all elements for computation and are not considered local operators, but global operators.
- Block-wise local operator. This class of operators require a block at the corresponding position in the input tensor and are widely used in layers associated with convolution, such as convolution [35], maxpool [48]. The size of the input blocks is determined by the parameters set by the corresponding layer [42], including the size of the convolution kernel, padding, and dilation.
- Row-wise local operator. This class of operators requires rows of the input tensor and are widely used in layers associated with matrix operations, such as addition [64] and multiplication [12]. The rows required for computation, ensuring that the correct input is obtained for each local operator to perform its respective calculation, are determined by the matrix calculation principles as following:

$$\begin{pmatrix} a_1 \\ \vdots \\ a_m \end{pmatrix} \times (b_1 \dots b_n) = \begin{pmatrix} c_{11} & \dots & c_{1n} \\ \vdots & \ddots & \vdots \\ c_{m1} & \dots & c_{mn} \end{pmatrix}$$

Row-wised local operators split the input matrix and keep a copy of the layer parameter matrix on different devices, reducing transmission volume and avoiding the synchronization needed to combine calculation results from different devices. This is in contrast to TP, which splits the layer parameter matrix and transfers a copy of the input matrix to different devices. The calculation result of row a_1 is $(c_{11} \dots c_{1n})$, which is also a row and can be directly computed by the next matrix operation layer. And LOP treats matrices with only one row as global operators.

After obtaining the required input for each local operator and the corresponding input position in the previous layer through the above analysis, LOP determines which local operators need to be computed in the previous layer to obtain the input for the current local operator. This establishes the dependency between local operators, which should be considered when scheduling local operators in LOSS. Notice that when the result of an operator is used by several local operators in the following layer, especially for block-wise local operators, LOP allows some operators to be repeatedly computed by the

robot and the GPU server, which avoids synchronization with high transmission costs in robotic IoT by introducing a small amount of redundant calculation. And we leave the support for additional types of local operators as future work.

4.2 Local Operator Scheduling Strategy

LOSS formulates the problem of scheduling local operators as a nonlinear optimization problem, modeled as follows:

First, we define OP_i as the set of operators in the i_{th} layer, including both local and global operators. When the i_{th} layer is a global operator layer, $|OP_i|$ is 1, as it only has one operator. We then define $X_i \subseteq OP_i$ as the set of operators executed on robots and $Y_i \subseteq OP_i$ as the set of operators executed on the GPU server, where $X_i \cup Y_i = OP_i$. $X_i \cap Y_i \neq \emptyset$ when the result of an operator is used by several local operators in the following layer, especially for block-wise local operators; otherwise, $X_i \cap Y_i = \emptyset$.

Next, we denote the completion time of the i_{th} layer on robots as T_{robot}^i and that on the GPU server as T_{server}^i , as shown in Fig. 3. We define $compute(X)$ as the estimated computation time of X and $transmit(X)$ as the estimated transmission time of X under the given bandwidth, leading to the following formula:

$$T_{robot}^i = \begin{cases} compute(X_0) & i = 0 \\ T_{robot}^{i-1} + compute(X_i) & i > 0, M_i = \emptyset \\ MAX(T_{robot}^{i-1}, T_{server}^{i-1} + \\ transmit(M_i)) + compute(X_i) & i > 0, M_i \neq \emptyset \end{cases}$$

$$T_{server}^i = \begin{cases} transmit(Y_0) + compute(Y_0) & i = 0 \\ T_{server}^{i-1} + compute(Y_i) & i > 0, N_i = \emptyset \\ MAX(T_{server}^{i-1}, T_{robot}^{i-1} + \\ transmit(N_i)) + compute(Y_i) & i > 0, N_i \neq \emptyset \end{cases}$$

Here, $M_i = parent(X_i) - X_{i-1}$ and $N_i = parent(Y_i) - Y_{i-1}$, where $parent(X_i)$ is the set of operators found by LOP in the $i-1_{th}$ layer, whose outputs form the inputs of X_i . The MAX function is used to minimize the idle time when combining the input tensor of this layer, and $transmit(X)$ includes not only its own transmission time but also the wait time for the previous transmission to complete.

Next, we present the corresponding objective function and constraints of a DNN model with N layers as a nonlinear optimization problem:

$$\min T_{robot}^N \quad (1)$$

$$\text{s.t. } M_i = \emptyset, \forall i \in \Pi \quad (2)$$

$$N_i = \emptyset, \forall i \in \Pi \quad (3)$$

Here, the layers in Π are those whose output data amounts are larger than the raw input data. Constraints are inspired by the key observation used in existing layer partitioning methods to

limit the transmission overhead, which states that the output data amounts in some intermediate layers of a DNN model are significantly smaller than that of its raw input data [17].

LOSS solves the above nonlinear optimization problem using the differential evolution algorithm [43] and schedules the computation and transmission of each local operator based on the obtained solution. It is important to note that when applying Intra-DP to a special model without any local operator layers, LOSS will degrade to the existing layer partitioning method. When applying Intra-DP to DNN models with complex structures as directed acyclic graphs (DAGs) (e.g., MobileNet [47], ResNet [49]), rather than simple chain-like DNN models (e.g., VGG19 [46]) as above, the $i-1_{th}$ layer in the above modeling process should be replaced by the parent layer in the corresponding DAG.

4.3 Algorithms of Intra-DP

Algorithm 1: Intra-DP_client

```

Input: Data input for inference input; DNN model model
Output: The inference result ret
Data: input of  $i_{th}$  layer  $Z_i$ ; schedule plan of  $i_{th}$  layer under
       the  $b$  bandwidth  $X_i^b, M_i^b, N_i^b$ 
       // profile phase on robot
1 info_robot = ProfileModel(model)
2 SendToServer(model, info_robot)
3  $X, M, N = ReceiveFromServer()$ 
       // inference phase on robot
4  $b = PredictsBandwidth()$ 
5  $Z_0 = input$ 
6 foreach  $i_{th}$  layer in model do
7   if  $M_i^b \neq \emptyset$  then
8     |  $Z_i = combine(Z_i, ReceiveFromServer())$ 
9   end
10  if  $N_i^b \neq \emptyset$  then
11    | SendToServer( $Z_i, N_i^b$ )
12  end
13  if  $X_i^b \neq \emptyset$  and  $Z_i \neq \emptyset$  then
14    |  $Z_{i+1} = compute(Z_i, X_i^b)$ 
15  end
16  else
17    |  $Z_{i+1} = \emptyset$ 
18  end
19 end
20 ret =  $Z_{N+1}$ 
21 return ret

```

Here, we present the algorithm of Intra-DP for both the client side on the robot and the server side on the GPU server, as shown in Fig. 4. The client part is given in Alg. 1 and the server part is given in Alg. 2. Both sides must first enter the profile phase and provide the basic information for

Algorithm 2: Intra-DP_server

Data: input of i_{th} layer Z_i ; schedule plan of i_{th} layer under the b bandwidth Y_i^b, M_i^b, N_i^b

```

1 // profile phase on server
2 model,info_robot = ReceiveFromClient()
3 info_server = ProfileModel(model)
4 X,Y,M,N = LOSS(info_robot,info_server)
5 SendToClient(X,M,N)
6 Z0 = ∅
7 foreach  $i_{th}$  layer in model do
8   if  $M_i^b \neq \emptyset$  then
9     | SendToClient( $Z_i, M_i^b$ )
10  end
11  if  $N_i^b \neq \emptyset$  then
12    |  $Z_i = \text{combine}(Z_i, \text{ReceiveFromClient}())$ 
13  end
14  if  $Y_i^b \neq \emptyset$  and  $Z_i \neq \emptyset$  then
15    |  $Z_{i+1} = \text{compute}(Z_i, Y_i^b)$ 
16  end
17  else
18    |  $Z_{i+1} = \emptyset$ 
19  end
20 end

```

LOSS (including the model structure with local and global operators, and the estimate functions *compute* and *transmit*), represented in *info_robot* and *info_server*. Then, Intra-DP generates the schedule plan of local operators under various bandwidths. Compared with the inference time, the solution of LOSS takes longer to complete, but these profiles are completed in advance, and only need to select the corresponding schedule plan according to the actual bandwidth during actual use. The copy of the model on the GPU server (Fig. 4, line 1 in Alg. 2) allows Intra-DP to switch among various schedule plans flexibly.

5 Implementation

We implement Intra-DP on Python and PyTorch. Intra-DP is easy to use and requires only three lines of code to apply to existing ML applications, as shown in Fig. 5. This is achieved by hooking around the forward method of the model, and in the first forward call we profile the model using the default PyTorch profiler and schedule; then we intercept and parallelize all the following forward calls as scheduled.

```

165 # Import package of Intra-DP
166 import intraDP
167 # Define a VGG19 model as usual
168 vgg19 = VGG19().to(device)
169 # Apply Intra-DP
170 IDP = intraDP(ip = "192.168.50.1")
171 IDP.start_client(model = vgg19)
172 # Run model for inference as usual
173 result = vgg19(input)

```

Figure 5: An example of applying Intra-DP to a VGG19 [46] model, where “192.168.50.1” is the IP address of the GPU server.

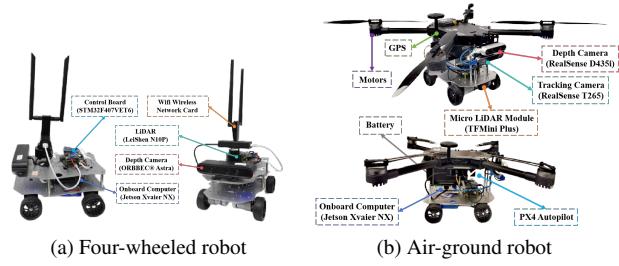


Figure 6: The detailed composition of the robot platforms

	inference	transmission	standby
Power (W)	13.35	4.25	4.04

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

6 Evaluation

Testbed. The evaluation was conducted on a custom four-wheeled robot (Fig 6a), and a custom air-ground robot(Fig 6b). They are equipped with a Jetson Xavier NX [38] 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 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

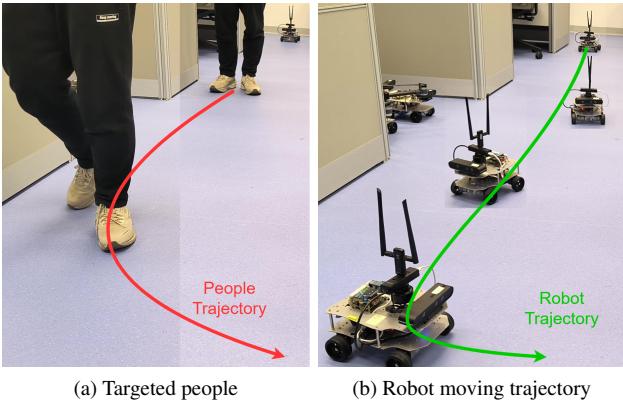


Figure 7: A real-time people-tracking robotic application on our robot based on a well-known human pose estimation ML model, Kapao [34].

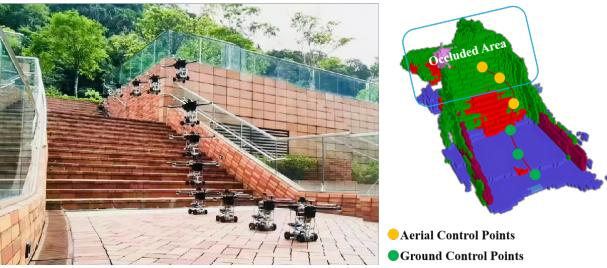


Figure 8: By predicting occlusions in advance, AGRNav [51] gains an accurate perception of the environment and avoids collisions, resulting in efficient and energy-saving paths.

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.

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 7), and AGRNav, an autonomous navigation application on our air-ground robot (Fig 8). 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 distributed inference. And we have verified several models common to mobile devices on a larger scale to further corroborate our observations and findings: DenseNet [19], VGGNet [46], ConvNeXt [53], RegNet [58].

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.

Baselines. We selected two SOTA pipeline parallelism methods as baselines: DSCCS [27], 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.

The evaluation questions are as follows:

- RQ1: How does Intra-DP benefit real-world robotic applications compared to baseline systems in terms of inference time and energy consumption?
- RQ2: How sensitive is Intra-DP to bandwidth fluctuation in robotic IoT?
- RQ3: How does Intra-DP perform on models common to mobile devices on a larger scale?
- RQ4: What are the limitations and potentials of Intra-DP?

6.1 Inference Time

The upper part of Tab.4 demonstrates that Intra-DP significantly reduced Kapao’s inference time compared to SPSO-GA and DSCCS, both indoors and outdoors. Specifically, Intra-DP achieved a 16.3% and 14.9% reduction indoors, and a 23.7% and 15% reduction outdoors, compared to SPSO-GA and DSCCS, respectively. For AGRNav, the performance gain of Intra-DP and baselines varied, as shown in the lower part of Tab.4. Intra-DP reduced AGRNav’s inference time by 36.2% and 27.8% indoors, and 41.1% and 30.8% outdoors, compared to SPSO-GA and DSCCS, respectively.

Transmission time accounts for up to 70.45% of the total inference time in SPSO-GA and DSCCS, highlighting the significant transmission bottlenecks faced by existing methods based on the PP paradigm, even with state-of-the-art layer partitioning. 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. Intra-DP’s transmission time cannot reach close to 100% of its inference time because it can only overlap the execution of local operators, and not all layers in the models of Kapao and AGRNav are local operator layers. Consequently, Intra-DP can only parallelize execution in some layers, not all layers.

Model(number of parameters)	Local computation time/s	System	Transmission time/s		Inference time/s		Percentage(%)	
			indoors	outdoors	indoors	outdoors	indoors	outdoors
kapao(77M)	1.01(± 0.03)	SPSO-GA	0.25(± 0.14)	0.31(± 0.15)	0.37(± 0.24)	0.44(± 0.25)	67.56	70.45
		DSCCS	0.21(± 0.1)	0.24(± 0.12)	0.36(± 0.2)	0.40(± 0.17)	58.33	60.21
		Intra-DP	0.24(± 0.15)	0.28(± 0.13)	0.31(± 0.14)	0.34(± 0.12)	77.42	82.35
agrnav(0.84M)	0.60(± 0.04)	SPSO-GA	0.25(± 0.11)	0.35(± 0.24)	0.47(± 0.21)	0.56(± 0.05)	53.19	62.49
		DSCCS	0.10(± 0.05)	0.15(± 0.05)	0.41(± 0.11)	0.47(± 0.12)	24.39	31.91
		Intra-DP	0.24(± 0.08)	0.26(± 0.07)	0.30(± 0.09)	0.33(± 0.07)	78.65	79.47

Table 4: Average transmission time, inference time, percentage that transmission time accounts for of the total inference time and their standard deviation ($\pm n$) of Kapao and AGRNav in different environments with different systems. “Local computation” refers to inference the entire model locally on the robot.

Model(number of parameters)	System	Power consumption(W)		Energy consumption(J) per inference	
		indoors	outdoors	indoors	outdoors
kapao(77M)	Local SPSO-GA	10.61(± 0.49)	10.61(± 0.49)	9.79(± 0.03)	9.79(± 0.03)
		5.49(± 1.52)	5.35(± 1.37)	2.03(± 0.82)	2.35(± 1.04)
		6.38(± 2.21)	6.63(± 2.38)	2.30(± 0.55)	2.65(± 0.55)
		7.05(± 1.63)	6.94(± 0.98)	2.19(± 0.62)	2.35(± 0.42)
agrnav(0.84M)	Local SPSO-GA	8.11(± 0.25)	8.11(± 0.25)	4.86(± 0.01)	4.86(± 0.01)
		5.86(± 1.60)	7.25(± 1.54)	2.75(± 0.22)	4.06(± 0.57)
		6.21(± 1.50)	7.29(± 1.55)	2.55(± 0.19)	3.43(± 0.18)
		7.52(± 0.51)	8.04(± 0.45)	2.26(± 0.15)	2.63(± 0.15)

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

The large standard deviation in transmission time outdoors for all systems indicates that bandwidth fluctuated more frequently and more fiercely outdoors compared to indoors, which is consistent with the observations in Fig. 2. Furthermore, the lower average bandwidth for outdoor scenarios (see Sec. 2.1) results in increased transmission and inference times relative to indoor scenarios.

6.2 Energy Consumption

Table 5 presents the power consumption over time and energy consumption per inference for Kapao and AGRNav using Intra-DP and baseline methods. Compared to SPSO-GA and DSCCS, Intra-DP exhibits higher power consumption over time. This can be attributed to Intra-DP’s computation at the operator granularity, where finer granularity results in lower GPU resource utilization and enables repeated computation of certain operators to avoids synchronization in LOSS.

Despite the higher power consumption over time, Intra-DP achieves the lowest energy consumption per inference for both Kapao and AGRNav, primarily due to its shortest inference time. Intra-DP avoids the need for synchronization with high transmission costs in robotic IoT by introducing a small amount of redundant computation. The additional energy consumed during Intra-DP’s computation phase is

significantly lower than the energy wasted by the prolonged inference times of SPSO-GA and DSCCS. Although SPSO-GA aims to optimize energy consumption, its advantages in power consumption over time diminish when considering energy consumption per inference due to extended inference times. This is because SPSO-GA solely focuses on minimizing power consumption over time, potentially at the cost of prolonged inference time.

6.3 Micro-Event Analysis

To further investigate Intra-DP’s sensitivity to bandwidth fluctuations in robotic IoT, we recorded the real-time bandwidth and Intra-DP’s corresponding inference time response, as depicted in Fig. 9. When the bandwidth fluctuates, Intra-DP’s inference time also fluctuates due to the changing bandwidth. However, the amplitude of inference time fluctuation is significantly smaller than that of bandwidth fluctuation, thanks to Intra-DP’s ability to flexibly switch between scheduling plans. Intra-DP generates scheduling plans for different bandwidths based on possible bandwidth fluctuation ranges in advance and adopts the corresponding scheduling plan based on the predicted bandwidth during the inference phase. This approach enables Intra-DP to adapt to varying network conditions and maintain stable performance, even in the presence

Model(number of parameters)	Local computation time/ms	System	Transmission time/ms		Inference time/ms		Percentage(%)	
			indoors	outdoors	indoors	outdoors	indoors	outdoors
DenseNet121(7M)	74.5(± 18.7)	SPSO-GA	55.0(± 33.7)	66.5(± 33.2)	76.7(± 36.3)	89.7(± 35.5)	71.76	74.15
		DSCCS	16.2(± 40.9)	20.8(± 51.9)	81.4(± 27.2)	86.6(± 27.7)	19.95	24.07
		Intra-DP	53.4(± 34.5)	52.9(± 23.9)	74.5(± 850.7)	55.1(± 15.6)	71.70	96.05
RegNet(54M)	175.0(± 23.6)	SPSO-GA	55.1(± 33.6)	66.7(± 33.6)	73.5(± 35.9)	86.5(± 35.3)	74.90	77.04
		DSCCS	47.6(± 47.8)	60.5(± 54.0)	77.8(± 39.3)	86.2(± 37.9)	61.22	70.22
		Intra-DP	49.6(± 21.7)	59.9(± 23.4)	55.0(± 24.8)	64.2(± 25.2)	90.18	93.34
ConvNeXt(88M)	160.2(± 21.0)	SPSO-GA	55.4(± 33.9)	66.9(± 34.9)	73.8(± 35.4)	86.7(± 36.3)	75.13	77.15
		DSCCS	46.9(± 43.1)	56.7(± 52.1)	72.4(± 35.7)	84.7(± 36.3)	64.78	66.95
		Intra-DP	50.4(± 32.2)	61.9(± 34.8)	53.9(± 26.2)	65.7(± 27.7)	93.51	94.23
VGG19(143M)	118.0(± 18.9)	SPSO-GA	55.7(± 33.5)	67.2(± 35.0)	68.8(± 33.9)	80.6(± 35.0)	80.84	83.33
		DSCCS	38.9(± 47.1)	41.6(± 53.8)	65.2(± 28.1)	75.5(± 27.1)	59.75	55.09
		Intra-DP	44.8(± 20.9)	51.5(± 15.0)	47.6(± 18.1)	53.6(± 14.7)	94.15	96.07
ConvNeXt(197M)	316.7(± 31.0)	SPSO-GA	57.1(± 38.9)	67.1(± 34.5)	80.5(± 40.8)	90.9(± 35.0)	70.87	73.89
		DSCCS	56.0(± 36.1)	67.0(± 37.6)	79.2(± 35.9)	90.6(± 35.4)	70.72	73.98
		Intra-DP	56.4(± 34.7)	66.5(± 33.7)	59.7(± 26.6)	68.0(± 26.6)	94.43	97.88

Table 6: Average transmission time, inference time, percentage that transmission time accounts for of the total inference time and their standard deviation ($\pm n$) of common AI models in different environments with different systems.

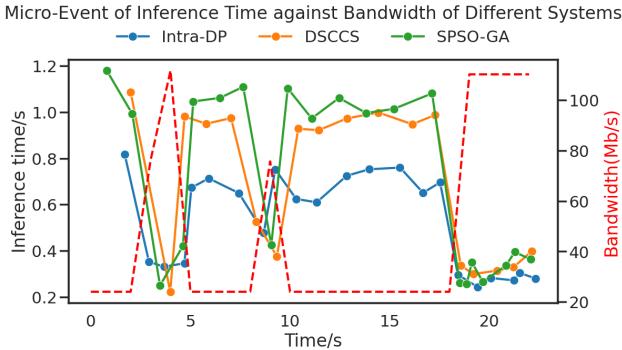


Figure 9: Real-time bandwidth and inference time of Intra-DP and baselines.

of bandwidth fluctuations.

6.4 Validation on a larger range of models

We evaluated Intra-DP and baselines on a wide range of models commonly used in mobile devices, with parameter counts varying (detailed in Tab. 6 and Tab. 7). Our results confirm that transmission time constitutes a significant portion of the total inference time in both DSCCS and SPSO-GA, leading to wasteful inference time and energy consumption compared to Intra-DP. Although Intra-DP’s outperformance remains consistent across various models, we observed that the performance gain is relatively smaller on models with fewer parameters. This is because Intra-DP’s performance improvement is primarily achieved through the parallel execution of local operators. When a model employs more global operator

layers or has fewer parameters, the number of local operators available for parallel execution is reduced, limiting the optimization potential for Intra-DP to enhance performance.

6.5 Lessons learned

Global optimal solution. The performance of Intra-DP heavily depends on the quality of the solution obtained by LOSS, with higher solution quality (closer to the global optimal solution) leading to better performance. However, as DNN models become increasingly complex with more layers, the nonlinear optimization problem in LOSS faces exponentially higher parameter dimensions and complexity, resulting in unacceptable profile times. Solving the global optimal solution of nonlinear optimization problems in finite time remains an open issue, and finding a fast and high-quality solution algorithm for Intra-DP is left for future work.

Model structure. During the implementation and evaluation of Intra-DP, we discovered that the presence of more local operator layers allows for increased parallel execution during model inference, thereby enhancing the performance improvement of Intra-DP. Future work should focus on supporting additional types of local operators and exploring the possibility of transforming global operators into local ones through lightweight synchronization techniques, based on their computational characteristics (e.g., synchronize the sum results in softmax instead of directly transferring the full input tensor and re-calculating).

Future work. It is of interest to explore further improvements of Intra-DP, such as a distributed inference system for multi-robot to minimize overall inference time and energy consumption. Such advancements could enable faster

Model(number of parameters)	System	Power consumption(W) indoors	Power consumption(W) outdoors	Energy consumption(J) per inference indoors	Energy consumption(J) per inference outdoors
DenseNet121(7M)	Local	8.2(± 0.27)	8.2(± 0.27)	0.46(± 0.04)	0.46(± 0.04)
	SPSO-GA	4.87(± 0.23)	4.74(± 0.19)	0.37(± 0.02)	0.42(± 0.02)
	DSCCS	6.91(± 0.45)	6.86(± 0.46)	0.56(± 0.04)	0.59(± 0.04)
	Intra-DP	5.36(± 0.79)	5.79(± 0.24)	0.4(± 0.06)	0.32(± 0.01)
RegNet(54M)	Local	9.0(± 0.3)	9.0(± 0.3)	1.37(± 0.02)	1.37(± 0.02)
	SPSO-GA	4.83(± 0.21)	4.8(± 0.18)	0.36(± 0.02)	0.41(± 0.02)
	DSCCS	5.84(± 1.79)	5.36(± 1.34)	0.45(± 0.14)	0.46(± 0.12)
	Intra-DP	5.24(± 1.43)	5.28(± 1.52)	0.29(± 0.08)	0.34(± 0.1)
ConvNeXt(88M)	Local	9.7(± 0.34)	9.7(± 0.34)	1.34(± 0.02)	1.34(± 0.02)
	SPSO-GA	4.93(± 0.25)	4.78(± 0.18)	0.36(± 0.02)	0.41(± 0.02)
	DSCCS	6.01(± 0.27)	5.71(± 1.56)	0.439(± 0.05)	0.48(± 0.13)
	Intra-DP	6.68(± 1.23)	6.68(± 1.21)	0.36(± 0.07)	0.44(± 0.08)
VGG19(143M)	Local	9.78(± 0.34)	9.78(± 0.34)	0.95(± 0.02)	0.95(± 0.02)
	SPSO-GA	4.9(± 0.25)	4.82(± 0.2)	0.34(± 0.02)	0.39(± 0.02)
	DSCCS	6.58(± 2.14)	6.93(± 2.35)	0.43(± 0.14)	0.52(± 0.18)
	Intra-DP	6.51(± 1.74)	7.32(± 1.52)	0.31(± 0.08)	0.39(± 0.08)
ConvNeXt(197M)	Local	10.72(± 0.38)	10.72(± 0.38)	3.12(± 0.03)	3.12(± 0.03)
	SPSO-GA	5.1(± 0.27)	4.99(± 0.2)	0.41(± 0.02)	0.45(± 0.02)
	DSCCS	5.06(± 0.31)	5.02(± 0.37)	0.4(± 0.02)	0.45(± 0.03)
	Intra-DP	4.57(± 0.23)	4.54(± 0.25)	0.27(± 0.01)	0.31(± 0.02)

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

and more robust wireless distributed inference in real-world robotic IoT.

7 Conclusion

In this paper, we present Intra-DP, a high-performance distributed inference system optimized for robotic IoT networks. By breaking up the granularity of model inference into local operators via LOP and applying adaptive scheduling to the computation and transmission of each local operator via LOSS, 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. We envision that the fast and energy-efficient inference of Intra-DP will foster the real-world deployment of diverse AI robotic tasks in the field.

References

- [1] iPerf - Download iPerf3 and original iPerf pre-compiled binaries.
- [2] Toni Adame, Marc Carrascosa-Zamacois, and Boris Bel-lalta. Time-sensitive networking in ieee 802.11 be: On the way to low-latency wifi 7. *Sensors*, 21(15):4954, 2021.
- [3] Majid Altamimi, Atef Abd Rabou, Kshirasagar Naik, and Amiya Nayak. Energy cost models of smartphones for task offloading to the cloud. *IEEE Transactions on Emerging Topics in Computing*, 3(3):384–398, 2015.
- [4] Anh-Quan Cao and Raoul de Charette. Monoscene: Monocular 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3991–4001, 2022.
- [5] Xing Chen, Jianshan Zhang, Bing Lin, Zheyi Chen, Katinka Wolter, and Geyong Min. Energy-efficient offloading for dnn-based smart iot systems in cloud-edge environments. *IEEE Transactions on Parallel and Distributed Systems*, 33(3):683–697, 2021.
- [6] Daniel Crankshaw, Gur-Eyal Sela, Xiangxi Mo, Corey Zumar, Ion Stoica, Joseph Gonzalez, and Alexey Tumanov. Inferline: latency-aware provisioning and scaling for prediction serving pipelines. In *Proceedings of the 11th ACM Symposium on Cloud Computing*, pages 477–491, 2020.
- [7] Ingrid Daubechies, Ronald DeVore, Simon Foucart, Boris Hanin, and Guergana Petrova. Nonlinear approxi-

- mation and (deep) relu networks. *Constructive Approximation*, 55(1):127–172, 2022.
- [8] Alexandre Défossez, Yossi Adi, and Gabriel Synnaeve. 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. Performance impact of los and nlos transmissions in dense cellular networks. *IEEE Transactions on Wireless Communications*, 15(3):2365–2380, 2015.
- [10] Khalid Elgazzar, Patrick Martin, and Hossam S Hasanein. Cloud-assisted computation offloading to support mobile services. *IEEE Transactions on Cloud Computing*, 4(3):279–292, 2014.
- [11] Zhou Fang, Tong Yu, Ole J Mengshoel, and Rajesh K Gupta. 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*, pages 2067–2070, 2017.
- [12] Kayvon Fatahalian, Jeremy Sugerman, and Pat Hanrahan. Understanding the efficiency of gpu algorithms for matrix-matrix multiplication. In *Proceedings of the ACM SIGGRAPH/EUROGRAPHICS conference on Graphics hardware*, pages 133–137, 2004.
- [13] Stefan Gheorghe and Mihai Ivanovici. Model-based weight quantization for convolutional neural network compression. In *2021 16th International Conference on Engineering of Modern Electric Systems (EMES)*, pages 1–4. IEEE, 2021.
- [14] Cheng Gong, Yao Chen, Ye Lu, Tao Li, Cong Hao, and Deming Chen. Vecq: Minimal loss dnn model compression with vectorized weight quantization. *IEEE Transactions on Computers*, 70(5):696–710, 2020.
- [15] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021.
- [16] Vinayak Honkote, Dileep Kurian, Sriram Muthukumar, Dibyendu Ghosh, Satish Yada, Kartik Jain, Bradley Jackson, Ilya Klotchkov, Mallikarjuna Rao Nimmagadda, Shreela Dattawadkar, et al. 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)*, pages 48–50. IEEE, 2019.
- [17] Chuang Hu, Wei Bao, Dan Wang, and Fengming Liu. Dynamic adaptive dnn surgery for inference acceleration on the edge. In *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, pages 1423–1431. IEEE, 2019.
- [18] Yang Hu, Connor Imes, Xuanang Zhao, Souvik Kundu, Peter A Beerel, Stephen P Crago, and John Paul Walters. Pipeedge: Pipeline parallelism for large-scale model inference on heterogeneous edge devices. In *2022 25th Euromicro Conference on Digital System Design (DSD)*, pages 298–307. IEEE, 2022.
- [19] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks, 2018.
- [20] Glenn Jocher, Alex Stoken, Jirka Borovec, Liu Changyu, Adam Hogan, Ayush Chaurasia, Laurentiu Diaconu, Francisco Ingham, Adrien Colmagro, Hu Ye, et al. ultralytics/yolov5: v4. 0-nn. silu () activations, weights & biases logging, pytorch hub integration. *Zenodo*, 2021.
- [21] K J Joseph, Salman Khan, Fahad Shahbaz Khan, and Vineeth N Balasubramanian. Towards open world object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5830–5840, June 2021.
- [22] Yiping Kang, Johann Hauswald, Cao Gao, Austin Rovinski, Trevor Mudge, Jason Mars, and Lingjia Tang. Neurorurgeon: Collaborative intelligence between the cloud and mobile edge. *ACM SIGARCH Computer Architecture News*, 45(1):615–629, 2017.
- [23] Nam Sung Kim, Todd Austin, David Baauw, Trevor Mudge, Krisztián Flautner, Jie S Hu, Mary Jane Irwin, Mahmut Kandemir, and Vijaykrishnan Narayanan. Leakage current: Moore’s law meets static power. *computer*, 36(12):68–75, 2003.
- [24] Ang Li, Shuaiwen Leon Song, Jieyang Chen, Jiajia Li, Xu Liu, Nathan R Tallent, and Kevin J Barker. Evaluating modern gpu interconnect: Pcie, nvlink, nv-sli, nvswitch and gpudirect. *IEEE Transactions on Parallel and Distributed Systems*, 31(1):94–110, 2019.
- [25] Qingbiao Li, Fernando Gama, Alejandro Ribeiro, and Amanda Prorok. Graph neural networks for decentralized multi-robot path planning. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 11785–11792. IEEE, 2020.
- [26] Yiming Li, Zhiding Yu, Christopher Choy, Chaowei Xiao, Jose M Alvarez, Sanja Fidler, Chen Feng, and Anima Anandkumar. Voxformer: Sparse voxel transformer for camera-based 3d semantic scene completion.

- In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9087–9098, 2023.
- [27] Huanghuang Liang, Qianlong Sang, Chuang Hu, Dazhao Cheng, Xiaobo Zhou, Dan Wang, Wei Bao, and Yu Wang. Dnn surgery: Accelerating dnn inference on the edge through layer partitioning. *IEEE transactions on Cloud Computing*, 2023.
- [28] Bing Lin, Yinhao Huang, Jianshan Zhang, Junqin Hu, Xing Chen, and Jun Li. Cost-driven off-loading for dnn-based applications over cloud, edge, and end devices. *IEEE Transactions on Industrial Informatics*, 16(8):5456–5466, 2019.
- [29] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363, 2020.
- [30] Ruofeng Liu and Nakjung Choi. 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):1–25, 2023.
- [31] Shuai Liu, Xin Li, Huchuan Lu, and You He. Multi-object tracking meets moving uav. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8876–8885, June 2022.
- [32] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-margin softmax loss for convolutional neural networks. *arXiv preprint arXiv:1612.02295*, 2016.
- [33] Antoni Masiukiewicz. 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):79–84, 2019.
- [34] William McNally, Kanav Vats, Alexander Wong, and John McPhee. Rethinking keypoint representations: Modeling keypoints and poses as objects for multi-person human pose estimation. In *European Conference on Computer Vision*, pages 37–54. Springer, 2022.
- [35] Thaha Mohammed, Carlee Joe-Wong, Rohit Babbar, and Mario Di Francesco. Distributed inference acceleration with adaptive dnn partitioning and offloading. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*, pages 854–863. IEEE, 2020.
- [36] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostafa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al. Efficient large-scale language model training on gpu clusters using megatron-lm. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15, 2021.
- [37] Mohammad Noormohammadpour and Cauligi S Raghavendra. Datacenter traffic control: Understanding techniques and tradeoffs. *IEEE Communications Surveys & Tutorials*, 20(2):1492–1525, 2017.
- [38] NVIDIA. The world’s smallest ai supercomputer. <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-xavier-series/>, 2024.
- [39] Takeshi Ohkawa, Kazushi Yamashina, Hitomi Kimura, Kanemitsu Ootsu, and Takashi Yokota. Fpga components for integrating fpgas into robot systems. *IEICE TRANSACTIONS on Information and Systems*, 101(2):363–375, 2018.
- [40] Yuanteng Pei, Matt W Mutka, and Ning Xi. Connectivity and bandwidth-aware real-time exploration in mobile robot networks. *Wireless Communications and Mobile Computing*, 13(9):847–863, 2013.
- [41] pytorch. pytorch. <https://pytorch.org/>, 2024.
- [42] pytorch. pytorch. <https://pytorch.org/docs/stable/generated/torch.nn.Conv2d.html>, 2024.
- [43] A Kai Qin, Vicky Ling Huang, and Ponnuthurai N Suganthan. Differential evolution algorithm with strategy adaptation for global numerical optimization. *IEEE transactions on Evolutionary Computation*, 13(2):398–417, 2008.
- [44] Yi Ren, Chih-Wei Tung, Jyh-Cheng Chen, and Frank Y Li. Proportional and preemption-enabled traffic offloading for ip flow mobility: Algorithms and performance evaluation. *IEEE Transactions on Vehicular Technology*, 67(12):12095–12108, 2018.
- [45] Nurul I Sarkar and Osman Mussa. The effect of people movement on wi-fi link throughput in indoor propagation environments. In *IEEE 2013 Tencon-Spring*, pages 562–566. IEEE, 2013.
- [46] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015.
- [47] Debjyoti Sinha and Mohamed El-Sharkawy. Thin mobilenet: An enhanced mobilenet architecture. In *2019 IEEE 10th annual ubiquitous computing, electronics & mobile communication conference (UEMCON)*, pages 0280–0285. IEEE, 2019.

- [48] Luna Sun, Zhenxue Chen, QM Jonathan Wu, Hongjian Zhao, Weikai He, and Xinghe Yan. Ampnet: Average-and max-pool networks for salient object detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(11):4321–4333, 2021.
- [49] Sasha Targ, Diogo Almeida, and Kevin Lyman. Resnet in resnet: Generalizing residual architectures. *arXiv preprint arXiv:1603.08029*, 2016.
- [50] Hao Wang, Sreeram Potluri, Miao Luo, Ashish Kumar Singh, Sayantan Sur, and Dhabaleswar K Panda. Mvapich2-gpu: optimized gpu to gpu communication for infiniband clusters. *Computer Science-Research and Development*, 26(3):257–266, 2011.
- [51] Junming Wang, Zekai Sun, Xiuxian Guan, Tianxiang Shen, Zongyuan Zhang, Tianyang Duan, Dong Huang, Shixiong Zhao, and Heming Cui. Agrnav: Efficient and energy-saving autonomous navigation for air-ground robots in occlusion-prone environments. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [52] Lin Wang and Kuk-Jin Yoon. Knowledge distillation and student-teacher learning for visual intelligence: A review and new outlooks. *IEEE transactions on pattern analysis and machine intelligence*, 44(6):3048–3068, 2021.
- [53] Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16133–16142, 2023.
- [54] Huaming Wu, William J Knottenbelt, and Katinka Wolter. An efficient application partitioning algorithm in mobile environments. *IEEE Transactions on Parallel and Distributed Systems*, 30(7):1464–1480, 2019.
- [55] Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. A brief overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal of Automatica Sinica*, 10(5):1122–1136, 2023.
- [56] Zhaoyang Xia, Youquan Liu, Xin Li, Xinge Zhu, Yuexin Ma, Yikang Li, Yuenan Hou, and Yu Qiao. Scpnet: Semantic scene completion on point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17642–17651, 2023.
- [57] Yecheng Xiang and Hyoseung Kim. Pipelined data-parallel cpu/gpu scheduling for multi-dnn real-time inference. In *2019 IEEE Real-Time Systems Symposium (RTSS)*, pages 392–405. IEEE, 2019.
- [58] Jing Xu, Yu Pan, Xinglin Pan, Steven Hoi, Zhang Yi, and Zenglin Xu. Regnet: self-regulated network for image classification. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [59] Min Xue, Huaming Wu, Guang Peng, and Katinka Wolter. Ddpqn: An efficient dnn offloading strategy in local-edge-cloud collaborative environments. *IEEE Transactions on Services Computing*, 15(2):640–655, 2021.
- [60] Xinlei Yang, Hao Lin, Zhenhua Li, Feng Qian, Xingyao Li, Zhiming He, Xudong Wu, Xianlong Wang, Yunhao Liu, Zhi Liao, et al. Mobile access bandwidth in practice: Measurement, analysis, and implications. In *Proceedings of the ACM SIGCOMM 2022 Conference*, pages 114–128, 2022.
- [61] Yang Yang, Li Juntao, and Peng Lingling. Multi-robot path planning based on a deep reinforcement learning dqn algorithm. *CAAI Transactions on Intelligence Technology*, 5(3):177–183, 2020.
- [62] Xinyou Yin, JAN Goudriaan, Egbert A Lantinga, JAN Vos, and Huub J Spiertz. A flexible sigmoid function of determinate growth. *Annals of botany*, 91(3):361–371, 2003.
- [63] Chaoqun Yue, Ruofan Jin, Kyoungwon Suh, Yanyuan Qin, Bing Wang, and Wei Wei. Linkforecast: Cellular link bandwidth prediction in lte networks. *IEEE Transactions on Mobile Computing*, 17(7):1582–1594, 2017.
- [64] Anthony Zee. Law of addition in random matrix theory. *Nuclear Physics B*, 474(3):726–744, 1996.
- [65] Yonghao Zhuang, Hexu Zhao, Lianmin Zheng, Zhuohan Li, Eric Xing, Qirong Ho, Joseph Gonzalez, Ion Stoica, and Hao Zhang. On optimizing the communication of model parallelism. *Proceedings of Machine Learning and Systems*, 5, 2023.