

Fault-Tolerant Collaborative Inference through the Edge-PRUNE Framework

Jani Boutellier¹ Bo Tan² Jari Nurmi²

Abstract

Collaborative inference has received significant research interest in machine learning as a vehicle for distributing computation load, reducing latency, as well as addressing privacy preservation in communications. Recent collaborative inference frameworks have adopted dynamic inference methodologies such as early-exit and run-time partitioning of neural networks. However, as machine learning frameworks scale in the number of inference inputs, e.g., in surveillance applications, fault tolerance related to device failure needs to be considered. This paper presents the Edge-PRUNE distributed computing framework, built on a formally defined model of computation, which provides a flexible infrastructure for fault tolerant collaborative inference. The experimental section of this work shows results on achievable inference time savings by collaborative inference, presents fault tolerant system topologies and analyzes their cost in terms of execution time overhead.

1. Introduction

Since a few years already, the execution of machine learning workloads has been moving from servers and the cloud to less powerful platforms, such as embedded and mobile devices. A considerable hindrance to this progress has been the significant computation load of machine learning inference, especially related to deep neural network (DNN) architectures. In order to match neural network complexity with computation platform resources, several different approaches have been developed: lightweight architectures such as MobileNets (Howard et al., 2017) attempt to maintain high inference accuracy despite drastically reduced

¹School of Technology and Innovations, University of Vaasa, Vaasa, Finland ²Faculty of Information Technology and Communication Sciences, Tampere University, Tampere, Finland. Correspondence to: Jani Boutellier <jani.boutellier@uwasa.fi>.

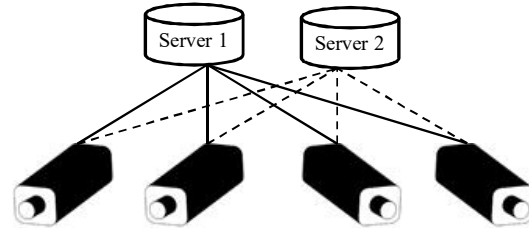


Figure 1. An example of a surveillance camera system with fault tolerance achieved by redundant nodes. The interconnection pattern equals to a $K_{2,4}$ complete bipartite graph, where 2 is the number of edge servers and 4 is the number of endpoint devices.

number of trainable parameters; post-hoc optimizations such as dense layer pruning (Zhu & Gupta, 2018), separable convolutions (Jaderberg et al., 2014) and weight quantization (Courbariaux et al., 2016) reduce inference time by approximating the original trained neural architecture, whereas neural accelerators (Skillman & Edsö, 2020; Han et al., 2016) leverage specialized hardware to speed up inference.

Orthogonal to the aforementioned techniques, distributed and collaborative inference have emerged as a notable branch of research. In these approaches, the neural network inference workload is distributed between low-resource *endpoint devices* and high performance *edge servers* (or the cloud); early milestone works of this direction are Neurosurgeon (Kang et al., 2017) and DDNN (Teerapittayanon et al., 2017). In conjunction with collaborative inference, several dynamic neural network techniques have been successfully adopted: for instance, early-exit (Teerapittayanon et al., 2017) can terminate inference at intermediate layers saving on communication bandwidth, whereas dynamic on-loading (Almeida et al., 2021) decides at runtime the target execution platform of DNN layers.

Collaborative inference can be used to target a variety of optimization objectives: reducing endpoint device computation load, end-to-end latency optimization, energy reduction or server workload reduction. A less frequently mentioned by-product of collaborative inference is privacy preservation. This topic has been extensively studied in (He et al., 2019), where key observations point out that malicious *model inversion attacks* against collaborative inference are significantly

less effective if they cannot access the feature vectors produced by the early neural network layers – especially if inference has passed one or more dense layers of the DNN. In terms of collaborative inference this means that the endpoint device should perform the inference of as many early DNN layers as possible before transmitting the intermediate feature vectors over a network interface to server processing.

This paper addresses the topic of collaborative inference, especially from the point of multiple inference inputs and system fault tolerance. An example scenario for multi-input collaborative inference is a smart surveillance camera system (see Figure 1), where several smart cameras have been deployed across a site for performing object detection and/or tracking. For privacy preservation, the smart cameras perform the inference of early DNN layers, and transmit the intermediate feature vectors to a local edge server for completing the inference.

In safety critical areas, node (endpoint device or server) failures related to malicious actions or hardware faults needs to be taken in account. Concretely, the surveillance system should be able to continue its operation if one or more of the endpoint devices become disabled, or even in the more severe case of server failure. Figure 1 illustrates a highly redundant configuration, where a single failure of any kind of resource (endpoint device, server or connection) does not incapacitate the overall system. Such redundancy evidently comes with a price, which in this case is the redundant server and the related connectivity.

This study, related to fault tolerant collaborative inference, is realized around the open source¹ *Edge-PRUNE* framework (Boutellier et al., 2022b). *Edge-PRUNE* is based on a formally defined *model of computation*, and provides the necessary theoretical infrastructure for specifying and designing collaborative inference between one or more endpoint devices and servers. Besides the theoretical basis, the *Edge-PRUNE* framework includes a self-sustained runtime engine that is hardware and training framework agnostic, providing a software environment for both endpoint devices and servers. Although not detailed in this work, *Edge-PRUNE* also has inherent support for conditional computing.

2. Related Work

Significant early works on collaborative inference were DDNN (Teerapittayanon et al., 2017) and Neurosurgeon (Kang et al., 2017). DDNN proposed distributing inference across endpoint, edge and cloud resources, also introducing early exits for reducing communication. Neurosurgeon, on the other hand, presented a scheduler for intelligently distributing neural network layers across endpoint and server resources. Edgent (Li et al., 2018a) further developed Neu-

rosurgeon’s concept by introducing DNN right-sizing, joint optimization of early exits and DNN partitioning. Edgent later evolved into Boomerang (Zeng et al., 2019) inspired by the early exit mechanism of BranchyNet (Teerapittayanon et al., 2016). IONN (Jeong et al., 2018) also continued in the vein of Neurosurgeon, however based on the offloading concept: the endpoint device can upload DNN partitions to an edge server for optimizing mobile device energy consumption, among other optimization goals. Similar to Neurosurgeon, also IONN is based on Caffe (Jia et al., 2014). Recently, SplitNets (Dong et al., 2022) combined neural architecture search with multi-input partition point search.

JointDNN (Eshratifar et al., 2019) introduced a directed acyclic graph (DAG) based model for DNN partitioning, optimizing for energy and latency. Besides partitioning, JointDNN also considers layer compression, similar to the preceding work JALAD (Li et al., 2018b), and the recent *supervised compression* work (Matsubara et al., 2022). A graph-based modeling approach is also adopted by DADS (Hu et al., 2019), the industrial effort Auto-Split (Banitalebi-Dehkordi et al., 2021) and D^3 (Zhang et al., 2021), enabling capturing of branched DNN topologies as opposed to simpler chain-like DNN structures. Finally, SPINN (Laskaridis et al., 2020) and DynO (Almeida et al., 2021) contribute to dynamic DNN partitioning, which is useful under, e.g., varying wireless network conditions.

An orthogonal approach to endpoint-server computation partitioning is taken by (Mao et al., 2017b;a; Zhao et al., 2018; Gao et al., 2021) that propose partitioning DNN inference across multiple endpoint devices.

3. The Edge-PRUNE Framework

The *Edge-PRUNE* framework used in this study significantly differs from the related works in the sense that it is based on a formal model of computation. The overall computation scheme of *Edge-PRUNE* is *dataflow*, similar to that of TensorFlow (Abadi et al., 2016). However, *Edge-PRUNE* goes further in computation modeling, by formalizing concepts such as data packaging, data rates, triggering of computations, and necessary conditions for deadlock-free conditional computations.

3.1. Model of Computation

The *Edge-PRUNE* framework relies on the VR-PRUNE model of computation (Boutellier et al., 2022a). In VR-PRUNE, a neural network is described as a directed graph $G = (A, F)$, where the *actors* A represent vertices that perform computation, such as inference of a DNN layer. The links F of graph G represent first-in-first-out (FIFO) buffers that carry data between actors. For each link $f \in F$, data is quantized into fixed-size *tokens* that in the neural network

¹Available at <https://gitlab.com/jboutell/vprf/-/tree/edge-prune>

context can be understood as feature vectors between layers. A FIFO $f \in F$ is connected to an actor $a \in A$ through a port p_a , such that $fifo(p_a) = f$.

An actor $a \in A$ can *fire* (perform a computation) when it has a sufficient number of input tokens available. To specify the required number of tokens, for each input port p_a of actor a , a non-negative integer-valued *token rate* $atr(p_a)$ is defined; once each input port of a has at least $atr(p_a)$ number of tokens in the associated FIFO buffer $fifo(p_a)$, actor a becomes *enabled*, that is, ready to fire. The exact moment in time when an enabled actor fires can depend, for instance, on availability of compute resources. As dataflow models in general, also VR-PRUNE is inherently *concurrent*: individual actors can execute in parallel, independent of others (Lee & Messerschmitt, 1987).

The VR-PRUNE model of computation balances between expressiveness and analyzability: while being expressive enough for allowing conditional computations, the model simultaneously provides means for analyzing graph consistency. The following subsection illustrates how conditional computation is expressed using VR-PRUNE concepts.

3.2. Conditional Computation

In order to maintain analyzability against graph deadlock and/or buffer overflow, the VR-PRUNE model restricts conditional computation to take place within *dynamic processing (sub)graphs*, DPGs: within a DPG, conditional computation can be realized between two *dynamic actors*. Figure 2 shows a minimal case of a DPG: the dynamic actor x provides a control signal (dashed connection) that at run time sets the input and output token rates of the dynamic actor y and the *dynamic processing actor* a . The range of allowed token rates and associated actor ports is expressed in the *control table* T of Figure 2: port p_{xc} dynamically sets the token rate of ports p_{x1} , p_{a1} , p_{a2} and p_{y1} to either 0 or to 1. With token rates set to 1, a becomes enabled, whereas token rate 0 disables execution of a . Actor b does not receive such a control signal and thus maintains static token rates.

3.3. Distributed Computing and Fault Tolerance

The Edge-PRUNE framework enables concurrent execution of actors both within a computing platform, and between computing platforms using a *mapping specification*. Within a platform, each actor $a \in A$ is mapped for execution to a specific CPU core, or to the local GPU. The same mapping specification is also used to set the execution platform (endpoint device or server) of each actor. To this extent, Edge-PRUNE provides mapping exploration functionality, which auto-generates endpoint-server mapping alternatives for discovering the best DNN partition point for collaborative inference.

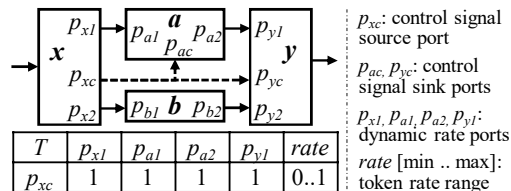


Figure 2. An Edge-PRUNE dynamic processing graph example.

The Edge-PRUNE mapping specification allows generalizing distributed computing also to system configurations that consist of one or more endpoint devices and/or servers, e.g., for achieving fault tolerance. From the dataflow model viewpoint, computing node malfunction means that either a subgraph of G ceases to produce tokens (endpoint failure), or a subgraph of G stops consuming tokens (server failure). Maintaining operation for the remaining system requires that the dataflow application needs to overcome such unexpected subgraph failures without ending up in deadlock. Currently, Edge-PRUNE implements fault tolerance on the dataflow graph level, whereas extending the fault tolerance behavior to model-based graph consistency analysis remains a prospective for future work.

3.4. Inference Engine

The Edge-PRUNE computing functionalities have been implemented in the C language to a lightweight runtime library for Linux-based platforms, enabling deployment to embedded devices as well as to servers. This inference engine is independent of neural network training frameworks (e.g., TensorFlow), but allows leveraging DNN acceleration libraries such as Intel oneDNN or ARM CL. Edge-PRUNE has deeply inbuilt support for GPU leverage, but can equally well operate on GPU-less platforms. Distributed computing functionality has been implemented using Linux Sockets, such that the endpoint devices are expected to establish an `ssh` connection to the server(s), delegating data security issues to the level of `ssh` connections. Fault tolerance functionality is implemented on the actor port level: ports responsible for inter-platform communication monitor and adjust to remote platform liveness by Linux socket error conditions: broken pipe, connection reset, no data sent. Finally, Edge-PRUNE does not constrain the nature (wired or cable) or number (shared or point-to-point) of connections between endpoints and servers.

4. Experiments

4.1. Collaborative Inference Partition Point Exploration

Figure 3 shows endpoint device inference time for the SSD-MobileNet v1 (Howard et al., 2017; Liu et al., 2016) ob-

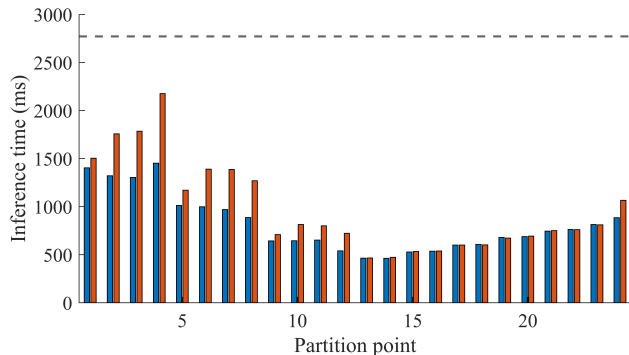


Figure 3. Collaborative inference partition point exploration for SSD-MobileNet v1 (300×300 input) under Edge-PRUNE. Bars reflect endpoint inference + communication time (blue: 100 Mbit Ethernet, red: 16 Mbit WiFi); horizontal dashed line is endpoint-only inference time (all per 1 frame).

ject detection CNN in a collaborative inference scenario where an ODROID N2 single-board computer (Hexacore ARM+GPU) acts as the endpoint device, and an Intel Core i7-8650U based platform as the edge server, interconnected by 100 Mbit Ethernet or 16 Mbit WiFi. The full-precision CNN (32 bit float) is implemented as an Edge-PRUNE application such that each Conv-BNorm-ReLU layer triplet is wrapped inside a dedicated actor, forming a graph of 53 actors and 69 links. Endpoint device inference time was explored by shifting the endpoint/server partition point actor-by-actor from inference input towards inference output, resulting in Figure 3. The graph shows that partition points 13 and 14 provide minimal endpoint inference time, $6.0\times$ higher than endpoint-only inference. On the endpoint device, the CNN layers inside Edge-PRUNE actors were implemented using ARM Compute Library functions.

4.2. Inference Time and Fault Tolerance

Both virtual and physically distributed system configurations were used to validate Edge-PRUNE collaborative inference and system behavior on computing node failures.

Virtual environment. Table 1 illustrates performance scaling of Edge-PRUNE collaborative inference for single and dual-server configurations with $1 \leq n \leq 6$ endpoint devices and complete bipartite graph $K_{m,n}$ interconnection topology. Per-frame processing time was measured in a virtual distributed environment: each endpoint process and each server process was assigned to a dedicated core on the 8-core Intel Core i7-8650U processor, and connections between endpoints and servers were handled over the Linux *loopback* network interface, which allowed using exactly the same server and endpoint software configurations as in a physically distributed system. Each endpoint performed

Table 1. Vehicle image classification collaborative inference time per frame in milliseconds for $1 \leq m \leq 2$ servers and $1 \leq n \leq 6$ endpoint devices. Upper half: endpoint, lower half: server.

N. OF ENDPOINTS	1	2	3	4	5	6
SINGLE-SERVER	4.7	4.9	5.2	5.4	6.2	7.0
DUAL-SERVER	4.9	5.1	5.2	5.5	6.3	7.1
SINGLE-SERVER	4.8	5.7	6.5	7.1	8.7	9.7
DUAL-SERVER	4.9	5.8	6.5	7.1	8.9	9.6

the inference of 7 initial layers (Conv2D-ReLU-MaxPool-Conv2D-ReLU-MaxPool-Dense) of a vehicle classification CNN (Xie et al., 2016), whereas each server process performed the inference of the last 5 layers of the same CNN. The 7 CNN layers of each endpoint subgraph were realized into 4 actors, whereas the 5 layers of the server subgraph were wrapped inside a single actor; therefore the largest $K_{2,6}$ configuration consisted of 26 actors and 30 FIFOs. On the endpoint side, Conv2D layer inference was implemented by the Intel oneDNN library. Table 1 shows: a) adding a second server for fault tolerance causes insignificant processing time overhead, and b) adding an endpoint increases per-node inference time by 11% on average. Endpoint and server fault behavior was observed by abrupt termination of endpoint/server processes, showing that the remaining system continues collaborative inference as expected.

Heterogeneous distributed system. To confirm Edge-PRUNE collaborative inference functionality and behavior on node fault on a physically distributed system, $K_{1,2}$ and $K_{2,1}$ configurations were established using a heterogeneous set of nodes: an Intel Core i7-8650U workstation, an ODROID N2, and an Intel Atom N270 based platform, all running Ubuntu Linux, and interconnected over 100 Mbit Ethernet. Using the same vehicle classification CNN (Xie et al., 2016), the heterogeneous configuration was used to validate that a) Edge-PRUNE recovers from communication breaks (temporarily disconnected cable), and b) maintains operation on permanent node failure (server power-down for $K_{2,1}$ and endpoint power-down for $K_{1,2}$).

5. Conclusion

In this work the topic of collaborative inference fault tolerance was studied for configurations consisting of one or more network-connected edge servers and one or more endpoint devices. The experimental study was done using the Edge-PRUNE framework, which was shown to scale successfully with an increasing number of endpoint devices and/or edge servers, and furthermore was shown to be capable of continuing operation after computing node failure.

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. Tensorflow: A system for large-scale machine learning. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pp. 265–283, 2016.
- Almeida, M., Laskaridis, S., Venieris, S. I., Leontiadis, I., and Lane, N. D. DynO: Dynamic onloading of deep neural networks from cloud to device. *ACM Transactions on Embedded Computing Systems (TECS)*, 2021.
- Banitalebi-Dehkordi, A., Vedula, N., Pei, J., Xia, F., Wang, L., and Zhang, Y. Auto-split: a general framework of collaborative edge-cloud ai. In *ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021.
- Boutellier, J., Ma, Y., Wu, J., Khan, M., and Bhattacharyya, S. S. VR-PRUNE: Decidable variable-rate dataflow for signal processing systems. *IEEE Transactions on Signal Processing*, 70:1819 – 1833, 2022a.
- Boutellier, J., Tan, B., and Nurmi, J. Edge-prune: Flexible distributed deep learning inference. *arXiv preprint arXiv:2204.12947*, 2022b.
- Courbariaux, M., Hubara, I., Soudry, D., El-Yaniv, R., and Bengio, Y. Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1. *arXiv preprint arXiv:1602.02830*, 2016.
- Dong, X., De Salvo, B., Li, M., Liu, C., Qu, Z., Kung, H., and Li, Z. Splitnets: Designing neural architectures for efficient distributed computing on head-mounted systems. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12559–12569, 2022.
- Eshratifar, A. E., Abrishami, M. S., and Pedram, M. JointDNN: an efficient training and inference engine for intelligent mobile cloud computing services. *IEEE Transactions on Mobile Computing*, 20(2):565–576, 2019.
- Gao, Z., Sun, S., Zhang, Y., Mo, Z., and Zhao, C. EdgeSP: Scalable multi-device parallel DNN inference on heterogeneous edge clusters. In *International Conference on Algorithms and Architectures for Parallel Processing*, pp. 317–333. Springer, 2021.
- Han, S., Liu, X., Mao, H., Pu, J., Pedram, A., Horowitz, M. A., and Dally, W. J. EIE: efficient inference engine on compressed deep neural network. In *ACM/IEEE International Symposium on Computer Architecture (ISCA)*, pp. 243–254. IEEE, 2016.
- He, Z., Zhang, T., and Lee, R. B. Model inversion attacks against collaborative inference. In *Annual Computer Security Applications Conference*, 2019.
- Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., and Adam, H. MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Hu, C., Bao, W., Wang, D., and Liu, F. Dynamic adaptive DNN surgery for inference acceleration on the edge. In *IEEE Conference on Computer Communications*, 2019.
- Jaderberg, M., Vedaldi, A., and Zisserman, A. Speeding up convolutional neural networks with low rank expansions. In *British Machine Vision Conference*, 2014.
- Jeong, H.-J., Lee, H.-J., Shin, C. H., and Moon, S.-M. IONN: Incremental offloading of neural network computations from mobile devices to edge servers. In *ACM Symposium on Cloud Computing*, 2018.
- Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., and Darrell, T. Caffe: Convolutional architecture for fast feature embedding. In *ACM international conference on Multimedia*, 2014.
- Kang, Y., Hauswald, J., Gao, C., Rovinski, A., Mudge, T., Mars, J., and Tang, L. Neurosurgeon: Collaborative intelligence between the cloud and mobile edge. *ACM SIGARCH Computer Architecture News*, 45(1):615–629, 2017.
- Laskaridis, S., Venieris, S. I., Almeida, M., Leontiadis, I., and Lane, N. D. SPINN: synergistic progressive inference of neural networks over device and cloud. In *Annual International Conference on Mobile Computing and Networking*, 2020.
- Lee, E. A. and Messerschmitt, D. G. Synchronous data flow. *Proceedings of the IEEE*, 75(9):1235–1245, 1987. ISSN 0018-9219.
- Li, E., Zhou, Z., and Chen, X. Edge intelligence: On-demand deep learning model co-inference with device-edge synergy. In *Workshop on Mobile Edge Communications*, 2018a.
- Li, H., Hu, C., Jiang, J., Wang, Z., Wen, Y., and Zhu, W. JALAD: Joint accuracy-and latency-aware deep structure decoupling for edge-cloud execution. In *IEEE International conference on parallel and distributed systems (ICPADS)*, 2018b.
- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.-Y., and Berg, A. C. SSD: Single shot multibox detector. In *European Conference on Computer Vision*, 2016.
- Mao, J., Chen, X., Nixon, K. W., Krieger, C., and Chen, Y. MoDNN: Local distributed mobile computing system for

- deep neural network. In *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2017a.
- Mao, J., Yang, Z., Wen, W., Wu, C., Song, L., Nixon, K. W., Chen, X., Li, H., and Chen, Y. MeDNN: A distributed mobile system with enhanced partition and deployment for large-scale DNNs. In *IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, 2017b.
- Matsubara, Y., Yang, R., Levorato, M., and Mandt, S. Supervised compression for resource-constrained edge computing systems. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2685–2695, 2022.
- Skillman, A. and Edsö, T. A technical overview of Cortex-M55 and Ethos-U55: Arm’s most capable processors for endpoint AI. In *IEEE Hot Chips Symposium (HCS)*, 2020.
- Teerapittayanon, S., McDanel, B., and Kung, H.-T. Branchynet: Fast inference via early exiting from deep neural networks. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pp. 2464–2469. IEEE, 2016.
- Teerapittayanon, S., McDanel, B., and Kung, H.-T. Distributed deep neural networks over the cloud, the edge and end devices. In *IEEE International conference on distributed computing systems (ICDCS)*, 2017.
- Xie, R., Huttunen, H., Lin, S., Bhattacharyya, S. S., and Takala, J. Resource-constrained implementation and optimization of a deep neural network for vehicle classification. In *European Signal Processing Conference*, pp. 1862–1866, 2016.
- Zeng, L., Li, E., Zhou, Z., and Chen, X. Boomerang: On-demand cooperative deep neural network inference for edge intelligence on the industrial internet of things. *IEEE Network*, 33(5), 2019.
- Zhang, B., Xiang, T., Zhang, H., Li, T., Zhu, S., and Gu, J. Dynamic DNN decomposition for lossless synergistic inference. In *IEEE International Conference on Distributed Computing Systems Workshops (ICDCSW)*, 2021.
- Zhao, Z., Barijough, K. M., and Gerstlauer, A. DeepThings: Distributed adaptive deep learning inference on resource-constrained IoT edge clusters. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 37(11):2348–2359, 2018.
- Zhu, M. and Gupta, S. To prune, or not to prune: exploring the efficacy of pruning for model compression. In *International Conference on Learning Representations (ICLR) Workshops*, 2018.