

New Problems in Distributed Inference for DNN Models on Robotic IoT

(Anonymous Authors)

Abstract—The abstract goes here.

I. INTRODUCTION

The rapid progress in machine learning (ML) techniques has led to remarkable achievements in various fundamental robotic tasks, such as object detection[6], [11], [13], robotic control[8], [21], [25], and environmental perception[3], [9], [23]. However, deploying these ML approaches 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 limited battery capacity of robots. Relying solely on additional computing accelerators on robots (e.g., GPU[15], FPGA[16], SoC[5]) increases energy consumption (e.g., 62% for [13] in our experiments) due to the computationally-intensive nature of DNN models, while offloading the entire model to the cloud cannot meet the latency requirements of robotic applications because of limited bandwidth, described in Sec.II-A.

Distributed inference[?], which involves offloading a portion of the inference computation onto other powerful GPU devices, has emerged as a promising approach to meet the latency requirements of robotic applications and extend the battery lifetime of robots. By leveraging the high computing capabilities of powerful GPUs through Internet of Things for these robots (robotic IoT), distributed inference accelerates the inference process while alleviating the local computational burden, thereby reducing energy consumption. This paradigm has been widely adopted in data centers[?], where numerous GPUs are utilized for large model inference, such as in the case of ChatGPT[22].

However, all existing parallel methods for distributed inference in data center are ill-suited in robotic IoT. In data center, there are mainly three kinds of parallel methods: Data parallelism (DP) replicates the model across devices, and lets each replica handle one micro-batch (i.e., splits of a batch which is a set of training inputs for each device); Tensor parallelism (TP) splits a single DNN layer (often too large to fit in one device) over devices; Pipeline parallelism (PP) places different layers of a DNN model over devices (layer partition) and pipelines the inference to reduce devices' idling time (pipeline execution). In this paper, we show that several problems are hindering existing parallel methods applying on robotic IoT.

Problem 1 (DP). The small batch sizes inherent to robotic IoT applications (typically 1, 3, etc.) hinder the computation of micro-batches, rendering DP inapplicable for robotic IoT.

In data centers, DP is feasible due to the large batch sizes employed (e.g., 16 or 32 images), allowing for the division of inputs into micro-batches that still contain several complete inputs (e.g., 2 or 4 images). However, in the context of robotic IoT, real-time performance is crucial, necessitating immediate inference upon receiving inputs, which typically have smaller batch sizes (e.g., 1 or 3 images). Further splitting these inputs would result in micro-batches containing incomplete inputs (e.g., 1/4 or 1/2 of an image), which cannot be computed in parallel to speed up inference.

Problem 2 (TP). TP requires frequent synchronization among devices, leading to unacceptable communication overhead in robotic IoT. In TP, synchronizing and combining computation results from different devices after processing the split layer entails significant communication overhead, which necessitates dedicated high-speed interconnects (e.g., 400 Gbps for NVLink[7]) even within data centers. Conversely, robots must prioritize seamless mobility and primarily depend on wireless connections, which inherently possess limited bandwidth. todo: need more explain and experiment data

Consequently, existing distributed inference approach[?] in robotic IoT primarily adopts the PP paradigm, which consist of two main parts: layer partition (places different layers of a DNN model across robots and GPU devices) and pipeline execution. Existing methods for robotic IoT primarily focus on layer partition of PP to achieve fast and energy-efficient inference, because the pipeline execution of PP enhances inference throughput, rather than reducing the completion time of a single inference, which is the most critical requirement in robotic IoT, and automatically partition a model over multiple devices through transforming the parallelization optimization problem into a cost minimization problem. By capitalizing on the fact that the number of parameters in DNN models typically decreases layer by layer, Allocating more layers to robots reduces transmission volume, while assigning more layers to GPU devices accelerates the computation phase of inference. These methods strike a balance on layer partition based on application-specific inference speed requirements and energy consumption demands.

Problem 3 (PP). Existing methods based on PP faces significant challenges due to transmission bottlenecks on robotic IoT (e.g., up to 70% of inference time and 50% of energy consumption in our experiments), which cannot be effectively mitigated through PP scheduling. Firstly, PP is unable to

overlap the transmission and computation phases within the same inference[?] PP can only overlap these phases across multiple inferences, which increases throughput but not the speed of a single inference. Secondly, reducing transmission volume is not an ideal solution for robotic IoT. Allocating more DNN model layers to robots does reduce the transmission volume and the associated transmission time. However, this approach increases energy consumption and under-utilizes the powerful GPU devices for accelerating inference, and may not necessarily reduce overall inference time.

In this paper, we take the first step to reveal and evaluate the problems hindering existing parallel methods for distributed inference applying on robotic IoT. These findings aim to raise research effort to find a new parallel method to speedup distributed inference on robotic IoT, so that the DNN models deployed on real-world robots can achieve fast and energy-efficient inference and it will nurture diverse ML applications deployed on mobile robots in the field.

The rest of the paper is organized as follows: the chapter II provides background; the chapter III describes in detail about the problems; the chapter IV; provides with evaluation results; the chapter V concludes.

II. BACKGROUND

A. Characteristics of Robotic IoT Networks

Wireless transmission of Robotic IoT. In real-world robotic IoT scenarios, devices often need to navigate and move around for tasks such as search and exploration. While wireless networks provide the required high mobility, they also come with 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[10]. Furthermore, the actual available bandwidth of wireless networks is often reduced in practice due to factors such as movement of the devices[12], [17], occlusion from by physical barriers[4], [19], and preemption of the wireless channel by other devices[2], [18].

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). The hardware and wireless network settings are as described in Sec. IV. We believe our setup represents the state-of-the-art (SOTA) computation and communication capabilities of robotic IoT devices. We saturated the wireless network connection with iperf [1] and recorded the average bandwidth capacity between these robots every 0.1s for 5 minutes. For comparison, we also recorded the actual maximum bandwidth in our lab without movement or occlusion (stable) under the same experimental equipment.

The results, shown in Fig.1, indicate a 44% and 56% degradation in average bandwidth capacity for indoor and outdoor scenarios, respectively, when compared to stable records without movement or occlusion. The outdoor environment exhibited higher instability, with bandwidth frequently dropping

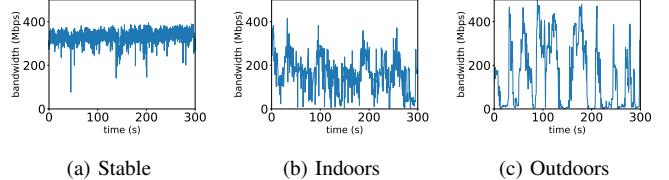


Fig. 1: The instability of wireless transmission in robotic IoT networks.

to extremely low values around 0 Mbps. This can be attributed to the lack of walls in open outdoor areas to reflect wireless signals, and the presence of obstacles, such as trees, between communicating robots. Consequently, fewer signals can be received in outdoor areas compared to indoor environments.

In summary, the wireless transmission of robotic IoT systems is constrained by limited bandwidth, both due to the theoretical upper limit of wireless transmission technologies and the practical instability of wireless networks.

Comparison with Datacenter Networks. Compared to robotic IoT networks, data center networks (used for large model inference, e.g., ChatGPT[22]) are wired and typically exhibit higher bandwidth capacity and lower bandwidth fluctuation. In data center networks, GPU devices are connected using high-speed networking technologies, such as InfiniBand[20] or PCIe[7], which offer bandwidths ranging from 40 Gbps to 500 Gbps. Bandwidth fluctuation in these networks is primarily caused by congestion on intermediate switches and can be mitigated through traffic scheduling on switches[14]. However, when distributed inference involves GPU devices deployed in commodity clouds (data centers), robotic IoT systems must access these clouds via the internet, which is subject to bandwidth fluctuations caused by variations in overall traffic volume. Reports from large-scale bandwidth testing services [24] reveal that the median bandwidth for accessing commodity clouds is only 137 Mbps in the US and 153 Mbps in China.

B. Existing parallel inference strategies in data center

Data parallelism. Data parallelism in distributed inference is an approach that divides the input data into smaller subsets (micro-batch) and processes them concurrently across multiple devices, thus leveraging their combined computational power. This parallel processing of data samples allows for faster inference speeds, . When data parallelism introduces increased throughput, efficient resource utilization, and the ability to handle large-scale data sets, it also introduces challenges such as communication overhead for synchronizing and potential difficulties in maintaining effectiveness across devices with various available computational resources. Its scalability [8] is bounded by communication on low-end networks and the size of a total batch (a set of data for producing each parameter update).

Tensor parallelism. Tensor parallelism in distributed inference is an approach that partitions input and parameter tensors of a layer (e.g., transformer multi-head self-attention layer [30]) across GPUs, allowing for concurrent computation on different parts of this tensor. TP requires an all-reduce communication (i.e., TP.sync) to aggregate the tensors between repeated blocks, which typically lies on the critical path and is network-intense. Therefore, TP is usually deployed across GPUs within the same server to use fast intra-server GPU-to-GPU links (e.g., NVLink [42]). TP is mainly adopted as a complementary technique to help existing parallel training systems [8], [19], [25], [41] to support larger transformer layers.

Pipeline parallelism.

C. Related Work

Compressed communication. compression greatly reduces the communication traffic volume and is indeed essential for practical distributed inference over wireless networks. Some lossy gradient compression methods [38] (information is lost during inference and cannot be recovered) achieve up to 0.1% compression rate (i.e., size after compression divided by original size), but they cannot provide convergence guarantee [38]. In this paper, we only consider lossless compression methods (e.g., the lost information during compression is compensated with error compensation [22]). Even with gradient compression, communication still takes a major time portion in distributed training on robotic IoT devices for two reasons. With the rapid advancement of SOTA robotic IoT devices [39], the computation time on each device is also decreasing. As a result, the communication time is typically comparable to the computation time. Consequently, even with gradient compression, the straggler effect, which severely prolongs the communication time, still has a major impact on the distributed training process. In our experiments, a Jetson Xavier NX [40] device out of a fourdevice team computed gradients in 2.18s and ideally needed to wait for 1.47s upon the synchronization barrier in BSP (four devices push and pull the compressed gradients sized 2.1MByte, summing up to 134.4Mbit), which is comparable to (equal to 67.4% of) the computation time. Meanwhile, the straggler effect in the above indoor scenario caused each device to on average stall for 2.23 s in each iteration, equal to 102.2% of the computation time, severely degrading the training throughput.

Job scheduling. Many research efforts have been devoted to the offloading problems in cloud computing. For example, Altamimiet al. [8] pointed out that it would be feasible to extend the battery life-time and capabilities of smartphones by offloading tasks directly to the cloud. They designed a mathematical model to estimate the energy consumption during the task offloading process. Elgazzar et al. [9] proposed a framework of cloud-assisted mobile service provisioning for reliable deliveries. This framework can support online offloading based on the current network and resource status of mobile systems, while meeting user-defined energy constraints. Fang et al. [10] focused on QoS-aware scheduling

of heterogeneous servers for DNN inference workloads in cloud computing. They designed a deep reinforcement learning (DRL) based scheduler to maximize the Quality-of-Service (QoS) including the inference accuracy and response delay. Kumar and Lu [11] argued that IoT systems may not be time-saving or energy-efficient if their computation tasks are offloaded to the cloud directly. Since it would increase the traffic load of core network and cause high latency due to the massive data transmission between IoT devices and the cloud.

III. PROBLEMS IN PARALLEL INFERENCE STRATEGIES FOR DNN MODELS ON ROBOTIC IOT

A. Parallel inference strategies

- 1) Data parallelism:
- 2) Tensor parallelism:
- 3) Pipeline parallelism:

B. Existing distributed inference on robotic IoT

1) *Layer partition:* Hu et al. [23] designed a dynamic DNN surgery strategy to partition DNN inference between the cloud and edge at the granularity of the DNN layers. This strategy reduced the system latency and improved throughput by limiting data transmission, but it paid less attention to the offloading problem for DNN layers. Mohammed et al. [19] proposed an adaptive DNN partition scheme and a distributed algorithm based on the matching game method, where the DNN layers were offloaded to fog nodes. Neurosurgeon [4] claimed that excessive latency and energy consumption were generated when uploading massive data of DNNs to the cloud via the wireless network. To cope with this problem, a lightweight scheduler was designed to partition DNN-based applications automatically between end devices and the cloud at the granularity of DNN layers.

2) *Energy consumption:* Most of the aforementioned work tried to reduce the system latency in cloud/edge environments [16], [18], [23], but they did not consider reducing the system energy consumption of offloading DNN layers with deadline constraints. An application partitioning algorithm was presented in [20] for pursuing a trade-off between energy consumption and data transmission in dynamic mobile environments. Teerapittayanon et al. [21] proposed distributed DNNs over the cloud, edge, and IoT devices. They considered the data transmission cost but not the layer execution consumption, while deploying distributed DNNs in the cloud-edge environments. In the previous work [6], a cost-driven offloading scheme was designed for DNN-based smart IoT systems with deadline constraints over the cloud, edge, and IoT devices, where a discrete PSO algorithm was developed to reduce the system cost of executing DNN layers and transferring data. Different from this work, we further consider the energy consumption of each participating server and IoT device, and introduce the layer partition operations into the offloading decision-making process for DNN layers.

In general, most of these work focused on the offloading problem in distributed inference. However, it is still an open

issue to optimize the system energy consumption when offloading DNN layers with deadline constraints in robotic IoT.

C. Dilemma on Inference Time

D. Dilemma on Energy Consumption

IV. EVALUATION

Testbed. The evaluation was conducted on a custom four-wheeled robot (Fig 2a), and a custom air-ground robot (Fig 2b). They are equipped with a Jetson Xavier NX[15] 8G onboard computer serving as the ROS master. The system runs Ubuntu 18.04 and utilizes a SanDisk 256G memory card, with ROS2 Galactic installed for application development 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 host computer processes environmental information in ROS2 Galactic, performing path planning, navigation, and obstacle avoidance before transmitting velocity and control data to corresponding ROS topics. The controller then subscribes to these topics, executing robot control tasks.

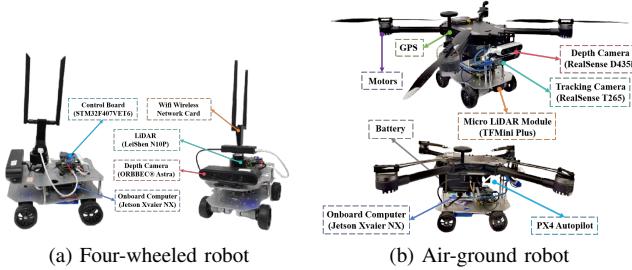


Fig. 2: The detailed composition of the robot platforms

Real-world Robotic Applications. We evaluated two kinds of typical real-world applications on robots: a real-time people-tracking robotic application on our four-wheeled robot as depicted in Fig 3 and a autonomous navigation on our air-ground robot as depicted in Fig 4.

V. CONCLUSION

The conclusion goes here.

ACKNOWLEDGMENT

The authors would like to thank...

REFERENCES

- [1] iPerf - Download iPerf3 and original iPerf pre-compiled binaries.
- [2] ADAME, T., CARRASCOSA-ZAMACOIS, M., AND BELLALTA, B. Time-sensitive networking in ieee 802.11 be: On the way to low-latency wifi 7. *Sensors* 21, 15 (2021), 4954.
- [3] CAO, A.-Q., AND DE CHARETTE, R. Monoscene: Monocular 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 3991–4001.
- [4] DING, M., WANG, P., LÓPEZ-PÉREZ, D., MAO, G., AND LIN, Z. Performance impact of los and nlos transmissions in dense cellular networks. *IEEE Transactions on Wireless Communications* 15, 3 (2015), 2365–2380.
- [5] HONKOTE, V., KURIAN, D., MUTHUKUMAR, S., GHOSH, D., YADA, S., JAIN, K., JACKSON, B., KLOTCOV, I., NIMMAGADDA, M. R., DATTAWADKAR, S., 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)* (2019), IEEE, pp. 48–50.
- [6] JOSEPH, K. J., KHAN, S., KHAN, F. S., AND BALASUBRAMANIAN, V. N. Towards open world object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2021), pp. 5830–5840.
- [7] LI, A., SONG, S. L., CHEN, J., LI, J., LIU, X., TALLENT, N. R., AND BARKER, K. J. Evaluating modern gpu interconnect: Pcie, nvlink, nv-sli, nvswitch and gpudirect. *IEEE Transactions on Parallel and Distributed Systems* 31, 1 (2019), 94–110.
- [8] LI, Q., GAMA, F., RIBEIRO, A., AND PROROK, A. Graph neural networks for decentralized multi-robot path planning. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (2020), IEEE, pp. 11785–11792.
- [9] LI, Y., YU, Z., CHOY, C., XIAO, C., ALVAREZ, J. M., FIDLER, S., FENG, C., AND ANANDKUMAR, A. Voxformer: Sparse voxel transformer for camera-based 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 9087–9098.
- [10] LIU, R., AND CHOI, N. A first look at wi-fi 6 in action: Throughput, latency, energy efficiency, and security. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* 7, 1 (2023), 1–25.

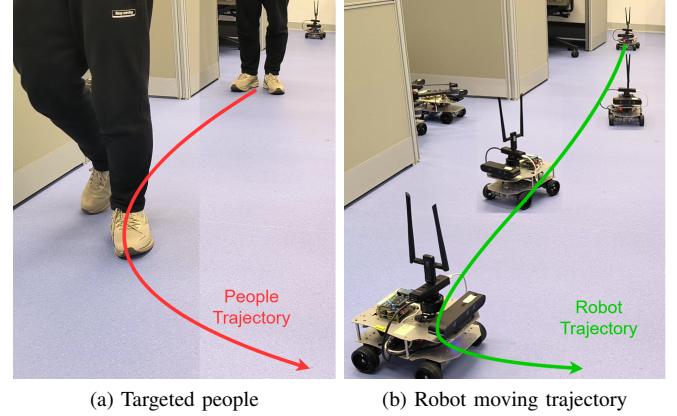


Fig. 3: A real-time people-tracking robotic application on our robot based on a well-known human pose estimation ML model, KAPAO[13].

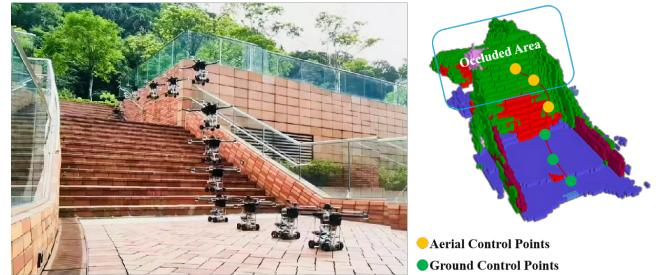


Fig. 4: By predicting occlusions in advance, AGRNav[21] gains an accurate perception of the environment and avoids collisions, resulting in efficient and energy-saving paths.

- [11] LIU, S., LI, X., LU, H., AND HE, Y. Multi-object tracking meets moving uav. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2022), pp. 8876–8885.
- [12] MASIUKIEWICZ, A. Throughput comparison between the new hew 802.11 ax standard and 802.11 n/ac standards in selected distance windows. *International Journal of Electronics and Telecommunications* 65, 1 (2019), 79–84.
- [13] MCNALLY, W., VATS, K., WONG, A., AND MCPHEE, J. Rethinking keypoint representations: Modeling keypoints and poses as objects for multi-person human pose estimation. In *European Conference on Computer Vision* (2022), Springer, pp. 37–54.
- [14] NOORMOHAMMADPOUR, M., AND RAGHAVENDRA, C. S. Datacenter traffic control: Understanding techniques and tradeoffs. *IEEE Communications Surveys & Tutorials* 20, 2 (2017), 1492–1525.
- [15] NVIDIA. The world's smallest ai supercomputer. <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-xavier-series/>, 2024.
- [16] OHKAWA, T., YAMASHINA, K., KIMURA, H., OOTSU, K., AND YOKOTA, T. Fpga components for integrating fpgas into robot systems. *IEICE TRANSACTIONS on Information and Systems* 101, 2 (2018), 363–375.
- [17] PEI, Y., MUTKA, M. W., AND XI, N. Connectivity and bandwidth-aware real-time exploration in mobile robot networks. *Wireless Communications and Mobile Computing* 13, 9 (2013), 847–863.
- [18] REN, Y., TUNG, C.-W., CHEN, J.-C., AND LI, F. Y. Proportional and preemption-enabled traffic offloading for ip flow mobility: Algorithms and performance evaluation. *IEEE Transactions on Vehicular Technology* 67, 12 (2018), 12095–12108.
- [19] SARKAR, N. I., AND MUSSA, O. The effect of people movement on wi-fi link throughput in indoor propagation environments. In *IEEE 2013 Tencon-Spring* (2013), IEEE, pp. 562–566.
- [20] WANG, H., POTLURI, S., LUO, M., SINGH, A. K., SUR, S., AND PANDA, D. K. Mvapich2-gpu: optimized gpu to gpu communication for infiniband clusters. *Computer Science-Research and Development* 26, 3 (2011), 257–266.
- [21] WANG, J., SUN, Z., GUAN, X., SHEN, T., ZHANG, Z., DUAN, T., HUANG, D., ZHAO, S., AND CUI, H. Agnav: Efficient and energy-saving autonomous navigation for air-ground robots in occlusion-prone environments. In *IEEE International Conference on Robotics and Automation (ICRA)* (2024).
- [22] WU, T., HE, S., LIU, J., SUN, S., LIU, K., HAN, Q.-L., AND TANG, Y. A brief overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal of Automatica Sinica* 10, 5 (2023), 1122–1136.
- [23] XIA, Z., LIU, Y., LI, X., ZHU, X., MA, Y., LI, Y., HOU, Y., AND QIAO, Y. Sepnet: Semantic scene completion on point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 17642–17651.
- [24] YANG, X., LIN, H., LI, Z., QIAN, F., LI, X., HE, Z., WU, X., WANG, X., LIU, Y., LIAO, Z., ET AL. Mobile access bandwidth in practice: Measurement, analysis, and implications. In *Proceedings of the ACM SIGCOMM 2022 Conference* (2022), pp. 114–128.
- [25] YANG, Y., JUNTAO, L., AND LINGLING, P. Multi-robot path planning based on a deep reinforcement learning dqn algorithm. *CAAI Transactions on Intelligence Technology* 5, 3 (2020), 177–183.